

1 Genetic control of mRNA splicing as a potential  
2 mechanism for incomplete penetrance of rare  
3 coding variants

4  
5 Jonah Einson<sup>1,2</sup>, Dafni Glinos<sup>2</sup>, Eric Boerwinkle<sup>3</sup>, Peter Castaldi<sup>4</sup>, Dawood Darbar<sup>5</sup>,  
6 Mariza de Andrade<sup>6</sup>, Patrick Ellinor<sup>7</sup>, Myriam Fornage<sup>8</sup>, Stacey Gabriel<sup>9</sup>, Soren  
7 Germer<sup>2</sup>, Richard Gibbs<sup>10</sup>, Craig P. Hersh<sup>11</sup>, Jill Johnsen<sup>12</sup>, Robert Kaplan<sup>13</sup>, Barbara A.  
8 Konkle<sup>12</sup>, Charles Kooperberg<sup>14</sup>, Rami Nassir<sup>15</sup>, Ruth J.F. Loos<sup>16</sup>, Deborah A Meyers<sup>17</sup>,  
9 Braxton D. Mitchell<sup>18,19</sup>, Bruce Psaty<sup>20</sup>, Ramachandran S. Vasan<sup>21</sup>, Stephen S. Rich<sup>22</sup>,  
10 Michael Rienstra<sup>23</sup>, Jerome I. Rotter<sup>24</sup>, Aabida Saferali<sup>11</sup>, M. Benjamin Shoemaker<sup>25</sup>,  
11 Edwin Silverman<sup>26</sup>, Albert Vernon Smith<sup>27</sup>, NHLBI Trans-Omics for Precision Medicine  
12 (TOPMed) Consortium, Pejman Mohammadi<sup>29</sup>, Stephane E. Castel<sup>2,30</sup>, Ivan Iossifov<sup>2,31</sup>,  
13 Tuuli Lappalainen<sup>32,33,2</sup>

14

15

16 <sup>1</sup>Department of Biomedical Informatics, Columbia University, <sup>2</sup>New York Genome Center, <sup>3</sup>University of  
17 Texas Health at Houston, <sup>4</sup>Department of Medicine, Brigham & Women's Hospital, <sup>5</sup>Department of  
18 Cardiology, University of Illinois at Chicago, <sup>6</sup>Department of Quantitative Health Sciences, Mayo Clinic,  
19 <sup>7</sup>Corrigan Minehan Heart Center, Massachusetts General Hospital, <sup>8</sup>Brown Foundation Institute of  
20 Molecular Medicine, McGovern Medical School, University of Texas Health at Houston, <sup>9</sup>Broad Institute,  
21 <sup>10</sup>Department of Molecular and Human Genetics, Baylor College of Medicine Human Genome  
22 Sequencing Center, <sup>11</sup>Channing Division of Network Medicine and Division of Pulmonary and Critical  
23 Care Medicine, Brigham and Women's Hospital, <sup>12</sup>Department of Hematology, University of Washington,  
24 <sup>13</sup>Department of Epidemiology & Population Health, Albert Einstein College of Medicine, <sup>14</sup>Fred  
25 Hutchinson Cancer Research Center, <sup>15</sup>Department of Pathology, School of Medicine, Umm Al-Qura  
26 University, <sup>16</sup>Environmental Medicine & Public Health, Icahn School of Medicine at Mount Sinai,  
27 <sup>17</sup>Department of Medicine, University of Arizona, <sup>18</sup>Department of Medicine, University of Maryland School  
28 of Medicine, <sup>19</sup>Geriatrics Research and Education Clinical Center, Baltimore Veterans Administration  
29 Medical Center, <sup>20</sup>Cardiovascular Health Research Unit, Departments of Medicine, Epidemiology, and  
30 Health Systems and Population Health, University of Washington, <sup>21</sup>Department of Medicine, Boston  
31 University, <sup>22</sup>Public Health Sciences, University of Virginia, <sup>23</sup>Clinical Cardiology, UMGC Cardiology, <sup>24</sup>The  
32 Institute for Translational Genomics and Population Sciences, Department of Pediatrics, The Lundquist  
33 Institute for Biomedical Innovation at Harbor-UCLA Medical Center, <sup>25</sup>Department of Medicine, Vanderbilt  
34 University, <sup>26</sup>Channing Division of Network Medicine and Division of Pulmonary and Critical Care  
35 Medicine, Brigham & Women's Hospital, <sup>27</sup>Department of Biostatistics, University of Michigan,  
36 <sup>29</sup>Department of Integrative Structural and Computational Biology, The Scripps Research Institute,  
37 <sup>30</sup>Variant Bio, <sup>31</sup>Cold Spring Harbor Laboratory, <sup>32</sup>Department of Systems Biology, Columbia University,  
38 <sup>33</sup>Department of Gene Technology, KTH Royal Institute of Technology

39

## 40 Abstract

41 Exonic variants present some of the strongest links between genotype and phenotype.  
42 However, these variants can have significant inter-individual pathogenicity differences, known  
43 as variable penetrance. In this study, we propose a model where genetically controlled mRNA  
44 splicing modulates the pathogenicity of exonic variants. By first cataloging exonic inclusion from  
45 RNA-seq data in GTEx v8, we find that pathogenic alleles are depleted on highly included  
46 exons. Using a large-scale phased WGS data from the TOPMed consortium, we observe that  
47 this effect may be driven by common splice-regulatory genetic variants, and that natural  
48 selection acts on haplotype configurations that reduce the transcript inclusion of putatively  
49 pathogenic variants, especially when limiting to haploinsufficient genes. Finally, we test if this  
50 effect may be relevant for autism risk using families from the Simons Simplex Collection, but  
51 find that splicing of pathogenic alleles has a penetrance reducing effect here as well. Overall,  
52 our results indicate that common splice-regulatory variants may play a role in reducing the  
53 damaging effects of rare exonic variants.

## 54 Introduction

55 Incomplete penetrance is a well known phenomenon, where an individual carries a disease-  
56 associated allele, but develops no symptoms of the disease themselves (Forrest et al. 2022;  
57 Gettler et al. 2021; Shawky 2014). Similarly, variable expressivity refers to analogous gradual  
58 differences in disease severity; here we refer to both as variable penetrance. These instances  
59 are likely underreported in the literature due to ascertainment bias, when many studies are  
60 based on sequencing due to a prior genetic condition (Cooper et al. 2013; Dewey et al. 2016).  
61 Even amongst Mendelian disease variants, which are typically thought of as having strong  
62 effects on phenotype, differing levels of severity have been observed between carriers (Chen et  
63 al. 2016). These changes have been attributed to epistatic or additive effects of genetic  
64 modifiers, as well as environmental modifiers of penetrance, which can be difficult to control in  
65 an experimental setting (Maya et al. 2018). When looking at incomplete penetrance in specific  
66 diseases, genetic modifiers have been mapped, for example, to BRCA in breast cancer (Milne  
67 and Antoniou 2011), and RET in Hirschsprung's disease (Emison et al. 2005). Modified  
68 penetrance has also been studied in the context of polygenic risk scores, where multiple  
69 common risk variants increase the expected pathogenicity of a disease-relevant variant (Fahed  
70 et al. 2020). However, genome-wide patterns underlying modified penetrance are still poorly

71 known. One potential mechanism for incomplete penetrance are cis-regulatory mechanisms that  
72 affect the regulation of a gene carrying a pathogenic variant. This model has been tested with  
73 expression quantitative trait loci (eQTLs) acting as modifiers of penetrance (Castel et al. 2018),  
74 but can be expanded to other types of gene regulatory processes, such as mRNA splicing.  
75 While eQTLs control the dosage of their target genes, splicing alters inclusion of variant-carrying  
76 exons in transcripts, which could potentially have a large effect on the overall pathogenicity of a  
77 damaging variant.

78  
79 Alternative splicing is responsible for the great diversity of isoform structures observed across  
80 human tissues and cell types (Keren et al. 2010). With regard to coding variant interpretation,  
81 exons with lower expression have been shown to be less likely to harbor pathogenic variants,  
82 while ubiquitously included exons can be prioritized for gene disrupting rare variants (Cummings  
83 et al. 2020). Autistic individuals with variants on the same exons have been shown to have  
84 remarkably similar disease phenotypes, putatively due to the variants having similar effects on  
85 gene dosage or function, a notable finding given the extreme heterogeneity of the condition  
86 (Chiang et al. 2021). Additionally, splicing can be influenced by common genetic variation, as  
87 evidenced by the many studies that use large scale WGS and transcriptomic datasets to map  
88 splicing quantitative trait loci (sQTLs) (Alasoo et al. 2019; Consortium 2020; Garrido-Martín et  
89 al. 2021; Kerimov et al. 2020). sQTLs in general have been implicated in disease risk and other  
90 genetic traits (Li et al. 2016; Noble et al. 2020; Ongen and Dermitzakis 2015).

91  
92 In this study, we build upon the finding that transcript usage of genes containing alleles  
93 contributes to the allele's pathogenicity, and ask if common splice-regulatory variants may  
94 partially drive this phenomenon and affect inter-individual variation in penetrance. Expanding on  
95 previous methodology (Castel et al. 2018), we look for non-random haplotype combinations of  
96 sQTL variants and putatively pathogenic rare variants in population scale datasets. Such an  
97 observation could indicate that haplotype combinations have an effect on fitness, and by proxy,  
98 disease risk. In doing so, we develop a general framework for modeling common and rare  
99 variant haplotypes in a population, with a corresponding test to detect deviations from the null  
100 (Figure 1, Supplemental Figure 1). These analyses will improve our understanding of how  
101 variants across the annotation and allele frequency spectrum act together to shape human traits  
102 and could ultimately aid our interpretation of rare variants in a clinical context.

## 103 Results

### 104 **Deleterious rare alleles accumulate at lowly spliced exons with respect to the** 105 **population**

106 We first tested the hypothesis that rare pathogenic alleles (CADD > 15) (Rentzsch et al. 2019)  
107 are more likely to occur at less spliced-in exons (Figure 1). To accomplish this, we used bulk  
108 RNA-sequencing (RNA-seq) and whole genome sequencing (WGS) data from the Genotype  
109 Tissue Expression Project (GTEx) v8 release, which is representative of a general population  
110 free of severe genetic disease. We defined variants as rare if their variant frequency in gnomAD  
111 (Karczewski et al. 2020) was less than 0.5% and they appeared 5 or fewer times among the 838  
112 GTEx WGS donors.

113 To begin, we calculated percent spliced in (PSI) scores for all annotated protein-coding gene  
114 exons across 18 GTEx tissues, and only kept exons with sufficient splicing variability across  
115 individuals (Methods, Supplemental Table 1, Supplemental Figure S2A). We extracted rare  
116 alleles that fell on variably spliced exons, separating alleles within 10bp of a splice junction to  
117 avoid cases where the allele is more likely to directly affect splicing. To compare the splicing of  
118 each donor with a deleterious allele to the population distribution per exon, we calculated PSI Z-  
119 scores across all tissues with available data (Supplemental Figure S2B, Methods). We found  
120 that PSI Z-scores were significantly different between exons carrying deleterious (N = 19,178)  
121 and non-deleterious (N = 49,575) rare alleles (Mann-Whitney U-Test:  $p = 2.577 \times 10^{-4}$ ). This rank  
122 difference was accounted for by a modest decrease in mean PSI Z-score among donors that  
123 carried deleterious alleles in a given exon, which was consistent across tissues and across  
124 variant consequence annotations (Figure 2, Supplemental Figure S3). Notably, stop-gained  
125 variants had the strongest association with low PSI Z-scores - even stronger than the signal for  
126 variants close to splice junction - but the overall result was present for multiple annotation  
127 categories (Supplemental Figure S3). This suggests that the signal is not solely driven by the  
128 most pathogenic variants nor direct rare variant effects on splicing. These results extend the  
129 previous work, comparing different exons and showing accumulation of stop-gained variants on  
130 those with lower inclusion (Cummings et al. 2020). Here, observe a similar pattern when  
131 comparing different individuals within a given exon, consistent with the hypothesis that the  
132 penetrance of coding alleles is reduced when they fall on more lowly included exons. However,  
133 this approach does not discern the underlying reasons for splicing differences between

134 individuals, including alleles that may drive a decrease in splicing and their haplotype  
135 combinations with rare alleles.

### 136 **A general model for coding allele-QTL haplotype configurations**

137 We next sought to test if regulatory alleles on the same haplotype as rare coding alleles  
138 contribute to this phenomenon, using phased whole genome sequencing (WGS) data. Since  
139 directly quantifying the penetrance of coding alleles is difficult, our approach was to observe  
140 modified penetrance through the lens of purifying selection, where high-penetrance haplotype  
141 combinations would be depleted from the general population. Advantageously, this technique  
142 allows us to use large phased WGS datasets where individual gene expression data is not  
143 available.

144  
145 Initially, splice-regulatory alleles were cataloged in GTEx through quantitative trait locus (QTL)  
146 mapping, using the percent spliced in (PSI or  $\psi$ ) (Pervouchine et al. 2013) of each exon as a  
147 quantitative phenotype. These alleles are hence referred to as  $\psi$ QTLs. We use the “ $\psi$ ”  
148 nomenclature to differentiate from sQTLs, where the splicing phenotype can vary between  
149 studies and is often less interpretable for downstream applications.  $\psi$ QTL mapping and  
150 properties are described in (Einson et al. 2022). Briefly, we mapped  $\psi$ QTLs from GTEx v8 using  
151 the same filtered set of PSI scores across 18 tissues as in the previous analyses (see Methods).  
152 We compiled a set of 5,196 cross-tissue  $\psi$ QTL genes (one sVariant and one sExon per gene),  
153 and recorded which alleles led to higher or lower sExon inclusion. We also mapped secondary  
154 sExons across  $\psi$ QTL genes where the top sVariant was also associated with splicing in the  
155 same direction as the top sExon in the same gene, which were used to expand the amount of  
156 genic space where rare variants could be considered.

157  
158 Next, to robustly test for non-random haplotype combinations of rare exonic alleles and common  
159  $\psi$ QTL alleles, we describe an approach that quantifies the significance of deviations in  
160 haplotype combinations from the null in a dataset, taking variable  $\psi$ QTL allele frequencies into  
161 account: In most datasets,  $\psi$ QTL alleles that may have an effect on rare variant penetrance are  
162 non-uniformly distributed, and thus we expect an unequal number of high and low penetrance  
163 haplotypes under the null (Figure 3). To account for this, we model these data using the  
164 Poisson-Binomial distribution, a generalization of the Binomial distribution describing the sum of  
165  $n$  independent but non-identically distributed Bernoulli random variables. (González et al. 2016;  
166 Hong 2013; Wang 1993) When looking at counts of haplotype combinations, the probability of

167 observing a high-penetrance haplotype is assigned according to the relevant  $\psi$ QTL allele  
168 frequency, independently across QTL genes. To apply the model to haplotypes extracted from  
169 phased genetic data, we developed a bootstrapping procedure that approximates the  
170 cumulative distribution function of the Poisson-Binomial, constituting a convenient method for  
171 calculating the significance, enrichment/depletion effect sizes ( $\epsilon$ ) and confidence intervals when  
172 comparing enrichment scores between groups i.e. haplotypes with deleterious vs. non-  
173 deleterious rare alleles (see Methods for details). In simulations, our method was well powered  
174 to detect deviations from the null across all tested theoretical allele frequency distributions, and  
175 performed well against other methods that directly calculate and approximate the CDF of the  
176 Poisson-binomial. (Figure 4, Supplemental Figure S4). This approach is generalizable to other  
177 analyses of haplotype combinations; here we apply it to test nonrandom combinations of  $\psi$ QTL  
178 and rare coding alleles.

### 179 **High penetrance haplotypes are depleted in TOPMed and GTEx**

180 After defining a theoretical model that describes counts of common regulatory alleles and rare  
181 coding alleles in a given population, we tested three datasets for evidence of selection against  
182 high penetrance coding alleles driven by genetically regulated splicing.

#### 183 Enrichment in GTEx

184 We identified  $\psi$ QTL-rare allele haplotypes using population and read-backed phased (Castel et  
185 al. 2016) WGS data from GTEx V8, labeling haplotypes in putative high and low penetrance  
186 configurations according to whether the rare alternative allele was on the higher or lower  
187 inclusion  $\psi$ QTL haplotype, respectively (Figure 1 & 3). We limited our analysis to European-  
188 Americans, since the  $\psi$ QTL data is dominated by European ancestries, with rare variants  
189 annotated to potentially deleterious (CADD > 15) and non-deleterious (CADD < 15) variants as  
190 described in Methods. In total, 14,767 haplotypes were identified, spanning 714 individuals and  
191 2,475 genes (Supplemental Figure S5). We observed an overall depletion of putative high-  
192 penetrance haplotypes ( $\epsilon = -0.0156$ , Poisson-binomial test  $p = 1.006 \times 10^{-6}$ ), consistent with our  
193 hypothesis. However, we did not detect a stronger depletion for putatively deleterious rare  
194 alleles ( $p = 0.508$ , Figure 5), possibly due to the modest sample size of GTEx limiting our  
195 statistical power.

## 196 Enrichment in TOPMed

197 Next, we increased our power to detect evidence of selection against putative high penetrance  
198 haplotypes by using population-phased WGS data from 44,634 European-American ancestry  
199 individuals in 19 TOPMed cohorts, post-filtering (Methods, Supplemental Figure S5). The large  
200 sample size in TOPMed allowed us to limit the analysis to exonic variants with 10 or fewer  
201 occurrences (excluding singletons due to limitations of population-based phasing), or  $<0.0213\%$   
202 minor allele frequency. With the same set of  $\psi$ QTLs from GTEx, we identified the haplotype of  
203 38,869 rare alleles that fell in primary and secondary sExons. Across all protein-coding genes  
204 and rare alleles, we observed a modest but significant overall depletion of high penetrance  
205 haplotypes than expected ( $\epsilon = -0.0037$ , Poisson-binomial  $p = 3.43 \times 10^{-4}$ ). Haplotypes with  
206 putatively deleterious rare alleles had some indication of being more depleted than those with  
207 non-deleterious rare alleles, but not to a degree that reached statistical significance ( $p = .100$ ,  
208 Figure 5). However, we hypothesized that this result would be more pronounced in genes with  
209 stronger  $\psi$ QTLs, as well as genes known to be intolerant to loss of function variation. When  
210 focusing on genes with stronger  $\psi$ QTLs where the  $\Delta$ PSI score was in the top quartile ( $\Delta$ PSI  $>$   
211  $0.076$ ) the difference was again not significant ( $p = 0.248$ ). However, when quantifying gene  
212 constraint with LOEUF (Karczewski et al. 2020) and limiting to genes in the first quartile among  
213 sGenes (LOEUF  $< 0.460$ ), we detected a significant difference in high-penetrance haplotype  
214 depletion between the two groups ( $p = 0.048$ ), suggesting that splicing may play a greater role  
215 in modifying penetrance in genes known to be constrained. Finally, while we would expect to  
216 see the greatest effects of purifying selection among constrained genes with strong  $\psi$ QTLs, the  
217 small number of such genes limits our power and no significant association was detected ( $p =$   
218  $0.982$ ). We found that across genes in general,  $\Delta$ PSI and LOEUF were positively correlated, so  
219 genes with high  $\Delta$ PSI and low LOEUF were uncommon (Supplemental Figure S6C). While  
220 subtle, these results suggest that deleterious rare alleles are more likely to be carried on exons  
221 that are skipped due to the effects of common regulatory variants, especially in constrained  
222 genes.

223

224 Next, we wanted to explore if any genes or classes of genes drove our observation of high-  
225 penetrance haplotype depletion. To this end, using the same TOPMed data, we tested for  
226 nonrandom haplotype combinations on a gene-by-gene basis, instead of pooling haplotypes  
227 across all genes as in the previous approach. For 2,396 genes with more than 10  $\psi$ QTL-coding  
228 variant haplotypes across all available individuals, we ran a Poisson-binomial test for high-  
229 penetrance haplotype depletion (Supplemental Figure S7). We observed little signal, with



230 approximately equal numbers of genes with enrichment and depletion of high and low  
231 penetrance haplotypes. However, only 411 of the genes had more than 30 deleterious allele  
232 haplotypes, indicating that our power is quite limited. Thus, our results indicate that observing  
233 signals of modified penetrance at the gene level in population cohorts is very challenging.

234

### 235 **Genetically controlled splicing's contribution to disease gene variant penetrance**

236 In addition to studying the general population as above, we next turned to investigate  
237 nonrandom distribution of  $\psi$ QTL-coding allele haplotypes in a disease cohort: the Simons  
238 Simplex Collection (SSC) with 2,380 Autism Spectrum Disorder (ASD) simplex families. Rare  
239 coding variants are known to contribute to the etiology of ASD (Iossifov et al. 2014; Sanders et  
240 al. 2015; Sanders et al. 2012), and the large set of transmission-resolved WGS data available in  
241 the SSC make it a suitable dataset to search for haplotype patterns indicative of modified  
242 penetrance. While de novo variants also play an important role in autism risk (Iossifov et al.  
243 2014), their number is so low that we chose to focus on inherited variants.

244 First, we sought to replicate the depletion of potential high-penetrance haplotypes observed in  
245 TOPMed, using SSC parents, who are a cohort of unrelated individuals, phenotypically healthy  
246 but with potential enrichment of ASD risk variants due to having a child with ASD. We analyzed  
247 all genes with a  $\psi$ QTL in GTEx, limiting our analysis to coding alleles with 3 or fewer  
248 occurrences across all parents, and removing genes with an unusually high number of rare  
249 variant haplotypes (Supplemental Figure S5). Singleton variants were included, since their  
250 haplotype can be confidently resolved using phasing by transmission. We recapitulated the  
251 patterns observed in TOPMed, with a significant depletion of high-penetrance haplotypes with  
252 deleterious rare alleles ( $\epsilon = -0.019$ , Poisson-binomial  $p = 2.11 \times 10^{-8}$ ), with high-penetrance  
253 haplotypes carrying deleterious rare alleles more depleted than those carrying non-deleterious  
254 rare alleles (Comparison p-value = 0.042, Figure 5).

255

256 Next, we sought to analyze potential splicing modifiers of the penetrance of disease-causing  
257 alleles in SSC by focusing on rare inherited variants in ASD-implicated genes. These alleles,  
258 while potentially contributing to ASD in the proband, are also carried on the same haplotypic  
259 background by a healthy parent and often a healthy sibling. Thus, both increased or decreased  
260 penetrance  $\psi$ QTL configurations could be possible (Supplemental Figure S8) To test this, we  
261 analyzed deviation in haplotype frequencies in parents, probands, and siblings, among the 218  
262 out of the 1,010 genes implicated in ASD risk according to SFARI Gene (Banerjee-Basu and  
263 Packer) that also had a  $\psi$ QTL. No significant deviation was detected in SSC parents ( $\epsilon = -$

264 0.0278,  $p = 0.122$ ). Interestingly, across probands and unaffected siblings we found that  
265 putatively highly penetrant haplotypes with deleterious coding alleles were depleted ( $\epsilon = -0.055$   
266 &  $-0.047$ ,  $p = 0.020$  &  $0.088$  respectively). While it seems counterintuitive to see depletion of  
267 penetrant haplotypes in individuals with ASD, we reason that this penetrance reducing effect  
268 may be acting to protect parents from developing phenotypes of ASD. We find that the SFARI  
269 genes tend to be highly constrained, compared to all protein coding genes (Supplemental  
270 Figure S8B) (Neale et al. 2012), and that these same alleles were also highly depleted among  
271 unrelated individuals in TOPMed (Figure 6), further corroborating the overall observed pattern of  
272 selection for penetrance reducing haplotype combinations.

## 273 Discussion

274 In this study, we have expanded our model of *cis*-regulatory alleles as modifiers of penetrance  
275 of coding variants (Castel et al. 2018) to directly consider splice-regulatory alleles as potential  
276 additional drivers. We first show that individuals carrying potentially deleterious rare mutations  
277 at variably spliced exons tend to use those exons in transcripts less frequently. This observation  
278 could indicate that the penetrance of these rare alleles is reduced by their exclusion from  
279 transcripts. However, this approach does not reveal the reason. One approach to potentially  
280 shed light on this would be analysis of allele-specific transcript structure, but this is not possible  
281 with short read RNA-sequencing. However, our model could be tested in larger future studies  
282 with long-read sequencing technology (Glinos et al. 2021).

283  
284 Thus, we investigate common splice-regulatory variants ( $\psi$ QTLs) as potential modifiers of  
285 penetrance of rare alleles in their target exons. Across different datasets, we have  
286 demonstrated and replicated the result that high-penetrance haplotype configurations of rare  
287 alleles and  $\psi$ QTLs alleles are depleted. These findings emphasize the importance of alternative  
288 splicing as one of the many processes that regulate human traits, and suggest that splicing is  
289 involved in variable penetrance of coding variants.

290  
291 Through this research, we derived a novel approach for calculating the cumulative distribution  
292 function of the Poisson-binomial distribution, as well as a metric for evaluating a dataset's  
293 deviation from an expected distribution or difference between two data sets (the comparison  
294 test). This method is well suited for very large datasets, and has further applications in genetic  
295 and non-genetic analyses where data is expected to follow the Poisson-binomial.

296

297 While we were able to detect a genome-wide signal of nonrandom combinations of splice-  
298 associated and coding alleles, it must be noted that finding evidence of modified penetrance in  
299 population cohorts is difficult, and requires very large sample sizes. This is particularly true on  
300 an individual gene level: Even in a dataset as large as TOPMed, which contains tens of  
301 thousands of donors, few genes have reasonable statistical power to detect depletion of high-  
302 penetrance haplotype configurations individually. Furthermore, the biologically and medically  
303 important genes where variant penetrance is of most interest are also highly constrained and  
304 depleted of functional genetic variation overall, further limiting the data to test for haplotype  
305 combinations in the general population.

306

307 An alternative approach is to study regulatory variation underlying modified penetrance in  
308 disease cohorts with well annotated disease-causing variants, linking haplotype patterns with  
309 phenotype variation between and within families. The Simons Simplex Collection had some  
310 limitations in this respect: most ASD-contributing rare variants are not known and the trait is  
311 highly polygenic, making it difficult to compare penetrance of variants in the same gene between  
312 families. Furthermore, in simplex families many causal variants are *de-novo*, but their total  
313 number is small for statistical analysis. In the future, large ASD studies with multiplex families  
314 may better capture ASD instances with heritable variant etiology. Furthermore, experimental  
315 validation, for example with genome editing, may be a fruitful approach.

316

317 Overall these results suggest that depletion of high-penetrance  $\psi$ QTL - coding variant  
318 haplotypes is robust across many data sources and gene sets. However, the data does not  
319 sufficiently support the hypothesis that modified penetrance by genetically controlled splicing is  
320 a significant driver for ASD risk, but that may provide some protection in families with a known  
321 incidence of autism.

322

323 In conclusion, this study provides evidence that splice-regulatory alleles play a role in controlling  
324 the impact of rare coding alleles with putatively deleterious effects. Understanding the  
325 importance of these mechanisms will be crucial for building a holistic model of genetic  
326 contribution to human phenotypic variation. We hope that in the future the prognosis of  
327 individuals carrying rare variants will be informed by genomic context that extends beyond  
328 coding regions.

329

## 330 Methods

### 331 Data Sources

332 In this project, we utilize bulk RNA sequencing and WGS from the Genotype-Tissue Expression  
333 (GTEx) Project Version 8 (Consortium 2020), WGS from 19 cohorts included in the Trans-  
334 Omics for Precision Medicine Project freeze 8 ([https://topmed.nhlbi.nih.gov/topmed-whole-](https://topmed.nhlbi.nih.gov/topmed-whole-genome-sequencing-methods-freeze-8)  
335 [genome-sequencing-methods-freeze-8](https://topmed.nhlbi.nih.gov/topmed-whole-genome-sequencing-methods-freeze-8)) (Supplemental Table 2) and WGS from simplex families  
336 in the Simons Simplex Collection (SSC).

### 337 GTEx PSI quantification and filtering

338 Percent spliced in (PSI) was calculated from GTEx V8 RNA-seq data. We limited our analysis to  
339 18 tissues, which were chosen for their coverage of tissue diversity GTEx and their coverage of  
340 the most coding genes possible (Table S1). Exon PSI for protein-coding genes was quantified  
341 using the Integrative Pipeline for Splicing Analysis (IPSA),(Pervouchine et al. 2013; 2020) which  
342 was run on Google Cloud through Terra (<https://github.com/guigolab/ipsa-nf>). The  
343 ‘-unstranded’ flag was used during the sjcount process. Exons were defined by the modified  
344 version of Gencode annotation v26 used in GTEx V8, which collapses genes with multiple  
345 isoforms to a single isoform per gene  
346 ([https://storage.googleapis.com/gtex\\_analysis\\_v8/reference/gencode.v26.](https://storage.googleapis.com/gtex_analysis_v8/reference/gencode.v26.GRCh38.genes.gtf)  
347 [GRCh38.genes.gtf](https://storage.googleapis.com/gtex_analysis_v8/reference/gencode.v26.GRCh38.genes.gtf)).

348  
349 For downstream analyses, PSI data for each tissue was prepared by 1) removing exons with  
350 data available in less than 50% of donors and 2) removing exons with fewer than 10 unique  
351 values across all available donors (Table S1). These data were normalized for QTL mapping by  
352 randomly breaking any ties between two individuals with the same PSI at an exon, then  
353 applying inverse-normal transformation across all individuals. Filtered and normalized PSI calls  
354 were saved in BED format with start/end position corresponding to each gene’s transcription  
355 start side (TSS), which serves as a reference for where to define windows for QTL mapping.  
356 The gene containing each exon was included in the BED files for use with QTLtools’ group  
357 permutation mode.

## 358 **PSI Z-Score Analysis in GTEx**

359 We compiled a list of all exons with sufficiently variable splicing in at least one GTEx tissue, as  
360 defined in the previous step, and saved the genomic coordinates of these exons in BED format.  
361 Rare variants (gnomAD AF < .01) that fell on variably spliced exons were extracted from GTEx  
362 WGS VCFs, and were subsequently filtered to variants that appeared less than 6 and greater  
363 than 1 time. Rare variant CADD scores and annotations with respect to the relevant gene were  
364 extracted as well. Some rare variants were annotated as ‘intronic’ because CADD v1.5 uses a  
365 different annotation that in rare cases does not correspond to gencode v26. Rare variant calls  
366 from exons represented disproportionately, either due to length or to high number of variants at  
367 the exon, were removed. Threshold for removing an exon was defined as  $Q3 + 1.5 * IQR$ , where  
368  $Q3$  is the third quartile of the number of rare variants per exon, where  $IQR$  is the interquartile  
369 range of the number of rare variants per exon. For all remaining variants, we computed the PSI  
370 Z-score of the individual that carried the variant at that specific exon, across all tissues where  
371 the exon was expressed and sufficiently variable. The PSI-Z score for a particular individual  $i$  at  
372 an exon  $j$  in tissue  $k$  is calculated as  $(\psi_{ijk} - \mu_j) / \sigma_j$ , where  $\psi_{ijk}$  is an individual’s PSI level at a  
373 particular exon and tissue, and  $\mu_j$  and  $\sigma_j$  are the mean and standard deviation of PSI for an exon  
374  $j$  across all individuals with data available for that exon in tissue  $k$ . Importantly, we do not  
375 normalize PSI for this analysis, to preserve signal from exons with high PSI Z-scores.

## 376 **Primary $\psi$ QTL mapping, collapsing, and secondary $\psi$ QTL mapping**

377 For each of the 18 GTEx V8 tissue groups, QTL mapping was run on every exon that passed  
378 filtering, using all genetic variants with an allele frequency greater than 5% within 1Mb of the  
379 gene’s transcription start site. We used QTLtools (Delaneau et al. 2017) run in grouped  
380 permutation mode, with groups defined by gene. This strategy controls for correlation between  
381 exons that are part of the same gene. 15 PEER factors recalculated from normalized PSI, 5  
382 genetic principal components (PCs), as well as sex, WGS PCR batch, and sequencing platform  
383 were also included as covariates in the QTL model, as recommended in the GTEx V8 STAR  
384 methods.(Consortium 2020)

385  
386 For every exon, we selected the most significant variant, and for every gene the most significant  
387 exon. We then compiled the  $\psi$ QTL results across tissues to achieve a set of cross-tissue top  
388  $\psi$ QTLs. When a gene was significant across multiple tissues, we used the tissue where the  
389 effect size ( $\Delta$ PSI score) was the highest. This process ensured that a gene was only included

390 once in our final set of  $\psi$ QTLs, and was labeled by one variant that is associated to splicing  
391 (sVariant).

392

393 Since the splicing of multiple exons within a gene is often correlated, we implemented an  
394 approach to identify additional exons whose splicing the sVariant is associated with.

395 Consideration of multiple exons per gene is desirable because it increases the amount of  
396 genetic space where rare variant haplotypes can be identified. For each gene with a significant  
397  $\psi$ QTL, we ran a nominal QTLtools pass of just the sVariant against PSI of all other exons in the  
398 gene. We then considered secondary exons with a Bonferroni-corrected  $p < 0.05$  if QTL effect  
399 direction was the same as the top exon.

400

401 This procedure produced the final set of common variant-exon pairs used in all downstream  
402 analyses (10,901 sExons, across 5,198 sGenes). Haplotype calls from phased, filtered WGS  
403 datasets (see next section) were compiled by extracting rare variants that fell within sExons,  
404 and recording if the variant appeared on the same haplotype as the high inclusion or low  
405 inclusion  $\psi$ QTL allele. (Code available at [https://github.com/jeinson/mp\\_manuscript](https://github.com/jeinson/mp_manuscript))

## 406 **WGS filtering across datasets**

407 **Genotype Tissue Expression Project (GTEx):** Read-aware Phased WGS data was used from  
408 all 838 samples included in GTEx v8. (Consortium 2020), (Supplementary Information Section  
409 2.4) For use in haplotype calling, the following filters were applied 1) Variants were extracted with  
410 an allele frequency less than 0.005 in gnomAD, and singleton variants without read-backing to  
411 support their phase call were removed. 2) Samples from donors that did not self-identify as  
412 European American were removed. Since the  $\psi$ QTL data from GTEx is based on 85% European  
413 Americans, the sVariants selected from these data may not capture allele frequencies and  
414 haplotype structures in other ancestries, and differing numbers of rare variants across ancestries  
415 might bias the results. 3) Haplotype calls from genes represented disproportionately, either due  
416 to length or to high number of variants at the gene, were removed. Threshold for removing a gene  
417 was defined as  $Q3 + 1.5 * IQR$ , where Q3 is the third quartile of the number of haplotypes per  
418 gene, where IQR is the interquartile range of the number of haplotypes per gene.

419

420 **Trans-Omics for Precision Medicine Initiative (TOPMed):** Population-phased WGS data from  
421 donors of European-American ancestry were used from TOPMed, since this matches the  
422 population source of the sQTL data from GTEx (see above). To define individuals of European

423 ancestry, we used the approach outlined in (Morris et al. 2019). Briefly, TOPMed samples were  
424 projected onto the first 20 principal components estimated from the 1000 Genomes Phase 3  
425 (1000G) project (Auton et al. 2015) using FastPCA v2.0 (Galinsky et al. 2016). Only bi-allelic  
426 variants shared between the two datasets, and that passed a strict set of criteria (MAF >1%, minor  
427 allele count >5, genotyping call rate >95%, Hardy-Weinberg p-value >1x10<sup>-6</sup>) were used to  
428 calculate the principal components. Expectation Maximization (EM) (Chen and Maitra 2015)  
429 clustering was used to compute the probabilities of cluster membership and eigenvectors 1, 2, 5,  
430 6 and 8 were selected for efficiently separating the individuals of White European and American  
431 ancestry (subpopulation codes CEU, GBR, FIN, CEU, IBS and TSI) from other ancestry groups.  
432 Finally, eight predefined clusters were chosen for EM clustering based on sensitivity analyses.  
433 This resulted in 52,426 TOPMed individuals clustering together with the 1000G CEU, GBR, FIN,  
434 CEU, IBS and TSI subpopulation, and they were termed of White ancestry. We kept 19 cohorts  
435 ([Supplemental Table 2](#)), and 49,542 individuals, filtering out the remaining cohorts which  
436 collectively contained less than 5% of all haplotypes.

437

438 To define rare coding variants for downstream analysis, we extracted SNPs and small indels with  
439 more than 1 and 10 or fewer occurrences; singletons were removed due to unreliable population-  
440 based phasing. To account for unusually long genes, and genes with an unusually high number  
441 of rare variants, we applied the same filtering procedure as step 3 from the GTEx analysis to  
442 produce a final set of rare variant haplotypes.

443

444 **Simons Simplex Collection (SSC):** Phased WGS data was used from 2,380 families. Simplex  
445 families consist of a proband child diagnosed with Autism Spectrum Disorder (ASD), an  
446 unaffected sibling, and two unaffected parents (Turner et al. 2016). We genotype the SSC whole-  
447 genome data set (An et al. 2018; Ruzzo et al. 2019; Yoon et al. 2021) using the transmission  
448 mode of our Multinomial Genotyper (Iossifov et al. 2012) that produces only high-quality  
449 mendelian family genotypes. The whole-genome sequence and the genotype calls are available  
450 to qualified researchers through the Simons Foundation. In addition, we transmission-phased the  
451 heterozygous variants on a per-variant basis when possible, using the genotypes of both parents.  
452 Since this method is accurate for singleton variants in probands, these were included in  
453 downstream analysis.

454

455 We additionally removed genes that contained an unusually high number of rare coding variants  
 456 across parents, using the same outlier definition as in the previous two datasets. This set of  
 457 variants post-filtering were considered in siblings and probands in downstream analyses.

## 458 Haplotype calling from phased genetic data and filtering

459  $\psi$ QTL-coding allele haplotypes were generated using a similar procedure across all three  
 460 phased-resolved WGS datasets. First, all rare variants were extracted among sExons using the  
 461 filters described above, considering variants that fell in primary and secondary sExons, taking  
 462 account of the haplotype phase assignment. Then, the genotype of sVariants, and phase for  
 463 heterozygous cases, was extracted from VCFs and haplotypes were labeled as high-penetrance  
 464 ( $\beta = 1$ ) and low penetrance ( $\beta = 0$ ) according to our model for splice QTLs as a modifier of  
 465 penetrance (Figure 1).

466

## 467 Table 2: Properties of 3 WGS datasets used in this study

468 Across all datasets, we extract rare variants that fall on primary and secondary sExons.

	GTEEx	TOPMed	SSC - Parents
<b>N Donors</b>	714	44,634	4,731
<b>Phasing Method</b>	Population Based & Read backed phasing (SHAPEIT2(O'Connell et al. 2014) and PhASer (Castel et al. 2016))	Population Phasing (Eagle) (Loh et al. 2016)	Phasing by transmission
<b>Singletons included</b>	Yes, in calls with RNA-seq read backing. Otherwise, no	No	Yes
<b>Rare variant allele frequency cutoff</b>	0.5% MAF in gnomad. (No count cutoff due to the relative small size of the GTEEx WGS dataset)	Appears 10 or fewer times (i.e. 0.0257% MAF)	Appears $\leq$ 3 times (i.e. 0.126% MAF)

469



## 470 **Test for depletion of regulatory haplotypes that increase penetrance**

471 We sought to test the hypothesis that QTL-coding allele haplotype combinations are present in  
472 the population at frequencies that deviate from a baseline expectation, based on allele  
473 frequencies alone. Such a result could indicate high-penetrance haplotypes with deleterious  
474 variants being removed from the population by natural selection. The total number of high  
475 penetrance haplotypes arising from  $\psi$ QTLs with varying allele frequencies can be modeled by  
476 the Poisson-Binomial distribution, which is a generalization of the binomial distribution. While a  
477 binomial describes the sum of  $n$  independent identically distributed bernoulli random variables,  
478 the Poisson-binomial describes the sum of  $n$  independent but non-identically distributed  
479 bernoulli random variables. Therefore, the distribution must be parameterized by a vector of  
480 probabilities of length  $n$ . While we could calculate P-values using a variety of methods that  
481 obtain the CDF of the Poisson-binomial, (Hong 2013) these methods all lack a way to quantify  
482 the magnitude of the effect size. Furthermore, they measure deviation from the null but do not  
483 allow comparison of two data sets (in our case, haplotypes carrying non-deleterious and  
484 deleterious coding alleles) Therefore, we developed the following procedure that approximates  
485 the Poisson-binomial CDF. This has the advantage of generating a quantifiable effect size for  
486 deviation from the null model, as well as corresponding confidence intervals.

487  
488 Our procedure for approximating the Poisson-binomial, and subsequently testing for non-  
489 random occurrences of putative high-penetrance haplotypes, which we applied to each WGS  
490 dataset in this study, is as follows:

491  
492 For each observation of a heterozygous coding allele that falls in a sExon, let  $L$  and  $H$  represent  
493 the low and high exon inclusion  $\psi$ QTL haplotype respectively, and let  $B$  and  $b$  represent the  
494 coding variant reference and minor allele respectively. Here, we focus on rare variants, with our  
495 main interest being deleterious ones, and we here treat rare alleles as independent. Using  
496 variant phasing information, for a given haplotype  $g$ , we define an indicator function  $\beta$  which is  
497 set equal to 1, corresponding to putatively high-penetrance, if the coding allele falls on the  
498 highly included sExon, and 0 otherwise. The genotype of the major coding allele is irrelevant,  
499 and for rare variants  $b/b$  homozygotes are absent in practice.

500

501

$$\beta(g) = \begin{cases} 1 & \text{if } g \in (Hb/HB), (Hb/LB) \\ 0 & \text{if } g \in (Lb, LB), (Lb/HB) \end{cases}$$

502 Next, we define an expectation function on  $\beta$ , under the null model where observing a high-  
503 penetrance and low-penetrance haplotype are equally likely.  $\mathbb{E}[\beta(g)]$  is dependent on the  
504 heterozygosity of the  $\psi$ QTL variant in an individual. Assuming independence of rare variants, if  
505 an individual is heterozygous for a  $\psi$ QTL allele, the probability that an exonic variant will land in  
506 a high-penetrance configuration is 0.5. If an individual is homozygous for the  $\psi$ QTL allele, the  
507 probability that the exonic variant will land in a high-penetrance configuration is dependent on  
508 the  $\psi$ QTL's allele frequency.

509

$$\mathbb{E}[\beta(g)] = \begin{cases} 0.5 & \text{if } g \in (L/H) \\ (n(H/H) + 1)/(n(H/H) + n(L/L)) & \text{if } g \in (L/L), (H/H) \end{cases}$$

510  
511

512 We define the expectation of observing a homozygous  $\psi$ QTL allele as the proportion of high  
513 inclusion  $\psi$ QTL homozygotes in the dataset, plus a pseudo-count, to avoid getting an  
514 expectation of 0 in datasets where the low inclusion allele is much more common. This method  
515 does not assume Hardy-Weinberg equilibrium for the  $\psi$ QTL allele, but requires that the  
516 proportion of homozygotes for the two alleles be recalculated on each dataset. This approach  
517 was used for the GTEx and TOPMed analyses. Alternatively, the expectation of  $\beta$  under the null  
518 model can also be calculated as follows:

519

$$\mathbb{E}[\beta(g)] = \begin{cases} 0.5 & \text{if } g \in (L/H) \\ f(H)^2/(f(H)^2 + (1 - f(H))^2) & \text{if } g \in (H/H), (L/L) \end{cases}$$

520  
521

522 Where  $f(H)$  is the population frequency of the high exon inclusion  $\psi$ QTL allele. We took this  
523 approach for haplotypes from SSC, where counting alleles across the whole dataset was  
524 infeasible due to the structure of the dataset, and used  $\psi$ QTL allele frequencies from gnomad  
525 3.0 (Karczewski et al. 2020).

526

527 The function  $\beta$  is evaluated across all individuals, sGenes, and rare variants in sExons in a  
528 dataset. The average observed deviation from the expected totals of high and low penetrance  
529 haplotypes ( $\varepsilon$ ) is calculated as follows:

530

$$\varepsilon = \frac{1}{N} \sum_{n=1}^N \beta(g_n) - \mathbb{E}[\beta(g_n)]$$

531  
532

533 where  $N$  is the total number of considered haplotypes.  $\epsilon$  can be interpreted as the effect size of  
534 depletion/enrichment of high-penetrance haplotypes in the dataset such that  $\epsilon < 0$  would  
535 indicate a depletion of high-penetrance haplotypes.

536

537 We quantify the significance of  $\epsilon$  by bootstrapping all haplotypes, generating 95% confidence  
538 intervals and drawing two-sided empirical  $P$ -values as

539

$$P(H_0) = 2 \min \left[ \frac{\sum_{b=1}^B \epsilon_b < 0}{B}, \frac{\sum_{b=1}^B \epsilon_b > 0}{B} \right]$$

540

541

542 where  $B$  is the total number of bootstraps. In practice, we found that 1,000 bootstraps was  
543 enough to accurately approximate the Poisson-binomial distribution, while managing runtime.

544

545 Although the test was designed for counts of haplotypes, this approach is generalizable to any  
546 system that can be modeled by a Poisson-binomial distribution. Therefore, to benchmark our  
547 test, we simulated data from several theoretical allele frequency distributions by sampling from  
548 beta distributions with various shape parameters, including one distribution where its  
549 parameters were estimated direction from our set of  $\psi$ QTLs from GTEx using the method of  
550 moments estimator (Figure 3, Supplemental Figure 4). We found that our bootstrapping  
551 procedure accurately approximated the Poisson-binomial distribution for all inputs tested.

552 However, the magnitude of  $\epsilon$  - but not direction - is dependent on the shape of the theoretical  
553 allele frequency distribution, so comparing magnitudes of  $\epsilon$  across distinct datasets should be  
554 done with caution. The accuracy of our method increased with larger sample sizes. Therefore,  
555 we recommend using this approach when handling data where  $N > 1,000$  (Supplemental Figure  
556 S4).

557

558 As an extension to this procedure, we can also conveniently calculate the significance of a  
559 difference in  $\epsilon$  between two similar datasets  $A$  and  $B$ , for example, between haplotypes where  
560 the rare variant is putatively deleterious vs. haplotypes where the rare variant is non-deleterious:

561

$$\epsilon_{comp} = \left( \frac{1}{N_A} \sum_{n=1}^{N_A} \beta(g_{A_n}) - \mathbb{E}[\beta(g_{A_n})] \right) - \left( \frac{1}{N_B} \sum_{n=1}^{N_B} \beta(g_{B_n}) - \mathbb{E}[\beta(g_{B_n})] \right)$$

562

563 We then apply the bootstrapping procedure as in the standard case, and draw P-values  
564 accordingly. The corresponding P-value from this procedure is referred to as the “comparison  
565 test” in the main text.

566

567 This test is implemented in the S**T**atistic for Modified P**E**Netrance (STAMPEN) R package  
568 that is available to download here (<https://github.com/jeinson/stampen>)

## 569 Data Availability

570 All code used to perform analyses and generate figures is available at  
571 [https://github.com/jeinson/mp\\_manuscript](https://github.com/jeinson/mp_manuscript). Qualified researchers requiring data access can  
572 apply for GTEx, and TOPMed data through dbGaP, and SSC data through the Simons  
573 foundation. We include a function to generate simulated data in the stampen R package  
574 (<https://github.com/jeinson/stampen>). PSI and  $\psi$ QTLs from GTEx v8 can be download from the  
575 repository for (Einson et al. 2022) at <https://zenodo.org/record/7275062#.Y9gc0OzMJf0>

## 576 Acknowledgements

577 J.E. thanks members of the Lappalainen lab for thoughtful discussions and feedback throughout  
578 this project.

579 Molecular data for the Trans-Omics in Precision Medicine (TOPMed) program was supported by  
580 the National Heart, Lung and Blood Institute (NHLBI). Whole genome sequencing (WGS) for the  
581 Trans-Omics in Precision Medicine (TOPMed) program was supported by the National Heart,  
582 Lung and Blood Institute (NHLBI). Core support including centralized genomic read mapping  
583 and genotype calling, along with variant quality metrics and filtering were provided by the  
584 TOPMed Informatics Research Center (3R01HL-117626-02S1; contract HHSN268201800002I).  
585 Core support including phenotype harmonization, data management, sample-identity QC, and  
586 general program coordination were provided by the TOPMed Data Coordinating Center  
587 (R01HL-120393; U01HL-120393; contract HHSN268201800001I) and TOPMed MESA Multi-  
588 Omics (HHSN2682015000031/HSN26800004).

589 Cohort specific acknowledgements for the 19 TOPMed cohorts used in this study are included  
590 in [Supplemental Table 2](#). The content is solely the responsibility of the authors and does not  
591 necessarily represent the official views of the National Institutes of Health.

## 592 **Funding and Sequencing Center Information**

- 593 1. Genome Sequencing for NHLBI TOPMed: Women's Health Initiative (phs001237) was  
594 performed at Broad Institute Genomics Platform (HHSN268201500014C).
- 595 2. Genome Sequencing for NHLBI TOPMed: Genetic Epidemiology of COPD Study  
596 (phs000951) was performed at Northwest Genomics Center (3R01HL089856-08S1).
- 597 3. Genome Sequencing for NHLBI TOPMed: Atherosclerosis Risk in Communities Study  
598 VTE cohort (phs001211) was performed at Baylor College of Medicine Human Genome  
599 Sequencing Center (3U54HG003273-12S2 / HHSN268201500015C).
- 600 4. Genome Sequencing for NHLBI TOPMed: Framingham Heart Study (phs000974) was  
601 performed at Broad Institute Genomics Platform (HHSN268201600034I).
- 602 5. Genome Sequencing for NHLBI TOPMed: My Life, Our Future: Genotyping for Progress  
603 in Hemophilia (phs001515) was performed at Baylor College of Medicine Human  
604 Genome Sequencing Center (HHSN268201600033I).
- 605 6. Genome Sequencing for NHLBI TOPMed: Mount Sinai BioMe Biobank (phs001644) was  
606 performed at McDonnell Genome Institute (3UM1HG008853-01S2).
- 607 7. Genome Sequencing for NHLBI TOPMed: Cardiovascular Health Study (phs001368)  
608 was performed at Broad Institute Genomics Platform (HHSN268201600034I).
- 609 8. Genome Sequencing for NHLBI TOPMed: Multi-Ethnic Study of Atherosclerosis  
610 (phs001416) was performed at Broad Institute Genomics Platform  
611 (HHSN268201600034I, 3U54HG003067-13S1).
- 612 9. Genome Sequencing for NHLBI TOPMed: Coronary Artery Risk Development in Young  
613 Adults (phs001612) was performed at Baylor College of Medicine Human Genome  
614 Sequencing Center (HHSN268201600033I).
- 615 10. Genome Sequencing for NHLBI TOPMed: Mayo Clinic Venous Thromboembolism Study  
616 (phs001402) was performed at Baylor College of Medicine Human Genome Sequencing  
617 Center (3U54HG003273-12S2 / HHSN268201500015C).
- 618 11. Genome Sequencing for NHLBI TOPMed: Lung Tissue Research Consortium  
619 (phs001662) was performed at Broad Institute Genomics Platform  
620 (HHSN268201600034I).

- 621 12. Genome Sequencing for NHLBI TOPMed: The Vanderbilt University BioVU Atrial  
622 Fibrillation Genetics Study (phs001624) was performed at Baylor College of Medicine  
623 Human Genome Sequencing Center (3UM1HG008898-01S3).
- 624 13. Genome Sequencing for NHLBI TOPMed: Vanderbilt Genetic Basis of Atrial Fibrillation  
625 (phs001032) was performed at Broad Institute Genomics Platform (3R01HL092577-  
626 06S1).
- 627 14. Genome Sequencing for NHLBI TOPMed: Hispanic Community Health Study - Study of  
628 Latinos (phs001395) was performed at Baylor College of Medicine Human Genome  
629 Sequencing Center (HHSN268201600033I).
- 630 15. Genome Sequencing for NHLBI TOPMed: Severe Asthma Research Program  
631 (phs001446) was performed at New York Genome Center Genomics  
632 (HHSN268201500016C).
- 633 16. Genome Sequencing for NHLBI TOPMed: Massachusetts General Hospital Atrial  
634 Fibrillation Study (phs001062) was performed at Broad Institute Genomics Platform  
635 (3U54HG003067-12S2 / 3U54HG003067-13S1; 3U54HG003067-12S2 /  
636 3U54HG003067-13S1; 3UM1HG008895-01S2).
- 637 17. Genome Sequencing for NHLBI TOPMed: Heart and Vascular Health Study  
638 (phs000993) was performed at Broad Institute Genomics Platform (3R01HL092577-  
639 06S1).
- 640 18. Genome Sequencing for NHLBI TOPMed: Groningen Genetics of Atrial Fibrillation Study  
641 (phs001725) was performed at Baylor College of Medicine Human Genome Sequencing  
642 Center (3UM1HG008898-01S3).
- 643 19. Genome Sequencing for NHLBI TOPMed: Genetics of Cardiometabolic Health in the  
644 Amish (phs000956) was performed at Broad Institute Genomics Platform  
645 (3R01HL121007-01S1).

646 J.E and TL were supported by NIH grants R01GM122924, R01MH106842. P.M. was supported  
647 by NIGMS grant R01GM140287. I.I. was supported by the Simons Center for Quantitative  
648 Biology at Cold Spring Harbor Laboratory, SFARI Grants SF497800, SF677963, SF666590, and  
649 the Centers for Common Disease Genomics grant (UM1 HG008901). Support for title page  
650 creation and format was provided by AuthorArranger, a tool developed at the National Cancer  
651 Institute.

## 652 Conflict Statement

653 T.L. is a paid advisor to GSK, Pfizer, Goldfinch Bio and Variant Bio, and has equity in Variant  
654 Bio.

## 655 References

- 656 Alasoo K, Rodrigues J, Danesh J, Freitag DF, Paul DS, Gaffney DJ. Genetic effects on  
657 promoter usage are highly context-specific and contribute to complex traits. Parker S, McCarthy  
658 MI, editors. eLife. eLife Sciences Publications, Ltd; 2019 Jan 8;8:e41673.
- 659 An J-Y, Lin K, Zhu L, Werling DM, Dong S, Brand H, et al. Genome-wide de novo risk score  
660 implicates promoter variation in autism spectrum disorder. Science. American Association for  
661 the Advancement of Science; 2018 Dec 14;362(6420):eaat6576.
- 662 Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, et al. A global  
663 reference for human genetic variation. Nature. Nature Publishing Group; 2015  
664 Oct;526(7571):68–74.
- 665 Banerjee-Basu S, Packer A. SFARI Gene: an evolving database for the autism research  
666 community | Disease Models & Mechanisms | The Company of Biologists [Internet]. [cited 2022  
667 Aug 2]. Available from: [https://journals.biologists.com/dmm/article/3/3-4/133/2349/SFARI-Gene-](https://journals.biologists.com/dmm/article/3/3-4/133/2349/SFARI-Gene-an-evolving-database-for-the-autism)  
668 [an-evolving-database-for-the-autism](https://journals.biologists.com/dmm/article/3/3-4/133/2349/SFARI-Gene-an-evolving-database-for-the-autism)
- 669 Castel SE, Cervera A, Mohammadi P, Aguet F, Reverter F, Wolman A, et al. Modified  
670 penetrance of coding variants by cis-regulatory variation contributes to disease risk. Nat Genet.  
671 2018 Sep;50(9):1327–34.
- 672 Castel SE, Mohammadi P, Chung WK, Shen Y, Lappalainen T. Rare variant phasing and  
673 haplotypic expression from RNA sequencing with phASER. Nature Communications. 2016  
674 08/online;7:12817.
- 675 Chen R, Shi L, Hakenberg J, Naughton B, Sklar P, Zhang J, et al. Analysis of 589,306 genomes  
676 identifies individuals resilient to severe Mendelian childhood diseases. Nature Biotechnology.  
677 Nature Publishing Group; 2016 May;34(5):531–8.
- 678 Chen W-C, Maitra R. R: EM Algorithm for Model-Based Clustering of Finite Mixture Gaussian  
679 Distribution [Internet]. 2015 [cited 2022 Jun 24]. Available from: [https://search.r-](https://search.r-project.org/CRAN/refmans/EMCluster/html/00Index.html)  
680 [project.org/CRAN/refmans/EMCluster/html/00Index.html](https://search.r-project.org/CRAN/refmans/EMCluster/html/00Index.html)
- 681 Chiang AH, Chang J, Wang J, Vitkup D. Exons as units of phenotypic impact for truncating  
682 mutations in autism. Mol Psychiatry. 2021 May;26(5):1685–95.
- 683 Consortium TGte. The GTEx Consortium atlas of genetic regulatory effects across human  
684 tissues. Science. American Association for the Advancement of Science; 2020 Sep  
685 11;369(6509):1318–30.
- 686 Cooper DN, Krawczak M, Polychronakos C, Tyler-Smith C, Kehrer-Sawatzki H. Where  
687 genotype is not predictive of phenotype: towards an understanding of the molecular basis of  
688 reduced penetrance in human inherited disease. Hum Genet. 2013 Oct 1;132(10):1077–130.
- 689 Cummings BB, Karczewski KJ, Kosmicki JA, Seaby EG, Watts NA, Singer-Berk M, et al.  
690 Transcript expression-aware annotation improves rare variant interpretation. Nature. Nature  
691 Publishing Group; 2020 May;581(7809):452–8.



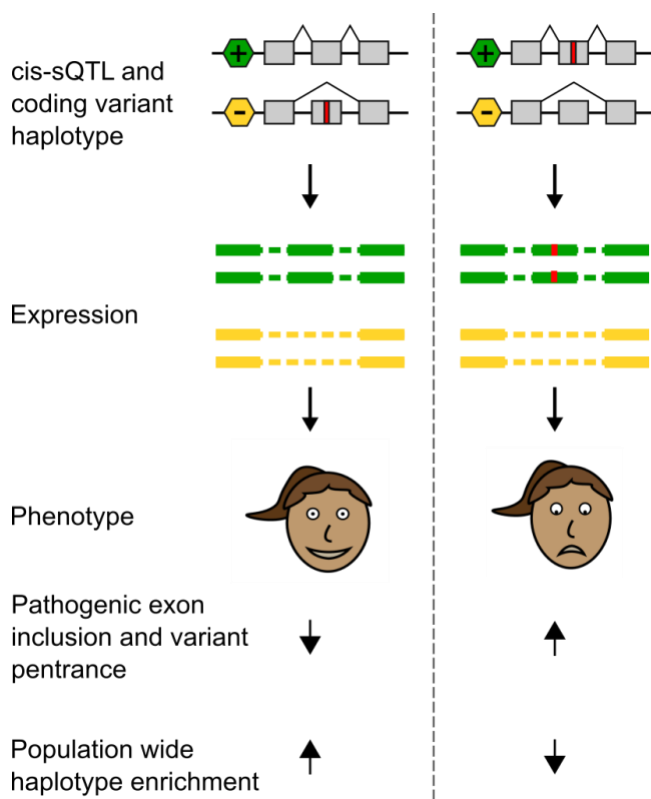
- 692 Delaneau O, Ongen H, Brown AA, Fort A, Panousis NI, Dermitzakis ET. A complete tool set for  
693 molecular QTL discovery and analysis. *Nat Commun*. 2017 May 18;8(1):1–7.
- 694 Dewey FE, Murray MF, Overton JD, Habegger L, Leader JB, Fetterolf SN, et al. Distribution and  
695 clinical impact of functional variants in 50,726 whole-exome sequences from the DiscovEHR  
696 study. *Science*. American Association for the Advancement of Science; 2016 Dec  
697 23;354(6319):aaf6814.
- 698 Einson J, Minaeva M, Rafi F, Lappalainen T. The impact of genetically controlled splicing on  
699 exon inclusion and protein structure [Internet]. *bioRxiv*; 2022 [cited 2022 Dec 22]. p.  
700 2022.12.05.518915. Available from:  
701 <https://www.biorxiv.org/content/10.1101/2022.12.05.518915v1>
- 702 Emison ES, McCallion AS, Kashuk CS, Bush RT, Grice E, Lin S, et al. A common sex-  
703 dependent mutation in a RET enhancer underlies Hirschsprung disease risk. *Nature*. 2005  
704 Apr;434(7035):857–63.
- 705 Fahed AC, Wang M, Homburger JR, Patel AP, Bick AG, Neben CL, et al. Polygenic background  
706 modifies penetrance of monogenic variants for tier 1 genomic conditions. *Nature*  
707 *Communications*. Nature Publishing Group; 2020 Aug 20;11(1):3635.
- 708 Forrest IS, Chaudhary K, Vy HMT, Petrazzini BO, Bafna S, Jordan DM, et al. Population-Based  
709 Penetrance of Deleterious Clinical Variants. *JAMA*. 2022 Jan 25;327(4):350–9.
- 710 Galinsky KJ, Bhatia G, Loh P-R, Georgiev S, Mukherjee S, Patterson NJ, et al. Fast Principal-  
711 Component Analysis Reveals Convergent Evolution of ADH1B in Europe and East Asia. *The*  
712 *American Journal of Human Genetics*. 2016 Mar 3;98(3):456–72.
- 713 Garrido-Martín D, Borsari B, Calvo M, Reverter F, Guigó R. Identification and analysis of  
714 splicing quantitative trait loci across multiple tissues in the human genome. *Nat Commun*. 2021  
715 Feb 1;12(1):727.
- 716 Gettler K, Levantovsky R, Moscati A, Giri M, Wu Y, Hsu N-Y, et al. Common and Rare Variant  
717 Prediction and Penetrance of IBD in a Large, Multi-ethnic, Health System-based Biobank  
718 Cohort. *Gastroenterology*. 2021 Apr;160(5):1546–57.
- 719 Glinos DA, Garborcauskas G, Hoffman P, Ehsan N, Jiang L, Gokden A, et al. Transcriptome  
720 variation in human tissues revealed by long-read sequencing [Internet]. *bioRxiv*; 2021 [cited  
721 2022 May 31]. p. 2021.01.22.427687. Available from:  
722 <https://www.biorxiv.org/content/10.1101/2021.01.22.427687v1>
- 723 González J, Wiberg M, von Davier AA. A Note on the Poisson's Binomial Distribution in Item  
724 Response Theory. *Applied Psychological Measurement*. SAGE Publications Inc; 2016 Jun  
725 1;40(4):302–10.
- 726 Hong Y. On computing the distribution function for the Poisson binomial distribution.  
727 *Computational Statistics & Data Analysis*. 2013 Mar;59:41–51.
- 728 Iossifov I, O’Roak BJ, Sanders SJ, Ronemus M, Krumm N, Levy D, et al. The contribution of de  
729 novo coding mutations to autism spectrum disorder. *Nature*. 2014 Nov 13;515(7526):216–21.

- 730 lossifov I, Ronemus M, Levy D, Wang Z, Hakker I, Rosenbaum J, et al. De Novo Gene  
731 Disruptions in Children on the Autistic Spectrum. *Neuron*. 2012 Apr 26;74(2):285–99.
- 732 Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The mutational  
733 constraint spectrum quantified from variation in 141,456 humans. *Nature*. Nature Publishing  
734 Group; 2020 May;581(7809):434–43.
- 735 Keren H, Lev-Maor G, Ast G. Alternative splicing and evolution: diversification, exon definition  
736 and function. *Nat Rev Genet*. Nature Publishing Group; 2010 May;11(5):345–55.
- 737 Kerimov N, Hayhurst JD, Manning JR, Walter P, Kolberg L, Peikova K, et al. eQTL Catalogue: a  
738 compendium of uniformly processed human gene expression and splicing QTLs. *bioRxiv*. 2020  
739 Jan 29;2020.01.29.924266.
- 740 Li YI, Geijn B van de, Raj A, Knowles DA, Petti AA, Golan D, et al. RNA splicing is a primary link  
741 between genetic variation and disease. *Science*. American Association for the Advancement of  
742 Science; 2016 Apr 29;352(6285):600–4.
- 743 Loh P-R, Danecek P, Palamara PF, Fuchsberger C, A Reshef Y, K Finucane H, et al.  
744 Reference-based phasing using the Haplotype Reference Consortium panel. *Nat Genet*. Nature  
745 Publishing Group; 2016 Nov;48(11):1443–8.
- 746 Maya I, Sharony R, Yacobson S, Kahana S, Yeshaya J, Tenne T, et al. When genotype is not  
747 predictive of phenotype: implications for genetic counseling based on 21,594 chromosomal  
748 microarray analysis examinations. *Genet Med*. 2018 Jan;20(1):128–31.
- 749 Milne RL, Antoniou AC. Genetic modifiers of cancer risk for BRCA1 and BRCA2 mutation  
750 carriers. *Annals of Oncology*. Elsevier; 2011 Jan 1;22:i11–7.
- 751 Morris JA, Kemp JP, Youlten SE, Laurent L, Logan JG, Chai RC, et al. An atlas of genetic  
752 influences on osteoporosis in humans and mice. *Nat Genet*. Nature Publishing Group; 2019  
753 Feb;51(2):258–66.
- 754 Neale BM, Kou Y, Liu L, Ma'ayan A, Samocha KE, Sabo A, et al. Patterns and rates of exonic  
755 de novo mutations in autism spectrum disorders. *Nature*. 2012 May;485(7397):242–5.
- 756 Noble JD, Balmant KM, Dervinis C, de los Campos G, Resende MFRJ, Kirst M, et al. The  
757 Genetic Regulation of Alternative Splicing in *Populus deltoides*. *Front. Plant Sci.* [Internet].  
758 *Frontiers*; 2020 [cited 2020 Sep 14];11. Available from:  
759 <https://www.frontiersin.org/articles/10.3389/fpls.2020.00590/full>
- 760 O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, et al. A General Approach  
761 for Haplotype Phasing across the Full Spectrum of Relatedness. *PLOS Genetics*. Public Library  
762 of Science; 2014 Apr 17;10(4):e1004234.
- 763 Ongen H, Dermitzakis ET. Alternative Splicing QTLs in European and African Populations. *Am J*  
764 *Hum Genet*. 2015 Oct 1;97(4):567–75.
- 765 Pervouchine DD, Knowles DG, Guigó R. Intron-centric estimation of alternative splicing from  
766 RNA-seq data. *Bioinformatics*. 2013 Jan 15;29(2):273–4.

- 767 Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. CADD: predicting the  
768 deleteriousness of variants throughout the human genome. *Nucleic Acids Research*. 2019 Jan  
769 8;47(D1):D886–94.
- 770 Ruzzo EK, Pérez-Cano L, Jung J-Y, Wang L, Kashef-Haghighi D, Hartl C, et al. Inherited and  
771 De Novo Genetic Risk for Autism Impacts Shared Networks. *Cell*. 2019 Aug 8;178(4):850-  
772 866.e26.
- 773 Sanders SJ, He X, Willsey AJ, Ercan-Sencicek AG, Samocha KE, Cicek AE, et al. Insights into  
774 Autism Spectrum Disorder Genomic Architecture and Biology from 71 Risk Loci. *Neuron*. 2015  
775 Sep 23;87(6):1215–33.
- 776 Sanders SJ, Murtha MT, Gupta AR, Murdoch JD, Raubeson MJ, Willsey AJ, et al. De novo  
777 mutations revealed by whole-exome sequencing are strongly associated with autism. *Nature*.  
778 2012 May;485(7397):237–41.
- 779 Shawky RM. Reduced penetrance in human inherited disease. *Egyptian Journal of Medical*  
780 *Human Genetics*. 2014 Apr 1;15(2):103–11.
- 781 Turner TN, Hormozdiari F, Duyzend MH, McClymont SA, Hook PW, Iossifov I, et al. Genome  
782 Sequencing of Autism-Affected Families Reveals Disruption of Putative Noncoding Regulatory  
783 DNA. *The American Journal of Human Genetics*. 2016 Jan 7;98(1):58–74.
- 784 Wang YH. ON THE NUMBER OF SUCCESSES IN INDEPENDENT TRIALS. *Statistica Sinica*.  
785 Institute of Statistical Science, Academia Sinica; 1993;3(2):295–312.
- 786 Yoon S, Munoz A, Yamrom B, Lee Y, Andrews P, Marks S, et al. Rates of contributory de novo  
787 mutation in high and low-risk autism families. *Commun Biol*. Nature Publishing Group; 2021 Sep  
788 1;4(1):1–10.
- 789 IPSA-nf [Internet]. Guigo Lab; 2020 [cited 2021 Aug 3]. Available from:  
790 <https://github.com/guigolab/ipsa-nf>
- 791
- 792
- 793

## 794 Figures

795



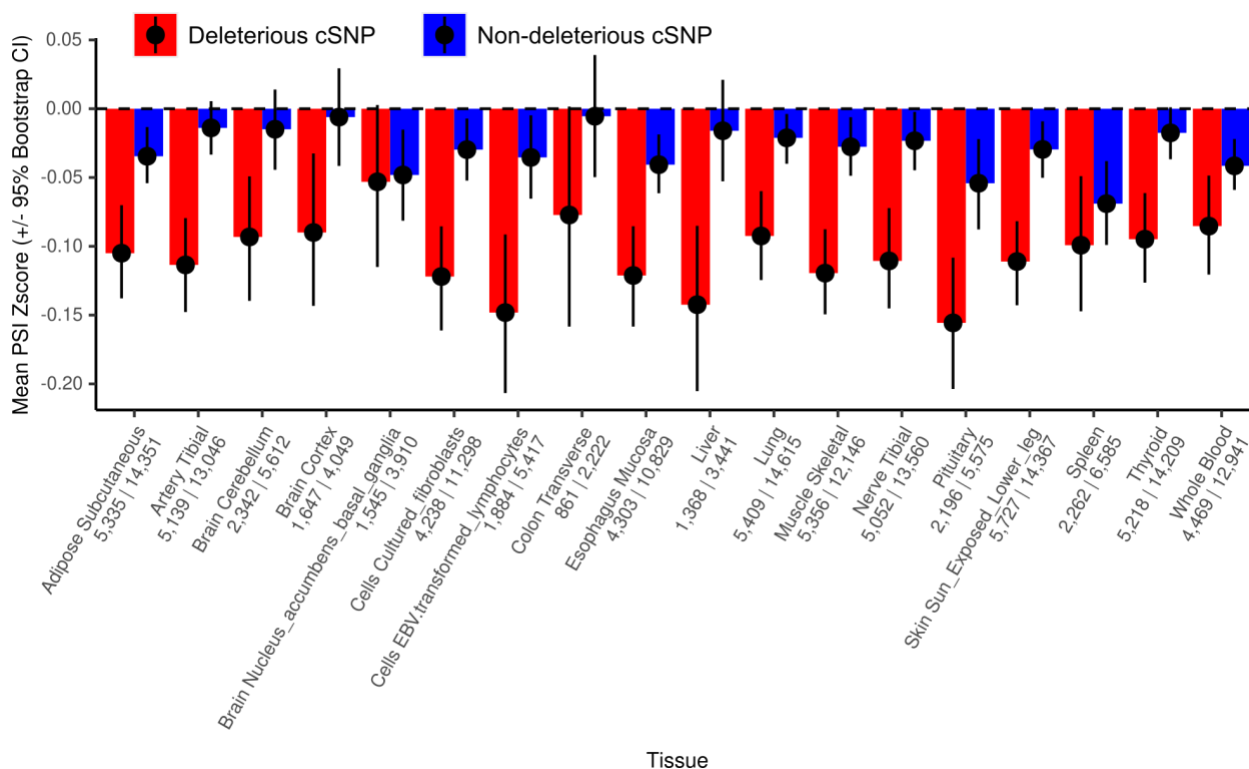
796

### 797 **Figure 1. Splice-regulatory variants as modifiers of penetrance hypothesis**

798 The hypothesis of this study is illustrated with an example of an individual who is heterozygous  
799 for both a  $\psi$ QTL and a coding variant. The two possible haplotype configurations result in either  
800 a reduced or increased penetrance state of the coding allele, depending if the allele is on the  
801 more lowly or highly included exon respectively. We predict that natural selection would deplete  
802 those that fall in a high penetrance configuration in the general population. See Supplementary  
803 Figure S1 for a quantitative description of the model.

804

805



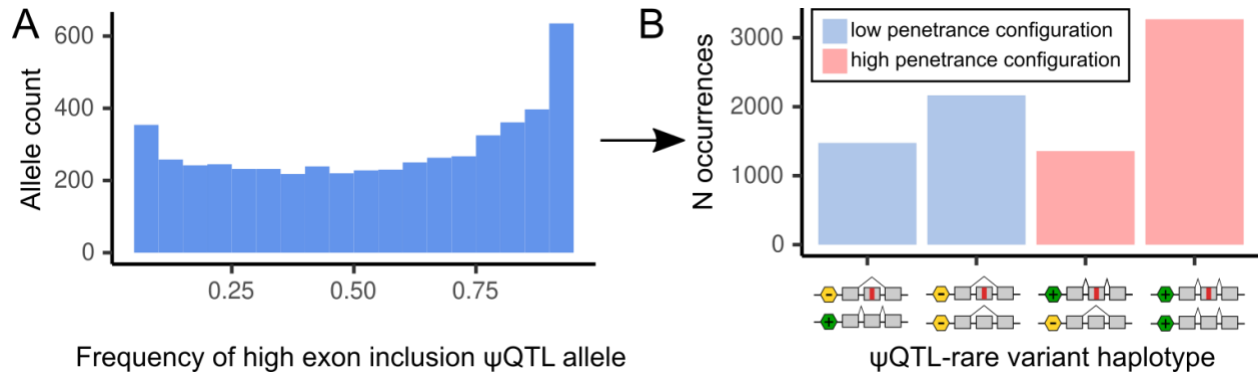
806

807 **Figure 2: Mean PSI Z-scores across tissues**

808 Mean decrease in PSI Z-scores among individuals carrying rare alleles at variably spliced exons  
 809 across 18 GTEx tissues, split by deleterious (CADD > 15) and non-deleterious (CADD < 15)  
 810 rare variants. The number of deleterious and non-deleterious alleles respectively are printed  
 811 below each tissue name. Error bars represent 95% bootstrapped confidence intervals.

812

813



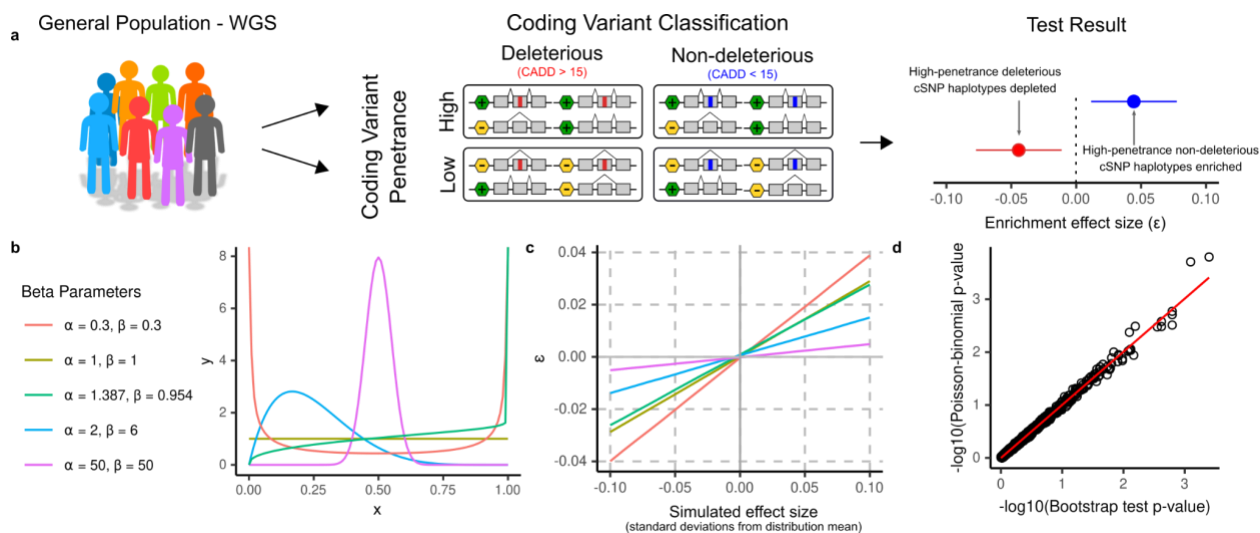
814

Frequency of high exon inclusion  $\psi$ QTL allele

815 **Figure 3:  $\psi$ QTL high inclusion allele frequencies and haplotype counts in GTEx.**

816 A. Distribution of allele frequencies for  $\psi$ QTLs that lead to higher exon inclusion. High inclusion  
817  $\psi$ QTL allele frequencies are skewed to the right, meaning  $\psi$ QTLs that include their target exon  
818 are more common in the general population. B. As a result of the nonuniform frequency  
819 distribution of high inclusion sQTL alleles, we expect to see more high penetrance haplotype  
820 configurations in general. This motivates the necessity to design a test that accounts for this  
821 difference.

822

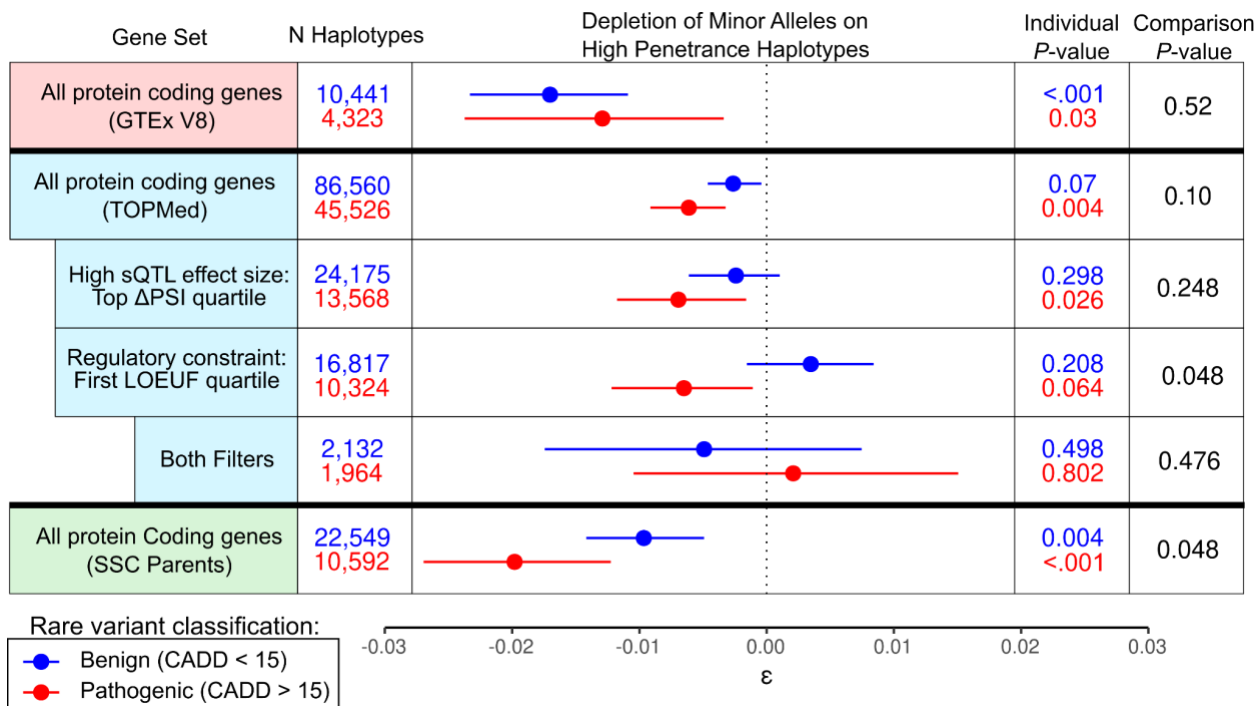


823

824 **Figure 4: The Poisson-binomial distribution models haplotype configuration counts**

825 **a.** We use phased variant calls from WGS across large populations to test for deviation in the  
 826 frequencies of  $\psi$ QTL-coding variant haplotype configurations. The magnitude and effect  
 827 direction of deviation, which we call  $\epsilon$ , is calculated using a procedure described in Methods.  
 828 The magnitude of  $\epsilon$  - but importantly not its direction - depends on the underlying  $\psi$ QTL allele  
 829 frequency distribution, as the probability of observing a high penetrance haplotype is dependent  
 830 on the  $\psi$ QTL allele frequency at each gene. Counts of highly penetrant haplotypes are modeled  
 831 by the Poisson-Binomial distribution. When running our test, we frequently divide haplotypes  
 832 into those with deleterious ( $CADD > 15$ ) and non-deleterious ( $CADD < 15$ ) coding variants,  
 833 which serve as a negative control where we do not expect to see evidence of purifying  
 834 selection. **b.** To verify that our test captures deviations from the null under any theoretical allele  
 835 frequency distribution, we simulated datasets by drawing samples from various Beta  
 836 distributions with different parameters. The Beta is defined by shape parameters  $\alpha$  and  $\beta$ . The  
 837 parameters  $\alpha = 1.387$  and  $\beta = 0.954$  were estimated from the high-inclusion  $\psi$ QTL allele  
 838 frequency distribution in GTEx using the Method of Moments estimator. **c.** We benchmarked our  
 839 test by simulating data from distributions with increasingly larger deviations from the expected  
 840 mean, in order to test how the magnitude of  $\epsilon$  differs depending on the input distribution. This  
 841 diagram can be used as a reference for how to interpret the magnitude of epsilon, given a  
 842 dataset's underlying probability distribution **d.** P-values from a simulated dataset of haplotypes  
 843 from 1,000 individuals across 1,000 genes, with  $\psi$ QTL allele frequencies matching those in  
 844 GTEx. We find that our method accurately replicates the results from the Poisson-binomial  
 845 distribution, calculated using the 'poibin' (Hong 2013) R package.

846



847

848 **Figure 5: Rare alleles carried in predicted high penetrance  $\psi$ QTL configurations in GTEx,**  
 849 **TOPMed, and SSC Parents**

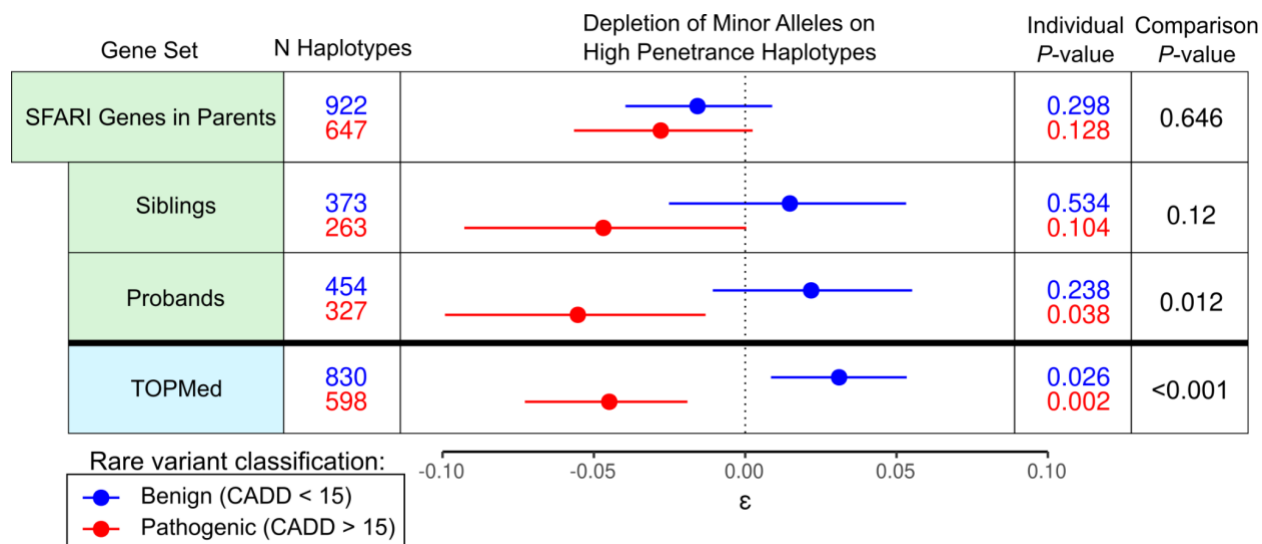
850 We tested for deviation in the frequencies of coding allele -  $\psi$ QTL configurations across all  
 851 protein coding genes with a significant  $\psi$ QTL. A negative value of  $\epsilon$  indicates fewer haplotypes  
 852 than expected given the population's  $\psi$ QTL allele frequencies. Individual  $p$ -values and 95%  
 853 confidence intervals were generated using our approximation of the Poisson-binomial cdf, with  
 854 1,000 bootstraps. Comparison  $P$ -values were generated with 1,000 bootstraps.

855

856

857





858

859 **Figure 6:  $\psi$ QTL haplotype configurations in Autism Spectrum Disorder implicated genes**  
 860 **in ASD families.**

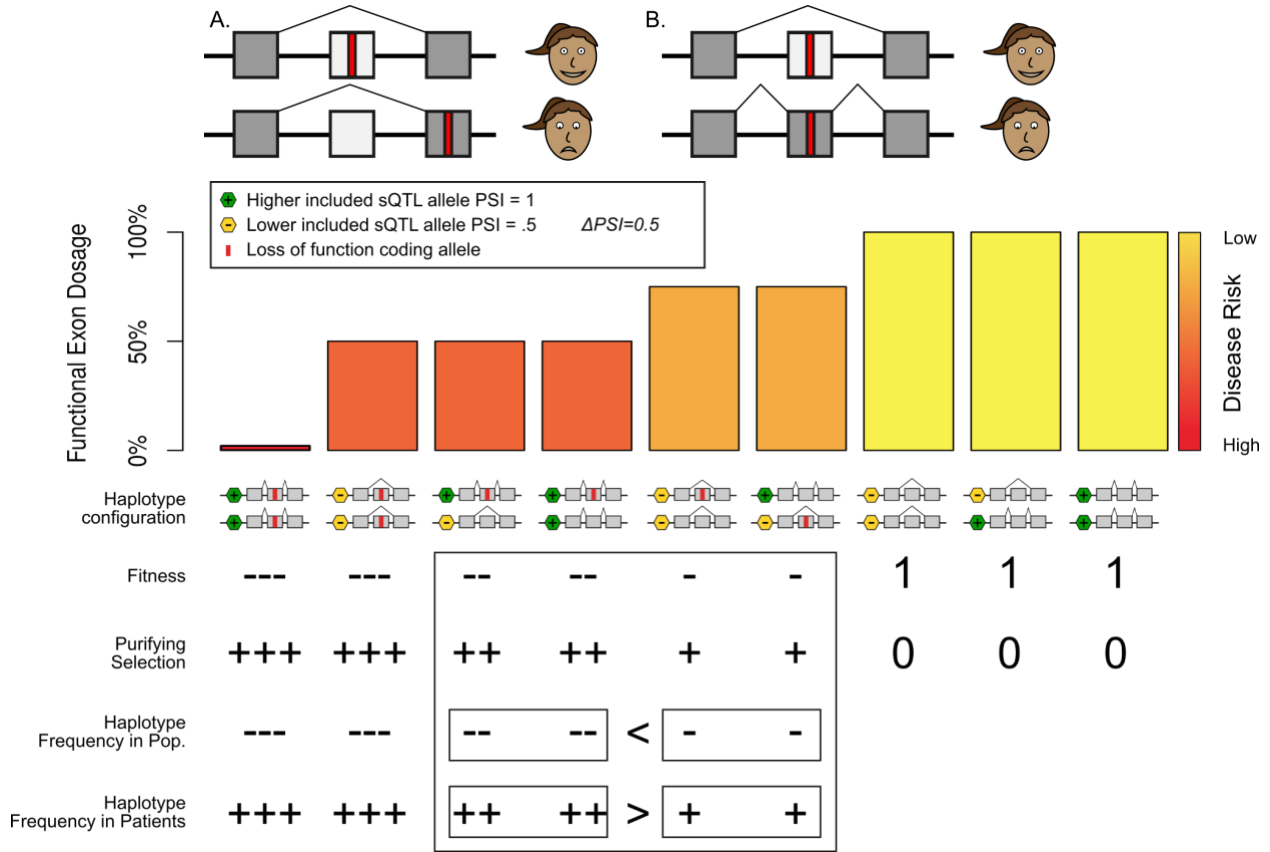
861 We tested for deviation in the frequencies of high penetrance variant -  $\psi$ QTL configurations in  
 862 ASD-implicated genes in probands and unaffected siblings in SSC families.

863

864

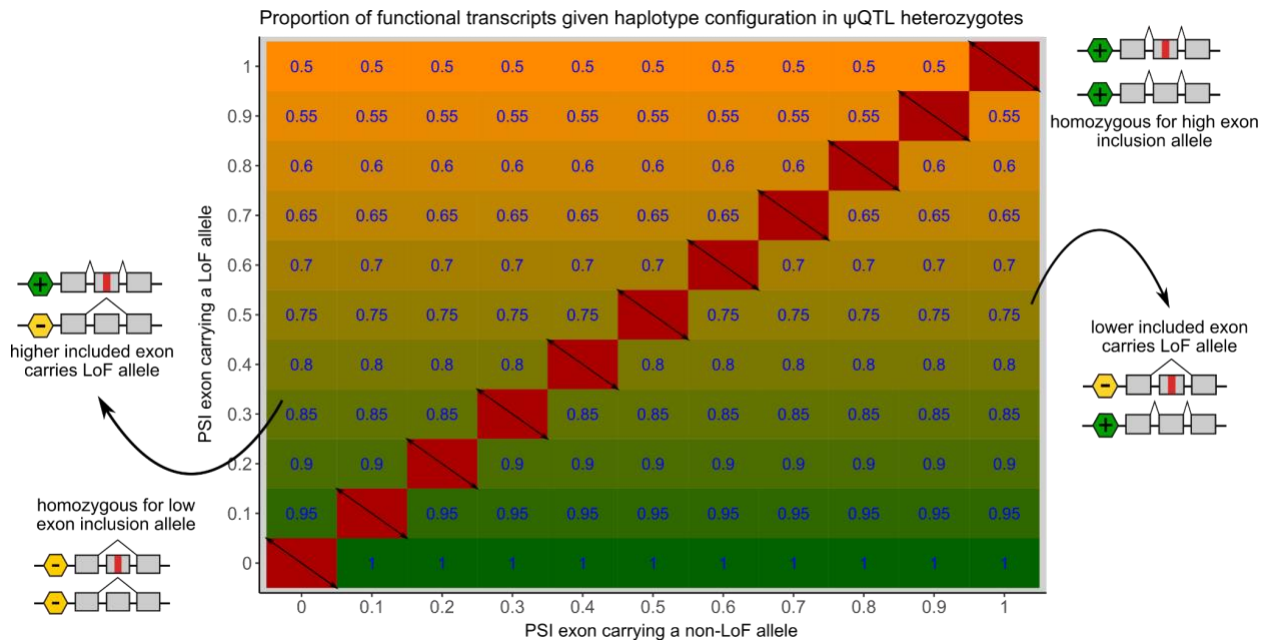
865 Supplemental Figures

866



867

868



869

870

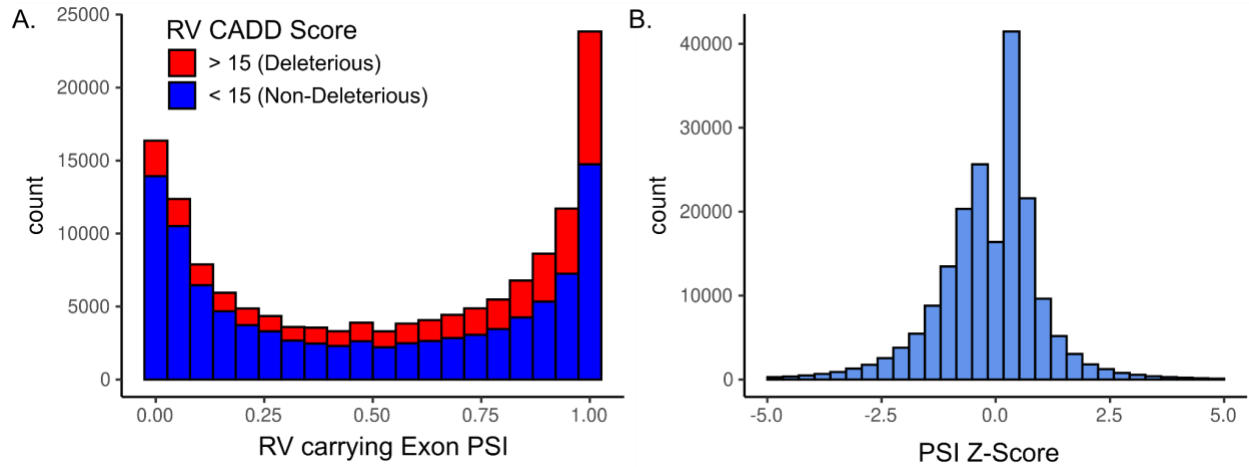
871 **Figure S1 - Splicing as a modifier of penetrance, in detail:**

872 Under the modified penetrance model that we consider here, regulatory variation that alters the  
873 dosage effect that a loss-of-function variant has on a gene is the primary driver of incomplete  
874 penetrance of the LoF variant. In this project, we focus specifically on exon splicing as a driver  
875 of this phenomenon. Generally, we consider regulatory alleles in this model to be selectively  
876 neutral, which is likely to be the case for most common regulatory variants.

877 In the top example, we present a scenario where two splicing isoforms for a particular  
878 gene exist, and their ratio is controlled by a  $\psi$ QTL where one allele causes a target exon to be  
879 included 100% of the time the gene is expressed, and the other allele causes the exon to be  
880 skipped 50% of the time it is expressed. If by chance, an individual carries a loss-of-function  
881 allele on the target exon (either by transmission or *de-novo* mutation), functional exon dosage is  
882 reduced to 75% if the loss-of-function variant lands on lower included haplotype. The functional  
883 dosage is further reduced to 50% if the loss-of-function variant lands on the higher included  
884 haplotype. In this example, it is important to note that loss of functional gene dosage is driven  
885 only by the haplotype carrying a loss-of-function allele, and its  $\psi$ QTL allele being a potential  
886 modifier of this. The other haplotype is fully functional and its sQTL allele is irrelevant. This is a  
887 subtle but pertinent distinction between the eQTL as a modifier of penetrance hypothesis  
888 (Castel et al. 2018), where the LoF haplotype is considered non-functional, and the the non-LoF  
889 haplotype is responsible for maintaining normal downstream function as modified by its eQTL  
890 allele. All assumptions about haplotype frequency in the population and haplotype frequency in  
891 diseased patients are the same across the two models.

892 In the lower figure, we generalize the model to include  $\psi$ QTLs of all effect sizes. For  
893 heterozygotes, the upper left corner of the plot represents putative high-penetrance haplotypes,  
894 and the lower right corner represents putative low-penetrance haplotypes. For  $\psi$ QTL  
895 homozygotes, deleterious or non-deleterious haplotype designations depend on the PSI of the  
896 alternative  $\psi$ QTL allele. At the population level, we hypothesize that purifying selection acts  
897 more strongly against high-penetrance haplotype combinations. However, we do not account for  
898 quantitative changes in functional dosage as they are likely to be highly gene-specific and  
899 mostly unknown.

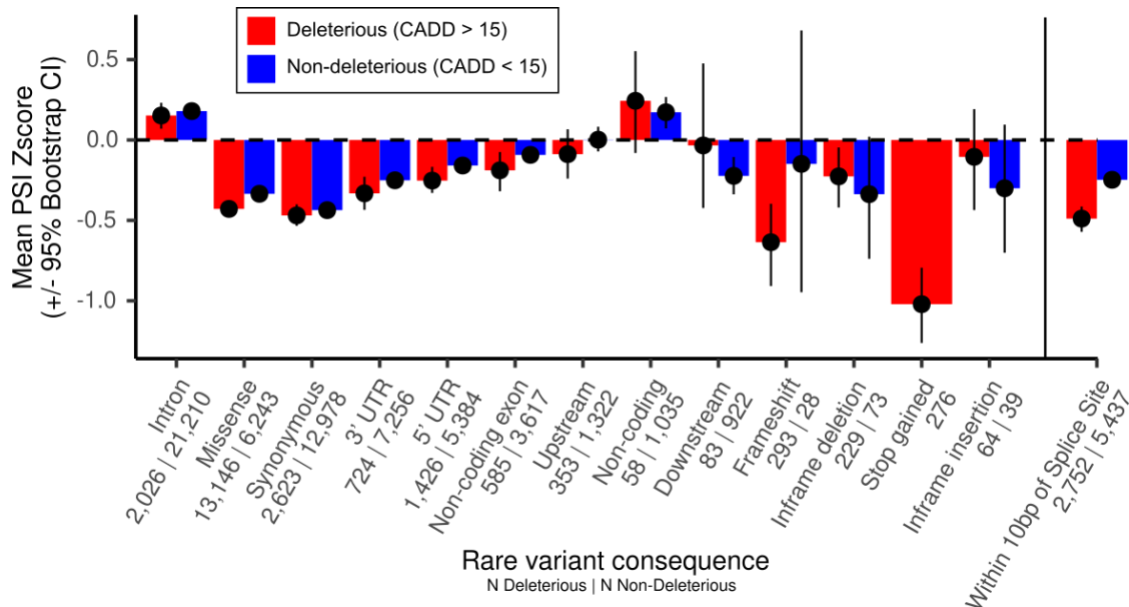
900



901  
902  
903  
904  
905  
906  
907  
908  
909

**Figure S2: PSI among exons carrying a rare variant and PSI Z scores.**

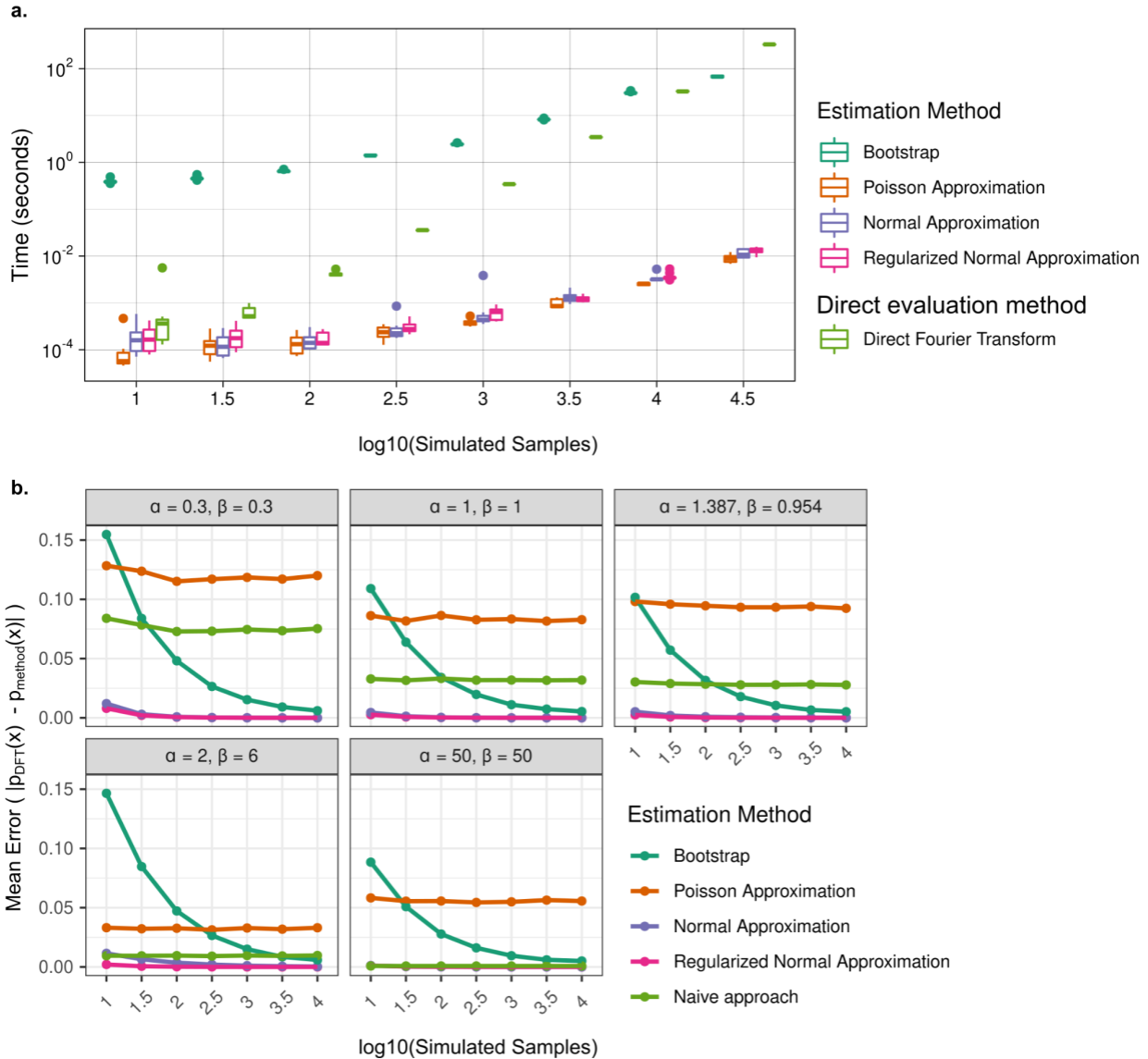
**A.** Distribution of percent spliced in (PSI) scores of all exons with sufficient variability across individuals in GTEx that carry rare variants. Colors indicate CADD score of the rare variant. In general, variants on more highly included exons are assigned a higher CADD score. **B.** PSI Z-scores are generated by fitting a normal distribution to PSI levels across GTEx individuals for a particular exon. For each exon, the PSI Z-score is in reference to the splicing of the same exon in the same tissue across all other donors with RNA-seq data available for that tissue.



910

911 **Figure S3: Mean Z-score (+/- 95% bootstrap CI) across annotations**

912 The number of rare alleles with deleterious and non-deleterious CADD designations  
 913 respectively are printed beneath each rare variant annotation. When data is available for an  
 914 individual in multiple tissues, we calculated the mean Z-score. When collapsing across tissues  
 915 and viewing by annotation, we see that deleterious alleles are depleted in most annotation  
 916 classes as well. Some variants may be annotated as “intronic” even though their loci are labeled  
 917 as exonic in the annotation used in the rest of the study.



918

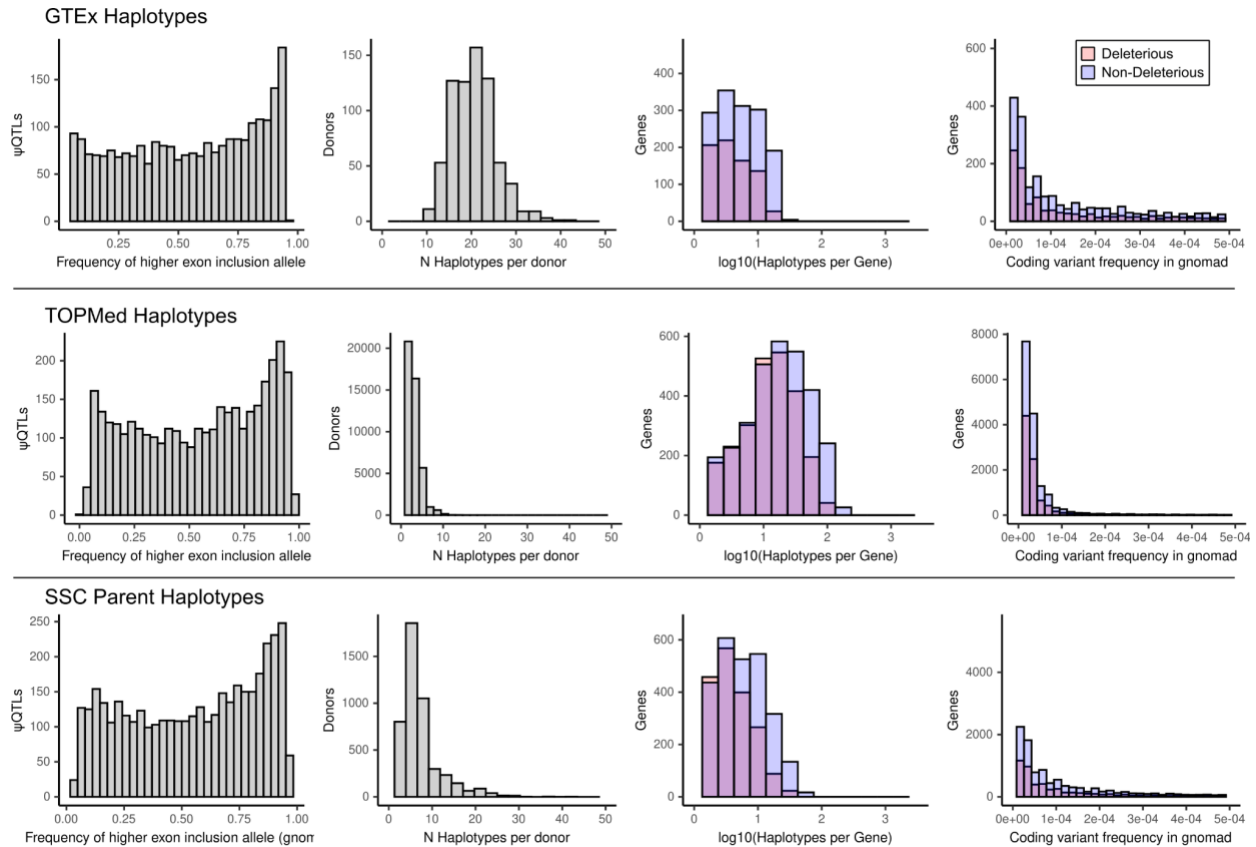
919 **Figure S4: Runtime and Accuracy benchmark of the Bootstrapped Poisson-Binomial**

920 For all benchmarking analyses, we compare our method, which approximates the cumulative  
 921 distribution function (CDF) of the Poisson-binomial distribution using a bootstrapping procedure,  
 922 to four other methods included in the 'poibin' R package. (Hong 2013) We use 5,000 bootstrap  
 923 samples here, but we found that in general 1,000 bootstrap samples balanced accuracy and  
 924 runtime. **a.** We measured the runtime to calculate the CDF of simulated datasets with uniform  
 925 probability distributions. We found that the bootstrap method outperformed the Direct Fourier  
 926 Transform (DFT) method for datasets with  $N > 10,000$ . DFT exceeded allocated memory for  
 927 more than 10,000 samples, which we frequently encounter when analyzing real data. **b.** The  
 928 bootstrap method performed more accurately with larger sample sizes, measured as the  
 929 absolute difference between the estimation method and the DFT method. We tested across

930 datasets with different distributions of  $p_j$ , the vector of probabilities that define each binary  
931 observation.  $p_j$ s were sampled from various beta distributions. The “naive approach,” for  
932 comparison, is a binomial test where  $p$  is the mean of  $p_j$ .

933

934



935

936 **Figure S5: Summaries of haplotype calls across the 3 WGS datasets**

937 In GTEx, TOPMed, and parents in the Simons Simplex Collection, we balanced sample size

938 and allele frequency cutoffs to compile the best set of haplotype configurations. Across each

939 dataset, we plot from left to right 1) the distribution of high exon inclusion  $\psi$ QTL allele

940 frequencies; 2) The number of haplotypes identified per donor, given the rare variant allele

941 frequency cutoffs (see Table 2). The larger the dataset, the more stringent we can be for

942 defining a 'rare' variant; 3) The number of haplotypes identified per gene; 4) The minor allele

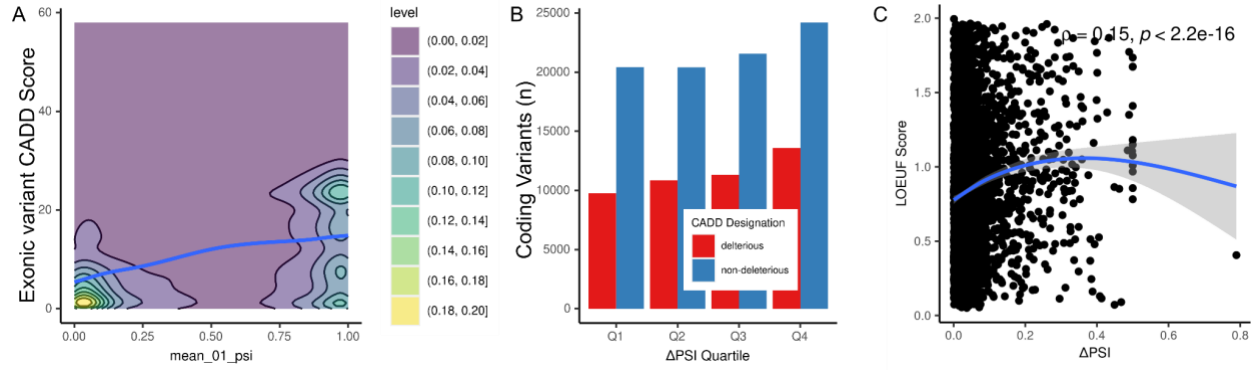
943 frequency in gnomad of all rare variants considered in the haplotype frequency analysis.

944 Deleterious and non-deleterious refer to the CADD score designation (less than and greater

945 than 15 respectively).

946





947

948 **Figure S6: cSNP annotation counts in TOPMed**

949 **A.** More deleterious (higher CADD) variants tend to fall on exons with higher baseline PSI. **B.**

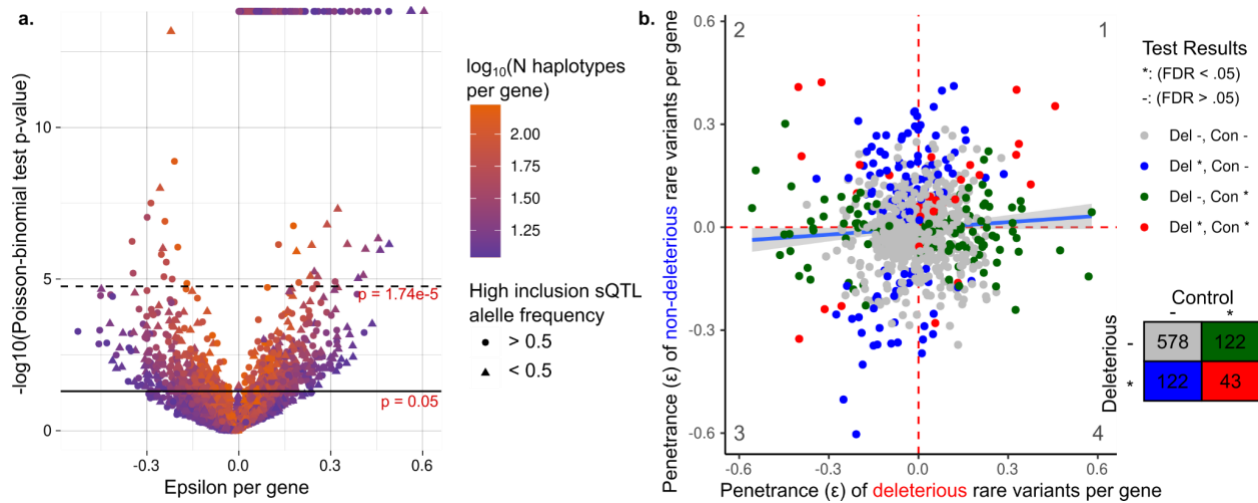
950 Haplotypes grouped by  $\Delta$ PSI Quartile. More rare variants, both deleterious and non-

951 deleterious, appear at exons with larger effect size sQTLs. **C.** Genes that are tolerant to loss-of-

952 function variants (high LOEUF) have  $\psi$ QTLs with a higher effect size ( $\Delta$ PSI).

953

954



955

956 **Figure S7: Gene by gene analysis in TOPMed**

957 **a.** Depletion of high-penetrance haplotype observations on a gene-by-gene basis in TOPMed.

958 For each gene with more than 10 observed haplotypes across donors, we test if any genes or  
959 classes of genes are driving the overall pattern of high-penetrance haplotype depletion. Each

960 point represents a single gene. **b.** Comparison of haplotype deviation between deleterious and

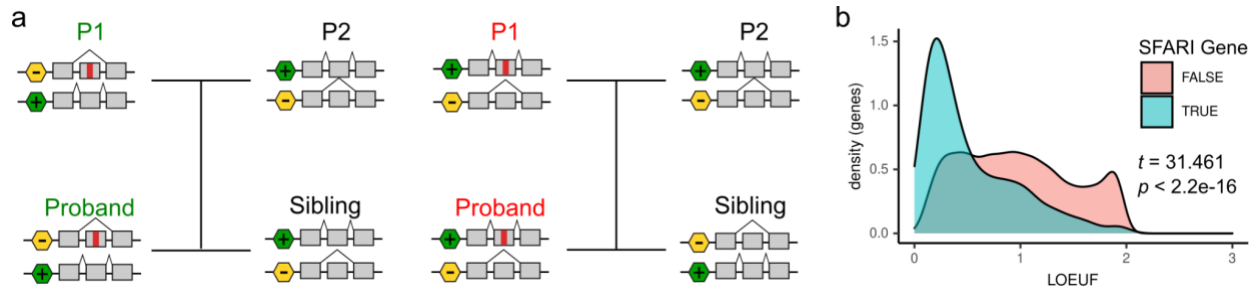
961 non-deleterious rare coding variants, among genes with greater than 10 haplotypes in both

962 categories. Under a model where highly penetrant deleterious cSNPs are depleted in the

963 population, we expect more blue-labeled genes in the third quadrant.

964

965



966

967 **Figure S8: Transmission patterns of splicing haplotypes**

968 **a.** When a parent carries a exonic variant in a putative low (green text) or high (red text)

969 penetrance haplotype configuration, they will almost always transmit it to a child in the same

970 haplotype configuration. **b.** Distribution of LOEUF scores among genes identified as relevant to

971 Autism Spectrum Disorder, by SFARI Gene. ASD genes have significant depletion of predicted

972 loss-of-function variants in general.

973

## 974 Supplemental Tables

### 975 **Supplementary Table 1: GTEx Tissues utilized for $\psi$ QTL calling, and the number of exons** 976 **pre and post filtering.**

Tissue	N Exons per tissue pre-filtering	N Exons per tissue post-filtering	Percent usable	Genes covered per tissue
Adipose_Subcutaneous	260,800	29,180	11.19%	8,585
Artery_Tibial	253,109	27,453	10.85%	8,127
Brain_Cerebellum	239,928	36,095	15.04%	8,605
Brain_Cortex	240,439	26,121	10.86%	7,857
Brain_Nucleus_accumbens_basal_ganglia	247,074	26,372	10.67%	7,998
Cells_Cultured_fibroblasts	230,752	28,486	12.34%	8,479
Cells_EBV.transformed_lymphocytes	220,547	37,837	17.16%	9,291
Colon_Transverse	231,647	29,066	12.55%	8,630
Esophagus_Mucosa	245,627	26,721	10.88%	7,984
Liver	224,469	21,605	9.62%	6,283
Lung	265,555	34,585	13.02%	9,387
Muscle_Skeletal	240,921	22,664	9.41%	6,788
Nerve_Tibial	261,375	30,771	11.77%	8,783
Pituitary	259,310	32,795	12.65%	8,774
Skin_Sun_Exposed_Lower_leg	259,438	29,570	11.40%	8,588
Spleen	241,122	30,379	12.60%	8,277
Thyroid	266,364	30,035	11.28%	8,586
Whole_Blood	236,866	23,135	9.77%	6,039

### 977 978 **Supplementary Table 2: TOPMed cohorts utilized and number of samples from each** 979 **cohort**

980

981 [Supplemental Table 2.xlsx](#)

982