

# Genetic control of mRNA splicing as a potential mechanism for incomplete penetrance of rare coding variants

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> Mariza de Andrade,<sup>6</sup> Patrick Ellinor,<sup>7</sup> Myriam Fornage,<sup>8</sup> Stacey Gabriel,<sup>9</sup> Soren 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. 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> Braxton D. Mitchell,<sup>18,19</sup> Bruce Psaty,<sup>20</sup> Ramachandran S. Vasani,<sup>21</sup> Stephen S. Rich,<sup>22</sup> Michael Rienstra,<sup>23</sup> Jerome I. Rotter,<sup>24</sup> Aabida Saferali,<sup>11</sup> Moore Benjamin Shoemaker,<sup>25</sup> Edwin Silverman,<sup>26</sup> Albert Vernon Smith,<sup>27</sup> NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium, Pejman Mohammadi,<sup>29</sup> Stephane E. Castel,<sup>2,30</sup> Ivan Iossifov,<sup>2,31</sup> Tuuli Lappalainen<sup>2,32,33\*</sup>

<sup>1</sup>Department of Biomedical Informatics, Columbia University, New York, NY 10027, USA

<sup>2</sup>New York Genome Center, New York, NY 10013, USA

<sup>3</sup>School of Public Health, University of Texas Health at Houston, Houston, TX 77030, USA

<sup>4</sup>Department of Medicine, Brigham & Women's Hospital, Boston, MA 02115, USA

<sup>5</sup>Department of Cardiology, University of Illinois at Chicago, Chicago, IL 60607, USA

<sup>6</sup>Department of Quantitative Health Sciences, Mayo Clinic, Rochester, MN 55905, USA

<sup>7</sup>Corrigan Minehan Heart Center, Massachusetts General Hospital, Boston, MA 02114, USA

<sup>8</sup>Brown Foundation Institute of Molecular Medicine, McGovern Medical School, University of Texas Health at Houston, Houston, TX 77030, USA

<sup>9</sup>Broad Institute, Cambridge, MA 02142, USA

<sup>10</sup>Department of Molecular and Human Genetics, Baylor College of Medicine Human Genome Sequencing Center, Houston, TX 77030, USA

<sup>11</sup>Channing Division of Network Medicine and Division of Pulmonary and Critical Care Medicine, Brigham and Women's Hospital, Boston, MA 02115, USA

<sup>12</sup>Department of Hematology, University of Washington, Seattle, WA 98195, USA

<sup>13</sup>Department of Epidemiology & Population Health, Albert Einstein College of Medicine, Bronx, NY 10461, USA

<sup>14</sup>Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA

<sup>15</sup>Department of Pathology, School of Medicine, Umm Al-Qura University, Mecca 24382, Saudi Arabia

<sup>16</sup>Environmental Medicine & Public Health, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA

<sup>17</sup>Department of Medicine, University of Arizona, Tucson, AZ 85721, USA

<sup>18</sup>Department of Medicine, University of Maryland School of Medicine, Baltimore, MD 21201, USA

<sup>19</sup>Geriatrics Research and Education Clinical Center, Baltimore Veterans Administration Medical Center, Baltimore, MD 21201, USA

<sup>20</sup>Cardiovascular Health Research Unit, Departments of Medicine, Epidemiology, and Health Systems and Population Health, University of Washington, Seattle, WA 98195, USA

<sup>21</sup>Department of Medicine, Boston University, Boston, MA 02118, USA

<sup>22</sup>Public Health Sciences, University of Virginia, Charlottesville, VA 22903, USA

<sup>23</sup>Clinical Cardiology, UMCG Cardiology, Groningen 09713, the Netherlands

<sup>24</sup>The Institute for Translational Genomics and Population Sciences, Department of Pediatrics, The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance, CA 90502, USA

<sup>25</sup>Department of Medicine, Vanderbilt University, Nashville, TN 37235, USA

<sup>26</sup>Channing Division of Network Medicine and Division of Pulmonary and Critical Care Medicine, Brigham & Women's Hospital, Boston, MA 02115, USA

<sup>27</sup>Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA

<sup>29</sup>Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, CA 92037, USA

<sup>30</sup>Variant Bio, Seattle, WA 98102, USA

<sup>31</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA

<sup>32</sup>Department of Systems Biology, Columbia University, New York, NY 10027, USA

<sup>33</sup>Department of Gene Technology, KTH Royal Institute of Technology, Stockholm 114 28, Sweden

\*Corresponding author: Department of Biomedical Informatics, Columbia University, New York, NY 10027, USA. Email: jeinson@nygenome.org (J.E.);

\*Corresponding author: Department of Systems Biology, Columbia University, New York, NY 10027, USA. Email: tlappalainen@nygenome.org (T.L.)

## Abstract

Exonic variants present some of the strongest links between genotype and phenotype. However, these variants can have significant inter-individual pathogenicity differences, known as variable penetrance. In this study, we propose a model where genetically controlled mRNA splicing modulates the pathogenicity of exonic variants. By first cataloging exonic inclusion from RNA-sequencing data in GTEx V8, we find that pathogenic alleles are depleted on highly included exons. Using a large-scale phased whole genome sequencing data from the TOPMed consortium, we observe that this effect may be driven by common splice-regulatory genetic variants, and that natural selection acts on haplotype configurations that reduce the transcript inclusion of putatively pathogenic variants, especially when limiting to haploinsufficient genes. Finally, we test if this effect may be relevant for autism risk using families from the Simons Simplex Collection, but find that splicing of pathogenic alleles has a penetrance reducing effect here as well. Overall, our results indicate that common splice-regulatory variants may play a role in reducing the damaging effects of rare exonic variants.

**Keywords:** incomplete penetrance, QTLs, alternative splicing, functional genomics, GTEx, TOPMed, Simons Simplex Collection, statistical genetics

Received: February 02, 2023. Accepted: April 18, 2023

© The Author(s) 2023. Published by Oxford University Press on behalf of The Genetics Society of America. All rights reserved. For permissions, please e-mail: journals.permissions@oup.com

## Introduction

Incomplete penetrance is a well-known phenomenon, where an individual carries a disease-associated allele, but develops no symptoms of the disease themselves (Shawky 2014; Gettler et al. 2021; Forrest et al. 2022). Similarly, variable expressivity refers to analogous gradual differences in disease severity; here, we refer to both as variable penetrance. These instances are likely under-reported in the literature due to ascertainment bias, when many studies are based on sequencing due to a prior genetic condition (Cooper et al. 2013; Dewey et al. 2016). Even amongst Mendelian disease variants, which are typically thought of as having strong effects on phenotype, differing levels of severity have been observed between carriers (Chen et al. 2016). These changes have been attributed to epistatic or additive effects of genetic modifiers, as well as environmental modifiers of penetrance, which can be difficult to control in an experimental setting (Maya et al. 2018). When looking at incomplete penetrance in specific diseases, genetic modifiers have been mapped, for example, to BRCA in breast cancer (Milne and Antoniou 2011), and RET in Hirschsprung's disease (Emison et al. 2005). Modified penetrance has also been studied in the context of polygenic risk scores, where multiple common risk variants increase the expected pathogenicity of a disease-relevant variant (Fahed et al. 2020). However, genome-wide patterns underlying modified penetrance are still poorly known. One potential mechanism for incomplete penetrance is cis-regulatory mechanisms that affect the regulation of a gene carrying a pathogenic variant. This model has been tested with expression quantitative trait loci (eQTLs) acting as modifiers of penetrance (Castel et al. 2018), but can be expanded to other types of gene regulatory processes, such as mRNA splicing. While eQTLs control the dosage of their target genes, splicing alters inclusion of variant-carrying exons in transcripts, which could potentially have a large effect on the overall pathogenicity of a damaging variant.

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

In this study, we build upon the finding that transcript usage of genes containing alleles contributes to the allele's pathogenicity, and ask if common splice-regulatory variants may partially drive this phenomenon and affect inter-individual variation in penetrance. Expanding on previous methodology (Castel et al. 2018), we look for nonrandom haplotype combinations of sQTL variants and putatively pathogenic rare variants in population scale datasets. Such an observation could indicate that haplotype combinations have an effect on fitness, and by proxy, disease risk. In doing so, we develop a general framework for modeling common and rare variant haplotypes in a population, with a corresponding

test to detect deviations from the null (Fig. 1 and Supplementary Fig. 1). These analyses will improve our understanding of how variants across the annotation and allele frequency spectrum act together to shape human traits and could ultimately aid our interpretation of rare variants in a clinical context.

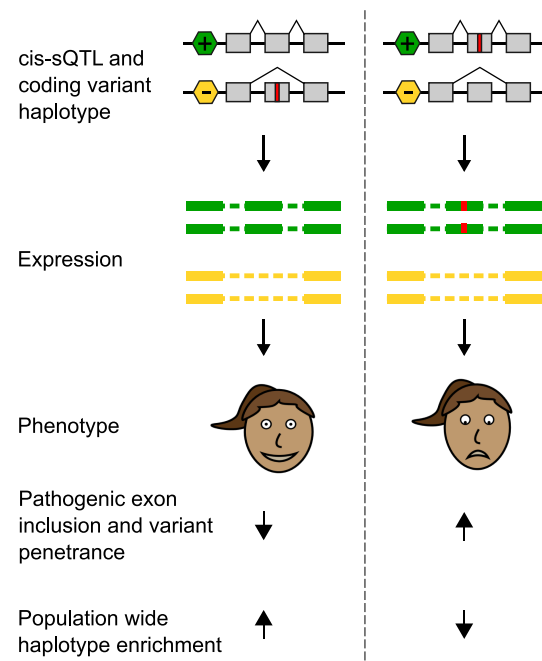
## Methods

### Data sources

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

### GTEx percent spliced in quantification and filtering

Percent spliced in (PSI) was calculated from GTEx V8 RNA-seq data. We limited our analysis to 18 tissues, which were chosen for their coverage of tissue diversity GTEx and their coverage of the most coding genes possible (Supplementary Table 1). Exon PSI for protein-coding genes was quantified using the Integrative Pipeline for Splicing Analysis (IPSA) (Pervouchine et al. 2013; IPSA-nf 2020), which was run on Google Cloud through Terra (<https://github.com/guigolab/ipsa-nf>). The “-unstranded” flag was used during the sjcount process. Exons were defined by the modified version of Gencode annotation v26 used in GTEx V8, which collapses genes with multiple isoforms to a single isoform



**Fig. 1.** Splice-regulatory variants as modifiers of penetrance hypothesis. The hypothesis of this study is illustrated with an example of an individual who is heterozygous for both a sQTL and a coding variant. The 2 possible haplotype configurations result in either a reduced or increased penetrance state of the coding allele, depending if the allele is on the more lowly or highly included exon, respectively. We predict that natural selection would deplete those that fall in a high-penetrance configuration in the general population. See Supplementary Fig. 1 for a quantitative description of the model.

per gene ([https://storage.googleapis.com/gtex\\_analysis\\_v8/reference/genencode.v26.GRCh38.genes.gtf](https://storage.googleapis.com/gtex_analysis_v8/reference/genencode.v26.GRCh38.genes.gtf)).

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

### Percent spliced in Z-score analysis in GTEx

We compiled a list of all exons with sufficiently variable splicing in at least one GTEx tissue, as defined in the previous step, and saved the genomic coordinates of these exons in BED format. Rare variants (gnomAD AF < 0.01) that fell on variably spliced exons were extracted from GTEx WGS VCFs, and were subsequently filtered to variants that appeared less than 6 and greater than 1 time. Rare variant Combined Annotation Dependent Depletion (CADD) (Rentzsch et al. 2019) scores and annotations with respect to the relevant gene were extracted as well. Because CADD v1.5 uses a different VEP annotation that in some cases does not correspond the exon annotations used previously, we re-annotated rare variants using VEP v93.2 and gencode v28, taking the most deleterious annotation when a rare variant covered multiple transcripts. Rare variant calls from exons represented disproportionately, either due to length or to high number of variants at the exon, were removed. Threshold for removing an exon was defined as  $Q3 + 1.5 * IQR$ , where  $Q3$  is the third quartile of the number of rare variants per exon and where  $IQR$  is the interquartile range of the number of rare variants per exon. For all remaining variants, we computed the PSI Z-score of the individual that carried the variant at that specific exon, across all tissues where the exon was expressed and sufficiently variable. The PSI Z-score for a particular individual  $i$  at 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 particular exon and tissue, and  $\mu_j$  and  $\sigma_j$  are the mean and standard deviation of PSI for an exon  $j$  across all individuals with data available for that exon in tissue  $k$ . Importantly, we do not normalize PSI for this analysis, to preserve signal from exons with high PSI Z-scores.

### Primary quantitative trait locus mapping, collapsing, and secondary $\psi$ quantitative trait locus mapping

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

For every exon, we selected the most significant variant, and for every gene the most significant exon. We then compiled the QTL results across tissues to achieve a set of cross-tissue top QTLs. From here forward we refer to these QTLs as  $\psi$ QTLs, using

the  $\psi$  to distinguish from convention sQTLs in order to emphasize that the quantitative trait is exon PSI. When a gene was significant across multiple tissues, we used the tissue where the effect size ( $\Delta$ PSI score) was the highest. This process ensured that a gene was only included once in our final set of  $\psi$ QTLs, and was labeled by one variant that is associated to splicing (sVariant).

Since the splicing of multiple exons within a gene is often correlated, we implemented an approach to identify additional exons whose splicing the sVariant is associated with. Consideration of multiple exons per gene is desirable because it increases the amount of genetic space where rare variant haplotypes can be identified. For each gene with a significant  $\psi$ QTL, we ran a nominal QTLtools pass of just the sVariant against PSI of all other exons in the gene. We then considered secondary exons with a Bonferroni-corrected  $P < 0.05$  if QTL effect direction was the same as the top exon.

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

### Whole genome sequencing filtering across datasets

#### Genotype-Tissue Expression (GTEx) project

Read-aware phased WGS data were used from all 838 samples included in GTEx V8 (Consortium 2020) (Supplementary Information Section 2.4). For use in haplotype calling, the following filters were applied as follows: (1) Variants were extracted with an allele frequency less than 0.005 in gnomAD, and singleton variants without read backing to support their phase call were removed. (2) Samples from donors that did not self-identify as European American were removed. Since the  $\psi$ QTL data from GTEx is based on 85% European Americans, the sVariants selected from these data may not capture allele frequencies and haplotype structures in other ancestries, and differing numbers of rare variants across ancestries might bias the results. (3) Haplotype calls from genes represented disproportionately, either due to length or to high number of variants at the gene, were removed. Threshold for removing a gene was defined as  $Q3 + 1.5 * IQR$ , where  $Q3$  is the third quartile of the number of haplotypes per gene and where  $IQR$  is the interquartile range of the number of haplotypes per gene.

#### Trans-Omics for Precision Medicine (TOPMed) initiative

Population-phased WGS data from donors of European-American ancestry were used from TOPMed, since this matches the population source of the  $\psi$ QTL data from GTEx (see above). To define individuals of European ancestry, we used the approach outlined in Morris et al. (2019). Briefly, TOPMed samples were projected onto the first 20 principal components estimated from the 1000 Genomes Phase 3 (1000G) project (Auton et al. 2015) using FastPCA v2.0 (Galinsky et al. 2016). Only biallelic variants shared between the 2 datasets, and that passed a strict set of criteria (minor allele frequency [MAF] > 1%, minor allele count > 5, genotyping call rate > 95%, and Hardy-Weinberg  $P$ -value >  $1 \times 10^{-6}$ ) were used to calculate the principal components. Expectation Maximization (EM) (Chen and Maitra 2015) clustering was used to compute the probabilities of cluster membership, and eigenvectors 1, 2, 5, 6, and 8 were selected for efficiently separating

the individuals of White European and American ancestry (subpopulation codes CEU, GBR, FIN, CEU, IBS, and TSI) from other ancestry groups. Finally, 8 predefined clusters were chosen for EM clustering based on sensitivity analyses. This resulted in 52,426 TOPMed individuals clustering together with the 1000G CEU, GBR, FIN, CEU, IBS, and TSI subpopulation, and they were termed of White ancestry. We kept 19 cohorts (Supplementary Table 2), and 49,542 individuals, filtering out the remaining cohorts, which collectively contained less than 5% of all haplotypes.

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

### Simons Simplex Collection (SSC)

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

We additionally removed genes that contained an unusually high number of rare coding variants across parents, using the same outlier definition as in the previous 2 datasets. This set of variants post-filtering was considered in siblings and probands in downstream analyses.

### Haplotype calling from phased genetic data and filtering

$\psi$ QTL-coding allele haplotypes were generated using a similar procedure across all 3 phased-resolved WGS datasets. First, all rare variants were extracted among sExons using the filters described above, considering variants that fell in primary and secondary sExons, taking account of the haplotype phase assignment. Then, the genotype of sVariants, and phase for heterozygous cases, was extracted from VCFs, and haplotypes were labeled as high-penetrance ( $\beta = 1$ ) and low penetrance ( $\beta = 0$ ) according to our model for splice QTLs as a modifier of penetrance (Fig. 1).

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

We sought to test the hypothesis that QTL-coding allele haplotype combinations are present in the population at frequencies that deviate from a baseline expectation, based on allele frequencies alone. Such a result could indicate high-penetrance haplotypes with deleterious variants being removed from the population by natural selection. The total number of high-penetrance haplotypes arising from  $\psi$ QTLs with varying allele frequencies can be modeled by the Poisson-binomial distribution, which is a generalization of the binomial distribution. While a binomial describes the sum of  $n$  independent identically distributed Bernoulli random variables, the Poisson-binomial describes the sum of  $n$  independent but non-identically distributed Bernoulli random

variables. Therefore, the distribution must be parameterized by a vector of probabilities of length  $n$ . While we could calculate  $P$ -values using a variety of methods that obtain the cumulative distribution function (CDF) of the Poisson-binomial (Hong 2013), these methods all lack a way to quantify the magnitude of the effect size. Furthermore, they measure deviation from the null but do not allow comparison of 2 datasets (in our case, haplotypes carrying non-deleterious and deleterious coding alleles). Therefore, we developed the following procedure that approximates the Poisson-binomial CDF. This has the advantage of generating a quantifiable effect size for deviation from the null model, as well as corresponding confidence intervals.

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

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

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

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

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

We define the expectation of observing a homozygous  $\psi$ QTL allele as the proportion of high inclusion  $\psi$ QTL homozygotes in the dataset, plus a pseudo-count, to avoid getting an expectation of 0 in datasets where the low inclusion allele is much more common. This method does not assume Hardy-Weinberg equilibrium for the  $\psi$ QTL allele, but requires that the proportion of homozygotes for the 2 alleles be recalculated on each dataset. This approach was used for the GTEx and TOPMed analyses. Alternatively, the expectation of  $\beta$  under the null model can also be calculated as follows:

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

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



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

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

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

We quantify the significance of  $\epsilon$  by bootstrapping all haplotypes, generating 95% confidence intervals and drawing 2-sided empirical  $P$ -values as

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

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

Although the test was designed for counts of haplotypes, this approach is generalizable to any system that can be modeled by a Poisson-binomial distribution. Therefore, to benchmark our test, we simulated data from several theoretical allele frequency distributions by sampling from beta distributions with various shape parameters, including one distribution where its parameters were estimated direction from our set of  $\psi$ QTLs from GTEx using the method of moments estimator (Fig. 3 and Supplementary Fig. 4). We found that our bootstrapping procedure accurately approximated the Poisson-binomial distribution for all inputs tested. However, the magnitude of  $\epsilon$ —but not direction—is dependent on the shape of the theoretical allele frequency distribution, so comparing magnitudes of  $\epsilon$  across distinct datasets should be done with caution. The accuracy of our method increased with larger sample sizes. Therefore, we recommend using this approach when handling data where  $N > 1,000$  (Supplementary Fig. 4).

As an extension to this procedure, we can also conveniently calculate the significance of a difference in  $\epsilon$  between 2 similar datasets  $A$  and  $B$ , for example, between haplotypes where the rare variant is putatively deleterious vs haplotypes where the rare variant is non-deleterious:

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

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

This test is implemented in the STatistic for Modified PENetrance (STAMPEN) R package that is available to download here (<https://github.com/jeinson/stampen>).

## Results

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

We first tested the hypothesis that rare pathogenic alleles (CADD > 15) (Rentzsch et al. 2019) are more likely to occur at less

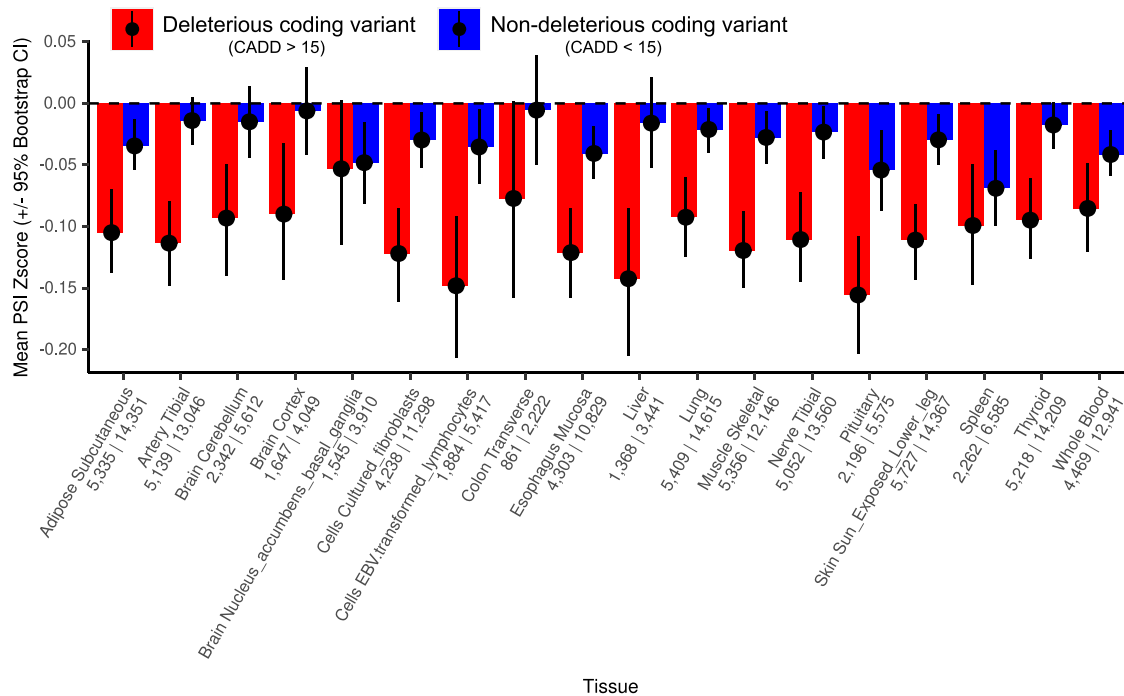
spliced-in exons (Fig. 1). To accomplish this, we used bulk RNA-sequencing (RNA-seq) and WGS data from the Genotype-Tissue Expression (GTEx) project V8 release, which is representative of a general population free of severe genetic disease. We defined variants as rare if their variant frequency in gnomAD (Karczewski et al. 2020) was less than 0.5%, and they appeared 5 or fewer times among the 838 GTEx WGS donors.

To begin, we calculated PSI scores for all annotated protein-coding gene exons across 18 GTEx tissues, and only kept exons with sufficient splicing variability across individuals (Methods, Supplementary Table 1, and Supplementary Fig. 2a). We extracted rare alleles that fell on variably spliced exons, separating alleles within 10 bp of a splice junction to avoid cases where the allele is more likely to directly affect splicing. To compare the splicing of each donor with a deleterious allele to the population distribution per exon, we calculated PSI Z-scores across all tissues with available data (Supplementary Fig. 2b, Methods). We found that PSI Z-scores were significantly different between exons carrying deleterious ( $N = 19,178$ ) and non-deleterious ( $N = 49,575$ ) rare alleles (Mann-Whitney U test:  $P = 2.577 \times 10^{-4}$ ). This rank difference was accounted for by a modest decrease in mean PSI Z-score among donors that carried deleterious alleles in a given exon, which was consistent across tissues and across variant consequence annotations (Fig. 2 and Supplementary Fig. 3). Notably, stop-gained variants had the strongest association with low PSI Z-scores—even stronger than the signal for variants close to splice junction—but the overall result was present for multiple annotation categories (Supplementary Fig. 3). This suggests that the signal is not solely driven by the most pathogenic variants nor direct rare variant effects on splicing. These results extend the previous work, comparing different exons and showing accumulation of stop-gained variants on those with lower inclusion (Cummings et al. 2020). Here, we observe a similar pattern when comparing different individuals within a given exon, consistent with the hypothesis that the penetrance of coding alleles is reduced when they fall on more lowly included exons. However, this approach does not discern the underlying reasons for splicing differences between individuals, including alleles that may drive a decrease in splicing and their haplotype combinations with rare alleles.

### A general model for coding allele-quantitative trait locus haplotype configurations

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

Initially, splice-regulatory alleles were cataloged in GTEx through quantitative trait locus (QTL) mapping, using the percent spliced in (PSI or  $\psi$ ) (Pervouchine et al. 2013) of each exon as a quantitative phenotype. These alleles are hence referred to as  $\psi$ QTLs. We use the “ $\psi$ ” nomenclature to differentiate from sQTLs, where the splicing phenotype can vary between studies and is often less interpretable for downstream applications.  $\psi$ QTL mapping and properties are described in Einson et al. (2022). Briefly, we mapped  $\psi$ QTLs from GTEx V8 using the same filtered set of PSI scores across 18 tissues as in the previous analyses (see Methods). We compiled a set of 5,196 cross-tissue  $\psi$ QTL genes [one significant  $\psi$ QTL variant (sVariant) and one exon associated



**Fig. 2.** Mean PSI Z-scores across tissues. Mean decrease in PSI Z-scores among individuals carrying rare alleles at variably spliced exons across 18 GTEx tissues, split by deleterious (CADD > 15) and non-deleterious (CADD < 15) rare variants. The number of deleterious and non-deleterious alleles, respectively, is printed below each tissue name. Error bars represent 95% bootstrapped confidence intervals.

with the sVariant (sExon) per gene], and recorded which alleles led to higher or lower sExon inclusion. We also mapped secondary sExons across  $\psi$ QTL genes where the most significant sVariant was also associated with splicing in the same direction as the top sExon in the same gene, which were used to expand the amount of genic space where rare variants could be considered.

Next, to robustly test for nonrandom haplotype combinations of rare exonic alleles and common  $\psi$ QTL alleles, we describe an approach that quantifies the significance of deviations in haplotype combinations from the null in a dataset, taking variable  $\psi$ QTL allele frequencies into account. In most datasets,  $\psi$ QTL alleles that may have an effect on rare variant penetrance are nonuniformly distributed, and thus, we expect an unequal number of high- and low-penetrance haplotypes under the null (Fig. 3). To account for this, we model these data using the Poisson-binomial distribution, a generalization of the binomial distribution describing the sum of  $n$  independent but non-identically distributed Bernoulli random variables. (Wang 1993; Hong 2013; González et al. 2016) When looking at counts of haplotype combinations, the probability of observing a high-penetrance haplotype is assigned according to the relevant  $\psi$ QTL allele frequency, independently across QTL genes. To apply the model to haplotypes extracted from phased genetic data, we developed a bootstrapping procedure that approximates the cumulative distribution function of the Poisson-binomial, constituting a convenient method for calculating the significance, enrichment/depletion effect sizes ( $\epsilon$ ), and confidence intervals when comparing enrichment scores between groups i.e. haplotypes with deleterious vs non-deleterious rare alleles (see Methods for details). In simulations, our method was well powered to detect deviations from the null across all tested theoretical allele frequency distributions, and performed well against other methods that directly calculate and approximate the CDF of the Poisson-binomial (Fig. 4 and Supplementary Fig. 4). We also found that the type I

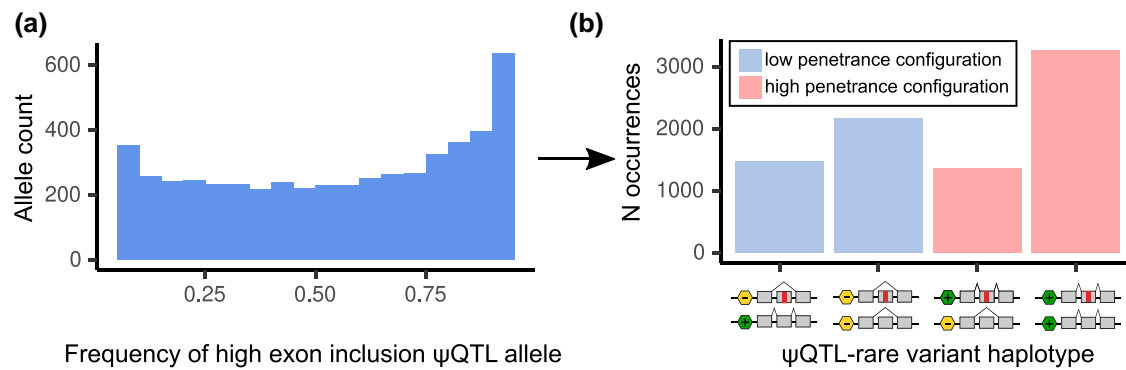
error rate was adequately controlled (Supplementary Fig. 5a), and that the test was well powered to detect differences between haplotype groups in our datasets, given their size (Supplementary Fig. 5b). Overall, this approach is generalizable to other analyses of haplotype combinations; here, we apply it to test nonrandom combinations of  $\psi$ QTL and rare coding alleles.

### High-penetrance haplotypes are depleted in TOPMed and GTEx

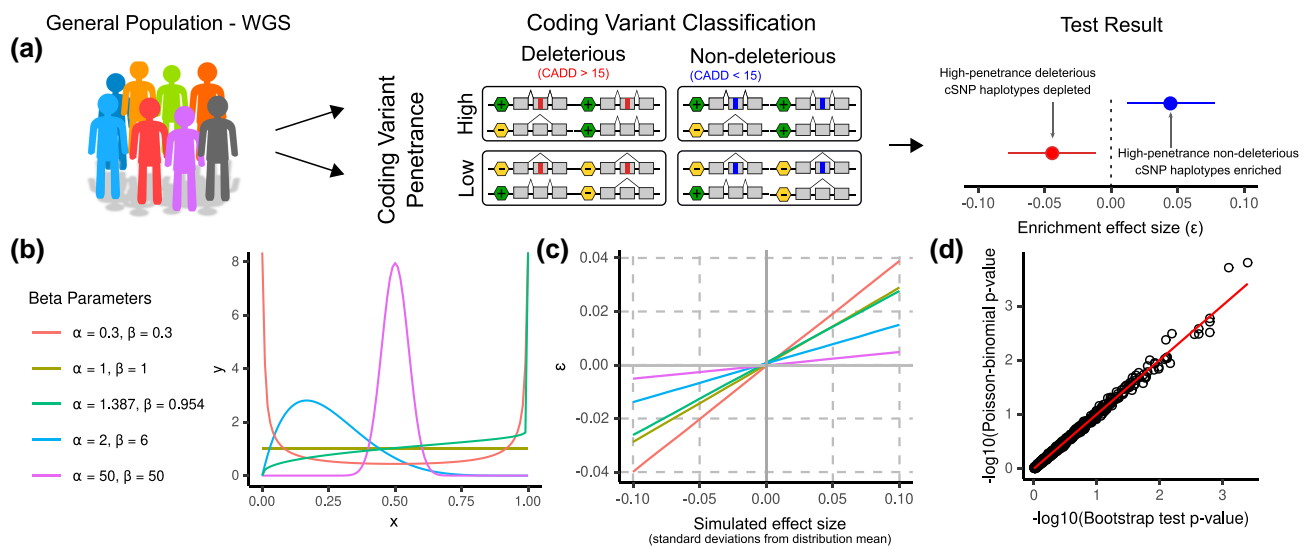
After defining a theoretical model that describes counts of common regulatory alleles and rare coding alleles in a given population, we tested 3 datasets for evidence of selection against high penetrance coding alleles driven by genetically regulated splicing (Table 1).

#### Enrichment in GTEx

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



**Fig. 3.**  $\psi$ QTL high inclusion allele frequencies and haplotype counts in GTEx. a) Distribution of allele frequencies for  $\psi$ QTLs that lead to higher exon inclusion. High inclusion  $\psi$ QTL allele frequencies are skewed to the right, meaning  $\psi$ QTLs that include their target exon are more common in the general population. b) As a result of the nonuniform frequency distribution of high inclusion  $\psi$ QTL alleles, we expect to see more high-penetrance haplotype configurations in general. This motivates the necessity to design a test that accounts for this difference.



**Fig. 4.** The Poisson-binomial distribution models haplotype configuration counts. a) We use phased variant calls from WGS across large populations to test for deviation in the frequencies of  $\psi$ QTL-coding variant haplotype configurations. The magnitude and effect direction of deviation, which we call  $\epsilon$ , are calculated using a procedure described in *Methods*. The magnitude of  $\epsilon$ —but importantly not its direction—depends on the underlying  $\psi$ QTL allele frequency distribution, as the probability of observing a high-penetrance haplotype is dependent on the  $\psi$ QTL allele frequency at each gene. Counts of highly penetrant haplotypes are modeled by the Poisson-binomial distribution. When running our test, we frequently divide haplotypes into those with deleterious (CADD > 15) and non-deleterious (CADD < 15) coding variants, which serve as a negative control where we do not expect to see evidence of purifying selection. b) To verify that our test captures deviations from the null under any theoretical allele frequency distribution, we simulated datasets by drawing samples from various beta distributions with different parameters. The beta is defined by shape parameters  $\alpha$  and  $\beta$ . The parameters  $\alpha = 1.387$  and  $\beta = 0.954$  were estimated from the high-inclusion  $\psi$ QTL allele frequency distribution in GTEx using the method of moments estimator. c) We benchmarked our test by simulating data from distributions with increasingly larger deviations from the expected mean, in order to test how the magnitude of  $\epsilon$  differs depending on the input distribution. This diagram can be used as a reference for how to interpret the magnitude of epsilon, given a dataset's underlying probability distribution. d) P-values from a simulated dataset of haplotypes from 1,000 individuals across 1,000 genes, with  $\psi$ QTL allele frequencies matching those in GTEx. We find that our method accurately replicates the results from the Poisson-binomial distribution, calculated using the “poibin” (Hong 2013) R package.

### Enrichment in TOPMed

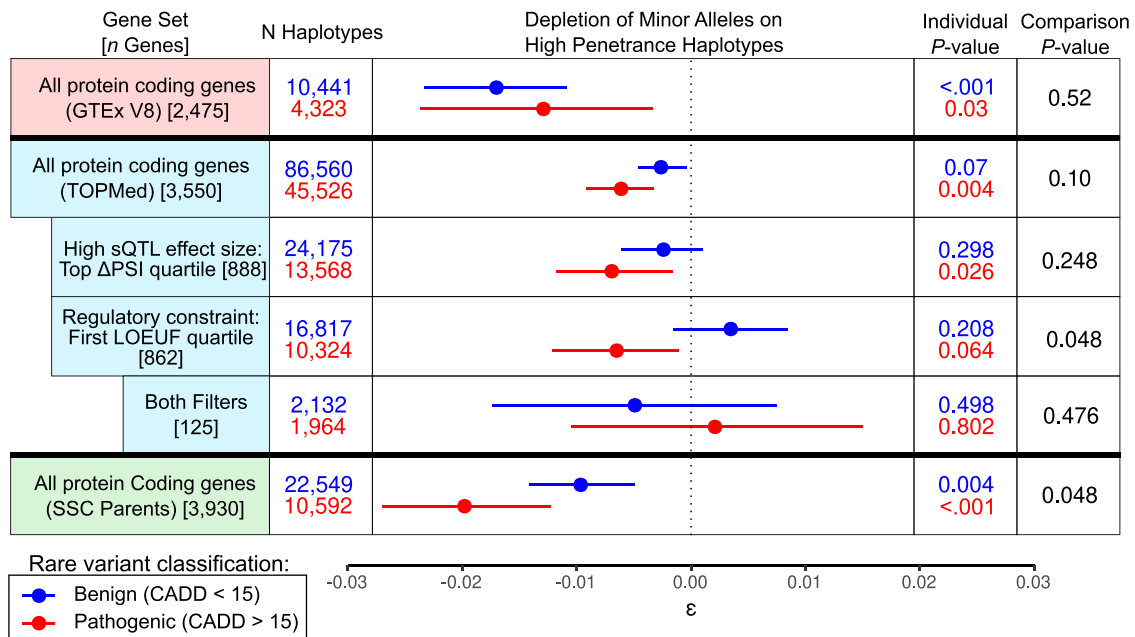
Next, we increased our power to detect evidence of selection against putative high-penetrance haplotypes by using population-phased WGS data from 44,634 European-American ancestry individuals in 19 TOPMed cohorts, post-filtering (*Methods*, Table 1, and Supplementary Fig. 6). The large sample size in TOPMed allowed us to limit the analysis to exonic variants with 10 or fewer occurrences (excluding singletons due to limitations of population-based phasing), or <0.0213% minor allele frequency. With the same set of  $\psi$ QTLs from GTEx, we identified the haplotype of 38,869 rare alleles that fell in primary and secondary sExons. Across all protein-coding

genes and rare alleles, we observed a modest but significant overall depletion of high-penetrance haplotypes than expected ( $\epsilon = -0.0037$ , Poisson-binomial  $P = 3.43 \times 10^{-4}$ ). Haplotypes with putatively deleterious rare alleles had some indication of being more depleted than those with non-deleterious rare alleles, but not to a degree that reached statistical significance ( $P = 0.100$ , Fig. 5). However, we hypothesized that this result would be more pronounced in genes with stronger  $\psi$ QTLs, as well as genes known to be intolerant to loss of function variation. When focusing on genes with stronger  $\psi$ QTLs where the  $\Delta$ PSI score was in the top quartile ( $\Delta$ PSI > 0.076), the difference was again not significant ( $P = 0.248$ ).

**Table 1.** Properties of 3 WGS datasets used in this study.

	GTE <sub>x</sub>	TOPMed	SSC—parents
N donors	714	44,634	4,731
Phasing method	Population based and 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
Singletons included	Yes, in calls with RNA-seq read backing. Otherwise, no	No	Yes
Rare variant allele frequency cutoff	0.5% MAF in gnomad. (No count cutoff due to the relative small size of the GTE <sub>x</sub> WGS dataset)	Appears 10 or fewer times (i.e. 0.0257% MAF)	Appears ≤ 3 times (i.e. 0.126% MAF)

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



**Fig. 5.** Rare alleles carried in predicted high penetrance  $\psi$ QTL configurations in GTE<sub>x</sub>, TOPMed, and SSC parents. We tested for deviation in the frequencies of coding allele— $\psi$ QTL configurations across all protein-coding genes with a significant  $\psi$ QTL. A negative value of  $\epsilon$  indicates fewer haplotypes than expected given the population's  $\psi$ QTL allele frequencies. Individual P-values and 95% confidence intervals were generated using our approximation of the Poisson-binomial CDF, with 1,000 bootstraps. Comparison P-values were generated with 1,000 bootstraps.

However, when quantifying gene constraint with the loss-of-function observed/expected upper bound fraction (LOEUF) score (Karczewski et al. 2020) and limiting to genes in the first quartile among sGenes (LOEUF < 0.460), we detected a significant difference in high-penetrance haplotype depletion between the 2 groups ( $P = 0.048$ ), suggesting that splicing may play a greater role in modifying penetrance in genes known to be constrained. Finally, while we would expect to see the greatest effects of purifying selection among constrained genes with strong  $\psi$ QTLs, the small number of such genes limits our power and no significant association was detected ( $P = 0.982$ ). We found that across genes in general,  $\Delta$ PSI and LOEUF were positively correlated, so genes with high  $\Delta$ PSI and low LOEUF were uncommon (Supplementary Fig. 7c). While subtle, these results suggest that deleterious rare alleles are more likely to be carried on exons that are skipped due to the effects of common regulatory variants, especially in constrained genes.

Next, we wanted to explore if any genes or classes of genes drove our observation of high-penetrance haplotype depletion. To this end, using the same TOPMed data, we tested for nonrandom haplotype combinations on a gene-by-gene basis, instead of pooling haplotypes across all genes as in the previous approach. For 2,396 genes with more than 10  $\psi$ QTL-coding variant haplotypes across all available individuals, we ran a Poisson-binomial test for high-penetrance haplotype depletion (Supplementary Fig. 8). We

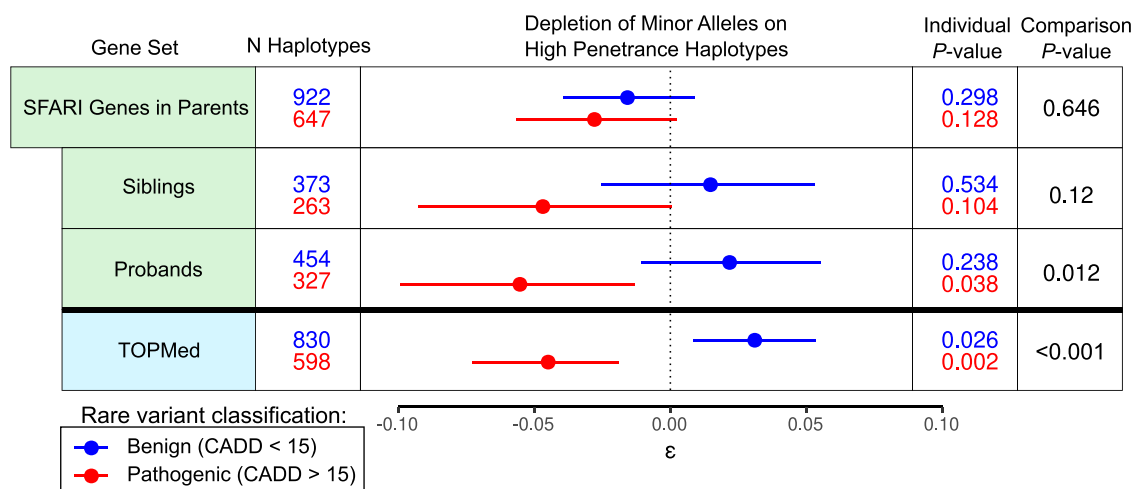
observed little signal, with approximately equal numbers of genes with enrichment and depletion of high- and low-penetrance haplotypes. However, only 411 of the genes had more than 30 deleterious allele haplotypes, indicating that our power is quite limited. Thus, our results indicate that observing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

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

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

First, we sought to replicate the depletion of potential high-penetrance haplotypes observed in TOPMed, using SSC parents, who are a cohort of unrelated individuals, phenotypically healthy but with potential enrichment of ASD risk variants due to having a





**Fig. 6.**  $\psi$ QTL haplotype configurations in Autism Spectrum Disorder implicated genes in ASD families. We tested for deviation in the frequencies of high penetrance variant— $\psi$ QTL configurations in ASD-implicated genes in parents, probands, and unaffected siblings in SSC families.

child with ASD. We analyzed all genes with a  $\psi$ QTL in GTEx, limiting our analysis to coding alleles with 3 or fewer occurrences across all parents, and removing genes with an unusually high number of rare variant haplotypes (Table 1, Supplementary Fig. 6). Singleton variants were included since their haplotype can be confidently resolved using phasing by transmission. We recapitulated the patterns observed in TOPMed, with a significant depletion of high-penetrance haplotypes with deleterious rare alleles ( $\epsilon = -0.019$ , Poisson-binomial  $P = 2.11 \times 10^{-8}$ ), with high-penetrance haplotypes carrying deleterious rare alleles more depleted than those carrying non-deleterious rare alleles (comparison  $P$ -value = 0.042, Fig. 5).

Next, we sought to analyze potential splicing modifiers of the penetrance of disease-causing alleles in SSC by focusing on rare inherited variants in ASD-implicated genes. These alleles, while potentially contributing to ASD in the proband, are also carried on the same haplotypic background by a healthy parent and often a healthy sibling. Thus, both increased or decreased penetrance  $\psi$ QTL configurations could be possible (Supplementary Fig. 9a). To test this, we analyzed deviation in haplotype frequencies in parents, probands, and siblings, among the 218 out of the 1,010 genes implicated in ASD risk according to SFARI Gene (Banerjee-Basu and Packer, 2010) that also had a  $\psi$ QTL. No significant deviation was detected in SSC parents ( $\epsilon = -0.0278$ ,  $P = 0.122$ ). Interestingly, across probands and unaffected siblings, we found that putatively highly penetrant haplotypes with deleterious coding alleles were depleted ( $\epsilon = -0.055$  and  $-0.047$ ,  $P = 0.038$  and  $0.104$ , respectively). While it seems counterintuitive to see depletion of penetrant haplotypes in individuals with ASD, we reason that this penetrance reducing effect may be acting to protect parents from developing phenotypes of ASD. We find that the SFARI genes tend to be highly constrained, compared to all protein-coding genes (Supplementary Fig. 9b) (Neale et al. 2012), and that these same alleles were also highly depleted among unrelated individuals in TOPMed (Fig. 6), further corroborating the overall observed pattern of selection for penetrance reducing haplotype combinations.

## Discussion

In this study, we have expanded our model of cis-regulatory alleles as modifiers of penetrance of coding variants (Castel et al. 2018) to directly consider splice-regulatory alleles as potential additional drivers. We first show that individuals carrying potentially deleterious rare mutations at variably spliced exons tend to use those exons in

transcripts less frequently. This observation could indicate that the penetrance of these rare alleles is reduced by their exclusion from transcripts. However, this approach does not reveal the reason. One approach to potentially shed light on this would be analysis of allele-specific transcript structure, but this is not possible with short read RNA-sequencing. However, our model could be tested in larger future studies with long-read sequencing technology (Glinos et al. 2021).

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

Through this research, we derived a novel approach for calculating the cumulative distribution function of the Poisson-binomial distribution, as well as a metric for evaluating a dataset's deviation from an expected distribution or difference between 2 datasets (the comparison test). This method is well suited for very large datasets, and has further applications in genetic and nongenetic analyses where data are expected to follow the Poisson-binomial.

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

An alternative approach is to study regulatory variation underlying modified penetrance in disease cohorts with well annotated disease-causing variants, linking haplotype patterns with phenotype variation between and within families. The Simons Simplex Collection had some limitations in this respect: Most ASD-contributing rare variants are not known, and the trait is highly polygenic, making it difficult to compare penetrance of variants in the same gene between families. Furthermore, in simplex families,

many causal variants are de novo, but their total number is small for statistical analysis. In the future, large ASD studies with multi-plex families may better capture ASD instances with heritable variant etiology. Furthermore, experimental validation, for example with genome editing, may be a fruitful approach.

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

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

## Data availability

All code used to perform analyses and generate figures is available at [https://github.com/jeinson/mp\\_manuscript](https://github.com/jeinson/mp_manuscript). Qualified researchers requiring data access can apply for GTEx, TOPMed data through dbGaP, and SSC data through the Simons Foundation. We include a function to generate simulated data in the stampen R package (<https://github.com/jeinson/stampen>). PSI and  $\psi$ QTLs from GTEx V8 can be downloaded from the repository for Einson *et al.* (2022) at <https://zenodo.org/record/7275062#.Y9gc0OzMjf0>.

[Supplemental material](#) available at GENETICS online.

## Acknowledgements

JE thanks members of the Lappalainen lab for the thoughtful discussions and feedback throughout this project.

Molecular data for the Trans-Omics for Precision Medicine (TOPMed) program were supported by the National Heart, Lung and Blood Institute (NHLBI). Whole genome sequencing (WGS) for the Trans-Omics for Precision Medicine (TOPMed) program was supported by the National Heart, Lung and Blood Institute (NHLBI). Core support including centralized genomic read mapping and genotype calling, along with variant quality metrics and filtering, were provided by the TOPMed Informatics Research Center (3R01HL-117626-02S1; contract HHSN268201800002I). Core support including phenotype harmonization, data management, sample-identity QC, and general program coordination were provided by the TOPMed Data Coordinating Center (R01HL-120393; U01HL-120393; contract HHSN268201800001I) and TOPMed MESA Multi-Omics (HHSN268201500003I/HSN26800004).

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

## Funding

- 1) Genome Sequencing for NHLBI TOPMed: Women's Health Initiative (phs001237) was performed at Broad Institute Genomics Platform (HHSN268201500014C).
- 2) Genome Sequencing for NHLBI TOPMed: Genetic Epidemiology of COPD Study (phs000951) was performed at Northwest Genomics Center (3R01HL089856-08S1).
- 3) Genome Sequencing for NHLBI TOPMed: Atherosclerosis Risk in Communities Study VTE cohort (phs001211) was performed at Baylor College of Medicine Human Genome Sequencing Center (3U54HG003273-12S2/HHSN268201500015C).
- 4) Genome Sequencing for NHLBI TOPMed: Framingham Heart Study (phs000974) was performed at Broad Institute Genomics Platform (HHSN268201600034I).
- 5) Genome Sequencing for NHLBI TOPMed: My Life, Our Future: Genotyping for Progress in Hemophilia (phs001515) was performed at Baylor College of Medicine Human Genome Sequencing Center (HHSN268201600033I).
- 6) Genome Sequencing for NHLBI TOPMed: Mount Sinai BioMe Biobank (phs001644) was performed at McDonnell Genome Institute (3UM1HG008853-01S2).
- 7) Genome Sequencing for NHLBI TOPMed: Cardiovascular Health Study (phs001368) was performed at Broad Institute Genomics Platform (HHSN268201600034I).
- 8) Genome Sequencing for NHLBI TOPMed: Multi-Ethnic Study of Atherosclerosis (phs001416) was performed at Broad Institute Genomics Platform (HHSN268201600034I, 3U54HG003067-13S1).
- 9) Genome Sequencing for NHLBI TOPMed: Coronary Artery Risk Development in Young Adults (phs001612) was performed at Baylor College of Medicine Human Genome Sequencing Center (HHSN268201600033I).
- 10) Genome Sequencing for NHLBI TOPMed: Mayo Clinic Venous Thromboembolism Study (phs001402) was performed at Baylor College of Medicine Human Genome Sequencing Center (3U54HG003273-12S2/HHSN268201500015C).
- 11) Genome Sequencing for NHLBI TOPMed: Lung Tissue Research Consortium (phs001662) was performed at Broad Institute Genomics Platform (HHSN268201600034I).
- 12) Genome Sequencing for NHLBI TOPMed: The Vanderbilt University BioVU Atrial Fibrillation Genetics Study (phs001624) was performed at Baylor College of Medicine Human Genome Sequencing Center (3UM1HG008898-01S3).
- 13) Genome Sequencing for NHLBI TOPMed: Vanderbilt Genetic Basis of Atrial Fibrillation (phs001032) was performed at Broad Institute Genomics Platform (3R01HL092577-06S1).
- 14) Genome Sequencing for NHLBI TOPMed: Hispanic Community Health Study—Study of Latinos (phs001395) was performed at Baylor College of Medicine Human Genome Sequencing Center (HHSN268201600033I).
- 15) Genome Sequencing for NHLBI TOPMed: Severe Asthma Research Program (phs001446) was performed at New York Genome Center Genomics (HHSN268201500016C).
- 16) Genome Sequencing for NHLBI TOPMed: Massachusetts General Hospital Atrial Fibrillation Study (phs001062) was performed at Broad Institute Genomics Platform (3U54HG003067-12S2/3U54HG003067-13S1; 3U54HG003067-12S2/3U54HG003067-13S1; 3UM1HG008895-01S2).
- 17) Genome Sequencing for NHLBI TOPMed: Heart and Vascular Health Study (phs000993) was performed at Broad Institute Genomics Platform (3R01HL092577-06S1).
- 18) Genome Sequencing for NHLBI TOPMed: Groningen Genetics of Atrial Fibrillation Study (phs001725) was performed at Baylor College of Medicine Human Genome Sequencing Center (3UM1HG008898-01S3).
- 19) Genome Sequencing for NHLBI TOPMed: Genetics of Cardiometabolic Health in the Amish (phs000956) was performed at Broad Institute Genomics Platform (3R01HL121007-01S1).

JE and TL were supported by National Institutes of Health (NIH) grants R01GM122924 and R01MH106842. PM was supported by National Institute of General Medical Sciences (NIGMS) grant R01GM140287. II was supported by the Simons Center for Quantitative Biology at Cold Spring Harbor Laboratory, SFARI grants SF497800, SF677963, and SF666590, and the Centers for Common Disease Genomics grant (UM1 HG008901). Support for title page creation and format was provided by AuthorArranger, a tool developed at the National Cancer Institute.

## Conflicts of interest statement

TL is a paid advisor to GSK, Pfizer, Goldfinch Bio, and Variant Bio and has equity in Variant Bio.

## References

- Alasoo K, Rodrigues J, Danesh J, Freitag DF, Paul DS, Gaffney DJ. Genetic effects on promoter usage are highly context-specific and contribute to complex traits. *eLife*. 2019;8:e41673. doi:10.7554/eLife.41673.
- An J-Y, Lin K, Zhu L, Werling DM, Dong S, Brand H, Wang HZ, Zhao X, Schwartz GB, Collins RL, et al. Genome-wide de novo risk score implicates promoter variation in autism spectrum disorder. *Science*. 2018;362(6420):eaat6576. doi:10.1126/science.aat6576.
- Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, Chakravarti A, Clark AG, Donnelly P, Eichler EE, et al. A global reference for human genetic variation. *Nature*. 2015;526(7571):68–74. doi:10.1038/nature15393.
- Banerjee-Basu S, Packer A. SFARI gene: an evolving database for the autism research community | disease models & mechanisms. *Dis Model Mech*. 2010;3(3–4):133–135. <https://journals.biologists.com/dmm/article/3/3-4/133/2349/SFARI-Gene-an-evolving-database-for-the-autism>
- Castel SE, Cervera A, Mohammadi P, Aguet F, Reverter F, Wolman A, Guigo R, Iossifov I, Vasileva A, Lappalainen T. Modified penetrance of coding variants by cis-regulatory variation contributes to disease risk. *Nat Genet*. 2018;50(9):1327–1334. doi:10.1038/s41588-018-0192-y.
- Castel SE, Mohammadi P, Chung WK, Shen Y, Lappalainen T. Rare variant phasing and haplotypic expression from RNA sequencing with phASER. *Nat Commun*. 2016;7(1):12817. 08/online. doi:10.1038/ncomms12817.
- Chen W-C, Maitra R. R: EM algorithm for model-based clustering of finite mixture Gaussian distribution [Internet]. 2015 (accessed 2022 Jun 24). <https://search.r-project.org/CRAN/refmans/EMCluster/html/00Index.html>.
- Chen R, Shi L, Hakenberg J, Naughton B, Sklar P, Zhang J, Zhou H, Tian L, Prakash O, Lemire M, et al. Analysis of 589,306 genomes identifies individuals resilient to severe Mendelian childhood diseases. *Nat Biotechnol*. 2016;34(5):531–538. doi:10.1038/nbt.3514.
- Chiang AH, Chang J, Wang J, Vitkup D. Exons as units of phenotypic impact for truncating mutations in autism. *Mol Psychiatry*. 2021;26(5):1685–1695. doi:10.1038/s41380-020-00876-3.
- Consortium TGte. The GTEx consortium atlas of genetic regulatory effects across human tissues. *Science*. 2020;369(6509):1318–1330. doi:10.1126/science.aaz1776.
- Cooper DN, Krawczak M, Polychronakos C, Tyler-Smith C, Kehrer-Sawatzki H. Where genotype is not predictive of phenotype: towards an understanding of the molecular basis of reduced penetrance in human inherited disease. *Hum Genet*. 2013;132(10):1077–1130. doi:10.1007/s00439-013-1331-2.
- Cummings BB, Karczewski KJ, Kosmicki JA, Seaby EG, Watts NA, Singer-Berk M, Mudge JM, Karjalainen J, Satterstrom FK, O'Donnell-Luria AH, et al. Transcript expression-aware annotation improves rare variant interpretation. *Nature*. 2020;581(7809):452–458. doi:10.1038/s41586-020-2329-2.
- Delaneau O, Ongen H, Brown AA, Fort A, Panousis NI, Dermitzakis ET. A complete tool set for molecular QTL discovery and analysis. *Nat Commun*. 2017;8(1):1–7. doi:10.1038/ncomms15452.
- Dewey FE, Murray MF, Overton JD, Habegger L, Leader JB, Fetterolf SN, O'Dushlaine C, Van Hout CV, Staples J, Gonzaga-Jauregui C, et al. Distribution and clinical impact of functional variants in 50,726 whole-exome sequences from the DiscovEHR study. *Science*. 2016;354(6319):aaf6814. doi:10.1126/science.aaf6814.
- Einson J, Minaeva M, Rafi F, Lappalainen T. The impact of genetically controlled splicing on exon inclusion and protein structure. *bioRxiv*. 2022: 518915. <https://doi.org/10.1101/2022.12.05.518915>, 22 December 2022, preprint: not peer reviewed.
- Emison ES, McCallion AS, Kashuk CS, Bush RT, Grice E, Lin S, Portnoy ME, Cutler DJ, Green ED, Chakravarti A. A common sex-dependent mutation in a RET enhancer underlies Hirschsprung disease risk. *Nature*. 2005;434(7035):857–863. doi:10.1038/nature03467.
- Fahed AC, Wang M, Homburger JR, Patel AP, Bick AG, Neben CL, Lai C, Brockman D, Philippakis A, Ellinor PT. Polygenic background modifies penetrance of monogenic variants for tier 1 genomic conditions. *Nat Commun*. 2020;11(1):3635. doi:10.1038/s41467-020-17374-3.
- Forrest IS, Chaudhary K, Vy HMT, Petrazzini BO, Bafna S, Jordan DM, Rocheleau G, Loos RJF, Nadkarni GN, Cho JH, et al. Population-based penetrance of deleterious clinical variants. *JAMA*. 2022;327(4):350–359. doi:10.1001/jama.2021.23686.
- Galinsky KJ, Bhatia G, Loh P-R, Georgiev S, Mukherjee S, Patterson NJ. Fast principal-component analysis reveals convergent evolution of ADH1B in Europe and East Asia. *The American Journal of Human Genetics*. 2016;98(3):456–472. doi:10.1016/j.ajhg.2015.12.022.
- Garrido-Martín D, Borsari B, Calvo M, Reverter F, Guigó R. Identification and analysis of splicing quantitative trait loci across multiple tissues in the human genome. *Nat Commun*. 2021;12(1):727. doi:10.1038/s41467-020-20578-2.
- Gettler K, Levantovsky R, Moscati A, Giri M, Wu Y, Hsu N-Y, Chuang L-S, Sazonovs A, Venkateswaran S, Korie U, et al. Common and rare variant prediction and penetrance of IBD in a large, multi-ethnic, health system-based biobank cohort. *Gastroenterology*. 2021;160(5):1546–1557. doi:10.1053/j.gastro.2020.12.034.
- Glinos DA, Garborcauskas G, Hoffman P, Ehsan N, Jiang L, Gokden A, Dai X, Aguet F, Brown KL, Garimella K, et al. Transcriptome variation in human tissues revealed by long-read sequencing. *bioRxiv*. 2021:427687. <https://doi.org/10.1101/2021.01.22.427687>, 31 May 2022, preprint: not peer reviewed.
- González J, Wiberg M, von Davier AA. A note on the Poisson's binomial distribution in item response theory. *Appl Psychol Meas*. 2016;40(4):302–310. doi:10.1177/0146621616629380.
- Hong Y. On computing the distribution function for the Poisson binomial distribution. *Comput Stat Data Anal*. 2013;59:41–51. doi:10.1016/j.csda.2012.10.006.
- Iossifov I, O'Roak BJ, Sanders SJ, Ronemus M, Krumm N, Levy D, Stessman HA, Witherspoon KT, Vives L, Patterson KE, et al. The contribution of de novo coding mutations to autism spectrum disorder. *Nature*. 2014;515(7526):216–221. doi:10.1038/nature13908.
- Iossifov I, Ronemus M, Levy D, Wang Z, Hakker I, Rosenbaum J, Yamrom B, Lee Y-h, Narzisi G, Leotta A, et al. De novo gene disruptions in children on the autistic spectrum. *Neuron*. 2012;74(2):285–299. doi:10.1016/j.neuron.2012.04.009.

- IPSA-nf [Internet]. Guigo Lab; 2020 [accessed 2021 Aug 3]. <https://github.com/guigolab/ipsa-nf>.
- Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, Collins RL, Laricchia KM, Ganna A, Birnbaum DP, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*. 2020;581(7809):434–443. doi:10.1038/s41586-020-2308-7.
- Keren H, Lev-Maor G, Ast G. Alternative splicing and evolution: diversification, exon definition and function. *Nat Rev Genet*. 2010; 11(5):345–355. doi:10.1038/nrg2776.
- Kerimov N, Hayhurst JD, Peikova K, Manning JR, Walter P, Kolberg L, Samoviča M, Sakthivel MP, Kuzmin I, Trevanion SJ, et al. eQTL catalogue: a compendium of uniformly processed human gene expression and splicing QTLs. *bioRxiv*. 2020:924266. <https://doi.org/10.1101/2020.01.29.924266>, 29 January 2020, preprint: not peer reviewed.
- Li YI, van de Geijn B, Raj A, Knowles DA, Petti AA, Golan D, Gilad Y, Pritchard JK. RNA splicing is a primary link between genetic variation and disease. *Science*. 2016;352(6285):600–604. doi:10.1126/science.aad9417.
- Loh P-R, Danecek P, Palamara PF, Fuchsberger C, A Reshef Y, K Finucane H, Schoenherr S, Forer L, McCarthy S, Abecasis GR, et al. Reference-based phasing using the haplotype reference consortium panel. *Nat Genet*. 2016;48(11):1443–1448. doi:10.1038/ng.3679.
- Maya I, Sharony R, Yacobson S, Kahana S, Yeshaya J, Tenne T, Agmon-Fishman I, Cohen-Vig L, Goldberg Y, Berger R, et al. When genotype is not predictive of phenotype: implications for genetic counseling based on 21,594 chromosomal microarray analysis examinations. *Genet Med*. 2018;20(1):128–131. doi:10.1038/gim.2017.89.
- Milne RL, Antoniou AC. Genetic modifiers of cancer risk for BRCA1 and BRCA2 mutation carriers. *Ann Oncol*. 2011;22(Suppl 1): i11–i17. doi:10.1093/annonc/mdq660.
- Morris JA, Kemp JP, Youlten SE, Laurent L, Logan JG, Chai RC, Vulpescu NA, Forgetta V, Kleinman A, Mohanty ST, et al. An atlas of genetic influences on osteoporosis in humans and mice. *Nat Genet*. 2019;51(2):258–266. doi:10.1038/s41588-018-0302-x.
- Neale BM, Kou Y, Liu L, Ma'ayan A, Samocha KE, Sabo A, Lin C-F, Stevens C, Wang L-S, Makarov V, et al. Patterns and rates of exonic de novo mutations in autism spectrum disorders. *Nature*. 2012; 485(7397):242–245. doi:10.1038/nature11011.
- Noble JD, Balmant KM, Dervinis C, de los Campos G, Resende MFRJ, Kirst M, et al. The genetic regulation of alternative splicing in *Populus deltoides*. *Front Plant Sci*. 2020;11:590. doi:10.3389/fpls.2020.00590.
- O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, Traglia M, Huang J, Huffman JE, Rudan I, et al. A general approach for haplotype phasing across the full spectrum of relatedness. *PLoS Genet*. 2014;10(4):e1004234. doi:10.1371/journal.pgen.1004234.
- Ongen H, Dermitzakis ET. Alternative splicing QTLs in European and African populations. *Am J Hum Genet*. 2015;97(4):567–575. doi:10.1016/j.ajhg.2015.09.004.
- Pervouchine DD, Knowles DG, Guigó R. Intron-centric estimation of alternative splicing from RNA-seq data. *Bioinformatics*. 2013; 29(2):273–274. doi:10.1093/bioinformatics/bts678.
- Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. CADD: predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids Res*. 2019;47(D1):D886–D894. doi:10.1093/nar/gky1016.
- Ruzzo EK, Pérez-Cano L, Jung J-Y, Wang L, Kashef-Haghighi D, Hartl C, Singh C, Xu J, Hoekstra JN, Leventhal O, et al. Inherited and de novo genetic risk for autism impacts shared networks. *Cell*. 2019; 178(4):850–866.e26. doi:10.1016/j.cell.2019.07.015.
- Sanders SJ, He X, Willsey AJ, Ercan-Sencicek AG, Samocha KE, Cicek AE, Murtha MT, Bal VH, Bishop SL, Dong S, et al. Insights into autism spectrum disorder genomic architecture and biology from 71 risk loci. *Neuron*. 2015;87(6):1215–1233. doi:10.1016/j.neuron.2015.09.016.
- Sanders SJ, Murtha MT, Gupta AR, Murdoch JD, Raubeson MJ, Willsey AJ, Ercan-Sencicek AG, DiLullo NM, Parikshak NN, Stein JL, et al. De novo mutations revealed by whole-exome sequencing are strongly associated with autism. *Nature*. 2012;485(7397): 237–241. doi:10.1038/nature10945.
- Shawky RM. Reduced penetrance in human inherited disease. *Egyptian Journal of Medical Human Genetics*. 2014;15(2): 103–111. doi:10.1016/j.ejmhg.2014.01.003.
- Turner TN, Hormozdiari F, Duyzend MH, McClymont SA, Hook PW, Iossifov I, Raja A, Baker C, Hoekzema K, Stessman HA, et al. Genome sequencing of autism-affected families reveals disruption of putative noncoding regulatory DNA. *The American Journal of Human Genetics*. 2016;98(1):58–74. doi:10.1016/j.ajhg.2015.11.023.
- Wang YH. On the number of successes in independent trials. *Statistica Sinica*. 1993;3(2):295–312.
- Yoon S, Munoz A, Yamrom B, Lee Y, Andrews P, Marks S, Wang Z, Reeves C, Winterkorn L, Krieger AM, et al. Rates of contributory de novo mutation in high and low-risk autism families. *Commun Biol*. 2021;4(1):1–10. doi:10.1038/s42003-021-02533-z.

Editor: H. Zhao