

# Quantitative prior probabilities for disease-causing variants reveal the top genetic contributors in inborn errors of immunity

Dylan Lawless<sup>\*1</sup>

<sup>1</sup>Department of Intensive Care and Neonatology, University Children's Hospital Zürich,  
University of Zürich, Switzerland.

April 8, 2025

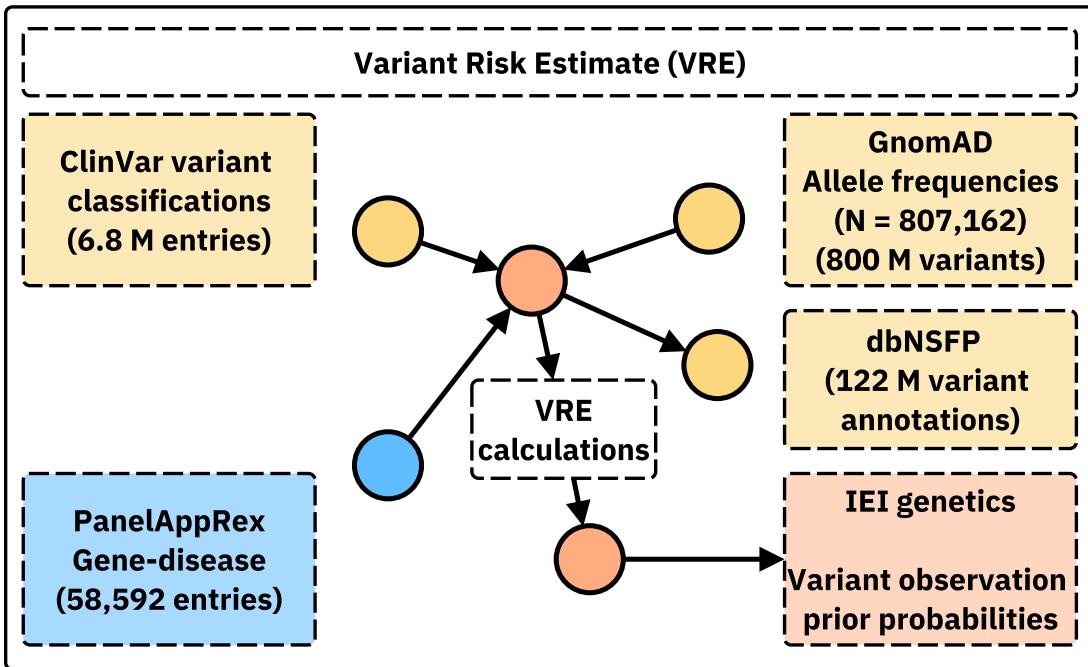
## Abstract

We present a novel framework for quantifying the prior probability of observing disease-associated variants in any gene for a given phenotype. By integrating large-scale genomic annotations, including population allele frequencies and ClinVar variant classifications, with Hardy-Weinberg-based calculations, our method estimates per-variant observation probabilities under autosomal dominant (AD), autosomal recessive (AR), and X-linked modes of inheritance. Applied to 557 genes implicated in primary immunodeficiency and inflammatory disease, our approach generated 54,814 variant probabilities. First, these detailed, pre-calculated results provide robust priors for any gene-disease combination. Second, a score positive total metric summarises the aggregate pathogenic burden, serving as an indicator of the likelihood of observing a patient with the disease and reflecting genetic constraint. Validation in *NFKB1* (AD) and *CFTR* (AR) disorders confirmed close concordance between predicted and observed case counts. The resulting datasets, available in both machine-readable and human-friendly formats, support Bayesian variant interpretation and clinical decision-making.<sup>1</sup>

---

<sup>\*</sup>Addresses for correspondence: [Dylan.Lawless@kispi.uzh.ch](mailto:Dylan.Lawless@kispi.uzh.ch)

<sup>1</sup> **Availability:** This data is integrated in public panels at <https://iei-genetics.github.io>. The source code and data are accessible as part of the variant risk estimation project at [https://github.com/DylanLawless/var\\_risk\\_est](https://github.com/DylanLawless/var_risk_est). The variant-level data is available from the Zenodo repository: <https://doi.org/10.5281/zenodo.15111583> (VarRiskEst PanelAppRex ID 398 gene variants.tsv). VarRiskEst is available under the MIT licence.



18

<sup>19</sup> **Acronyms**

<sup>20</sup> <b>ACMG</b> American College of Medical Genetics and Genomics.....	<sup>24</sup>
<sup>21</sup> <b>ACAT</b> Aggregated Cauchy Association Test .....	<sup>24</sup>
<sup>22</sup> <b>AD</b> Autosomal Dominant.....	<sup>4</sup>
<sup>23</sup> <b>ANOVA</b> Analysis of Variance .....	<sup>11</sup>
<sup>24</sup> <b>AR</b> Autosomal Recessive .....	<sup>4</sup>
<sup>25</sup> <b>BMF</b> Bone Marrow Failure.....	<sup>18</sup>
<sup>26</sup> <b>CD</b> Complement Deficiencies .....	<sup>19</sup>
<sup>27</sup> <b>CI</b> Confidence Interval.....	<sup>14</sup>
<sup>28</sup> <b>CF</b> Cystic Fibrosis .....	<sup>10</sup>
<sup>29</sup> <b>CFTR</b> Cystic Fibrosis Transmembrane Conductance Regulator.....	<sup>5</sup>
<sup>30</sup> <b>CVID</b> Common Variable Immunodeficiency .....	<sup>8</sup>
<sup>31</sup> <b>dbNSFP</b> database for Non-Synonymous Functional Predictions .....	<sup>5</sup>
<sup>32</sup> <b>GE</b> Genomics England .....	<sup>5</sup>
<sup>33</sup> <b>gnomAD</b> Genome Aggregation Database .....	<sup>5</sup>
<sup>34</sup> <b>HGVS</b> Human Genome Variation Society.....	<sup>5</sup>
<sup>35</sup> <b>HPC</b> High-Performance Computing.....	<sup>8</sup>
<sup>36</sup> <b>HWE</b> Hardy-Weinberg Equilibrium .....	<sup>4</sup>
<sup>37</sup> <b>IEI</b> Inborn Errors of Immunity.....	<sup>4</sup>
<sup>38</sup> <b>InDel</b> Insertion/Deletion .....	<sup>5</sup>
<sup>39</sup> <b>IUIS</b> International Union of Immunological Societies .....	<sup>5</sup>
<sup>40</sup> <b>LD</b> Linkage Disequilibrium .....	<sup>21</sup>
<sup>41</sup> <b>LOEUF</b> Loss-Of-function Observed/Expected Upper bound Fraction .....	<sup>11</sup>
<sup>42</sup> <b>LOF</b> Loss-of-Function .....	<sup>18</sup>
<sup>43</sup> <b>MOI</b> Mode of Inheritance .....	<sup>4</sup>
<sup>44</sup> <b>NFKB1</b> Nuclear Factor Kappa B Subunit 1 .....	<sup>5</sup>
<sup>45</sup> <b>OMIM</b> Online Mendelian Inheritance in Man .....	<sup>22</sup>
<sup>46</sup> <b>PID</b> Primary Immunodeficiency .....	<sup>4</sup>
<sup>47</sup> <b>PPI</b> Protein-Protein Interaction .....	<sup>5</sup>
<sup>48</sup> <b>SNV</b> Single Nucleotide Variant .....	<sup>4</sup>
<sup>49</sup> <b>SKAT</b> Sequence Kernel Association Test.....	<sup>24</sup>
<sup>50</sup> <b>STRINGdb</b> Search Tool for the Retrieval of Interacting Genes/Proteins.....	<sup>5</sup>
<sup>51</sup> <b>HSD</b> Honestly Significant Difference .....	<sup>11</sup>
<sup>52</sup> <b>UMAP</b> Uniform Manifold Approximation and Projection .....	<sup>18</sup>
<sup>53</sup> <b>UniProt</b> Universal Protein Resource .....	<sup>5</sup>
<sup>54</sup> <b>VEP</b> Variant Effect Predictor.....	<sup>5</sup>
<sup>55</sup> <b>XL</b> X-Linked .....	<sup>4</sup>

## 92 1 Introduction

93 In this study, we focused on reporting the probability of disease observation through  
94 genome-wide assessments of gene-disease combinations. Our central hypothesis was  
95 that by using highly curated annotation data including population allele frequen-  
96 cies, disease phenotypes, Mode of Inheritance (MOI) patterns, and variant classi-  
97 fications and by applying rigorous calculations based on Hardy-Weinberg Equilib-  
98 rium (HWE), we could accurately estimate the expected probabilities of observing  
99 disease-associated variants. Among other benefits, this knowledge can be used to  
100 derive genetic diagnosis confidence by incorporating these new priors.

101 In this report, we focused on known Inborn Errors of Immunity (IEI) genes, also re-  
102 ferred to as the Primary Immunodeficiency (PID) or Monogenic Inflammatory Bowel  
103 Disease genes (1–3) to validate our approach and demonstrate its clinical relevance.  
104 This application to a well-established genotype-phenotype set, comprising over 500  
105 gene-disease associations, underscores its utility (1).

106 Quantifying the risk that a newborn inherits a disease-causing variant is a fun-  
107 damental challenge in genomics. Classical statistical approaches grounded in HWE  
108 (4; 5) have long been used to calculate genetic MOI probabilities for Single Nucleotide  
109 Variant (SNV)s. However, applying these methods becomes more complex when ac-  
110 counting for different MOI, such as Autosomal Recessive (AR) versus Autosomal  
111 Dominant (AD) or X-Linked (XL) disorders. In AR conditions, for example, the  
112 occurrence probability must incorporate both the homozygous state and compound  
113 heterozygosity, whereas for AD and XL disorders, a single pathogenic allele is suffi-  
114 cient to cause disease. Advances in genetic research have revealed that MOI can be  
115 even more complex (6). Mechanisms such as dominant negative effects, haploinsuffi-  
116 ciency, mosaicism, and digenic or epistatic interactions can further modulate disease  
117 risk and clinical presentation, underscoring the need for nuanced approaches in risk  
118 estimation. Karczewski et al. (7) made significant advances; however, the remain-  
119 ing challenge lay in applying the necessary statistical genomics data across all MOI  
120 for any gene-disease combination Similar approaches have been reported for disease  
121 such Wilson disease, Mucopolysaccharidoses, Primary ciliary dyskinesia, and treat-  
122 able metabolic diseases, (8; 9), as reviewed by Hannah et al. (10).

123 To our knowledge all approaches to date have been limited to single MOI, specific  
124 to the given disease, or restricted to a small number of genes. We argue that our  
125 integrated approach is highly powerful because the resulting probabilities can serve  
126 as informative priors in a Bayesian framework for variant and disease probability  
127 estimation; a perspective that is often overlooked in clinical and statistical genetics.  
128 Such a framework not only refines classical HWE-based risk estimates but also has  
129 the potential to enrich clinicians' understanding of what to expect in a patient and to  
130 enhance the analytical models employed by bioinformaticians. The dataset also holds  
131 value for AI and reinforcement learning applications, providing an enriched version of  
132 the data underpinning frameworks such as AlphaFold (11) and AlphaMissense (12).

133 We introduced PanelAppRex to aggregate gene panel data from multiple sources,

including Genomics England (GE) PanelApp, ClinVar, and Universal Protein Resource (UniProt), thereby enabling advanced natural searches for clinical and research applications (2; 3; 13; 14). It automatically retrieves expert-curated panels, such as those from the NHS National Genomic Test Directory and the 100,000 Genomes Project, and converts them into machine-readable formats for rapid variant discovery and interpretation. We used PanelAppRex to label disease-associated variants. We also integrate key statistical genomic resources. The gnomAD v4 dataset compiles data from 807,162 individuals, encompassing over 786 million SNVs and 122 million Insertion/Deletion (InDel)s with detailed population-specific allele frequencies (7). database for Non-Synonymous Functional Predictions (dbNSFP) provides functional predictions for over 120 million potential non-synonymous and splicing-site SNVs, aggregating scores from 33 sources alongside allele frequencies from major populations (15). ClinVar offers curated variant classifications such as “Pathogenic”, “Likely pathogenic” and “Benign” mapped to HGVS standards and incorporating expert reviews (13).

## 2 Methods

### 2.1 Dataset

Data from Genome Aggregation Database (gnomAD) v4 comprised 807,162 individuals, including 730,947 exomes and 76,215 genomes (7). This dataset provided 786,500,648 SNVs and 122,583,462 InDels, with variant type counts of 9,643,254 synonymous, 16,412,219 missense, 726,924 nonsense, 1,186,588 frameshift and 542,514 canonical splice site variants. ClinVar data were obtained from the variant summary dataset (as of: 16 March 2025) available from the NCBI FTP site, and included 6,845,091 entries, which were processed into 91,319 gene classification groups and a total of 38,983 gene classifications; for example, the gene *A1BG* contained four variants classified as likely benign and 102 total entries (13). For our analysis phase we also used dbNSFP which consisted of a number of annotations for 121,832,908 SNVs (15). The PanelAppRex core model contained 58,592 entries consisting of 52 sets of annotations, including the gene name, disease-gene panel ID, diseases-related features, confidence measurements. (2) A Protein-Protein Interaction (PPI) network data was provided by Search Tool for the Retrieval of Interacting Genes/Proteins (STRINGdb), consisting of 19,566 proteins and 505,968 interactions (16). The Human Genome Variation Society (HGVS) nomenclature is used with Variant Effect Predictor (VEP)-based codes for variant IDs. We carried out validations for disease cohorts with Nuclear Factor Kappa B Subunit 1 (*NFKB1*) (17–20) and Cystic Fibrosis Transmembrane Conductance Regulator (*CFTR*) (21–23) to demonstrate applications in AD and AR disease genes, respectively. **Box 2.1** list the definitions from the International Union of Immunological Societies (IUIS) IEI for the major disease categories used throughout this study (1).

Box 2.1 Definitions for IEI Major Disease Categories

Major Category	Description
1. CID	Immunodeficiencies affecting cellular and humoral immunity
2. CID+	Combined immunodeficiencies with associated or syndromic features
3. PAD -	Predominantly Antibody Deficiencies
4. PIRD -	Diseases of Immune Dysregulation
5. PD -	Congenital defects of phagocyte number or function
6. IID -	Defects in intrinsic and innate immunity
7. AID -	Autoinflammatory Disorders
8. CD -	Complement Deficiencies
9. BMF -	Bone marrow failure

173

174 **2.2 Variant Class Observation Probability**

As a starting point, we considered the classical HWE for a biallelic locus:

$$p^2 + 2pq + q^2 = 1,$$

175 where  $p$  is the allele frequency,  $q = 1 - p$ ,  $p^2$  represents the homozygous dominant,  
 176  $2pq$  the heterozygous, and  $q^2$  the homozygous recessive genotype frequencies. For disease  
 177 phenotypes, particularly under AR MOI, the risk is traditionally linked to the  
 178 homozygous state ( $p^2$ ); however, to account for compound heterozygosity across multiple  
 179 variants, we extend this by incorporating the contribution from other pathogenic  
 180 alleles.

181 Our computational pipeline estimated the probability of observing a disease-associated  
 182 genotype for each variant and aggregated these probabilities by gene and ClinVar  
 183 classification. This approach included all variant classifications, not limited solely to  
 184 those deemed “pathogenic”, and explicitly conditioned the classification on the given  
 185 phenotype, recognising that a variant could only be considered pathogenic relative to  
 186 a defined clinical context. The core calculations proceeded as follows:

**1. Allele Frequency and Total Variant Frequency.** For each variant  $i$  in a gene, the allele frequency was denoted as  $p_i$ . For each gene, we defined the total variant frequency (summing across all reported variants in that gene) as:

$$P_{\text{tot}} = \sum_{i \in \text{gene}} p_i.$$

If any of the possible SNV had no observed allele ( $p_i = 0$ ), we assigned a minimal risk:

$$p_i = \frac{1}{\max(AN) + 1},$$

187 where  $\max(AN)$  was the maximum allele number observed for that gene. This adjust-  
188 ment ensured that a nonzero risk was incorporated even in the absence of observed  
189 variants.

190 **2. Occurrence Probability Based on MOI.** The probability that an individual  
191 was affected by a variant depended on the mode of MOI relative to a specific pheno-  
192 type. Specifically, we calculated the occurrence probability  $p_{\text{disease},i}$  for each variant  
193 as follows:

- For **AD** and **XL** variants, a single copy was sufficient, so

$$p_{\text{disease},i} = p_i.$$

- For **AR** variants, disease manifested when two pathogenic alleles were present. In this case, we accounted for both the homozygous state and the possibility of compound heterozygosity:

$$p_{\text{disease},i} = p_i^2 + 2p_i(P_{\text{tot}} - p_i).$$

**3. Expected Case Numbers and Case Detection Probability.** Given a population with  $N$  births (e.g. as seen in our validation studies,  $N = 69\,433\,632$ ), the expected number of cases attributable to variant  $i$  was calculated as:

$$E_i = N \cdot p_{\text{disease},i}.$$

The probability of detecting at least one affected individual for that variant was computed as:

$$P(\geq 1)_i = 1 - (1 - p_{\text{disease},i})^N.$$

**4. Aggregation by Gene and ClinVar Classification.** For each gene and for each ClinVar classification (e.g. “Pathogenic”, “Likely pathogenic”, “Uncertain significance”, etc.), we aggregated the results across all variants. The total expected cases for a given group was:

$$E_{\text{group}} = \sum_{i \in \text{group}} E_i,$$

and the overall probability of observing at least one case within the group was calculated as:

$$P_{\text{group}} = 1 - \prod_{i \in \text{group}} (1 - p_{\text{disease},i}).$$

194 **5. Data Processing and Implementation.** We implemented the calculations  
195 within a High-Performance Computing (HPC) pipeline and provided an example  
196 for a single dominant disease gene, *TNFAIP3*, in the source code to enhance repro-  
197 ducibility. Variant data were imported in chunks from the annotation database for  
198 all chromosomes (1-22, X, Y, M).

199 For each data chunk, the relevant fields were gene name, position, allele number,  
200 allele frequency, ClinVar classification, and HGVS annotations. Missing classifica-  
201 tions (denoted by ".") were replaced with zeros and allele frequencies were converted  
202 to numeric values. We then retained only the first transcript allele annotation for sim-  
203 plicity, as the analysis was based on genomic coordinates. Subsequently, the variant  
204 data were merged with gene panel data from PanelAppRex to obtain the disease-  
205 related MOI mode for each gene. For each gene, if no variant was observed for a  
206 given ClinVar classification (i.e.  $p_i = 0$ ), a minimal risk was assigned as described  
207 above. Finally, we computed the occurrence probability, expected cases, and the  
208 probability of observing at least one case using the equations presented.

209 The final results were aggregated by gene and ClinVar classification and used to  
210 generate summary statistics that reviewed the predicted disease observation proba-  
211 bilities.

## 212 **2.3 Validation of Autosomal Dominant Estimates Using *NFKB1***

213 To validate our genome-wide probability estimates in an AD gene, we focused on  
214 *NFKB1*. Our goal was to compare the expected number of *NFKB1*-related Common  
215 Variable Immunodeficiency (CVID) cases, as predicted by our framework, with the  
216 reported case count in a well-characterised national-scale PID cohort.

217 **1. Reference Dataset.** We used a reference dataset reported by Tuijnenburg  
218 et al. (17) to build a validation model in an AD disease gene. This study performed  
219 whole-genome sequencing of 846 predominantly sporadic, unrelated PID cases from  
220 the NIHR BioResource-Rare Diseases cohort. There were 390 CVID cases in the  
221 cohort. The study identified *NFKB1* as one of the genes most strongly associated  
222 with PID. Sixteen novel heterozygous variants including truncating, missense, and  
223 gene deletion variants, were found in *NFKB1* among the CVID cases.

224 **2. Cohort Prevalence Calculation.** Within the cohort, 16 out of 390 CVID  
cases were attributable to *NFKB1*. Thus, the observed cohort prevalence was

$$\text{Prevalence}_{\text{cohort}} = \frac{16}{390} \approx 0.041,$$

224 with a 95% confidence interval (using Wilson's method) of approximately (0.0254, 0.0656).

**3. National Estimate Based on Literature.** Based on literature, the prevalence of CVID in the general population was estimated as

$$\text{Prevalence}_{\text{CVID}} = \frac{1}{25\,000}.$$

For a UK population of

$$N_{\text{UK}} \approx 69\,433\,632,$$

the expected total number of CVID cases was

$$E_{\text{CVID}} \approx \frac{69\,433\,632}{25\,000} \approx 2777.$$

Assuming that the proportion of CVID cases attributable to *NFKB1* is equivalent to the cohort estimate, the literature extrapolated estimate is

$$\text{Estimated } \text{NFKB1} \text{ cases} \approx 2777 \times 0.041 \approx 114,$$

with a median value of approximately 118 and a 95% confidence interval of 70 to 181 cases (derived from posterior sampling).

**4. Bayesian Adjustment.** Recognising that the clinical cohort likely represents nearly all CVID cases (besides first-second degree relatives), two Bayesian adjustments were performed:

1. **Weighted Adjustment (emphasising the cohort,  $w = 0.9$ ):**

$$\text{Adjusted Estimate} = 0.9 \times 16 + 0.1 \times 114 \approx 26,$$

with a corresponding 95% confidence interval of approximately 21 to 33 cases.

2. **Mixture Adjustment (equal weighting,  $w = 0.5$ ):** Posterior sampling of the cohort prevalence was performed assuming

$$p \sim \text{Beta}(16 + 1, 390 - 16 + 1),$$

which yielded a Bayesian mixture adjusted median estimate of 67 cases with a 95% credible interval of approximately 43 to 99 cases.

**5. Predicted Total Genotype Counts.** The predicted total synthetic genotype count (before adjustment) was 456, whereas the predicted total genotypes adjusted for *synth\_flag* was 0. This higher synthetic count was set based on a minimal risk threshold, ensuring that at least one genotype is assumed to exist (e.g. accounting for a potential unknown de novo variant) even when no variant is observed in gnomAD (as per [section 2.2](#)).

<sup>239</sup> **6. Validation Test.** Thus, the expected number of *NFKB1*-related CVID cases  
<sup>240</sup> derived from our genome-wide probability estimates was compared with the observed  
<sup>241</sup> counts from the UK-based PID cohort. This comparison validates our framework for  
<sup>242</sup> estimating disease incidence in AD disorders.

## <sup>243</sup> 2.4 Validation Study for Autosomal Recessive CF Using *CFTR*

<sup>244</sup> To validate our framework for AR diseases, we focused on Cystic Fibrosis (CF).  
<sup>245</sup> For comparability sizes between the validation studies, we analysed the most com-  
<sup>246</sup> mon SNV in the *CFTR* gene, typically reported as “p.Arg117His” (GRCh38 Chr  
<sup>247</sup> 7:117530975 G/A, MANE Select HGVS p.ENST00000003084.11: p.Arg117His). Our  
<sup>248</sup> goal was to validate our genome-wide probability estimates by comparing the ex-  
<sup>249</sup> pected number of CF cases attributable to the p.Arg117His variant in *CFTR* with  
<sup>250</sup> the nationally reported case count in a well-characterised disease cohort (21–23).

**1. Expected Genotype Counts.** Let  $p$  denote the allele frequency of the p.Arg117His variant and  $q$  denote the combined frequency of all other pathogenic *CFTR* variants, such that

$$q = P_{\text{tot}} - p \quad \text{with} \quad P_{\text{tot}} = \sum_{i \in \text{CFTR}} p_i.$$

Under Hardy–Weinberg equilibrium for an AR trait, the expected frequencies were:

$$f_{\text{hom}} = p^2 \quad (\text{homozygous for p.Arg117His})$$

and

$$f_{\text{comphet}} = 2pq \quad (\text{compound heterozygotes carrying p.Arg117His and another pathogenic allele}).$$

For a population of size  $N$  (here,  $N \approx 69\,433\,632$ ), the expected number of cases were:

$$E_{\text{hom}} = N \cdot p^2, \quad E_{\text{comphet}} = N \cdot 2pq, \quad E_{\text{total}} = E_{\text{hom}} + E_{\text{comphet}}.$$

**2. Mortality Adjustment.** Since CF patients experience increased mortality, we adjusted the expected genotype counts using an exponential survival model (21–23). With an annual mortality rate  $\lambda \approx 0.004$  and a median age of 22 years, the survival factor was computed as:

$$S = \exp(-\lambda \cdot 22).$$

Thus, the mortality-adjusted expected genotype count became:

$$E_{\text{adj}} = E_{\text{total}} \times S.$$

**3. Bayesian Uncertainty Simulation.** To incorporate uncertainty in the allele frequency  $p$ , we modelled  $p$  as a beta-distributed random variable:

$$p \sim \text{Beta}(p \cdot \text{AN}_{\text{eff}} + 1, \text{AN}_{\text{eff}} - p \cdot \text{AN}_{\text{eff}} + 1),$$

251 using a large effective allele count ( $\text{AN}_{\text{eff}}$ ) for illustration. By generating 10,000 poste-  
252 rior samples of  $p$ , we obtained a distribution of the literature-based adjusted expected  
253 counts,  $E_{\text{adj}}$ .

**4. Bayesian Mixture Adjustment.** Since the national registry may not capture all nuances (e.g., reduced penetrance) of *CFTR*-related disease, we further combined the literature-based estimate with the observed national count (714 cases from the UK Cystic Fibrosis Registry 2023 Annual Data Report) using a 50:50 weighting:

$$E_{\text{Bayes}} = 0.5 \times (\text{Observed Count}) + 0.5 \times E_{\text{adj}}.$$

254 **5. Validation test.** Thus, the expected number of *CFTR*-related CF cases de-  
255 rived from our genome-wide probability estimates was compared with the observed  
256 counts from the UK-based CF registry. This comparison validated our framework for  
257 estimating disease incidence in AD disorders.

## 258 2.5 Protein Network and Genetic Constraint Interpretation

259 A PPI network was constructed using protein interaction data from STRINGdb (16).  
260 We previously prepared and reported on this dataset consisting of 19,566 proteins and  
261 505,968 interactions (<https://github.com/DylanLawless/ProteoMCLustR>). Node  
262 attributes were derived from log-transformed score-positive-total values, which in-  
263 formed both node size and colour. Top-scoring nodes (top 15 based on score) were  
264 labelled to highlight prominent interactions. To evaluate group differences in score-  
265 positive-total across major disease categories, one-way Analysis of Variance (ANOVA)  
266 was performed followed by Tukey Honestly Significant Difference (HSD) post hoc tests  
267 (and non-parametric Dunn's test for confirmation). GnomAD v4.1 constraint metrics  
268 data was used for the PPI analysis and was sourced from Karczewski et al. (7). This  
269 provided transcript-level metrics, such as observed/expected ratios, Loss-Of-function  
270 Observed/Expected Upper bound Fraction (LOEUF), pLI, and Z-scores, quantifying  
271 loss-of-function and missense intolerance, along with confidence intervals and related  
272 annotations for 211,523 observations.

## 273 2.6 Gene Set Enrichment Test

274 To test for overrepresentation of biological functions, the prioritised genes were com-  
275 pared against gene sets from MsigDB (including hallmark, positional, curated, motif,  
276 computational, GO, oncogenic, and immunologic signatures) and WikiPathways using

277 hypergeometric tests with FUMA (24; 25). The background set consisted of 24,304  
278 genes. Multiple testing correction was applied per data source using the Benjamini-  
279 Hochberg method, and gene sets with an adjusted P-value  $\leq 0.05$  and more than one  
280 overlapping gene are reported.

281 

### 3 Results

282 

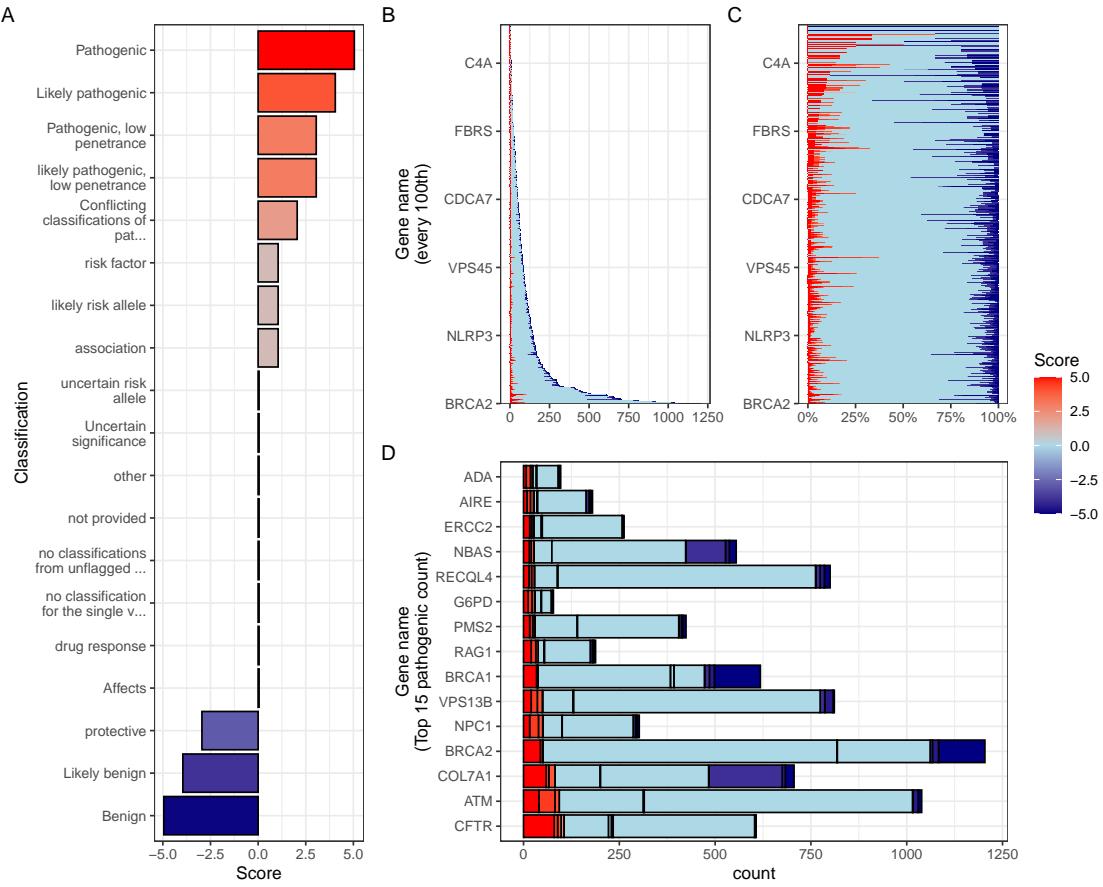
#### 3.1 Observation Probability Across Disease Genes

283 Our study integrated large-scale annotation databases with gene panels from Panel-  
284 elAppRex to systematically assess disease genes by MOI. By combining population  
285 allele frequencies with ClinVar clinical classifications, we computed an expected obser-  
286 vation probability for each SNV, representing the likelihood of encountering a variant  
287 of a specific pathogenicity for a given phenotype. We report these probabilities for  
288 54,814 ClinVar variant classifications across 557 genes (linked dataset (26)).

289 In practice, our approach computed a simple observation probability for every  
290 SNV across the genome and was applicable to any disease-gene panel. Here, we fo-  
291 cused on panels related to Primary Immunodeficiency or Monogenic Inflammatory  
292 Bowel Disease, using PanelAppRex panel ID 398 as a case study. **Figure 1** dis-  
293 plays all reported ClinVar variant classifications for this panel. The resulting natural  
294 scaling system (-5 to +5) accounts for the frequently encountered combinations of  
295 classification labels (e.g. benign to pathogenic). The resulting data set (26) is briefly  
296 shown in **Table 1** to illustrate that our method yielded estimations of the probability  
297 of observing a variant with a particular ClinVar classification.

Table 1: Example of the first several rows from our main results for 557 genes of Panel-  
elAppRex’s panel: (ID 398) Primary immunodeficiency or monogenic inflammatory  
bowel disease. “ClinVar Significance” indicates the pathogenicity classification as-  
signed by ClinVar, while “inVar Significance” indicates the pathogenicity classification  
assigned by ClinVar, while “Occurrence Prob” represents our calculated probability of  
observing the corresponding variant class for a given phenotype. Additional columns,  
such as population allele frequency, are not shown. (26)

Gene	Panel ID	ClinVar Clinical Significance	GRCh38 Pos	HGVSc (VEP)	HGVSp (VEP)	Inheritance	Occurrence Probability
ABI3	398	Uncertain significance	49210771	c.47G>A	p.Arg16Gln	AR	0.000000007
ABI3	398	Uncertain significance	49216678	c.265C>T	p.Arg89Cys	AR	0.000000005
ABI3	398	Uncertain significance	49217742	c.289G>A	p.Val97Met	AR	0.000000002
ABI3	398	Uncertain significance	49217781	c.328G>A	p.Gly110Ser	AR	0.000000002
ABI3	398	Uncertain significance	49217844	c.391C>T	p.Pro131Ser	AR	0.000000015
ABI3	398	Uncertain significance	49220257	c.733C>G	p.Pro245Ala	AR	0.000000022



**Figure 1: Summary of ClinVar clinical significance classifications in the PID gene panel.** (A) Shows the numeric score coding for each classification. Panels (B) and (C) display the tally of classifications per gene as absolute counts and as percentages, respectively. (D) Highlights the top 15 genes with the highest number of reported pathogenic classifications (score 5).

## 298 3.2 Validation studies

### 299 3.2.1 Validation of Dominant Disease Occurrence with *NFKB1*

300 To validate our genome-wide probability estimates for AD disorders, we focused  
 301 on *NFKB1*. We used a reference dataset from Tuijnenburg et al. (17), in which  
 302 whole-genome sequencing of 846 PID patients identified *NFKB1* as one of the genes  
 303 most strongly associated with the disease, with 16 *NFKB1*-related CVID cases at-  
 304 tributed to AD heterozygous variants. Our goal was to compare the predicted num-  
 305 ber of *NFKB1*-related CVID cases with the reported count in this well-characterised  
 306 national-scale cohort.

307 Our model calculated 0 known pathogenic variant *NFKB1*-related CVID cases  
 308 in the UK with a minimal risk of 456 unknown de novo variants. In the reference

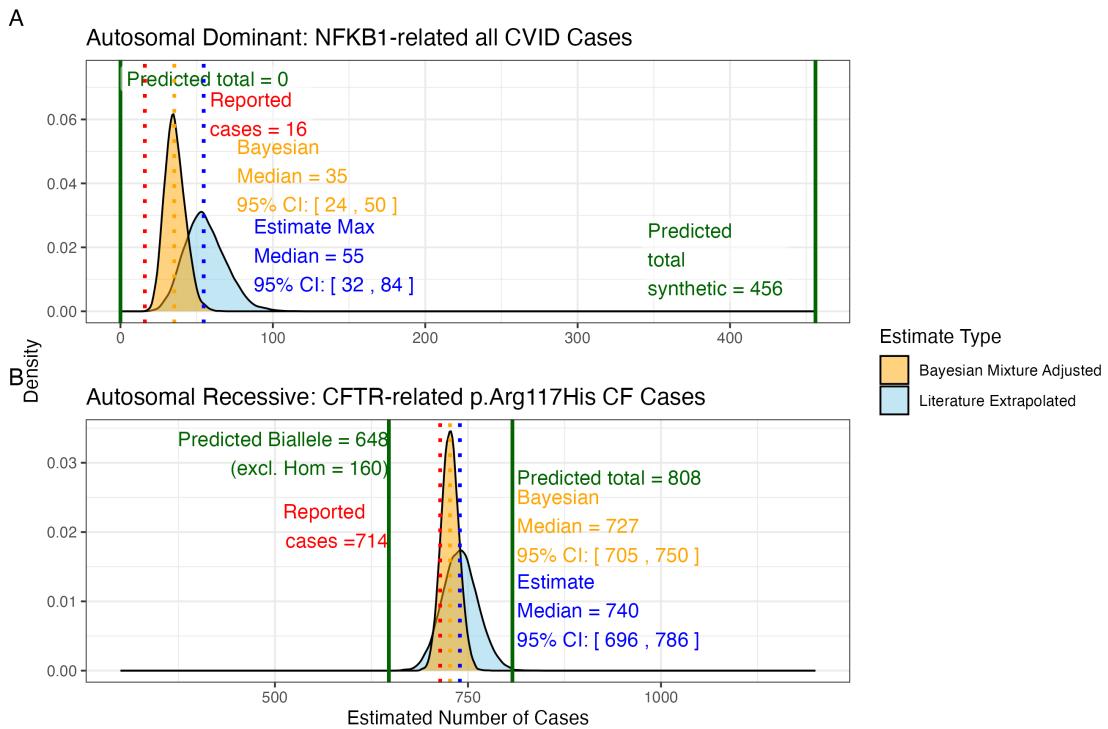
309 cohort, 16 *NFKB1* CVID cases were reported. We additionally wanted to account for  
310 potential under-reporting in the reference study. We used an extrapolated national  
311 CVID prevalence which yielded a median estimate of 118 cases (95% CI: 70–181),  
312 while a Bayesian-adjusted mixture estimate produced a median of 67 cases (95% CI:  
313 43–99). **Figure 2 (A)** illustrates that our predicted values reflect these ranges and  
314 are closer to the observed count. This case supports the validity of our integrated  
315 probability estimation framework for AD disorders, and represents a challenging ex-  
316 ample where pathogenic SNV are not reported in the reference population of gnomAD.  
317 Our min-max values successfully contained the true reported values.

318 **3.2.2 Validation of Recessive Disease Occurrence with *CFTR***

319 Our analysis predicted the number of CF cases attributable to carriage of the p.Arg117His  
320 variant (either as homozygous or as compound heterozygous with another pathogenic  
321 allele) in the UK. Based on HWE calculations and mortality adjustments, we pre-  
322 dicted approximately 648 cases arising from biallelic variants and 160 cases from  
323 homozygous variants, resulting in a total of 808 expected cases.

324 In contrast, the nationally reported number of CF cases was 714, as recorded in the  
325 UK Cystic Fibrosis Registry 2023 Annual Data Report (21). To account for factors  
326 such as reduced penetrance and the mortality-adjusted expected genotype, we derived  
327 a Bayesian-adjusted estimate via posterior simulation. Our Bayesian approach yielded  
328 a median estimate of 740 cases (95% Confidence Interval (CI): 696, 786) and a  
329 mixture-based estimate of 727 cases (95% CI: 705, 750). **Figure 2 (B)** illustrates  
330 the close concordance between the predicted values, the Bayesian-adjusted estimates,  
331 and the national report supports the validity of our approach for estimating disease.

332 **Figure S1** shows the final values for these genes of interest in a given population  
333 size and phenotype. It reveals that an allele frequency threshold of approximately  
334 0.000007 is required to observe a single heterozygous disease-causing variant carrier in  
335 the UK population for both genes. However, owing to the AR MOI pattern of *CFTR*,  
336 this threshold translates into more than 100,000 heterozygous carriers, compared to  
337 only 456 carriers for the AD gene *NFKB1*. Note that this allele frequency threshold,  
338 being derived from the current reference population, represents a lower bound that  
339 can become more precise as public datasets continue to grow. This marked difference  
340 underscores the significant impact of MOI patterns on population carrier frequencies  
341 and the observed disease prevalence.



**Figure 2: Prior probabilities compared to validation disease cohort metrics.**

(A) Density distributions for the number of *NFKB1*-related CVID cases in the UK. Our model (green) predicted 456 cases, which falls between the observed cohort count (red) of 390 and the upper extrapolated values. The blue curve represents maximum count of 1280, and the orange curve shows the Bayesian-adjusted mixture estimate of 835. (B) Density distributions for *CFTR*-related p.Arg117His CF cases. Our model (green) predicted 648 biallelic cases and 808 total cases. The nationally reported case count (red) was 714. The blue curve represents maximum extrapolated count of 740, and the orange curve shows the Bayesian-adjusted mixture estimate of 727. We observed close agreement among the reported disease cases and our integrated probability estimation framework.

### 342 3.2.3 Interpretation of ClinVar Variant Observations

343 **Figure 3** shows the two validation study PID genes, representing AR and dominant  
 344 MOI. **Figure 3 (A)** illustrates the overall probability of an affected birth by ClinVar  
 345 variant classification, whereas **Figure 3 (B)** depicts the total expected number of  
 346 cases per classification for an example population, here the UK, of approximately 69.4  
 347 million.

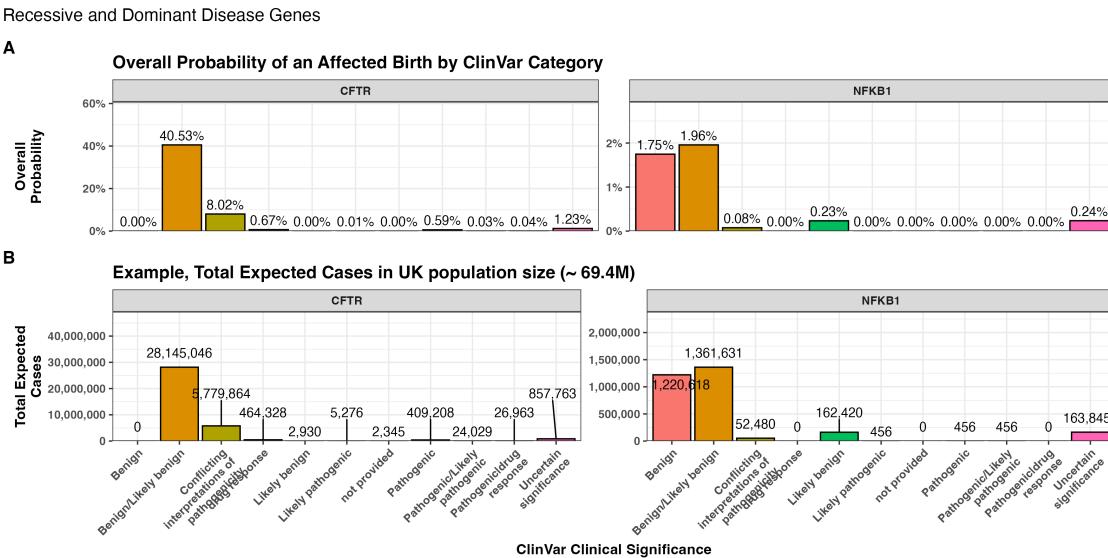


Figure 3: **Combined bar charts summarizing the genome-wide analysis of ClinVar clinical significance for the PID gene panel.** Panel (A) shows the overall probability of an affected birth by variant classification, and (B) displays the total expected number of cases per classification, both stratified by gene. These integrated results illustrate the variability in variant observations across genes and underpin our validation of the probability estimation framework.

### 348 3.3 Genetic constraint in high-impact protein networks

349 We next examined genetic constraint in high-impact protein networks across the whole  
 350 IEI gene set of over 500 known disease-gene phenotypes (1). By integrating ClinVar  
 351 variant classification scores with PPI data, we quantified the pathogenic burden per  
 352 gene and assessed its relationship with network connectivity and genetic constraint  
 353 (7; 16).

#### 354 3.3.1 Score-Positive-Total within IEI PPI network

355 The ClinVar classifications reported in **Figure 1** were scaled -5 to +5 based on their  
 356 pathogenicity. We were interested in positive (potentially damaging) but not negative

(benign) scoring variants, which are statistically incidental in this analysis. We tallied gene-level positive scores to give the score positive total metric. **Figure 4 (A)** shows the PPI network of disease-associated genes, where node size and colour encode the score positive total (log-transformed). The top 15 genes/proteins with the highest total prior probabilities of being observed with disease are labelled (as per **Figure 1**).

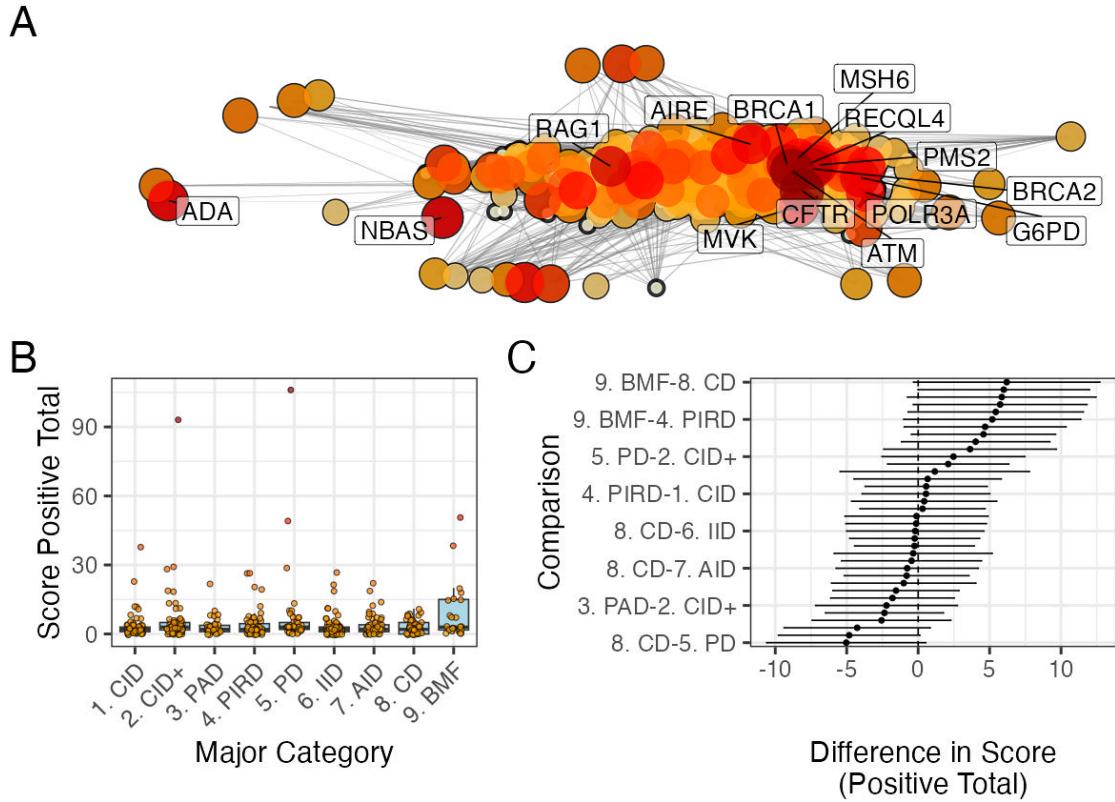


Figure 4: **PPI network and score positive total ClinVar significance variants.** (A) PPI network of disease-associated genes. Node size and colour represent the log-transformed score positive total, the top 15 genes/proteins with the highest probability of being observed in disease are labelled. (B) Distribution of score positive total across the major IEI disease categories. (C) Tukey HSD comparisons of mean differences in score positive total among all pairwise disease categories. Every 5th label is shown on y-axis.

### 3.3.2 Association Analysis of Score-Positive-Total across IEI Categories

We checked for any statistical enrichment in score positive totals, which represents the expected observation of pathogenicity, between the IEI categories. The one-way ANOVA revealed an effect of major disease category on score positive total ( $F(8, 500) = 2.82, p = 0.0046$ ), indicating that group means were not identical, which we observed in **Figure 4 (B)**. However, despite some apparent differences in median scores across

368 categories (i.e. 9. Bone Marrow Failure (BMF)), the Tukey HSD post hoc comparisons  
 369 **Figure 4 (C)** showed that all pairwise differences had 95% confidence intervals  
 370 overlapping zero, suggesting that individual group differences were not significant.

### 371 3.3.3 UMAP Embedding of the PPI Network

372 To address the density of the PPI network for the IEI gene panel, we applied Uniform  
 373 Manifold Approximation and Projection (UMAP) (**Figure 5**). Node sizes reflect  
 374 interaction degree, a measure of evidence-supported connectivity (16). We tested  
 375 for a correlation between interaction degree and score positive total. In **Figure**  
 376 **5**, gene names with degrees above the 95th percentile are labelled in blue, while  
 377 the top 15 genes by score positive total are labelled in yellow (as per **Figure 1**).  
 378 Notably, genes with high pathogenic variant loads segregated from highly connected  
 379 nodes, suggesting that Loss-of-Function (LOF) in hub genes is selectively constrained,  
 380 whereas damaging variants in lower-degree genes yield more specific effects. This  
 381 observation was subsequently tested empirically.

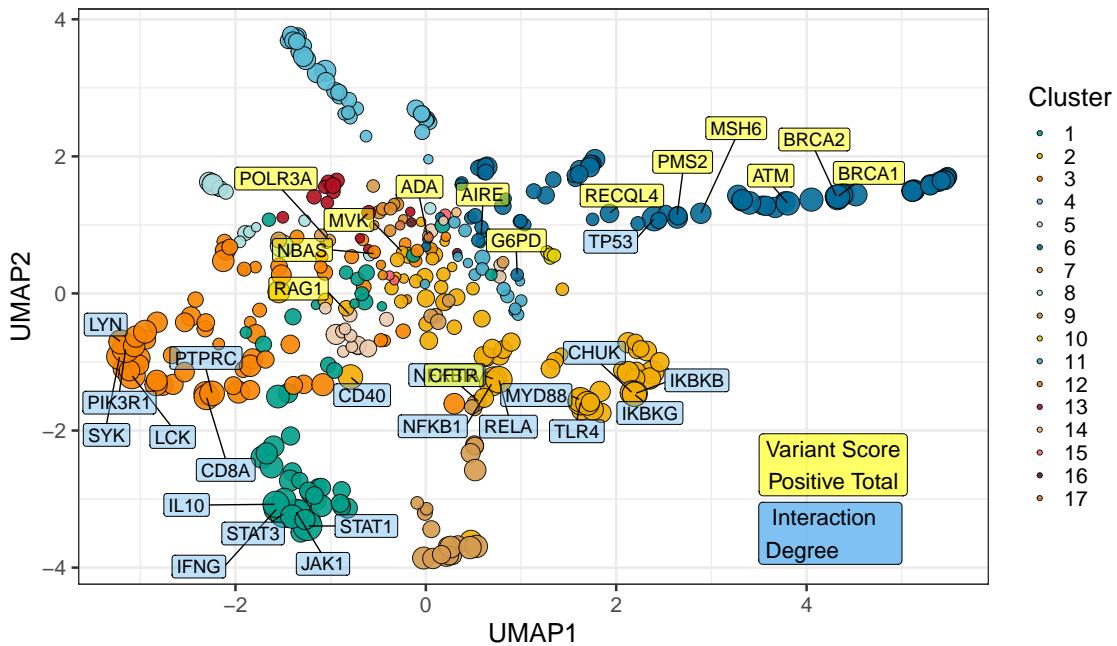


Figure 5: **UMAP embedding of the PPI network (p\_umap).** The plot projects the high-dimensional protein-protein interaction network into two dimensions, with nodes coloured by cluster and sized by interaction degree. Blue labels indicate hub genes (degree above the 95th percentile) and yellow labels mark the top 15 genes by score positive total (damaging ClinVar classifications). The spatial segregation suggests that genes with high pathogenic variant loads are distinct from highly connected nodes.

382    **3.3.4 Hierarchical Clustering of Enrichment Scores for Major Disease Cat-**  
383    **egories**

384    **Figure S2** presents a heatmap of standardised residuals for major disease categories  
385    across network clusters, as per **Figure 5**. A dendrogram clusters similar disease cate-  
386    gories, while the accompanying bar plot displays the maximum absolute standardised  
387    residual for each category. Notably, (8) Complement Deficiencies (CD) shows the  
388    highest maximum enrichment, followed by (9) BMF. While all maximum values  
389    exceed 2, the threshold for significance, this likely reflects the presence of protein  
390    clusters with strong damaging variant scores rather than uniform significance across  
391    all categories (i.e. genes from cluster 4 in 8 CD).

392    **3.3.5 PPI Connectivity, LOEUF Constraint and Enriched Network Clus-**  
393    **ter Analysis**

394    Based on the preliminary insight from **Figure S2**, we evaluated the relationship  
395    between network connectivity (PPI degree) and LOEUF constraint (LOEUF upper rank)  
396    Karczewski et al. (7) using Spearman's rank correlation. Overall, there was a weak  
397    but significant negative correlation ( $\rho = -0.181, p = 0.00024$ ) at the global scale,  
398    indicating that highly connected genes tend to be more constrained. A supplementary  
399    analysis (see **Figure 6**) did not reveal distinct visual associations between network  
400    clusters and constraint metrics, likely due to the high network density. However  
401    once stratified by gene clusters, the natural biological scenario based on quantitative  
402    PPI evidence (16), some groups showed strong correlations; for instance, cluster 2  
403    ( $\rho = -0.375, p = 0.000994$ ) and cluster 4 ( $\rho = -0.800, p < 0.000001$ ), while others did  
404    not. This indicated that shared mechanisms within pathway clusters may underpin  
405    genetic constraints, particularly for LOF intolerance. We observe that the score  
406    positive total metric effectively summarises the aggregate pathogenic burden across  
407    IEI genes, serving as a robust indicator of genetic constraint and highlighting those  
408    with elevated disease relevance.

409    **Figure 6 (C, D)** shows the re-plotted PPI networks for clusters with significant  
410    correlations between PPI degree and LOEUF upper rank. In these networks, node  
411    size is scaled by a normalised variant score, while node colour reflects the variant  
412    score according to a predefined palette.

413    **3.4 New Insight from Functional Enrichment**

414    To interpret the functional relevance of our prioritised IEI gene sets with the highest  
415    load of damaging variants (i.e. clusters 2 and 4 in **Figure 6**), we performed func-  
416    tional enrichment analysis for known disease associations using MsigDB with FUMA  
417    (i.e. GWAScatalog and Immunologic Signatures) (24). Composite enrichment pro-  
418    files (**Figure S4**) reveal that our enriched PPI clusters were associated with distinct  
419    disease-related phenotypes, providing functional insights beyond traditional IUIS IEI

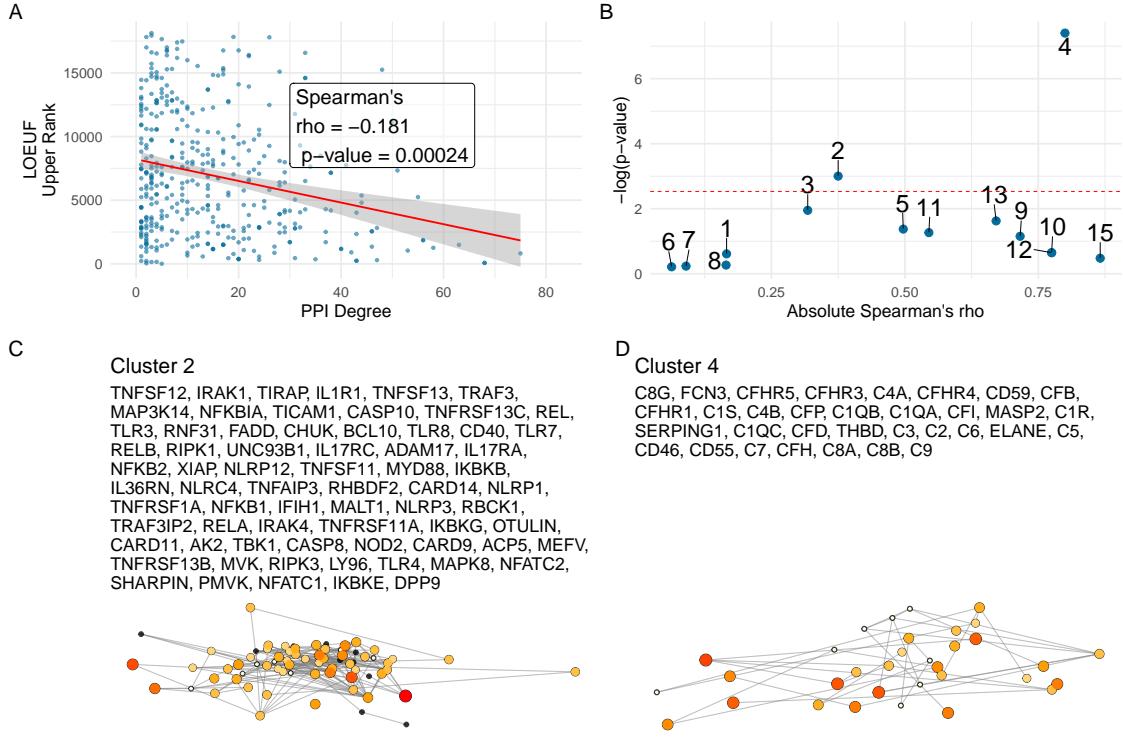


Figure 6: **Correlation between PPI degree and LOEUF upper rank.** (A) Ananlysis across all genes revealed a weak, significant negative correlation between PPI degree and LOEUF upper rank. (B) The cluster-wise analysis showed that clusters 2 and 4 exhibited moderate to strong correlations, while other clusters display weak or non-significant relationships. (C) and (D) Shows the new network plots for the significantly enriched clusters based on gnomAD constraint metrics.

groupings (1). The gene expression profiles shown in **Figure S5** (GTEx v8 54 tissue types) offer the tissue-specific context for these associations. Together, these results enable the annotation of IEI gene sets with established disease phenotypes, supporting a data-driven classification of IEI.

Based on these independent sources of interpretation, we observed that genes from cluster 2 were independently associated with specific inflammatory phenotypes, including ankylosing spondylitis, psoriasis, inflammatory bowel disease, and rheumatoid arthritis, as well as quantitative immune traits such as lymphocyte and neutrophil percentages and serum protein levels. In contrast, genes from Cluster 4 were linked to ocular and complement-related phenotypes, notably various forms of age-related macular degeneration (e.g. geographic atrophy and choroidal neovascularisation) and biomarkers of the complement system (e.g. C3, C4, and factor H-related proteins), with additional associations to nephropathy and pulmonary function metrics.

433    **3.5 Genome-wide Gene Distribution and Locus-specific Vari-**  
 434    **ant Occurrence**

435    **Figure 7 (A)** shows a genome-wide karyoplot of all IEI panel genes across GRCh38,  
 436    with colour-coding based on MOI. Figures (B) and (C) display zoomed-in locus plots  
 437    for *NFKB1* and *CFTTR*, respectively. In **Figure 7 (B)**, the probability of observing  
 438    variants with known classifications is high only for variants such as p.Ala475Gly,  
 439    which are considered benign in the AD *NFKB1* gene that is intolerant to LOF. In  
 440    **Figure 7 (C)**, high probabilities of observing patients with pathogenic variants in  
 441    *CFTTR* are evident, reproducing this well-established phenomenon. Furthermore, the  
 442    analysis of Linkage Disequilibrium (LD) using  $R^2$  shows that high LD regions can be  
 443    modelled effectively, allowing independent variant signals to be distinguished.

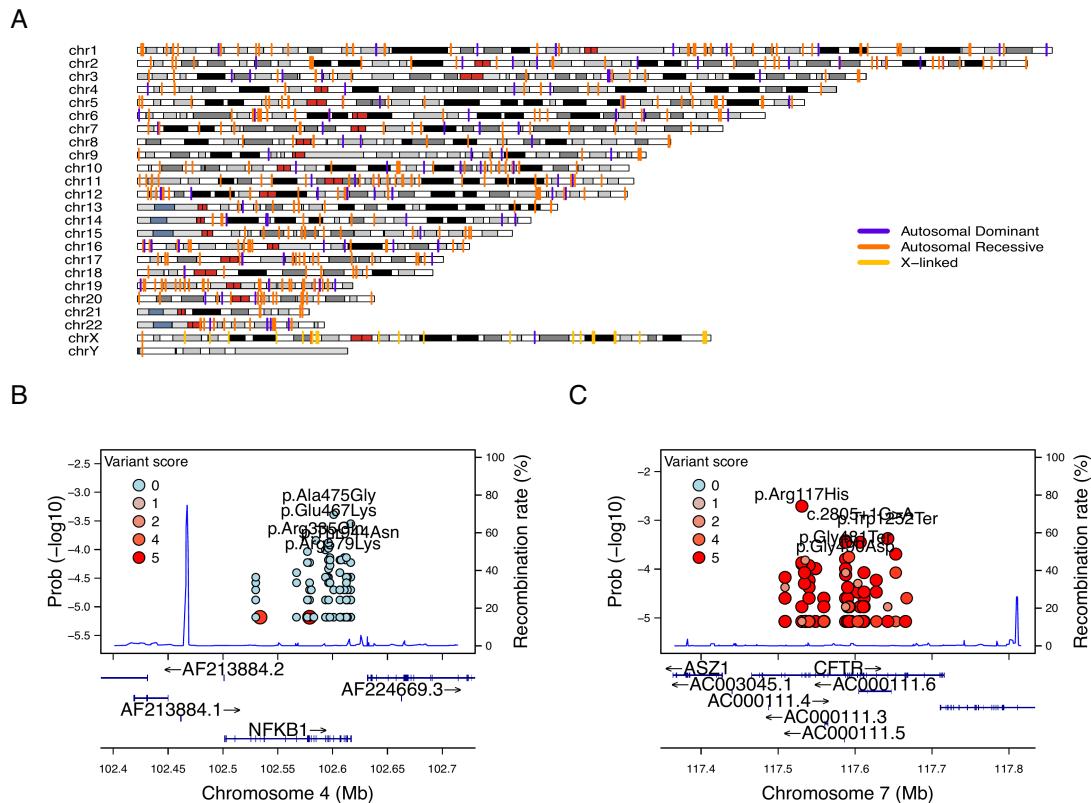


Figure 7: Genome-wide IEI, variant occurrence probability and LD by  $R^2$ . (A) Genome-wide karyoplot of all IEI panel genes mapped to GRCh38, with colours indicating MOI. (B) Zoomed-in locus plot for *NFKB1* showing variant observation probabilities; only benign variants such exhibit high probabilities in this AD gene intolerant to LOF. (C) Locus plot for *CFTTR* displaying high probabilities for pathogenic variants; due to the dense clustering of pathogenic variants, score filter  $>0$  was applied. Top five variant are labelled per gene.

### 444 3.5.1 Integration of Variant Probabilities into IEI Genetics Data

445 We integrated the computed prior probabilities for observing variants in all known  
 446 genes associated with a given phenotype (1), across AD, AR, and XL MOI, into  
 447 our IEI genetics framework. These calculations, derived from gene panels in Pan-  
 448 elAppRex, have yielded novel insights for the IEI disease panel. The final result  
 449 comprised of machine- and human-readable datasets, including the table of variant  
 450 classifications and priors available via a the linked repository (26), and a user-friendly  
 451 web interface that incorporates these new metrics.

452 **Figure 8** shows the interface summarising integrated variant data. Server-side  
 453 pre-calculation of summary statistics minimises browser load, while clinical signifi-  
 454 cance is converted to numerical metrics. Key quantiles (min, Q1, median, Q3, max)  
 455 for each gene are rendered as sparkline box plots, and dynamic URLs link table entries  
 456 to external databases (e.g. ClinVar, Online Mendelian Inheritance in Man (OMIM),  
 457 AlphaFold).

The screenshot shows a table titled "Viewer Zoom" with a search bar at the top. The table has 13 columns: Major category, Subcategory, Disease, Genetic defect, Inheritance, Gene score, Prior prob of observing pathogenic, ClinVar SNV classification, ClinVar all variant reports, OMIM, Alpha Missense / Uniprot ID, HPO combined, and HPO term. The "Gene score" column contains sparkline box plots. The "Prior prob of observing pathogenic" column contains numerical values. The "ClinVar SNV classification" and "ClinVar all variant reports" columns contain links. The "OMIM" column contains OMIM IDs. The "Alpha Missense / Uniprot ID" column contains AlphaMissense IDs. The "HPO combined" and "HPO term" columns contain HPO terms. The table is paginated at the bottom with "Previous" and "Next" buttons.

Figure 8: **Integration of variant probabilities into the IEI genetics framework.** The interface summarises the condensed variant data, with pre-calculated summary statistics and dynamic links to external databases. This integration enables immediate access to detailed variant classifications and prior probabilities for each gene.

## 458 4 Discussion

459 Our study presents, to our knowledge, the first comprehensive framework for calculat-  
 460 ing prior probabilities of observing disease-associated variants. By integrating large-  
 461 scale genomic annotations, including population allele frequencies from gnomAD (7),  
 462 variant classifications from ClinVar (13), and functional annotations from resources  
 463 such as dbNSFP, with classical Hardy-Weinberg-based calculations, we derived robust

<sup>464</sup> estimates for 54,814 ClinVar variant classifications across 557 IEI genes implicated in  
<sup>465</sup> PID and monogenic inflammatory bowel disease (1; 2).

<sup>466</sup> Our approach yielded two key results. First, our detailed, per-variant pre-calculated  
<sup>467</sup> results provide prior probabilities of observing disease-associated variants across all  
<sup>468</sup> MOI for any gene-disease combination. Second, the score positive total metric effec-  
<sup>469</sup> tively summarises the aggregate pathogenic burden across genes, serving as a robust  
<sup>470</sup> indicator of genetic constraint and highlighting those with elevated disease relevance.

Estimating disease risk in genetic studies is complicated by uncertainties in key parameters such as variant penetrance and the fraction of cases attributable to specific variants (6). In the simplest model, where a single, fully penetrant variant causes disease, the lifetime risk  $P(D)$  is equivalent to the genotype frequency  $P(G)$ . For an allele with frequency  $p$ , this translates to:

$$\text{Recessive: } P(D) = p^2,$$

$$\text{Dominant: } P(D) = 2p(1 - p) \approx 2p.$$

When penetrance is incomplete, defined as  $P(D | G)$ , the risk becomes:

$$P(D) = P(G) P(D | G).$$

In more realistic scenarios where multiple variants contribute to disease,  $P(G | D)$  denotes the fraction of cases attributable to a given variant. This leads to:

$$P(D) = \frac{P(G) P(D | G)}{P(G | D)}.$$

<sup>471</sup> Because both penetrance and  $P(G | D)$  are often uncertain, solving this equation  
<sup>472</sup> systematically poses a major challenge.

<sup>473</sup> Our framework addresses this challenge by combining variant classifications, pop-  
<sup>474</sup>ulation allele frequencies, and curated gene-disease associations. While imperfect on  
<sup>475</sup>an individual level, these sources exhibit predictable aggregate behaviour, supported  
<sup>476</sup>by James-Stein estimation principles (27). Curated gene-disease associations help  
<sup>477</sup>identify genes that explainable for most disease cases, allowing us to approximate  
<sup>478</sup> $P(G | D)$  close to one. In this way, we obtain robust estimates of  $P(G)$  (the fre-  
<sup>479</sup>quency of disease-associated genotypes), even when exact values of penetrance and  
<sup>480</sup>case attribution remain uncertain.

This approach allows us to pre-calculate priors and summarise the overall pathogenic burden using our *score positive total* metric. By focusing on a subset  $\mathcal{V}$  of variants that pass stringent filtering, where each  $P(G_i | D)$  is the probability that a case of disease  $D$  is attributable to variant  $i$ , we assume that, in aggregate,

$$\sum_{i \in \mathcal{V}} P(G_i | D) \approx 1.$$

<sup>481</sup> Even if the cumulative contribution is slightly less than one, the resultant risk esti-  
<sup>482</sup>mates remain robust within the broad confidence intervals typical of epidemiological

483 studies. By incorporating these pre-calculated priors into a Bayesian framework, our  
484 method refines risk estimates and enhances clinical decision-making despite inherent  
485 uncertainties.

486 Our results focused on IEI, but the genome-wide approach accommodates the  
487 distinct MOI patterns of AD, AR, and XL disorders. Whereas AD and XL conditions  
488 require only a single pathogenic allele, AR disorders necessitate the consideration of  
489 both homozygous and compound heterozygous states. These classical HWE-based  
490 estimates provide an informative baseline for predicting variant occurrence and serve  
491 as robust priors for Bayesian models of variant and disease risk estimation. This  
492 is an approach that has been underutilised in clinical and statistical genetics. As  
493 such, our framework refines risk calculations by incorporating MOI complexities and  
494 enhances clinicians' understanding of expected variant occurrences, thereby improving  
495 diagnostic precision.

496 Moreover, our method complements existing statistical approaches for aggregat-  
497 ing variant effects with methods like Sequence Kernel Association Test (SKAT) and  
498 Aggregated Cauchy Association Test (ACAT) (28–31)) and multi-omics integration  
499 techniques (32; 33), while remaining consistent with established variant interpretation  
500 guidelines from the American College of Medical Genetics and Genomics (ACMG)  
501 (34) and complementary frameworks (35; 36), as well as quality control protocols  
502 (37; 38). Standardised reporting for qualifying variant sets, such as ACMG Secondary  
503 Findings v3.2 (39), further contextualises the integration of these probabilities into  
504 clinical decision-making.

505 We acknowledge that our current framework is restricted to SNVs and does not in-  
506 corporate numerous other complexities of genetic disease, such as structural variants,  
507 de novo variants, hypomorphic alleles, overdominance, variable penetrance, tissue-  
508 specific expression, the Wahlund effect, pleiotropy, and others (6). In certain applica-  
509 tions, more refined estimates would benefit from including factors such as embryonic  
510 lethality, condition-specific penetrance, and age of onset (10). Our analysis also relies  
511 on simplifying assumptions of random mating, an effectively infinite population, and  
512 the absence of migration, novel mutations, or natural selection.

513 Future work will incorporate additional variant types and models to further refine  
514 these probability estimates. By continuously updating classical estimates with emerg-  
515 ing data and prior knowledge, we aim to enhance the precision of genetic diagnostics  
516 and ultimately improve patient care.

## 517 5 Conclusion

518 Our work generates prior probabilities for observing any variant classification in IEI  
519 genetic disease, providing a quantitative resource to enhance Bayesian variant inter-  
520 pretation and clinical decision-making.

521 **Acknowledgements**

522 We acknowledge Genomics England for providing public access to the PanelApp data.  
523 The use of data from Genomics England panelapp was licensed under the Apache  
524 License 2.0. The use of data from UniProt was licensed under Creative Commons  
525 Attribution 4.0 International (CC BY 4.0). ClinVar asks its users who distribute or  
526 copy data to provide attribution to them as a data source in publications and websites  
527 (13). dbNSFP version 4.4a is licensed under the Creative Commons Attribution-  
528 NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0); while we cite  
529 this dataset as used our research publication, it is not used for the final version which  
530 instead used ClinVar and gnomAD directly. GnomAD is licensed under Creative  
531 Commons Zero Public Domain Dedication (CC0 1.0 Universal). GnomAD request  
532 that usages cites the gnomAD flagship paper (7) and any online resources that include  
533 the data set provide a link to the browser, and note that tool includes data from the  
534 gnomAD v4.1 release.

535 **Competing interest**

536 We declare no competing interest.

537 **References**

- 538 [1] Stuart G. Tangye, Waleed Al-Herz, Aziz Bousfiha, Charlotte Cunningham-  
539 Rundles, Jose Luis Franco, Steven M. Holland, Christoph Klein, Tomohiro Morio,  
540 Eric Oksenhendler, Capucine Picard, Anne Puel, Jennifer Puck, Mikko R. J.  
541 Seppänen, Raz Somech, Helen C. Su, Kathleen E. Sullivan, Troy R. Torgerson,  
542 and Isabelle Meyts. Human Inborn Errors of Immunity: 2022 Update  
543 on the Classification from the International Union of Immunological Societies  
544 Expert Committee. *Journal of Clinical Immunology*, 42(7):1473–1507, October  
545 2022. ISSN 0271-9142, 1573-2592. doi: 10.1007/s10875-022-01289-3. URL  
546 <https://link.springer.com/10.1007/s10875-022-01289-3>.
- 547 [2] Dylan Lawless. PanelAppRex aggregates disease gene panels and facilitates  
548 sophisticated search. March 2025. doi: 10.1101/2025.03.20.25324319. URL  
549 <http://medrxiv.org/lookup/doi/10.1101/2025.03.20.25324319>.
- 550 [3] Antonio Rueda Martin, Eleanor Williams, Rebecca E. Foulger, Sarah Leigh,  
551 Louise C. Daugherty, Olivia Niblock, Ivone U. S. Leong, Katherine R. Smith,  
552 Oleg Gerasimenko, Eik Haraldsdottir, Ellen Thomas, Richard H. Scott, Emma  
553 Baple, Arianna Tucci, Helen Brittain, Anna De Burca, Kristina Ibañez, Dalia  
554 Kasperaviciute, Damian Smedley, Mark Caulfield, Augusto Rendon, and Ellen M.  
555 McDonagh. PanelApp crowdsources expert knowledge to establish consensus  
556 diagnostic gene panels. *Nature Genetics*, 51(11):1560–1565, November 2019.

- 557 ISSN 1061-4036, 1546-1718. doi: 10.1038/s41588-019-0528-2. URL <https://www.nature.com/articles/s41588-019-0528-2>.
- 558
- 559 [4] Oliver Mayo. A Century of Hardy–Weinberg Equilibrium. *Twin Research*  
560 and *Human Genetics*, 11(3):249–256, June 2008. ISSN 1832-4274, 1839–  
561 2628. doi: 10.1375/twin.11.3.249. URL [https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal\\_article](https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal_article).
- 562
- 563 [5] Nikita Abramovs, Andrew Brass, and May Tassabehji. Hardy-Weinberg Equi-  
564 librium in the Large Scale Genomic Sequencing Era. *Frontiers in Genetics*,  
565 11:210, March 2020. ISSN 1664-8021. doi: 10.3389/fgene.2020.00210. URL  
566 <https://www.frontiersin.org/article/10.3389/fgene.2020.00210/full>.
- 567
- 568 [6] Johannes Zschocke, Peter H. Byers, and Andrew O. M. Wilkie. Mendelian  
569 inheritance revisited: dominance and recessiveness in medical genetics. *Nature*  
570 *Reviews Genetics*, 24(7):442–463, July 2023. ISSN 1471-0056, 1471-0064.  
571 doi: 10.1038/s41576-023-00574-0. URL <https://www.nature.com/articles/s41576-023-00574-0>.
- 572
- 573 [7] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings,  
574 Jessica Alfoldi, Qingbo Wang, Ryan L Collins, Kristen M Laricchia, Andrea  
575 Ganna, Daniel P Birnbaum, et al. The mutational constraint spectrum quantified  
576 from variation in 141,456 humans. *Nature*, 581(7809):434–443, 2020.
- 577
- 578 [8] Sarah L. Bick, Aparna Nathan, Hannah Park, Robert C. Green, Monica H. Wo-  
579 jcik, and Nina B. Gold. Estimating the sensitivity of genomic newborn screen-  
580 ing for treatable inherited metabolic disorders. *Genetics in Medicine*, 27(1):  
101284, January 2025. ISSN 10983600. doi: 10.1016/j.gim.2024.101284. URL  
<https://linkinghub.elsevier.com/retrieve/pii/S1098360024002181>.
- 581
- 582 [9] Benjamin D. Evans, Piotr Słowiński, Andrew T. Hattersley, Samuel E. Jones,  
583 Seth Sharp, Robert A. Kimmitt, Michael N. Weedon, Richard A. Oram,  
584 Krasimira Tsaneva-Atanasova, and Nicholas J. Thomas. Estimating disease  
585 prevalence in large datasets using genetic risk scores. *Nature Communications*,  
586 12(1):6441, November 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-26501-7.  
URL <https://www.nature.com/articles/s41467-021-26501-7>.
- 587
- 588 [10] William B. Hannah, Mitchell L. Drumm, Keith Nykamp, Tiziano Prampano,  
589 Robert D. Steiner, and Steven J. Schrodi. Using genomic databases to de-  
590 termine the frequency and population-based heterogeneity of autosomal reces-  
591 sive conditions. *Genetics in Medicine Open*, 2:101881, 2024. ISSN 29497744.  
592 doi: 10.1016/j.gimo.2024.101881. URL <https://linkinghub.elsevier.com/retrieve/pii/S2949774424010276>.
- 593
- 594 [11] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov,  
595 Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek,  
Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J.

- 596 Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh  
597 Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy,  
598 Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer,  
599 Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Ko-  
600 ray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate pro-  
601 tein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August  
602 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL  
603 <https://www.nature.com/articles/s41586-021-03819-2>.
- 604 [12] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-  
605 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias  
606 Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-  
607 abis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense vari-  
608 ant effect prediction with AlphaMissense. *Science*, 381(6664):eadg7492, Septem-  
609 ber 2023. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.adg7492. URL  
610 <https://www.science.org/doi/10.1126/science.adg7492>.
- 611 [13] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen Chao,  
612 Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Wonhee  
613 Jang, Karen Karapetyan, Kenneth Katz, Chunlei Liu, Zenith Maddipatla, Adri-  
614 ana Malheiro, Kurt McDaniel, Michael Ovetsky, George Riley, George Zhou,  
615 J Bradley Holmes, Brandi L Kattman, and Donna R Maglott. ClinVar: im-  
616 proving access to variant interpretations and supporting evidence. *Nucleic Acids  
617 Research*, 46(D1):D1062–D1067, January 2018. ISSN 0305-1048, 1362-4962. doi:  
618 10.1093/nar/gkx1153. URL [http://academic.oup.com/nar/article/46/D1/  
619 D1062/4641904](http://academic.oup.com/nar/article/46/D1/D1062/4641904).
- 620 [14] The UniProt Consortium, Alex Bateman, Maria-Jesus Martin, Sandra Orchard,  
621 Michele Magrane, Aduragbemi Adesina, Shadab Ahmad, Bowler-Barnett, and  
622 Others. UniProt: the Universal Protein Knowledgebase in 2025. *Nucleic  
623 Acids Research*, 53(D1):D609–D617, January 2025. ISSN 0305-1048, 1362-4962.  
624 doi: 10.1093/nar/gkae1010. URL [https://academic.oup.com/nar/article/  
625 53/D1/D609/7902999](https://academic.oup.com/nar/article/53/D1/D609/7902999).
- 626 [15] Xiaoming Liu, Chang Li, Chengcheng Mou, Yibo Dong, and Yicheng Tu.  
627 dbNSFP v4: a comprehensive database of transcript-specific functional pre-  
628 dictions and annotations for human nonsynonymous and splice-site SNVs.  
629 *Genome Medicine*, 12(1):103, December 2020. ISSN 1756-994X. doi: 10.  
630 1186/s13073-020-00803-9. URL [https://genomemedicine.biomedcentral.  
631 com/articles/10.1186/s13073-020-00803-9](https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-020-00803-9).
- 632 [16] Damian Szklarczyk, Katerina Nastou, Mikaela Koutrouli, Rebecca Kirsch, Far-  
633 rokh Mehryary, Radja Hachilif, Dewei Hu, Matteo E Peluso, Qingyao Huang,  
634 Tao Fang, et al. The string database in 2025: protein networks with directional-  
635 ity of regulation. *Nucleic Acids Research*, 53(D1):D730–D737, 2025.

- 636 [17] Paul Tuijnneburg, Hana Lango Allen, Siobhan O. Burns, Daniel Greene,  
637 Machiel H. Jansen, and Others. Loss-of-function nuclear factor B subunit  
638 1 (NFKB1) variants are the most common monogenic cause of common vari-  
639 able immunodeficiency in Europeans. *Journal of Allergy and Clinical Im-*  
640 *munology*, 142(4):1285–1296, October 2018. ISSN 00916749. doi: 10.1016/  
641 j.jaci.2018.01.039. URL <https://linkinghub.elsevier.com/retrieve/pii/S0091674918302860>.
- 643 [18] WHO Scientific Group et al. Primary immunodeficiency diseases: report of a  
644 who scientific group. *Clin. Exp. Immunol.*, 109(1):1–28, 1997.
- 645 [19] Charlotte Cunningham-Rundles and Carol Bodian. Common variable immunod-  
646 eficiency: clinical and immunological features of 248 patients. *Clinical immunol-*  
647 *ogy*, 92(1):34–48, 1999.
- 648 [20] Eric Oksenhendler, Laurence Gérard, Claire Fieschi, Marion Malphettes, Gael  
649 Mouillot, Roland Jaussaud, Jean-François Viallard, Martine Gardembas, Lionel  
650 Galicier, Nicolas Schleinitz, et al. Infections in 252 patients with common variable  
651 immunodeficiency. *Clinical Infectious Diseases*, 46(10):1547–1554, 2008.
- 652 [21] Y Naito, F Adams, S Charman, J Duckers, G Davies, and S Clarke. Uk cystic  
653 fibrosis registry 2023 annual data report. *London: Cystic Fibrosis Trust*, 2023.
- 654 [22] Carlo Castellani, CFTR2 team, et al. Cftr2: how will it help care? *Paediatric*  
655 *respiratory reviews*, 14:2–5, 2013.
- 656 [23] Hartmut Grasemann and Felix Ratjen. Cystic fibrosis. *New England Journal*  
657 *of Medicine*, 389(18):1693–1707, 2023. doi: 10.1056/NEJMra2216474. URL  
658 <https://www.nejm.org/doi/full/10.1056/NEJMra2216474>.
- 659 [24] Kyoko Watanabe, Erdogan Taskesen, Arjen Van Bochoven, and Danielle  
660 Posthuma. Functional mapping and annotation of genetic associations with  
661 FUMA. *Nature Communications*, 8(1):1826, November 2017. ISSN 2041-1723.  
662 doi: 10.1038/s41467-017-01261-5. URL <https://www.nature.com/articles/s41467-017-01261-5>.
- 663 [25] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir,  
664 Pablo Tamayo, and Jill P. Mesirov. Molecular signatures database (MSigDB)  
665 3.0. *Bioinformatics*, 27(12):1739–1740, June 2011. ISSN 1367-4811, 1367-  
666 4803. doi: 10.1093/bioinformatics/btr260. URL <https://academic.oup.com/bioinformatics/article/27/12/1739/257711>.
- 667 [26] Dylan Lawless. Variant risk estimate probabilities for iei genes. March 2025. doi:  
668 10.5281/zenodo.15111584. URL <https://doi.org/10.5281/zenodo.15111584>.
- 669 [27] Bradley Efron and Carl Morris. Stein’s Estimation Rule and Its Competitors—  
670 An Empirical Bayes Approach. *Journal of the American Statistical Association*,  
671 68(341):117, March 1973. ISSN 01621459. doi: 10.2307/2284155. URL <https://www.jstor.org/stable/2284155?origin=crossref>.

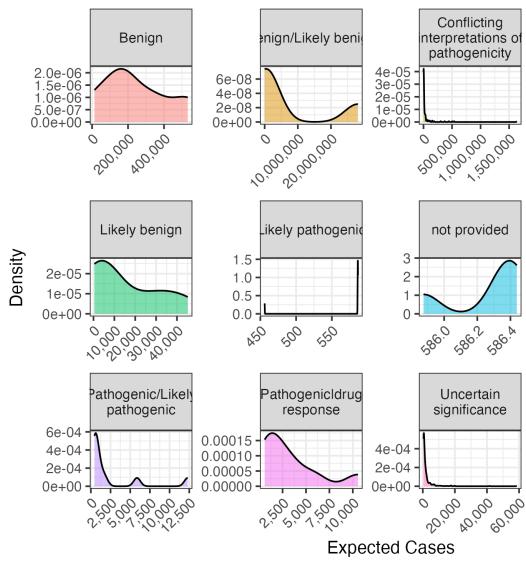
- 675 [28] Yaowu Liu, Sixing Chen, Zilin Li, Alanna C Morrison, Eric Boerwinkle, and  
676 Xihong Lin. Acat: a fast and powerful p value combination method for rare-  
677 variant analysis in sequencing studies. *The American Journal of Human Genetics*,  
678 104(3):410–421, 2019.
- 679 [29] Xiacao Li, Zilin Li, Hufeng Zhou, Sheila M Gaynor, Yaowu Liu, Han Chen, Ryan  
680 Sun, Rounak Dey, Donna K Arnett, Stella Aslibekyan, et al. Dynamic incorpora-  
681 tion of multiple in silico functional annotations empowers rare variant association  
682 analysis of large whole-genome sequencing studies at scale. *Nature genetics*, 52  
683 (9):969–983, 2020.
- 684 [30] Michael C Wu, Seunggeun Lee, Tianxi Cai, Yun Li, Michael Boehnke, and Xi-  
685 hong Lin. Rare-variant association testing for sequencing data with the sequence  
686 kernel association test. *The American Journal of Human Genetics*, 89(1):82–93,  
687 2011.
- 688 [31] Seunggeun Lee, Mary J Emond, Michael J Bamshad, Kathleen C Barnes, Mark J  
689 Rieder, Deborah A Nickerson, David C Christiani, Mark M Wurfel, and Xihong  
690 Lin. Optimal unified approach for rare-variant association testing with applica-  
691 tion to small-sample case-control whole-exome sequencing studies. *The American  
692 Journal of Human Genetics*, 91(2):224–237, 2012.
- 693 [32] Augustine Kong, Gudmar Thorleifsson, Michael L Frigge, Bjarni J Vilhjalmsson,  
694 Alexander I Young, Thorgeir E Thorgeirsson, Stefania Benonisdottir, Asmundur  
695 Oddsson, Bjarni V Halldorsson, Gisli Masson, et al. The nature of nurture:  
696 Effects of parental genotypes. *Science*, 359(6374):424–428, 2018.
- 697 [33] Laurence J Howe, Michel G Nivard, Tim T Morris, Ailin F Hansen, Humaira  
698 Rasheed, Yoonsoo Cho, Geetha Chittoor, Penelope A Lind, Teemu Palviainen,  
699 Matthijs D van der Zee, et al. Within-sibship gwas improve estimates of direct  
700 genetic effects. *BioRxiv*, pages 2021–03, 2021.
- 701 [34] Sue Richards, Nazneen Aziz, Sherri Bale, David Bick, Soma Das, Julie Gastier-  
702 Foster, Wayne W Grody, Madhuri Hegde, Elaine Lyon, Elaine Spector, et al.  
703 Standards and guidelines for the interpretation of sequence variants: a joint  
704 consensus recommendation of the american college of medical genetics and ge-  
705 nomics and the association for molecular pathology. *Genetics in medicine*, 17  
706 (5):405–423, 2015.
- 707 [35] Sean V Tavtigian, Steven M Harrison, Kenneth M Boucher, and Leslie G  
708 Biesecker. Fitting a naturally scaled point system to the acmg/amp variant  
709 classification guidelines. *Human mutation*, 41(10):1734–1737, 2020.
- 710 [36] Quan Li and Kai Wang. Intervar: clinical interpretation of genetic variants by  
711 the 2015 acmg-amp guidelines. *The American Journal of Human Genetics*, 100  
712 (2):267–280, 2017.

- 713 [37] Brent S Pedersen, Joe M Brown, Harriet Dashnow, Amelia D Wallace, Matt  
714 Velinder, Martin Tristani-Firouzi, Joshua D Schiffman, Tatiana Tvardik, Rong  
715 Mao, D Hunter Best, et al. Effective variant filtering and expected candidate  
716 variant yield in studies of rare human disease. *NPJ Genomic Medicine*, 6(1):1–8,  
717 2021.
- 718 [38] Carl A Anderson, Fredrik H Pettersson, Geraldine M Clarke, Lon R Cardon,  
719 Andrew P Morris, and Krina T Zondervan. Data quality control in genetic  
720 case-control association studies. *Nature protocols*, 5(9):1564–1573, 2010. URL  
721 <https://doi.org/10.1038/nprot.2010.116>.
- 722 [39] David T Miller, Kristy Lee, Noura S Abul-Husn, Laura M Amendola, Kyle Broth-  
723 ers, Wendy K Chung, Michael H Gollob, Adam S Gordon, Steven M Harrison,  
724 Ray E Hershberger, et al. Acmg sf v3. 2 list for reporting of secondary findings  
725 in clinical exome and genome sequencing: a policy statement of the american  
726 college of medical genetics and genomics (acmg). *Genetics in Medicine*, 25(8):  
727 100866, 2023.

## 6 Supplemental

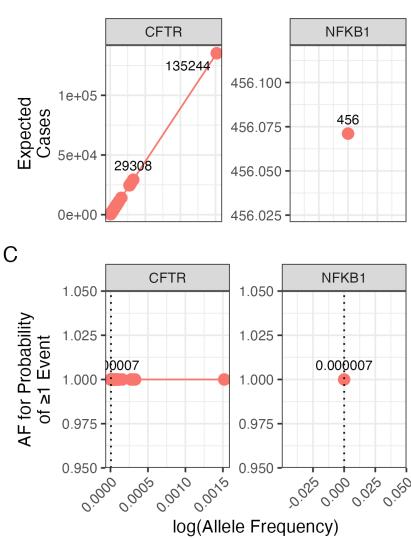
Condition: population size 69433632, phenotype PID-related, genes CFTR and NFKB1.

A

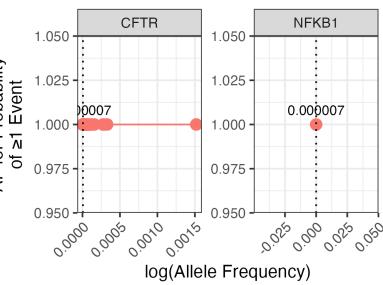


B

clinvar\_clnsig • Pathogenic



C



**Figure S1: Interpretation of probability of observing a variant classification.**  
The result from the chosen validation genes *CFTR* and *NFKB1* are shown. Case counts are dependant on the population size and phenotype. (A) The density plots of expected observations by ClinVar clinical significance. We then highlight the values for pathogenic variants specifically showing; (B) the allele frequency versus expected cases in this population size and (C) the probability of observing at least one event in this population size.

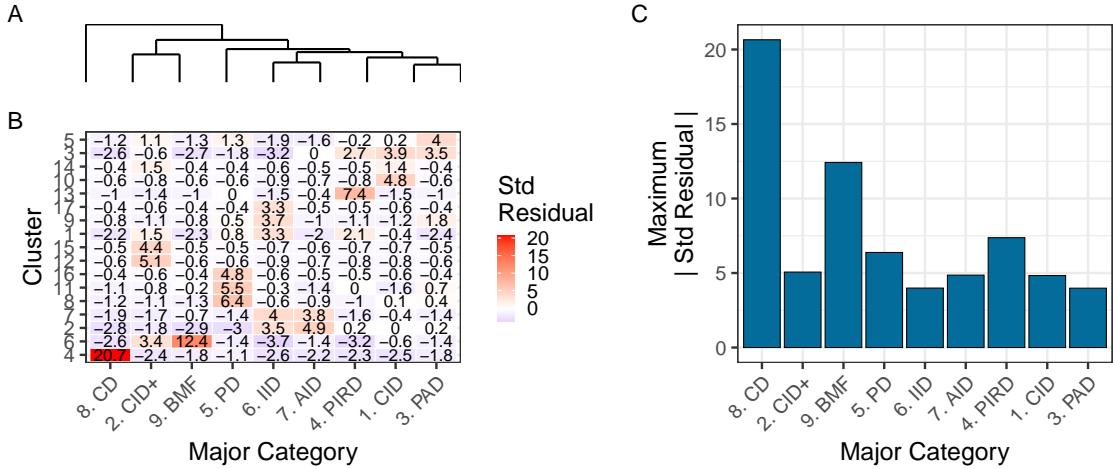


Figure S2: Hierarchical clustering of enrichment scores. The heatmap displays standardised residuals for major disease categories (x-axis) across network clusters (y-axis). A dendrogram groups similar disease categories, and the bar plot shows the maximum absolute residual per category. (8) CD and (9)BMF show the highest values, indicating significant enrichment or depletion (residuals  $> |2|$ ). Definitions in **Box 2.1**.

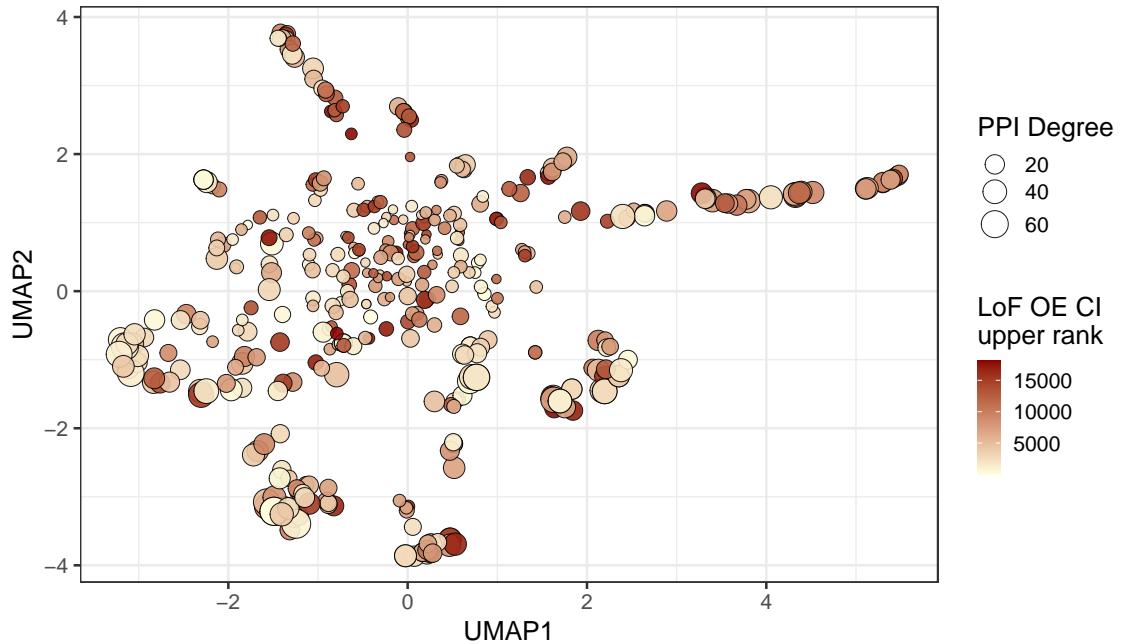
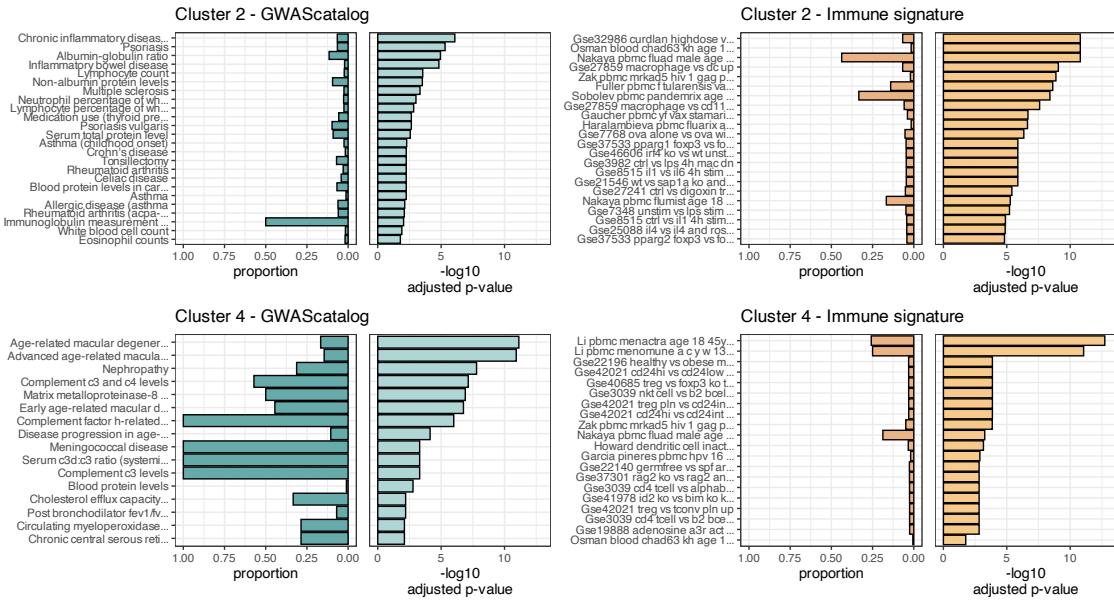


Figure S3: Supplementary analysis of PPI degree versus LOEUF upper rank with UMAP embedding of the PPI network. The relationship between PPI degree (size) and LOEUF upper rank (color) across gene clusters. No clear patterns are evident.



**Figure S4: Composite Enrichment Profiles for IEI Gene Sets.** We selected the top two enriched clusters (as per **Figure 6**) and performed functional enrichment analysis derived from known disease associations. For each gene set, the left panel displays the proportion of input genes overlapping with a curated gene set, and the right panel shows the  $-\log_{10}$  adjusted p-value from hypergeometric testing. These profiles, stratified by cluster (Cluster 2 and Cluster 4) and by gene set category (GWAScatalog and Immunologic Signatures), highlight distinct enrichment patterns that reflect differential pathogenic variant loads in the IEI gene panels.

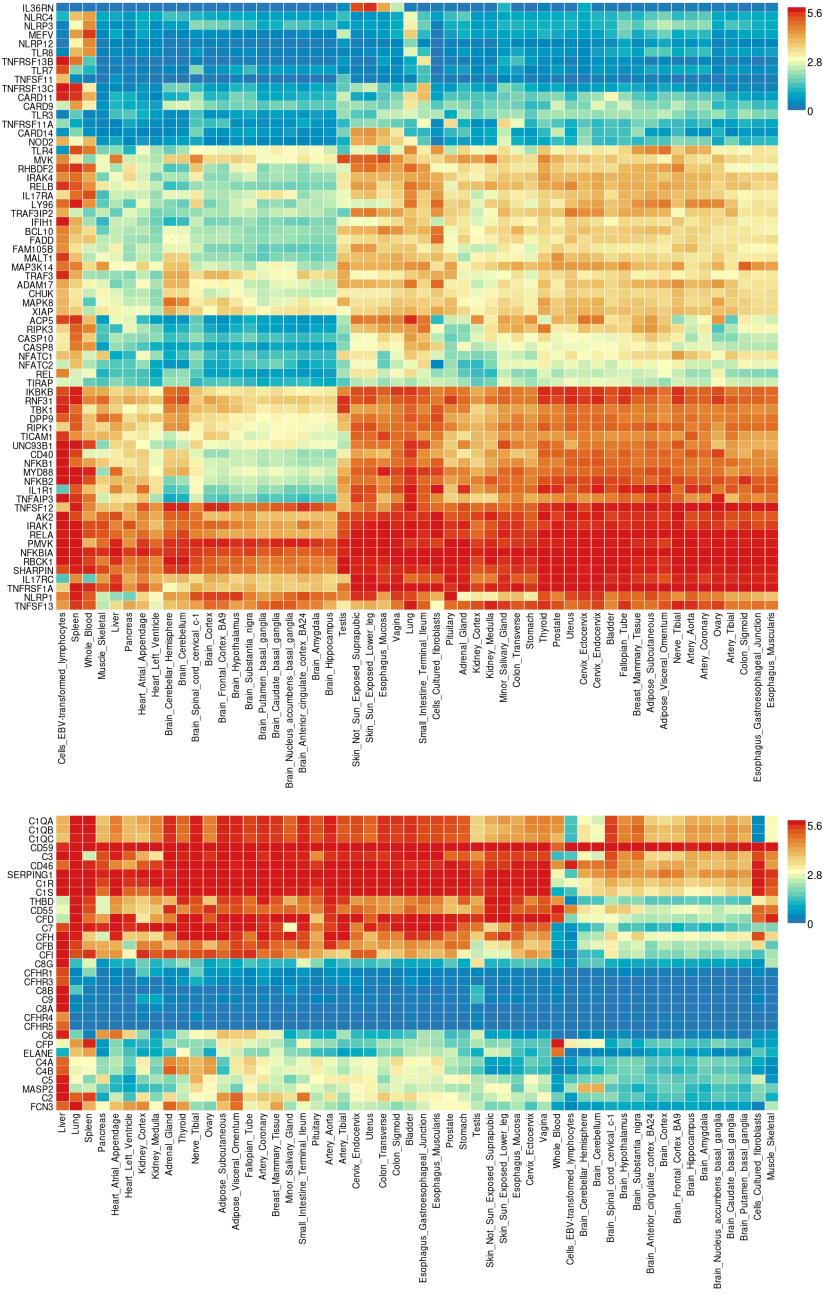


Figure S5: **Gene Expression Heatmaps for IEI Genes.** GTEx v8 data from 54 tissue types display the average expression per tissue label (log<sub>2</sub> transformed) for the IEI gene panels. Top: Cluster 2; Bottom: Cluster 4.

## 7 Clinical Genetics Application

In this section, we detail our approach to integrating sequencing data with prior pathogenicity evidence. Our method is designed to account for all possible outcomes of true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN), by first ensuring that every nucleotide corresponding to known pathogenic variants in a gene has been accurately sequenced. Only after confirming that these positions match the reference alleles (i.e. no unaccounted variant is present) do we calculate the probability that additional, alternative pathogenic variants (those not observed in the sequencing data) could be present. Our confidence interval (CI) for pathogenicity thus incorporates uncertainty from the entire process, including the tally of TP, FP, TN, and FN outcomes.

### 7.1 Methods

#### 7.1.1 Quality Control:

Before performing any probability calculations, we inspect the gVCF to confirm that all known pathogenic variant positions in the gene are adequately covered and appear as reference alleles. This step not only verifies true negatives (TN) but also flags instances where sequencing quality is insufficient, leading to missing sequence information, and prevents false confidence. For example, if a nucleotide position corresponding to a known pathogenic variant has low quality reads and fails QC, it is flagged as missing, thereby affecting the overall probability estimate for unobserved variants.

#### 7.1.2 Prior Probability Calculation:

For variants with an established ClinVar classification, the occurrence probability is derived directly from the allele frequency. For variants lacking a ClinVar label (i.e. variants of uncertain significance, VUS), we utilise an ACMG evidence score (0–100) to compute a prior probability as follows:

1. **Convert the ACMG Score:** The evidence score  $S$  is normalised to a fractional support level:

$$S_{\text{adj}} = \frac{S}{100}$$

This value reflects the strength of the pathogenic support.

2. **Assign a Minimal Risk ( $\epsilon$ ):** In the absence of a ClinVar classification, we assign a minimal risk based on the maximum observed allele number,  $\max(AN)$ , scaled by the evidence support:

$$\epsilon = \frac{1}{\max(AN) + 1} \times S_{\text{adj}}$$

756 This step ensures that even low-frequency variants receive a baseline risk pro-  
757 portional to the qualitative evidence.

- 758
- 759 3. **Adjust the Allele Frequency:** The observed allele frequency  $p_i$  is then in-  
creased by  $\epsilon$  to yield an adjusted frequency:

$$p_i^{\text{adj}} = p_i + \epsilon$$

758 This adjusted frequency reflects both the empirical observation and the ACMG  
759 evidence.

- 760 4. **Calculate the Prior Probability of Disease:**

- For **Autosomal Dominant (AD)** or **X-Linked (XL)** inheritance, the prior probability is:

$$p_{\text{disease}} = p_i^{\text{adj}}$$

- For **Autosomal Recessive (AR)** inheritance—which considers both homozygosity and compound heterozygosity—the probability is calculated as:

$$p_{\text{disease}} = \left( p_i^{\text{adj}} \right)^2 + 2 p_i^{\text{adj}} \left( P_{\text{tot}} - p_i^{\text{adj}} \right)$$

where

$$P_{\text{tot}} = \sum_{j \in \text{gene}} p_j^{\text{adj}}$$

761 **7.1.3 Deriving the Confidence Interval (CI)**

762 To capture uncertainty from all possible outcomes (TP, FP, TN, FN) in our sequencing  
763 and variant classification process, we propagate the variance arising from:

- 764
- 765
- 766
- The observed allele frequency and its adjustment via  $\epsilon$ .
  - The potential misclassification of variants (e.g. a VUS might be miscalled, contributing to FP or FN counts).
  - Missing sequence data at known pathogenic sites.

768 We demonstrate two methods for deriving the 95% CI of the final occurrence  
769 probability: (1) the Wilson score interval and (2) a Bayesian credible interval using  
770 a Beta distribution.

771 **1. Wilson Score Interval** Assume the adjusted occurrence probability is esti-  
 772 mated as  $\hat{p} = p_i^{\text{adj}}$  based on an effective sample size  $N$  (which reflects the number  
 773 of informative reads or quality-controlled observations). The Wilson score interval is  
 774 computed as:

$$\hat{p}_W = \frac{\hat{p} + \frac{z^2}{2N}}{1 + \frac{z^2}{N}}$$

$$\text{Margin} = \frac{z}{1 + \frac{z^2}{N}} \sqrt{\frac{\hat{p}(1 - \hat{p})}{N} + \frac{z^2}{4N^2}}$$

$$\text{CI}_{\text{Wilson}} = [\hat{p}_W - \text{Margin}, \hat{p}_W + \text{Margin}]$$

775 where  $z = 1.96$  for a 95% confidence level. This interval integrates uncertainty from  
 776 the adjusted allele frequency and any variability in the count data.

**2. Bayesian Credible Interval** Alternatively, we can model the uncertainty using a Bayesian framework. Suppose that, after accounting for TP, FP, TN, and FN outcomes, the posterior distribution of the pathogenic probability is approximated by a Beta distribution,  $\text{Beta}(\alpha, \beta)$ . Here, the parameters  $\alpha$  and  $\beta$  are chosen based on the effective counts of “successes” (e.g. detection or strong evidence of pathogenicity) and “failures” (e.g. absence or refutation), respectively. For example, if  $k$  is the effective number of positive events and  $N - k$  the negatives, then:

$$\alpha = k + 1, \quad \beta = N - k + 1.$$

The 95% credible interval is then given by the 2.5th and 97.5th percentiles of the Beta distribution:

$$\text{CI}_{\text{Bayesian}} = [\text{BetaInv}(0.025; \alpha, \beta), \text{BetaInv}(0.975; \alpha, \beta)],$$

777 where  $\text{BetaInv}(q; \alpha, \beta)$  denotes the quantile function of the Beta distribution at prob-  
 778 ability  $q$ .

779 Both methods integrate the uncertainty from the observed data, the adjustment  
 780 via  $\epsilon$  from the ACMG evidence score, and the potential misclassification or missing  
 781 sequence data. In our analysis, the resulting 95% CI for pathogenicity is derived from  
 782 such propagation of uncertainty, ensuring that all outcomes (TP, FP, TN, FN) are  
 783 reflected in the final confidence bounds.

## 784 7.2 Results

785 We illustrate our method with two examples:

786 **Example 1: Missing Sequence Information** In one case, a known pathogenic  
787 nucleotide position in *GENE\_XYZ* exhibited low quality reads and did not pass QC.  
788 This missing information prevents confirmation of the absence of the known variant (a  
789 potential false negative), thereby widening the uncertainty in our probability estimate.  
790 In such cases, the adjusted allele frequency is calculated with additional variance,  
791 leading to a broader CI. For instance, if the observed allele frequency is  $1.0 \times 10^{-5}$   
792 and after adjusting with the ACMG score the estimated occurrence probability is  
793  $1.0 \times 10^{-5}$ , the propagated uncertainty might yield a 95% CI of [0.70, 0.85]. This  
794 broader interval reflects the impact of missing sequence data on our confidence.

795 **Example 2: Heterozygous Variant in an Autosomal Recessive Gene** In  
796 another case, a patient carries a heterozygous variant in an autosomal recessive (AR)  
797 gene. In this scenario, there is also a second VUS in the same gene. Both variants  
798 are assessed using the ACMG evidence score adjustment. Their adjusted allele fre-  
799 quencies are used to compute the overall prior probability of disease, accounting for  
800 the possibility of compound heterozygosity. The two VUS are then ranked based on  
801 their evidence and the resulting 95% CIs. For instance, one variant may yield an  
802 occurrence probability of  $2.5 \times 10^{-4}$  with a 95% CI of [0.80, 0.88], while the other  
803 might have a lower probability of  $1.8 \times 10^{-4}$  with a CI of [0.75, 0.83]. The variant  
804 with the higher occurrence probability and narrower CI would be ranked as the more  
805 likely causal variant in the context of AR inheritance.

806 Table S1 shows the final variant results for a male patient carrying an X-linked  
807 loss-of-function (LOF) variant in *GENE\_XYZ* where all known pathogenic positions  
808 were confirmed as reference alleles. For the variant c. 1234del (p.Glu412Argfs\*5), the  
809 observed allele frequency is  $1.2 \times 10^{-5}$ . After applying the ACMG evidence score  
810 adjustment (for a VUS lacking a ClinVar classification), the adjusted allele frequency  
811 remains consistent with the observed data. The resulting occurrence probability is  
812  $1.2 \times 10^{-5}$ , and by propagating the uncertainty from the allele frequency, evidence  
813 score adjustment, and the full range of possible outcomes (TP, FP, TN, FN), we  
814 derive a 95% CI for causality of [0.92, 0.97]. This confirms the variant as the top  
815 causal variant in this patient, with no evidence of additional alternative pathogenic  
816 variants.

Table S1: Final Variant Results for Patient (XL LOF)

Parameter	Value
Gene	<i>GENE_XYZ</i>
Variant	c. 1234del (p.Glu412Argfs*5)
Variant Type	Loss-of-Function (LOF)
Inheritance	X-Linked (XL)
Patient Sex	Male (hemizygous)
Allele Frequency	$1.2 \times 10^{-5}$
Occurrence Probability	$1.2 \times 10^{-5}$
95% CI for Causality	[0.92, 0.97]
Clinical Interpretation	Top causal variant confirmed; no evidence of additional alternative pathogenic variants