

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

Dylan Lawless^{*1}

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

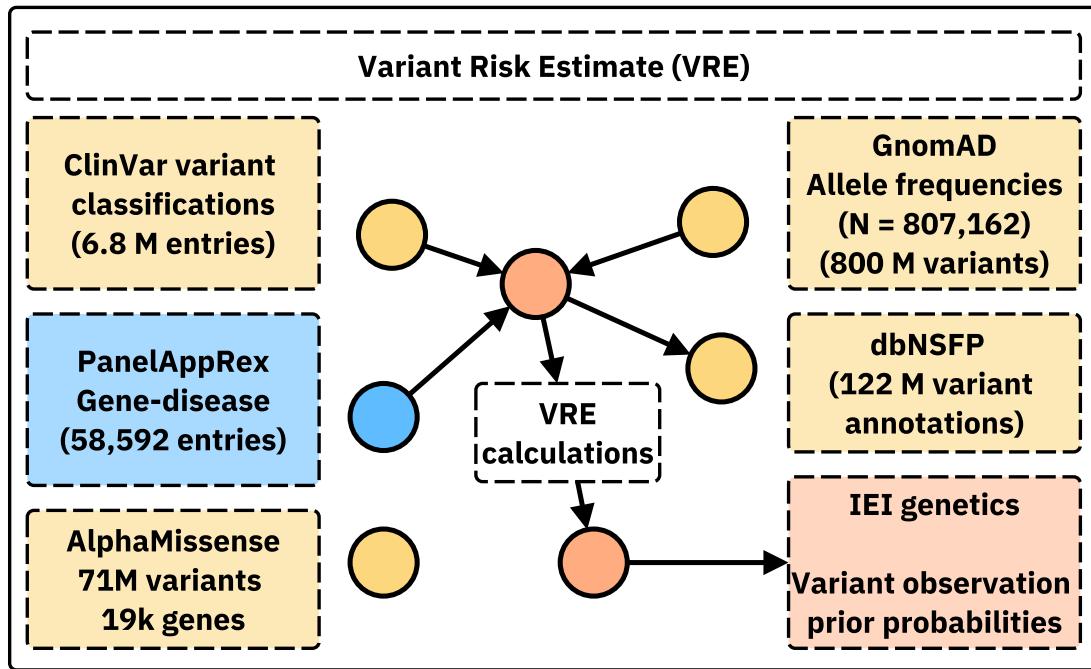
April 17, 2025

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.¹

^{*}Addresses for correspondence: Dylan.Lawless@kispi.uzh.ch

¹ **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. 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

¹⁹ Acronyms

²⁰ ACMG American College of Medical Genetics and Genomics.....	³⁰
²¹ ACAT Aggregated Cauchy Association Test	³⁰
²² AD Autosomal Dominant.....	⁴
²³ ANOVA Analysis of Variance	¹²
²⁴ AR Autosomal Recessive	⁴
²⁵ BMF Bone Marrow Failure.....	¹⁸
²⁶ CD Complement Deficiencies	¹⁹
²⁷ CI Confidence Interval.....	¹⁵
²⁸ CF Cystic Fibrosis	¹⁰
²⁹ CFTR Cystic Fibrosis Transmembrane Conductance Regulator.....	⁵
³⁰ CVID Common Variable Immunodeficiency	⁸
³¹ dbNSFP database for Non-Synonymous Functional Predictions	⁵
³² GE Genomics England	⁵
³³ gnomAD Genome Aggregation Database	⁵
³⁴ HGVS Human Genome Variation Society.....	⁵
³⁵ HPC High-Performance Computing.....	⁸
³⁶ HWE Hardy-Weinberg Equilibrium	⁴
³⁷ IEI Inborn Errors of Immunity.....	⁴
³⁸ Ig Immunoglobulin	²³
³⁹ InDel Insertion/Deletion	⁵
⁴⁰ IUIS International Union of Immunological Societies	⁶
⁴¹ LD Linkage Disequilibrium	²¹
⁴² LOEUF Loss-Of-function Observed/Expected Upper bound Fraction	¹²
⁴³ LOF Loss-of-Function	¹⁸
⁴⁴ MOI Mode of Inheritance	⁴
⁴⁵ NFKB1 Nuclear Factor Kappa B Subunit 1	⁵
⁴⁶ OMIM Online Mendelian Inheritance in Man	²⁷
⁴⁷ PID Primary Immunodeficiency	⁴
⁴⁸ PPI Protein-Protein Interaction	⁵
⁴⁹ SNV Single Nucleotide Variant	⁴
⁵⁰ SKAT Sequence Kernel Association Test.....	³⁰
⁵¹ STRINGdb Search Tool for the Retrieval of Interacting Genes/Proteins.....	⁵
⁵² HSD Honestly Significant Difference	¹²
⁵³ UMAP Uniform Manifold Approximation and Projection	¹⁸
⁵⁴ UniProt Universal Protein Resource.....	⁵
⁵⁵ VEP Variant Effect Predictor.....	⁵
⁵⁶ XL X-Linked	⁴

94 1 Introduction

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

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

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

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

133 value for AI and reinforcement learning applications, providing an enriched version of
134 the data underpinning frameworks such as AlphaFold (11) and AlphaMissense (12).

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

151 2 Methods

152 2.1 Dataset

153 Data from Genome Aggregation Database (gnomAD) v4 comprised 807,162 indi-
154 viduals, including 730,947 exomes and 76,215 genomes (7). This dataset provided
155 786,500,648 SNVs and 122,583,462 InDels, with variant type counts of 9,643,254 syn-
156 onymous, 16,412,219 missense, 726,924 nonsense, 1,186,588 frameshift and 542,514
157 canonical splice site variants. ClinVar data were obtained from the variant summary
158 dataset (as of: 16 March 2025) available from the NCBI FTP site, and included
159 6,845,091 entries, which were processed into 91,319 gene classification groups and a
160 total of 38,983 gene classifications; for example, the gene *A1BG* contained four vari-
161 ants classified as likely benign and 102 total entries (13). For our analysis phase
162 we also used dbNSFP which consisted of a number of annotations for 121,832,908
163 SNVs (15). The PanelAppRex core model contained 58,592 entries consisting of
164 52 sets of annotations, including the gene name, disease-gene panel ID, diseases-
165 related features, confidence measurements. (2) A Protein-Protein Interaction (PPI)
166 network data was provided by Search Tool for the Retrieval of Interacting Genes/Pro-
167 teins (STRINGdb), consisting of 19,566 proteins and 505,968 interactions (16). The
168 Human Genome Variation Society (HGVS) nomenclature is used with Variant Effect
169 Predictor (VEP)-based codes for variant IDs. We carried out validations for disease
170 cohorts with Nuclear Factor Kappa B Subunit 1 (*NFKB1*) (17–20) and Cystic Fibrosis
171 Transmembrane Conductance Regulator (*CFTR*) (21–23) to demonstrate applications

172 in AD and AR disease genes, respectively. AlphaMissense includes pathogenicity pre-
 173 diction classifications for 71 million variants in 19 thousand human genes (12; 26).
 174 We used these scores to compare against the probability of observing the same given
 175 variants. **Box 2.1** list the definitions from the International Union of Immunological
 176 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

177

2.2 Variant Class Observation Probability

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

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

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

185 Our computational pipeline estimated the probability of observing a disease-associated
 186 genotype for each variant and aggregated these probabilities by gene and ClinVar
 187 classification. This approach included all variant classifications, not limited solely to
 188 those deemed “pathogenic”, and explicitly conditioned the classification on the given
 189 phenotype, recognising that a variant could only be considered pathogenic relative to
 190 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:

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

191 where $\max(AN)$ was the maximum allele number observed for that gene. This adjustment
192 ensured that a nonzero risk was incorporated even in the absence of observed
193 variants.

194 **2. Occurrence Probability Based on MOI.** The probability that an individual
195 was affected by a variant depended on the mode of MOI relative to a specific pheno-
196 type. Specifically, we calculated the occurrence probability $p_{\text{disease},i}$ for each variant
197 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:

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

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

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

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

216 **2.3 Validation of Autosomal Dominant Estimates Using *NFKB1***

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

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

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

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

228 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

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

²²⁹ with a median value of approximately 118 and a 95% confidence interval of 70 to 181
²³⁰ cases (derived from posterior sampling).

²³¹ **4. Bayesian Adjustment.** Recognising that the clinical cohort likely represents
²³² nearly all CVID cases (besides first-second degree relatives), two Bayesian adjust-
²³³ ments were performed:

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

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

²³⁴ 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

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

²³⁵ which yielded a Bayesian mixture adjusted median estimate of 67 cases with a
²³⁶ 95% credible interval of approximately 43 to 99 cases.

²³⁷ **5. Predicted Total Genotype Counts.** The predicted total synthetic genotype
²³⁸ count (before adjustment) was 456, whereas the predicted total genotypes adjusted
²³⁹ for *synth_flag* was 0. This higher synthetic count was set based on a minimal risk
²⁴⁰ threshold, ensuring that at least one genotype is assumed to exist (e.g. accounting for
²⁴¹ a potential unknown de novo variant) even when no variant is observed in gnomAD
²⁴² (as per **section 2.2**).

243 **6. Validation Test.** Thus, the expected number of *NFKB1*-related CVID cases
244 derived from our genome-wide probability estimates was compared with the observed
245 counts from the UK-based PID cohort. This comparison validates our framework for
246 estimating disease incidence in AD disorders.

247 2.4 Validation Study for Autosomal Recessive CF Using *CFTR*

248 To validate our framework for AR diseases, we focused on Cystic Fibrosis (CF).
249 For comparability sizes between the validation studies, we analysed the most com-
250 mon SNV in the *CFTR* gene, typically reported as “p.Arg117His” (GRCh38 Chr
251 7:117530975 G/A, MANE Select HGVS p.ENST00000003084.11: p.Arg117His). Our
252 goal was to validate our genome-wide probability estimates by comparing the ex-
253 pected number of CF cases attributable to the p.Arg117His variant in *CFTR* with
254 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:

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

using a large effective allele count (AN_{eff}) for illustration. By generating 10,000 posterior samples of p , we obtained a distribution of the literature-based adjusted expected counts, E_{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:

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

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

2.5 Validation of SCID-specific Estimates Using PID–SCID Genes

To validate our genome-wide probability estimates for diagnosing a genetic variant in a patient with a PID phenotype, we focused on a subset of genes implicated in Severe Combined Immunodeficiency (SCID). Given that the overall panel corresponds to PID, but SCID represents a rarer subset, the probabilities were converted to values per million PID cases.

1. Incidence Conversion. Based on literature, PID occurs in approximately 1 in 1,000 births, whereas SCID occurs in approximately 1 in 100,000 births. Consequently, in a population of 1,000,000 births there are about 1,000 PID cases and 10 SCID cases. To express SCID-related variant counts on a per-million PID scale, the observed SCID counts were multiplied by 100. For example, if a gene is expected to cause SCID in 10 cases within the total PID population, then on a per-million PID basis the count is $10 \times 100 = 1,000$ cases (across all relevant genes).

2. Prevalence Calculation and Data Adjustment. For each SCID-associated gene (e.g. *IL2RG*, *RAG1*, *DCLRE1C*), the observed variant counts in the dataset were adjusted by multiplying by 100 so that the probabilities reflect the expected number of cases per 1,000,000 PID. In this manner, our estimates are directly comparable to known counts from SCID cohorts, rather than to national population counts as in previous validation studies.

282 **3. Integration with Prior Probability Estimates.** The predicted genotype
283 occurrence probabilities were derived from our framework across the PID gene panel.
284 These probabilities were then converted to expected case counts per million PID
285 cases by multiplying by 1,000,000. For instance, if the probability of observing a
286 pathogenic variant in *IL2RG* is p , the expected SCID-related count becomes $p \times 10^6$.
287 Similar conversions are applied for all relevant SCID genes.

288 **4. Bayesian Uncertainty and Comparison with Observed Data.** To address
289 uncertainty in the SCID-specific estimates, a Bayesian uncertainty simulation was
290 performed for each gene to generate a distribution of predicted case counts on a
291 per-million PID scale. The resulting median estimates and 95% credible intervals
292 were then compared against known national SCID counts compiled from independent
293 registries. This comparison permitted a direct evaluation of our framework’s accuracy
294 in predicting the occurrence of SCID-associated variants within a PID cohort.

295 **5. Validation Test.** Thus, by converting the overall probability estimates to a
296 per-million PID scale, our framework was directly validated against observed counts
297 for SCID.

298 **2.6 Protein Network and Genetic Constraint Interpretation**

299 A PPI network was constructed using protein interaction data from STRINGdb (16).
300 We previously prepared and reported on this dataset consisting of 19,566 proteins and
301 505,968 interactions (<https://github.com/DylanLawless/ProteoMCLustR>). Node
302 attributes were derived from log-transformed score-positive-total values, which in-
303 formed both node size and colour. Top-scoring nodes (top 15 based on score) were
304 labelled to highlight prominent interactions. To evaluate group differences in score-
305 positive-total across major disease categories, one-way Analysis of Variance (ANOVA)
306 was performed followed by Tukey Honestly Significant Difference (HSD) post hoc tests
307 (and non-parametric Dunn’s test for confirmation). GnomAD v4.1 constraint metrics
308 data was used for the PPI analysis and was sourced from Karczewski et al. (7). This
309 provided transcript-level metrics, such as observed/expected ratios, Loss-Of-function
310 Observed/Expected Upper bound Fraction (LOEUF), pLI, and Z-scores, quantifying
311 loss-of-function and missense intolerance, along with confidence intervals and related
312 annotations for 211,523 observations.

313 **2.7 Gene Set Enrichment Test**

314 To test for overrepresentation of biological functions, the prioritised genes were com-
315 pared against gene sets from MsigDB (including hallmark, positional, curated, motif,
316 computational, GO, oncogenic, and immunologic signatures) and WikiPathways using
317 hypergeometric tests with FUMA (24; 25). The background set consisted of 24,304

318 genes. Multiple testing correction was applied per data source using the Benjamini-
319 Hochberg method, and gene sets with an adjusted P-value ≤ 0.05 and more than one
320 overlapping gene are reported.

321 **2.8 Deriving novel PID classifications by genetic PPI and
322 clinical features**

323 We recategorised 315 immunophenotypic features from the original IUIS IEI annotations,
324 reducing the original multi-level descriptors (e.g. “decreased cd8, normal or
325 decreased cd4”) first to minimal labels (e.g.“low”) and second to binary outcomes (normal
326 vs. not-normal) for T cells, B cells, neutrophils, and immunoglobulins Each gene
327 was mapped to its PPI cluster derived from STRINGdb and UMAP embeddings from
328 previous steps. We first tested for non-random associations between these four binary
329 immunophenotypes and PPI clusters using χ^2 tests. To generate a data-driven PID
330 classification, we trained a decision tree (rpart) to predict PPI cluster membership
331 from the four immunophenotypic features plus the traditional IUIS Major and Subcat-
332 egory labels. Hyperparameters (complexity parameter = 0.001, minimum split = 10,
333 minimum bucket = 5, maximum depth = 30) were optimised via five-fold cross vali-
334 dation using the caret framework. Terminal node assignments were then relabelled
335 according to each group’s predominant abnormal feature profile.

336 **2.9 Probability of observing AlphaMissense pathogenicity**

337 We obtained the subset pathogenicity predictions from AlphaMissense via the Al-
338 phaFold database and whole genome data from the studies data repository(12; 26).
339 The AlphaMissense data (genome-aligned and amino acid substitutions) were merged
340 with the panel variants based on genomic coordinate and HGVS annotation. Occur-
341 rence probabilities were log-transformed and adjusted (y-axis displaying $\log_{10}(\text{occurrence}$
342 prob + 1e-5) + 5), to visualise the distribution of pathogenicity scores across the
343 residue sequence. A Kruskal-Wallis test was used to compare the observed disease
344 probability across clinical classification groups.

345 **3 Results**

346 **3.1 Observation Probability Across Disease Genes**

347 Our study integrated large-scale annotation databases with gene panels from Pan-
348 elAppRex to systematically assess disease genes by MOI. By combining population
349 allele frequencies with ClinVar clinical classifications, we computed an expected obser-
350 vation probability for each SNV, representing the likelihood of encountering a variant
351 of a specific pathogenicity for a given phenotype. We report these probabilities for
352 54,814 ClinVar variant classifications across 557 genes (linked dataset (27)).

353 In practice, our approach computed a simple observation probability for every
 354 SNV across the genome and was applicable to any disease-gene panel. Here, we fo-
 355 cused on panels related to Primary Immunodeficiency or Monogenic Inflammatory
 356 Bowel Disease, using PanelAppRex panel ID 398 as a case study. **Figure 1** dis-
 357 plays all reported ClinVar variant classifications for this panel. The resulting natural
 358 scaling system (-5 to +5) accounts for the frequently encountered combinations of
 359 classification labels (e.g. benign to pathogenic). The resulting data set (27) is briefly
 360 shown in **Table 1** to illustrate that our method yielded estimations of the probability
 361 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 PanelAppRex’s panel: (ID 398) Primary immunodeficiency or monogenic inflammatory bowel disease. “ClinVar Significance” indicates the pathogenicity classification assigned 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. (27)

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

362 3.2 Validation studies

363 3.2.1 Validation of Dominant Disease Occurrence with *NFKB1*

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

371 Our model calculated 0 known pathogenic variant *NFKB1*-related CVID cases
 372 in the UK with a minimal risk of 456 unknown de novo variants. In the reference
 373 cohort, 16 *NFKB1* CVID cases were reported. We additionally wanted to account for
 374 potential under-reporting in the reference study. We used an extrapolated national
 375 CVID prevalence which yielded a median estimate of 118 cases (95% CI: 70–181),
 376 while a Bayesian-adjusted mixture estimate produced a median of 67 cases (95% CI:
 377 43–99). **Figure S1 (A)** illustrates that our predicted values reflect these ranges and
 378 are closer to the observed count. This case supports the validity of our integrated

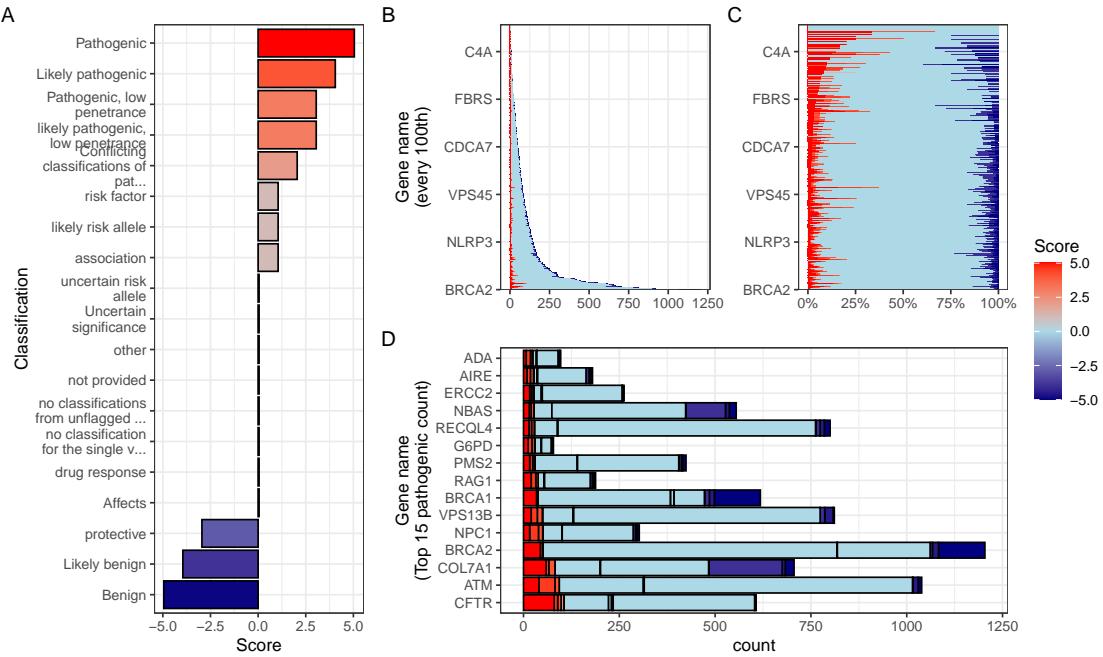


Figure 1: **Summary of ClinVar clinical significance classifications in the PID gene panel.** (A) Shows the numeric score coding for each classification. Panels (B) and (C) display the tally of classifications per gene as absolute counts and as percentages, respectively. (D) Highlights the top 15 genes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

379 probability estimation framework for AD disorders, and represents a challenging ex-
 380 ample where pathogenic SNV are not reported in the reference population of gnomAD.
 381 Our min-max values successfully contained the true reported values.

382 3.2.2 Validation of Recessive Disease Occurrence with *CFTR*

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

388 In contrast, the nationally reported number of CF cases was 714, as recorded in the
 389 UK Cystic Fibrosis Registry 2023 Annual Data Report (21). To account for factors
 390 such as reduced penetrance and the mortality-adjusted expected genotype, we derived
 391 a Bayesian-adjusted estimate via posterior simulation. Our Bayesian approach yielded
 392 a median estimate of 740 cases (95% Confidence Interval (CI): 696, 786) and a
 393 mixture-based estimate of 727 cases (95% CI: 705, 750). **Figure S1 (B)** illustrates
 394 the close concordance between the predicted values, the Bayesian-adjusted estimates,
 395 and the national report supports the validity of our approach for estimating disease.

396 **Figure S2** shows the final values for these genes of interest in a given population
397 size and phenotype. It reveals that an allele frequency threshold of approximately
398 0.000007 is required to observe a single heterozygous disease-causing variant carrier in
399 the UK population for both genes. However, owing to the AR MOI pattern of *CFTR*,
400 this threshold translates into more than 100,000 heterozygous carriers, compared to
401 only 456 carriers for the AD gene *NFKB1*. Note that this allele frequency threshold,
402 being derived from the current reference population, represents a lower bound that
403 can become more precise as public datasets continue to grow. This marked difference
404 underscores the significant impact of MOI patterns on population carrier frequencies
405 and the observed disease prevalence.

406 **3.2.3 Interpretation of ClinVar Variant Observations**

407 **Figure S9** shows the two validation study PID genes, representing AR and dominant
408 MOI. **Figure S9 (A)** illustrates the overall probability of an affected birth by ClinVar
409 variant classification, whereas **Figure S9 (B)** depicts the total expected number of
410 cases per classification for an example population, here the UK, of approximately 69.4
411 million.

412 **3.2.4 Validation of SCID-specific Disease Occurrence**

413 Given that SCID is a subset of PID, our probability estimates reflect the likelihood of
414 observing a genetic variant as a diagnosis when the phenotype is PID. However, we
415 additionally tested our results against SCID cohorts in **Figure S4**. The summarised
416 raw cohort data for SCID-specific gene counts are summarised and compared across
417 countries in **Figure S3**. True counts for *IL2RG* and *DCLRE1C* from ten distinct
418 locations yielded 95% confidence intervals surrounding our predicted values. For
419 *IL2RG*, the prediction was low (approximately 1 case per 1,000,000 PID), as expected
420 since loss-of-function variants in this X-linked gene are highly deleterious and rarely
421 observed in gnomAD. In contrast, the predicted value for *RAG1* was substantially
422 higher (553 cases per 1,000,000 PID) than the observed counts (ranging from 0 to
423 200). We attributed this discrepancy to the lower penetrance and higher background
424 frequency of *RAG1* variants in recessive inheritance, whereby reference studies may
425 underreport the true national incidence. Overall, we argued that agreement within
426 an order of magnitude was tolerable given the inherent uncertainties from reference
427 studies arising from variable penetrance and allele frequencies.

428 **3.3 Genetic constraint in high-impact protein networks**

429 We next examined genetic constraint in high-impact protein networks across the whole
430 IEI gene set of over 500 known disease-gene phenotypes (1). By integrating ClinVar
431 variant classification scores with PPI data, we quantified the pathogenic burden per

432 gene and assessed its relationship with network connectivity and genetic constraint
 433 (7; 16).

434 **3.3.1 Score-Positive-Total within IEI PPI network**

435 The ClinVar classifications reported in **Figure 1** were scaled -5 to +5 based on their
 436 pathogenicity. We were interested in positive (potentially damaging) but not negative
 437 (benign) scoring variants, which are statistically incidental in this analysis. We tallied
 438 gene-level positive scores to give the score positive total metric. **Figure 2 (A)** shows
 439 the PPI network of disease-associated genes, where node size and colour encode the
 440 score positive total (log-transformed). The top 15 genes with the highest total prior
 441 probabilities of being observed with disease are labelled (as per **Figure 1**).

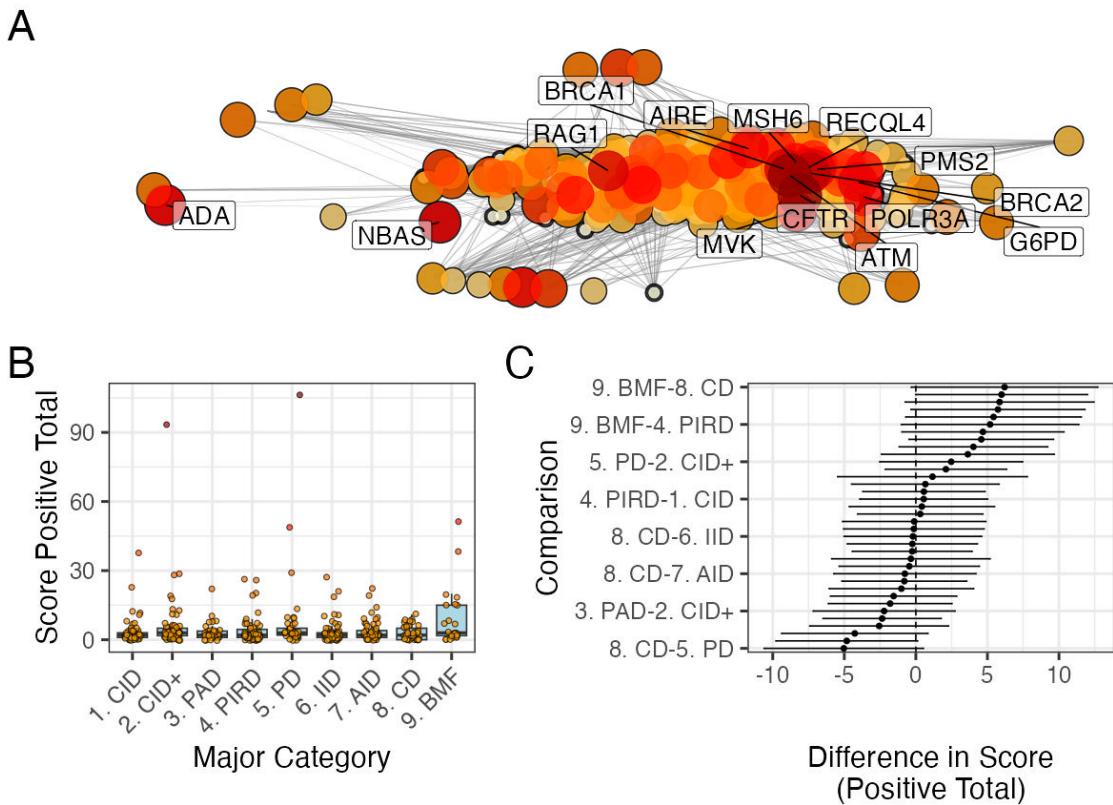


Figure 2: **PPI network and score positive total ClinVar significance variants.** (A) PPI network of disease-associated genes. Node size and colour represent the log-transformed score positive total, the top 15 genes/proteins with the highest probability of being observed in disease are labelled. (B) 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.

442 **3.3.2 Association Analysis of Score-Positive-Total across IEI Categories**

443 We checked for any statistical enrichment in score positive totals, which represents
444 the expected observation of pathogenicity, between the IEI categories. The one-way
445 ANOVA revealed an effect of major disease category on score positive total ($F(8, 500) =$
446 2.82, $p = 0.0046$), indicating that group means were not identical, which we observed
447 in **Figure 2 (B)**. However, despite some apparent differences in median scores across
448 categories (i.e. 9. Bone Marrow Failure (BMF)), the Tukey HSD post hoc compar-
449 isons **Figure 2 (C)** showed that all pairwise differences had 95% confidence intervals
450 overlapping zero, suggesting that individual group differences were not significant.

451 **3.3.3 UMAP Embedding of the PPI Network**

452 To address the density of the PPI network for the IEI gene panel, we applied Uniform
453 Manifold Approximation and Projection (UMAP) (**Figure 3**). Node sizes reflect
454 interaction degree, a measure of evidence-supported connectivity (16). We tested
455 for a correlation between interaction degree and score positive total. In **Figure**
456 **3**, gene names with degrees above the 95th percentile are labelled in blue, while
457 the top 15 genes by score positive total are labelled in yellow (as per **Figure 1**).
458 Notably, genes with high pathogenic variant loads segregated from highly connected
459 nodes, suggesting that Loss-of-Function (LOF) in hub genes is selectively constrained,
460 whereas damaging variants in lower-degree genes yield more specific effects. This
461 observation was subsequently tested empirically.

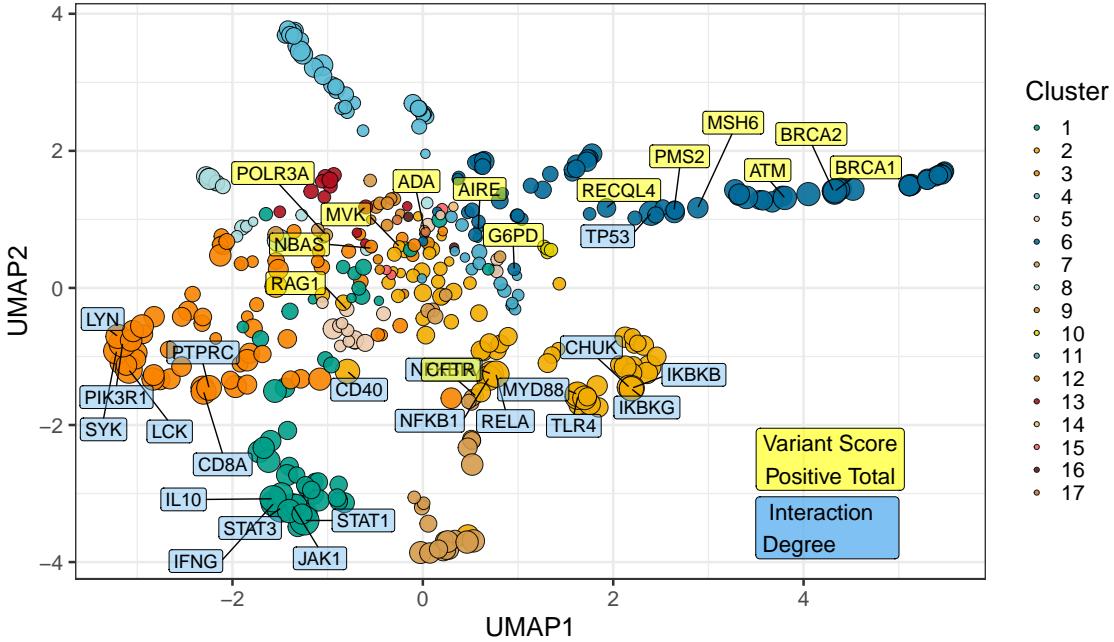


Figure 3: **UMAP embedding of the PPI network (p_umap).** The plot projects the high-dimensional protein-protein interaction network into two dimensions, with nodes coloured by cluster and sized by interaction degree. Blue labels indicate hub genes (degree above the 95th percentile) and yellow labels 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.

462 **3.3.4 Hierarchical Clustering of Enrichment Scores for Major Disease Cat-**
 463 **egories**

464 **Figure S5** presents a heatmap of standardised residuals for major disease categories
 465 across network clusters, as per **Figure 3**. A dendrogram clusters similar disease cate-
 466 gories, while the accompanying bar plot displays the maximum absolute standardised
 467 residual for each category. Notably, (8) Complement Deficiencies (CD) shows the
 468 highest maximum enrichment, followed by (9) BMF. While all maximum values
 469 exceed 2, the threshold for significance, this likely reflects the presence of protein
 470 clusters with strong damaging variant scores rather than uniform significance across
 471 all categories (i.e. genes from cluster 4 in 8 CD).

472 **3.3.5 PPI Connectivity, LOEUF Constraint and Enriched Network Clus-**
 473 **ter Analysis**

474 Based on the preliminary insight from **Figure S5**, we evaluated the relationship
 475 between network connectivity (PPI degree) and LOEUF constraint (LOEUF upper rank)
 476 Karczewski et al. (7) using Spearman's rank correlation. Overall, there was a weak

477 but significant negative correlation ($\rho = -0.181$, $p = 0.00024$) at the global scale,
 478 indicating that highly connected genes tend to be more constrained. A supplementary
 479 analysis (**Figure S6**) did not reveal distinct visual associations between network
 480 clusters and constraint metrics, likely due to the high network density. However
 481 once stratified by gene clusters, the natural biological scenario based on quantitative
 482 PPI evidence (16), some groups showed strong correlations; for instance, cluster 2
 483 ($\rho = -0.375$, $p = 0.000994$) and cluster 4 ($\rho = -0.800$, $p < 0.000001$), while others did
 484 not. This indicated that shared mechanisms within pathway clusters may underpin
 485 genetic constraints, particularly for LOF intolerance. We observe that the score
 486 positive total metric effectively summarises the aggregate pathogenic burden across
 487 IEI genes, serving as a robust indicator of genetic constraint and highlighting those
 488 with elevated disease relevance.

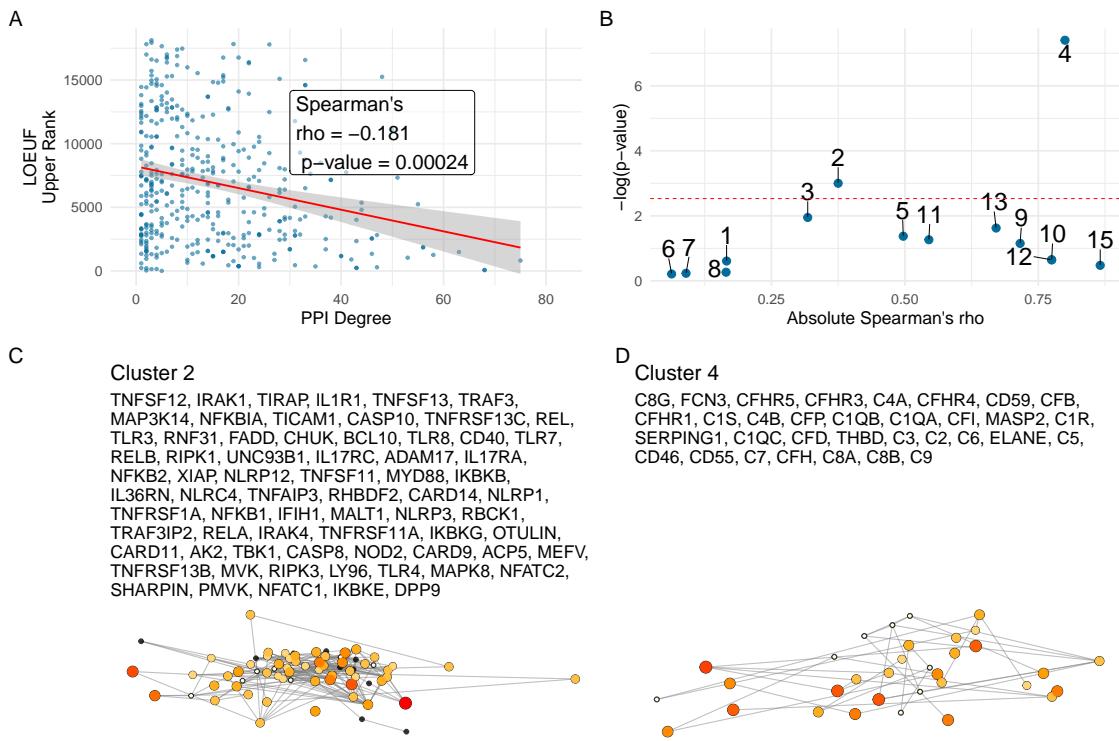


Figure 4: **Correlation between PPI degree and LOEUF upper rank.** (A) Ananlysis across all genes revealed a weak, significant negative correlation between PPI degree and LOEUF upper rank. (B) The cluster-wise analysis showed that clusters 2 and 4 exhibited moderate to strong correlations, while other 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.

489 **Figure 4 (C, D)** shows the re-plotted PPI networks for clusters with significant
 490 correlations between PPI degree and LOEUF upper rank. In these networks, node
 491 size is scaled by a normalised variant score, while node colour reflects the variant
 492 score according to a predefined palette.

493 3.4 New Insight from Functional Enrichment

494 To interpret the functional relevance of our prioritised IEI gene sets with the highest
495 load of damaging variants (i.e. clusters 2 and 4 in **Figure 4**), we performed func-
496 tional enrichment analysis for known disease associations using MsigDB with FUMA
497 (i.e. GWAScatalog and Immunologic Signatures) (24). Composite enrichment pro-
498 files (**Figure S7**) reveal that our enriched PPI clusters were associated with distinct
499 disease-related phenotypes, providing functional insights beyond traditional IUIS IEI
500 groupings (1). The gene expression profiles shown in **Figure S8** (GTEx v8 54 tissue
501 types) offer the tissue-specific context for these associations. Together, these results
502 enable the annotation of IEI gene sets with established disease phenotypes, supporting
503 a data-driven classification of IEI.

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

513 3.5 Genome-wide Gene Distribution and Locus-specific Vari- 514 ant Occurrence

515 **Figure 5 (A)** shows a genome-wide karyoplot of all IEI panel genes across GRCh38,
516 with colour-coding based on MOI. Figures **(B)** and **(C)** display zoomed-in locus plots
517 for *NFKB1* and *CFTR*, respectively. In **Figure 5 (B)**, the probability of observing
518 variants with known classifications is high only for variants such as p.Ala475Gly,
519 which are considered benign in the AD *NFKB1* gene that is intolerant to LOF. In
520 **Figure 5 (C)**, high probabilities of observing patients with pathogenic variants in
521 *CFTR* are evident, reproducing this well-established phenomenon. Furthermore, the
522 analysis of Linkage Disequilibrium (LD) using R^2 shows that high LD regions can be
523 modelled effectively, allowing independent variant signals to be distinguished.

524 3.6 Novel PID classifications derived from genetic PPI and 525 clinical features

526 We recategorised 315 immunophenotypic features from the original IUIS IEI annota-
527 tions, reducing detailed descriptions (e.g. “decreased cd8, normal or decreased cd4”),
528 first to minimal labels (e.g.“low”), and second to binary outcomes (normal vs. not-
529 normal) for T cells, B cells, neutrophils, and immunoglobulins (**Figure 6**). These

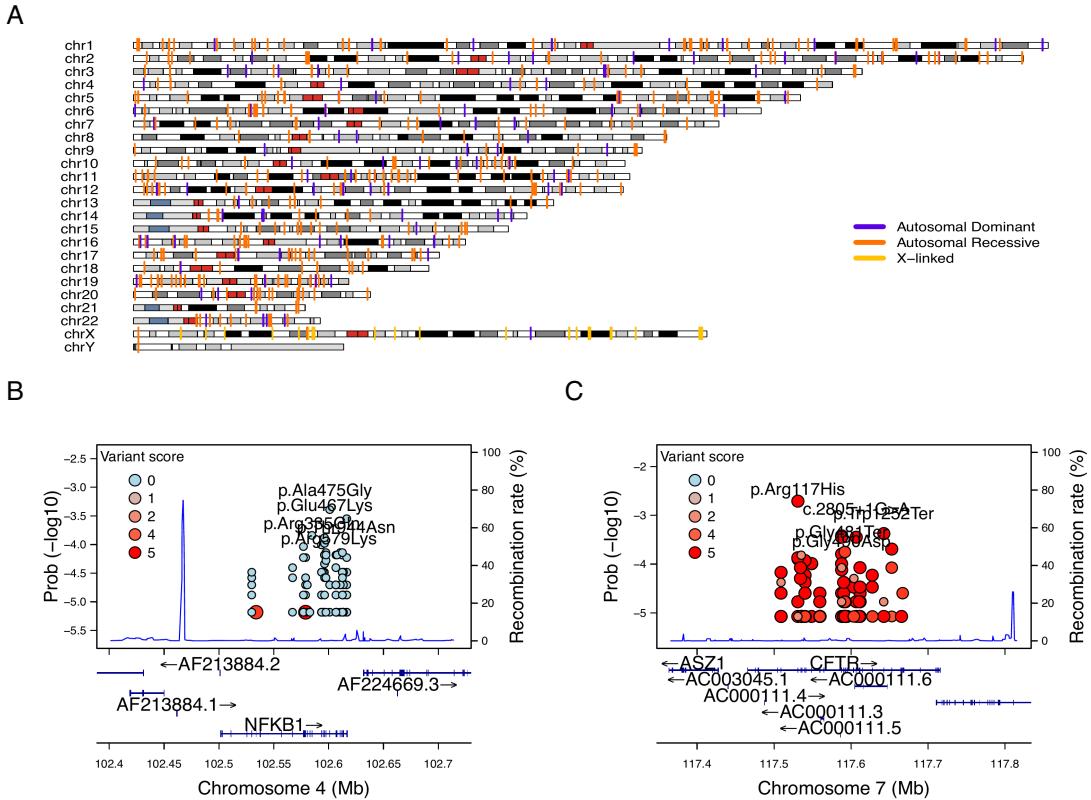


Figure 5: Genome-wide IEI, variant occurrence probability and LD by R^2 .
(A) Genome-wide karyoplot of all IEI panel genes mapped to GRCh38, with colours indicating MOI. (B) Zoomed-in locus plot for *NFKB1* showing variant observation probabilities; only benign variants such exhibit high probabilities in this AD gene intolerant to LOF. (C) Locus plot for *CFTR* 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.

530 simplified profiles were integrated with PPI network clustering from STRINGdb to
531 refine PID gene groupings. Chi-square analyses confirmed significant associations be-
532 between specific clinical abnormalities and PPI clusters (**Figure ??**). A decision tree
533 classifier, with hyperparameters optimised via 5-fold cross validation, demonstrated
534 high sensitivity and specificity, as shown in the confusion matrices and variable impor-
535 tance metrics (**Figure S10**). The resulting novel PID classifications, illustrated by
536 the decision tree and gene group distributions (**Figure 9**), provide a more coherent
537 and data-driven framework for categorising PID genes.

538 **3.7 Novel PID classifications derived from genetic PPI and**
539 **clinical features**

540 We recategorised 315 immunophenotypic features from the original IUIS IEI annotations,
541 reducing detailed descriptions (e.g. “decreased CD8, normal or decreased CD4”) to minimal labels
542 (e.g. “low”) and then binarising them (normal vs. not-normal) for T cells, B cells, Immunoglobulin (Ig) and neutrophils (**Figure 6**). These simplified profiles were mapped onto STRINGdb PPI clusters, revealing non-random distributions
544 ($\chi^2 < 1e-13$; **Figure 7**), indicating that network context captures key immunophenotypic variation.
546

547 We next compared four classifiers under 5-fold cross-validation to determine which
548 features predicted PPI clustering. As shown in **Figure 8**, the fully combined model
549 achieved the highest accuracy among the four: (i) phenotypes only (33 %) (i.e. T
550 cell, B cell, Ig, Neutrophil); (ii) phenotypes + IUIS major category (50 %) (e.g. CID.
551 See **Box 2.1** for more); (iii) IUIS major + subcategory only (59 %) (e.g. CID, T-B+
552 SCID); and (iv) phenotypes + IUIS major + subcategory (61 %). This demonstrated
553 that incorporating both traditional IUIS classifications and core immunophenotypic
554 markers into the PPI-based framework yielded the most robust discrimination of PID
555 gene clusters. Variable importance analysis highlighted abnormality status for Ig and
556 T cells were among the top ten features in addition to the other IUIS major and sub
557 categories. Per-class specificity remained uniform across the classes while sensitivity
558 dropped.

559 The PPI and immunophenotype model yielded 17 data-driven PID groups, whereas
560 incorporating the full complement of IUIS categories expanded this to 33 groups. For
561 clarity, we only demonstrate the decision tree from the smaller 17-group model in
562 **Figure 9**. Each terminal node is annotated by its predominant immunophenotypic
563 signature (for example, “group 65 with abnormal T cell and B cell features”), and the
564 full resulting gene counts per 33 class are plotted in **Figure 9**. Although, less user-
565 friendly, this data-driven taxonomy both aligns with and refines traditional IUIS IEI
566 classifications to provide a scaffold for large-scale computational analyses. Because
567 this framework is fully reproducible, alternative PPI embeddings or incorporate additional
568 molecular annotations can readily swapped to continue building on these PID
569 classification schemes.

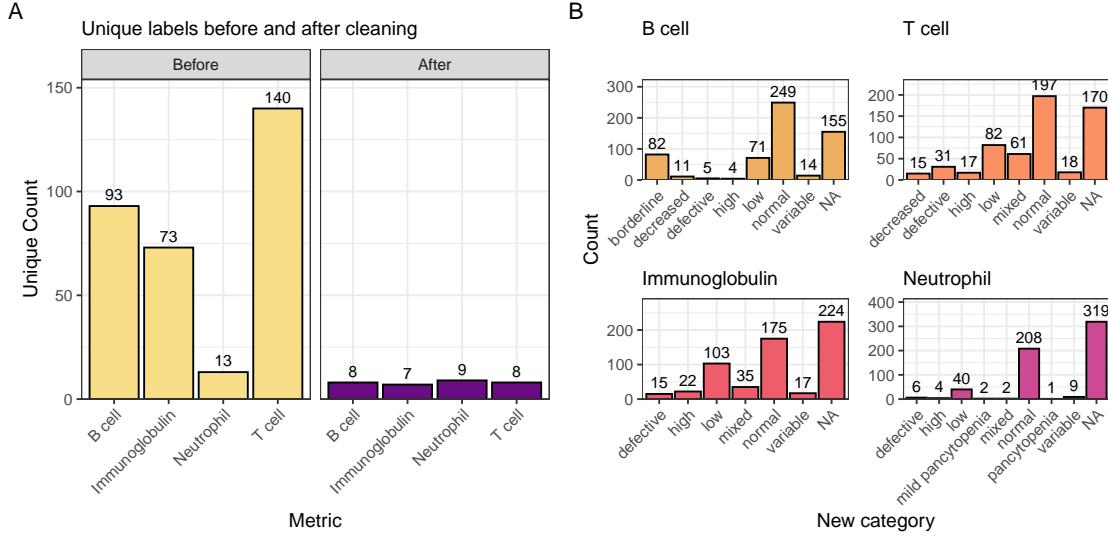


Figure 6: Distribution of immunophenotypic features before and after recategorisation. The original IUIS IEI descriptions contain information such as T cell-related “decreased cd8, normal or decreased cd4 cells” which we recategorise as “low”. The bar plot shows the count of unique labels for each status (normal, not_normal) across the T cell, B cell, Ig, and Neutrophil features.

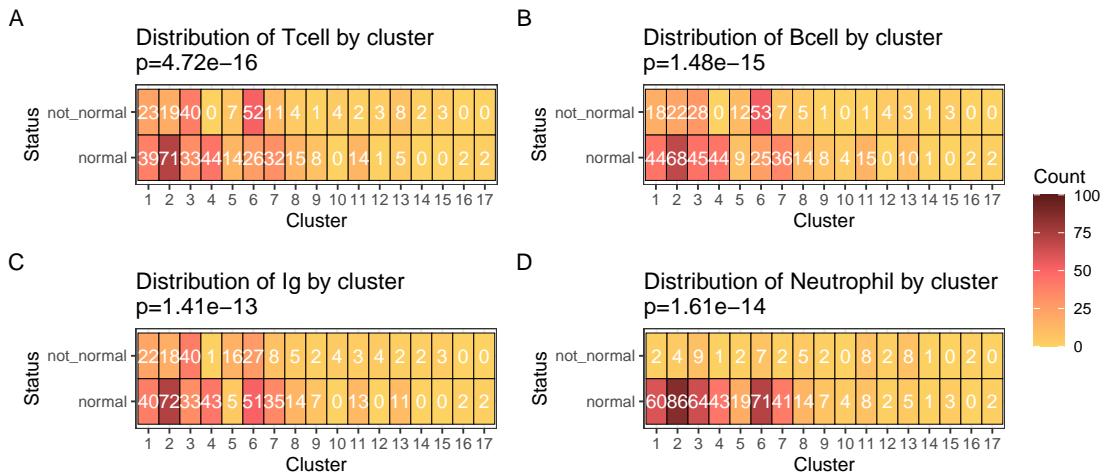


Figure 7: Heatmaps of clinical feature distributions by PPI cluster. The heatmaps display the count of observations for abnormality of each clinical feature (A) T cell, (B) B cell, (C) Immunoglobulin, (D) Neutrophil, in relation to the PPI clusters, with p-values from chi-square tests annotated in the titles.

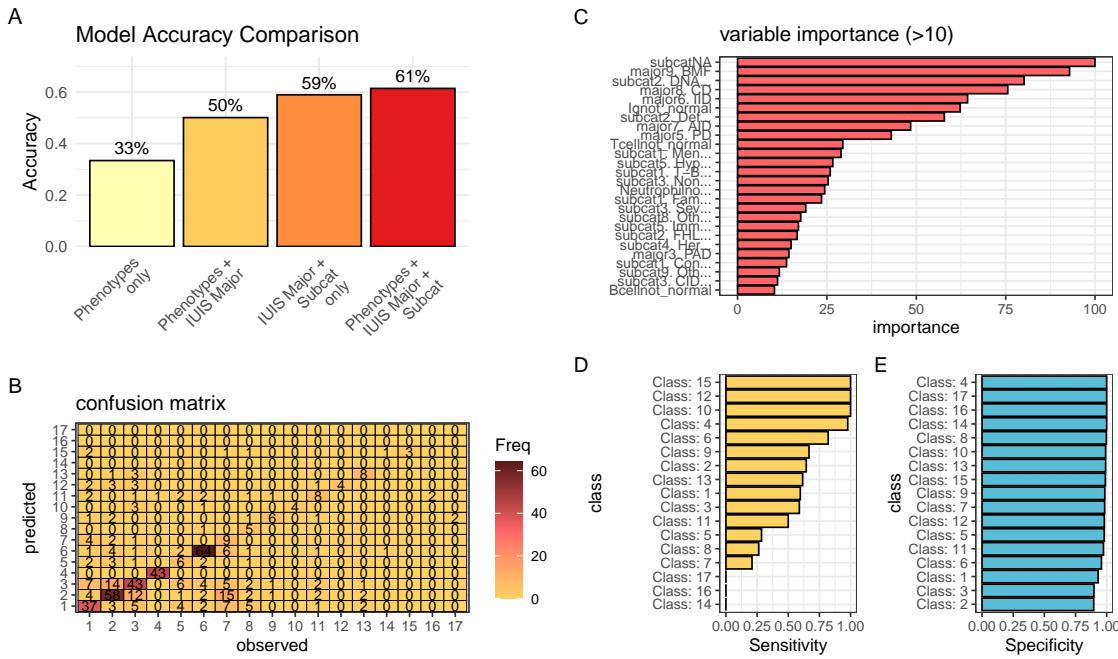


Figure 8: Performance comparison of PID classifiers. Classification predicting PPI cluster membership from IUIS major category, subcategory, and immunological features. (A) Overall accuracy for four rpart models used to predict PPI clustering. The combined model achieves 61.4 % accuracy, exceeding all simpler approaches. Nodes were split to minimize Gini impurity, pruned by cost-complexity ($cp = 0.001$), and validated via 5-fold cross-validation. (B-E) The summary statistics from the top model are detailed.

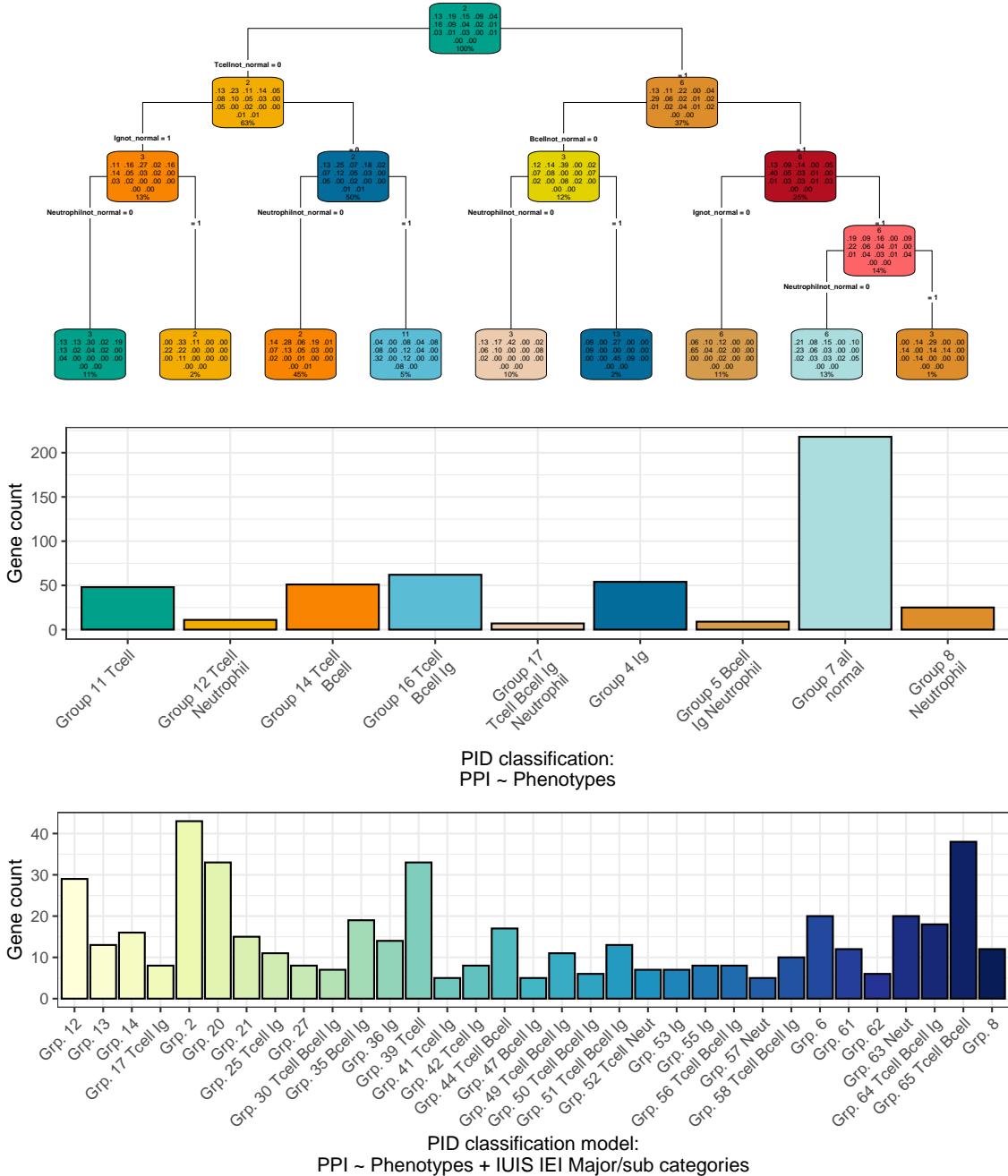


Figure 9: **Fine-tuned model for PID classification.** (Top) In each terminal node, the top block indicates the number of genes in the node; the middle block shows the fitted class probabilities (which sum to 1); and the bottom block displays the percentage of the total sample in that node. These metrics summarise the model’s assignment based on immunophenotypic and PPI features. (Middle) Bar plot presenting the distribution of novel PID classifications, where group labels denote the predominant abnormal clinical feature(s) (e.g. T cell, B cell, Ig, Neutrophil) characterising each group. (Bottom) The complete model including the traditional IUIS IEI categories.

570 **3.7.1 Integration of Variant Probabilities into IEI Genetics Data**

571 We integrated the computed prior probabilities for observing variants in all known
 572 genes associated with a given phenotype (1), across AD, AR, and XL MOI, into
 573 our IEI genetics framework. These calculations, derived from gene panels in Pan-
 574 elAppRex, have yielded novel insights for the IEI disease panel. The final result
 575 comprised of machine- and human-readable datasets, including the table of variant
 576 classifications and priors available via a the linked repository (27), and a user-friendly
 577 web interface that incorporates these new metrics.

578 **Figure 10** shows the interface summarising integrated variant data. Server-side
 579 pre-calculation of summary statistics minimises browser load, while clinical signifi-
 580 cance is converted to numerical metrics. Key quantiles (min, Q1, median, Q3, max)
 581 for each gene are rendered as sparkline box plots, and dynamic URLs link table entries
 582 to external databases (e.g. ClinVar, Online Mendelian Inheritance in Man (OMIM),
 583 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 pathogenicity, ClinVar SNV classification, ClinVar all variant reports, OMIM, Alpha Missense / Uniprot ID, HPO combined, and HPO term. The rows list various genetic conditions, such as SCID, CD3 deficiency, IL2RG deficiency, IL7Ra deficiency, ITPKB deficiency, JAK3 deficiency, and LAT deficiency, along with their respective details. Each row contains a sparkline box plot for the "Prior prob of pathogenicity" column. The "ClinVar all variant reports" column contains URLs. The "OMIM" column lists OMIM IDs. The "Alpha Missense / Uniprot ID" column lists UniProt IDs. The "HPO combined" and "HPO term" columns list Human Phenotype Ontology terms. The bottom of the table shows a page navigation bar with "Previous" and "Next" buttons, and a "Show 10" dropdown menu.

Figure 10: **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.

584 **3.8 Probability of observing AlphaMissense pathogenicity**

585 AlphaMissense provides pathogenicity scores for all possible amino acid substitutions;
 586 however, our results in **Figure 11** show that the most probable observations in pa-
 587 tients occur predominantly for benign or unknown variants. This finding places the
 588 likelihood of disease-associated substitutions into perspective and offers a data-driven
 589 foundation for future improvements in variant prediction. The values in **Figure 11**
 590 (**A**) can be directly compared to **Figure 1 (D)** to view the distribution of classifi-
 591 cations. A Kruskal-Wallis test was used to compare the observed disease probability

592 across clinical classification groups and no significant differences were detected. In
 593 general, most variants in patients are classified as benign or unknown, indicating
 594 limited discriminative power in the current classification, such that pathogenicity
 595 prediction does not infer observation prediction (**Figure S11**).

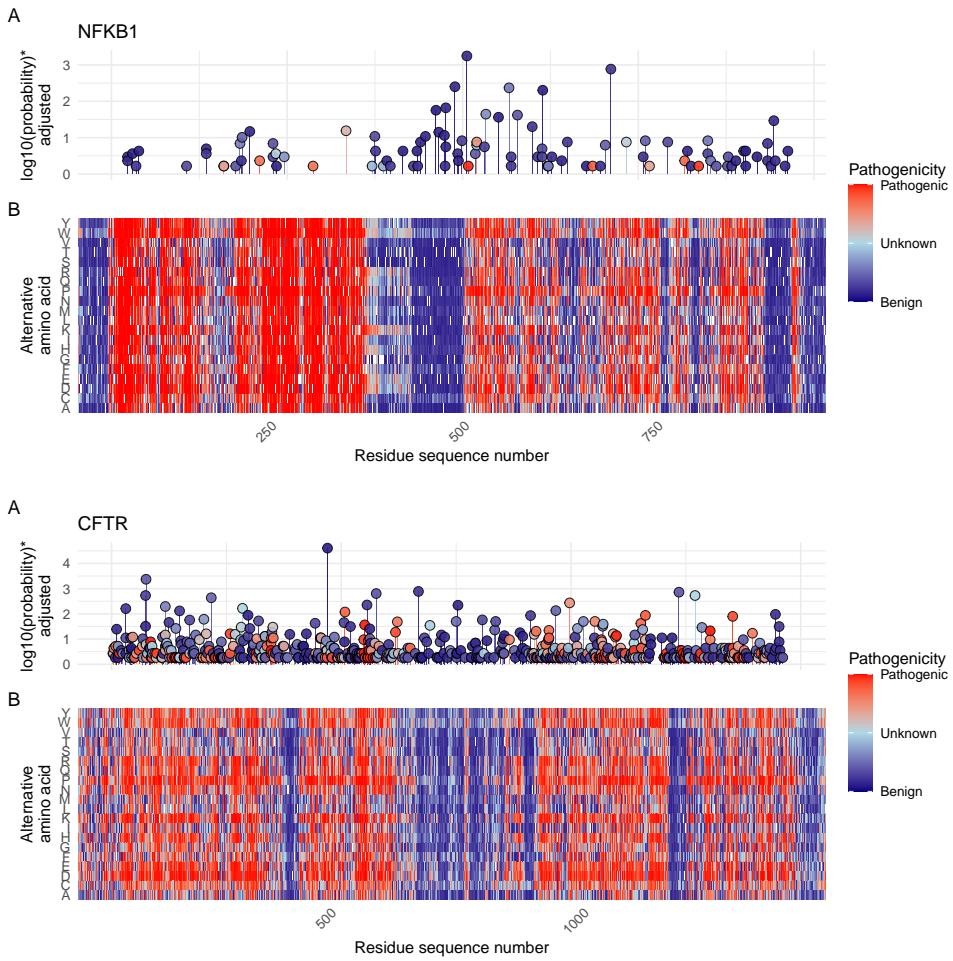


Figure 11: **(A) Probabilities of observing a patient with (B) AlphaMissense-derived pathogenicity scores.** Although AlphaMissense provides scores for every possible amino acid substitution, the most frequently observed variants in patients tend to be classified as benign or of unknown significance. This juxtaposition contextualises the likelihood of disease-associated substitutions and underlines prospects for refining predictive models. *Axis scaled for visibility near zero. Higher point indicates higher probability.

596 **4 Discussion**

597 Our study presents, to our knowledge, the first comprehensive framework for calculating
598 prior probabilities of observing disease-associated variants. By integrating large-
599 scale genomic annotations, including population allele frequencies from gnomAD (7),
600 variant classifications from ClinVar (13), and functional annotations from resources
601 such as dbNSFP, with classical Hardy-Weinberg-based calculations, we derived robust
602 estimates for 54,814 ClinVar variant classifications across 557 IEI genes implicated in
603 PID and monogenic inflammatory bowel disease (1; 2).

604 Our approach yielded two key results. First, our detailed, per-variant pre-calculated
605 results provide prior probabilities of observing disease-associated variants across all
606 MOI for any gene-disease combination. Second, the score positive total metric effec-
607 tively summarises the aggregate pathogenic burden across genes, serving as a robust
608 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:

$$\begin{aligned} \text{Recessive: } P(D) &= p^2, \\ \text{Dominant: } P(D) &= 2p(1 - p) \approx 2p. \end{aligned}$$

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:

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

609 Because both penetrance and $P(G | D)$ are often uncertain, solving this equation
610 systematically poses a major challenge.

611 Our framework addresses this challenge by combining variant classifications, pop-
612 ulation allele frequencies, and curated gene-disease associations. While imperfect on
613 an individual level, these sources exhibit predictable aggregate behaviour, supported
614 by James-Stein estimation principles (28). Curated gene-disease associations help
615 identify genes that explainable for most disease cases, allowing us to approximate
616 $P(G | D)$ close to one. In this way, we obtain robust estimates of $P(G)$ (the fre-
617 quency of disease-associated genotypes), even when exact values of penetrance and
618 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,

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

Even if the cumulative contribution is slightly less than one, the resultant risk estimates remain robust within the broad confidence intervals typical of epidemiological studies. By incorporating these pre-calculated priors into a Bayesian framework, our method refines risk estimates and enhances clinical decision-making despite inherent uncertainties.

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

Moreover, our method complements existing statistical approaches for aggregating variant effects with methods like Sequence Kernel Association Test (SKAT) and Aggregated Cauchy Association Test (ACAT) (29–32)) and multi-omics integration techniques (33; 34), while remaining consistent with established variant interpretation guidelines from the American College of Medical Genetics and Genomics (ACMG) (35) and complementary frameworks (36; 37), as well as quality control protocols (38; 39). Standardised reporting for qualifying variant sets, such as ACMG Secondary Findings v3.2 (40), further contextualises the integration of these probabilities into clinical decision-making.

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

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

655 5 Conclusion

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

659 Acknowledgements

660 We acknowledge Genomics England for providing public access to the PanelApp data.
661 The use of data from Genomics England panelapp was licensed under the Apache
662 License 2.0. The use of data from UniProt was licensed under Creative Commons
663 Attribution 4.0 International (CC BY 4.0). ClinVar asks its users who distribute or
664 copy data to provide attribution to them as a data source in publications and websites
665 (13). dbNSFP version 4.4a is licensed under the Creative Commons Attribution-
666 NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0); while we cite
667 this dataset as used our research publication, it is not used for the final version which
668 instead used ClinVar and gnomAD directly. GnomAD is licensed under Creative
669 Commons Zero Public Domain Dedication (CC0 1.0 Universal). GnomAD request
670 that usages cites the gnomAD flagship paper (7) and any online resources that include
671 the data set provide a link to the browser, and note that tool includes data from the
672 gnomAD v4.1 release. AlphaMissense asks to cite Cheng et al. (12) for usage in
673 research, with data available from Cheng et al. (26).

674 Competing interest

675 We declare no competing interest.

676 References

- 677 [1] Stuart G. Tangye, Waleed Al-Herz, Aziz Bousfiha, Charlotte Cunningham-
678 Rundles, Jose Luis Franco, Steven M. Holland, Christoph Klein, Tomohiro Morio,
679 Eric Oksenhendler, Capucine Picard, Anne Puel, Jennifer Puck, Mikko R. J.
680 Seppänen, Raz Somech, Helen C. Su, Kathleen E. Sullivan, Troy R. Torgerson,
681 and Isabelle Meyts. Human Inborn Errors of Immunity: 2022 Update
682 on the Classification from the International Union of Immunological Societies
683 Expert Committee. *Journal of Clinical Immunology*, 42(7):1473–1507, October
684 2022. ISSN 0271-9142, 1573-2592. doi: 10.1007/s10875-022-01289-3. URL
685 <https://link.springer.com/10.1007/s10875-022-01289-3>.
- 686 [2] Dylan Lawless. PanelAppRex aggregates disease gene panels and facilitates
687 sophisticated search. March 2025. doi: 10.1101/2025.03.20.25324319. URL
688 <http://medrxiv.org/lookup/doi/10.1101/2025.03.20.25324319>.

- 689 [3] Antonio Rueda Martin, Eleanor Williams, Rebecca E. Foulger, Sarah Leigh,
690 Louise C. Daugherty, Olivia Niblock, Ivone U. S. Leong, Katherine R. Smith,
691 Oleg Gerasimenko, Eik Haraldsdottir, Ellen Thomas, Richard H. Scott, Emma
692 Baple, Arianna Tucci, Helen Brittain, Anna De Burca, Kristina Ibañez, Dalia
693 Kasperaviciute, Damian Smedley, Mark Caulfield, Augusto Rendon, and Ellen M.
694 McDonagh. PanelApp crowdsources expert knowledge to establish consensus
695 diagnostic gene panels. *Nature Genetics*, 51(11):1560–1565, November 2019.
696 ISSN 1061-4036, 1546-1718. doi: 10.1038/s41588-019-0528-2. URL <https://www.nature.com/articles/s41588-019-0528-2>.
- 698 [4] Oliver Mayo. A Century of Hardy–Weinberg Equilibrium. *Twin Research
699 and Human Genetics*, 11(3):249–256, June 2008. ISSN 1832-4274, 1839-
700 2628. doi: 10.1375/twin.11.3.249. URL https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal_article.
- 702 [5] Nikita Abramovs, Andrew Brass, and May Tassabehji. Hardy-Weinberg Equi-
703 librium in the Large Scale Genomic Sequencing Era. *Frontiers in Genetics*,
704 11:210, March 2020. ISSN 1664-8021. doi: 10.3389/fgene.2020.00210. URL
705 <https://www.frontiersin.org/article/10.3389/fgene.2020.00210/full>.
- 706 [6] Johannes Zschocke, Peter H. Byers, and Andrew O. M. Wilkie. Mendelian
707 inheritance revisited: dominance and recessiveness in medical genetics. *Nature
708 Reviews Genetics*, 24(7):442–463, July 2023. ISSN 1471-0056, 1471-0064.
709 doi: 10.1038/s41576-023-00574-0. URL <https://www.nature.com/articles/s41576-023-00574-0>.
- 711 [7] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings,
712 Jessica Alfoldi, Qingbo Wang, Ryan L Collins, Kristen M Laricchia, Andrea
713 Ganna, Daniel P Birnbaum, et al. The mutational constraint spectrum quantified
714 from variation in 141,456 humans. *Nature*, 581(7809):434–443, 2020.
- 715 [8] Sarah L. Bick, Aparna Nathan, Hannah Park, Robert C. Green, Monica H. Wo-
716 jcik, and Nina B. Gold. Estimating the sensitivity of genomic newborn screen-
717 ing for treatable inherited metabolic disorders. *Genetics in Medicine*, 27(1):
718 101284, January 2025. ISSN 10983600. doi: 10.1016/j.gim.2024.101284. URL
719 <https://linkinghub.elsevier.com/retrieve/pii/S1098360024002181>.
- 720 [9] Benjamin D. Evans, Piotr Słowiński, Andrew T. Hattersley, Samuel E. Jones,
721 Seth Sharp, Robert A. Kimmitt, Michael N. Weedon, Richard A. Oram,
722 Krasimira Tsaneva-Atanasova, and Nicholas J. Thomas. Estimating disease
723 prevalence in large datasets using genetic risk scores. *Nature Communications*,
724 12(1):6441, November 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-26501-7.
725 URL <https://www.nature.com/articles/s41467-021-26501-7>.
- 726 [10] William B. Hannah, Mitchell L. Drumm, Keith Nykamp, Tiziano Prampano,
727 Robert D. Steiner, and Steven J. Schrodi. Using genomic databases to de-
728 termine the frequency and population-based heterogeneity of autosomal reces-
729 sive conditions. *Genetics in Medicine Open*, 2:101881, 2024. ISSN 29497744.

730 doi: 10.1016/j.gimo.2024.101881. URL <https://linkinghub.elsevier.com/retrieve/pii/S2949774424010276>.

- 731 [11] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov,
732 Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek,
733 Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J.
734 Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh
735 Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy,
736 Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer,
737 Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray
738 Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein
739 structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August
740 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL
741 <https://www.nature.com/articles/s41586-021-03819-2>.
- 742 [12] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Taylor
743 Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hassabis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Science*, 381(6664):eadg7492, September
744 2023. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.adg7492. URL
745 <https://www.science.org/doi/10.1126/science.adg7492>.
- 746 [13] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen Chao,
747 Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Wonhee Jang, Karen Karapetyan, Kenneth Katz, Chunlei Liu, Zenith Maddipatla, Adriana Malheiro, Kurt McDaniel, Michael Ovetsky, George Riley, George Zhou, J Bradley Holmes, Brandi L Kattman, and Donna R Maglott. ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Research*, 46(D1):D1062–D1067, January 2018. ISSN 0305-1048, 1362-4962. doi:
748 10.1093/nar/gkx1153. URL <http://academic.oup.com/nar/article/46/D1/D1062/4641904>.
- 749 [14] The UniProt Consortium, Alex Bateman, Maria-Jesus Martin, Sandra Orchard,
750 Michele Magrane, Aduragbemi Adesina, Shadab Ahmad, Bowler-Barnett, and Others. UniProt: the Universal Protein Knowledgebase in 2025. *Nucleic Acids Research*, 53(D1):D609–D617, January 2025. ISSN 0305-1048, 1362-4962. doi:
751 10.1093/nar/gkae1010. URL <https://academic.oup.com/nar/article/53/D1/D609/7902999>.
- 752 [15] Xiaoming Liu, Chang Li, Chengcheng Mou, Yibo Dong, and Yicheng Tu. dbNSFP v4: a comprehensive database of transcript-specific functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Genome Medicine*, 12(1):103, December 2020. ISSN 1756-994X. doi: 10.1186/s13073-020-00803-9. URL <https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-020-00803-9>.

- 771 [16] Damian Szklarczyk, Katerina Nastou, Mikaela Koutrouli, Rebecca Kirsch, Far-
772 rokh Mehryary, Radja Hachilif, Dewei Hu, Matteo E Peluso, Qingyao Huang,
773 Tao Fang, et al. The string database in 2025: protein networks with directional-
774 ity of regulation. *Nucleic Acids Research*, 53(D1):D730–D737, 2025.
- 775 [17] Paul Tuijnenburg, Hana Lango Allen, Siobhan O. Burns, Daniel Greene,
776 Machiel H. Jansen, and Others. Loss-of-function nuclear factor B subunit
777 1 (NFKB1) variants are the most common monogenic cause of common vari-
778 able immunodeficiency in Europeans. *Journal of Allergy and Clinical Im-*
779 *munology*, 142(4):1285–1296, October 2018. ISSN 00916749. doi: 10.1016/
780 j.jaci.2018.01.039. URL <https://linkinghub.elsevier.com/retrieve/pii/S0091674918302860>.
- 781 [18] WHO Scientific Group et al. Primary immunodeficiency diseases: report of a
782 who scientific group. *Clin. Exp. Immunol.*, 109(1):1–28, 1997.
- 783 [19] Charlotte Cunningham-Rundles and Carol Bodian. Common variable immunod-
784 eficiency: clinical and immunological features of 248 patients. *Clinical immunol-*
785 *ogy*, 92(1):34–48, 1999.
- 786 [20] Eric Oksenhendler, Laurence Gérard, Claire Fieschi, Marion Malphettes, Gael
787 Mouillot, Roland Jaussaud, Jean-François Viallard, Martine Gardembas, Lionel
788 Galicier, Nicolas Schleinitz, et al. Infections in 252 patients with common variable
789 immunodeficiency. *Clinical Infectious Diseases*, 46(10):1547–1554, 2008.
- 790 [21] Y Naito, F Adams, S Charman, J Duckers, G Davies, and S Clarke. Uk cystic
791 fibrosis registry 2023 annual data report. *London: Cystic Fibrosis Trust*, 2023.
- 792 [22] Carlo Castellani, CFTR2 team, et al. Cftr2: how will it help care? *Paediatric*
793 *respiratory reviews*, 14:2–5, 2013.
- 794 [23] Hartmut Grasemann and Felix Ratjen. Cystic fibrosis. *New England Journal*
795 *of Medicine*, 389(18):1693–1707, 2023. doi: 10.1056/NEJMra2216474. URL
796 <https://www.nejm.org/doi/full/10.1056/NEJMra2216474>.
- 797 [24] Kyoko Watanabe, Erdogan Taskesen, Arjen Van Bochoven, and Danielle
798 Posthuma. Functional mapping and annotation of genetic associations with
799 FUMA. *Nature Communications*, 8(1):1826, November 2017. ISSN 2041-1723.
800 doi: 10.1038/s41467-017-01261-5. URL <https://www.nature.com/articles/s41467-017-01261-5>.
- 801 [25] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir,
802 Pablo Tamayo, and Jill P. Mesirov. Molecular signatures database (MSigDB)
803 3.0. *Bioinformatics*, 27(12):1739–1740, June 2011. ISSN 1367-4811, 1367-
804 4803. doi: 10.1093/bioinformatics/btr260. URL <https://academic.oup.com/bioinformatics/article/27/12/1739/257711>.

- 808 [26] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-
809 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias
810 Sergeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-
811 abis, Pushmeet Kohli, and Žiga Avsec. Predictions for alphanonsense, September
812 2023. URL <https://doi.org/10.5281/zenodo.8208688>.
- 813 [27] Dylan Lawless. Variant risk estimate probabilities for iei genes. March 2025. doi:
814 10.5281/zenodo.15111584. URL <https://doi.org/10.5281/zenodo.15111584>.
- 815 [28] Bradley Efron and Carl Morris. Stein’s Estimation Rule and Its Competitors—
816 An Empirical Bayes Approach. *Journal of the American Statistical Association*,
817 68(341):117, March 1973. ISSN 01621459. doi: 10.2307/2284155. URL <https://www.jstor.org/stable/2284155?origin=crossref>.
- 818 [29] Yaowu Liu, Sixing Chen, Zilin Li, Alanna C Morrison, Eric Boerwinkle, and
819 Xihong Lin. Acat: a fast and powerful p value combination method for rare-
820 variant analysis in sequencing studies. *The American Journal of Human Genetics*,
821 104(3):410–421, 2019.
- 822 [30] Xihao Li, Zilin Li, Hufeng Zhou, Sheila M Gaynor, Yaowu Liu, Han Chen, Ryan
823 Sun, Rounak Dey, Donna K Arnett, Stella Aslibekyan, et al. Dynamic incorpora-
824 tion of multiple in silico functional annotations empowers rare variant association
825 analysis of large whole-genome sequencing studies at scale. *Nature genetics*, 52
826 (9):969–983, 2020.
- 827 [31] Michael C Wu, Seunggeun Lee, Tianxi Cai, Yun Li, Michael Boehnke, and Xi-
828 hong Lin. Rare-variant association testing for sequencing data with the sequence
829 kernel association test. *The American Journal of Human Genetics*, 89(1):82–93,
830 2011.
- 831 [32] Seunggeun Lee, Mary J Emond, Michael J Bamshad, Kathleen C Barnes, Mark J
832 Rieder, Deborah A Nickerson, David C Christiani, Mark M Wurfel, and Xihong
833 Lin. Optimal unified approach for rare-variant association testing with applica-
834 tion to small-sample case-control whole-exome sequencing studies. *The American
835 Journal of Human Genetics*, 91(2):224–237, 2012.
- 836 [33] Augustine Kong, Gudmar Thorleifsson, Michael L Frigge, Bjarni J Vilhjalmsson,
837 Alexander I Young, Thorgeir E Thorgeirsson, Stefania Benonisdottir, Asmundur
838 Oddsson, Bjarni V Halldorsson, Gisli Masson, et al. The nature of nurture:
839 Effects of parental genotypes. *Science*, 359(6374):424–428, 2018.
- 840 [34] Laurence J Howe, Michel G Nivard, Tim T Morris, Ailin F Hansen, Humaira
841 Rasheed, Yoonsu Cho, Geetha Chittoor, Penelope A Lind, Teemu Palviainen,
842 Matthijs D van der Zee, et al. Within-sibship gwas improve estimates of direct
843 genetic effects. *BioRxiv*, pages 2021–03, 2021.
- 844 [35] Sue Richards, Nazneen Aziz, Sherri Bale, David Bick, Soma Das, Julie Gastier-
845 Foster, Wayne W Grody, Madhuri Hegde, Elaine Lyon, Elaine Spector, et al.

- 847 Standards and guidelines for the interpretation of sequence variants: a joint
848 consensus recommendation of the american college of medical genetics and ge-
849 nomics and the association for molecular pathology. *Genetics in medicine*, 17
850 (5):405–423, 2015.
- 851 [36] Sean V Tavtigian, Steven M Harrison, Kenneth M Boucher, and Leslie G
852 Biesecker. Fitting a naturally scaled point system to the acmng/amp variant
853 classification guidelines. *Human mutation*, 41(10):1734–1737, 2020.
- 854 [37] Quan Li and Kai Wang. Intervar: clinical interpretation of genetic variants by
855 the 2015 acmng-amp guidelines. *The American Journal of Human Genetics*, 100
856 (2):267–280, 2017.
- 857 [38] Brent S Pedersen, Joe M Brown, Harriet Dashnow, Amelia D Wallace, Matt
858 Velinder, Martin Tristani-Firouzi, Joshua D Schiffman, Tatiana Tvrzik, Rong
859 Mao, D Hunter Best, et al. Effective variant filtering and expected candidate
860 variant yield in studies of rare human disease. *NPJ Genomic Medicine*, 6(1):1–8,
861 2021.
- 862 [39] Carl A Anderson, Fredrik H Pettersson, Geraldine M Clarke, Lon R Cardon,
863 Andrew P Morris, and Krina T Zondervan. Data quality control in genetic
864 case-control association studies. *Nature protocols*, 5(9):1564–1573, 2010. URL
865 <https://doi.org/10.1038/nprot.2010.116>.
- 866 [40] David T Miller, Kristy Lee, Noura S Abul-Husn, Laura M Amendola, Kyle Broth-
867 ers, Wendy K Chung, Michael H Gollob, Adam S Gordon, Steven M Harrison,
868 Ray E Hershberger, et al. Acmng sf v3. 2 list for reporting of secondary findings
869 in clinical exome and genome sequencing: a policy statement of the american
870 college of medical genetics and genomics (acmng). *Genetics in Medicine*, 25(8):
871 100866, 2023.

872 **6 Supplemental**

873 **6.1 Validation studies**

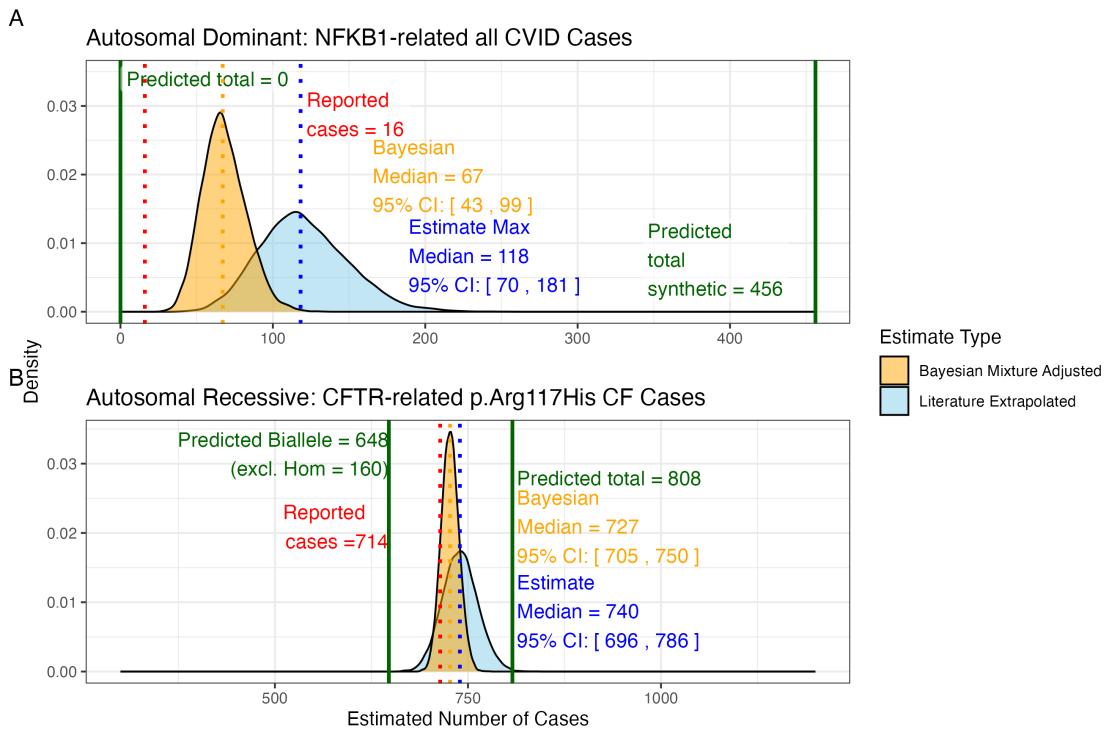


Figure S1: **Prior probabilities compared to validation disease cohort metrics.** (A) Density distributions for the number of *NFKB1*-related CVID cases in the UK. Our model (green) predicted 456 cases, which falls between the observed cohort count (red) of 390 and the upper extrapolated values. The 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.

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

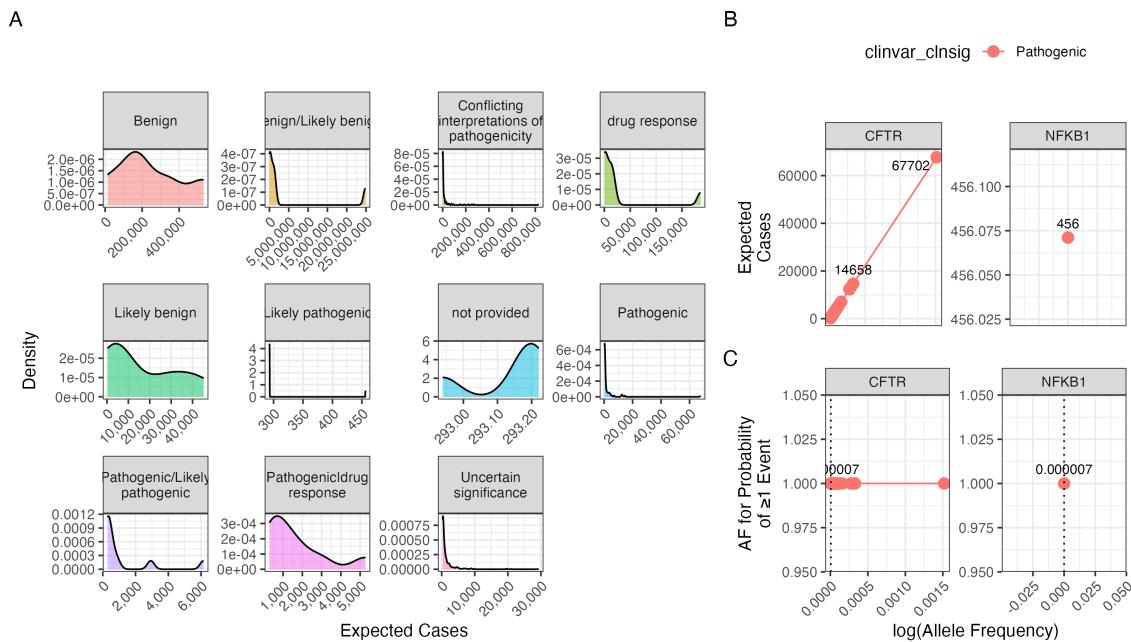


Figure S2: Interpretation of probability of observing a variant classification.
The result from the chosen validation genes *CFTR* and *NFKB1* are shown. Case counts are dependant on the population size and phenotype. (A) The density plots of expected observations by ClinVar clinical significance. We 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.

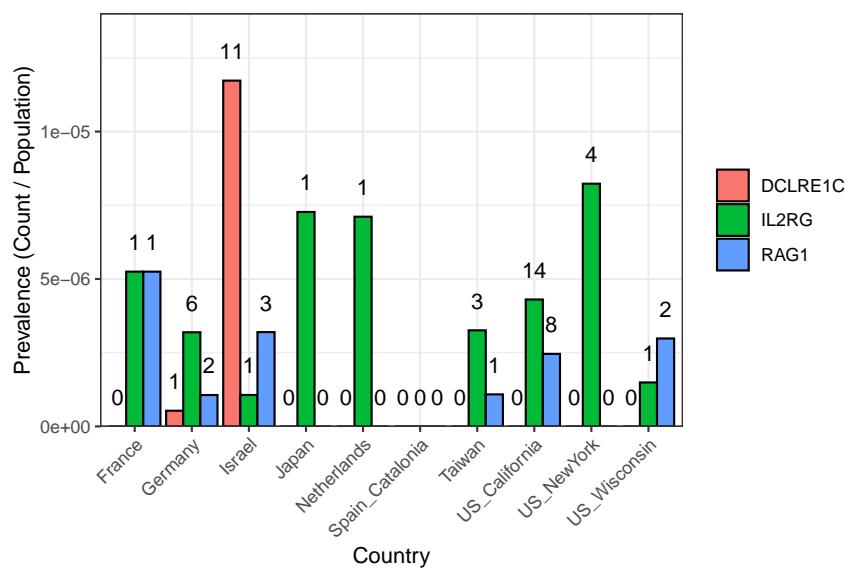


Figure S3: SCID-specific gene comparison across regions. The bar plot shows the prevalence of SCID-related cases (count divided by population) for each gene and country (or region), with numbers printed above the bars representing the actual counts in the original cohort (ranging from 0 to 11 per region and gene).

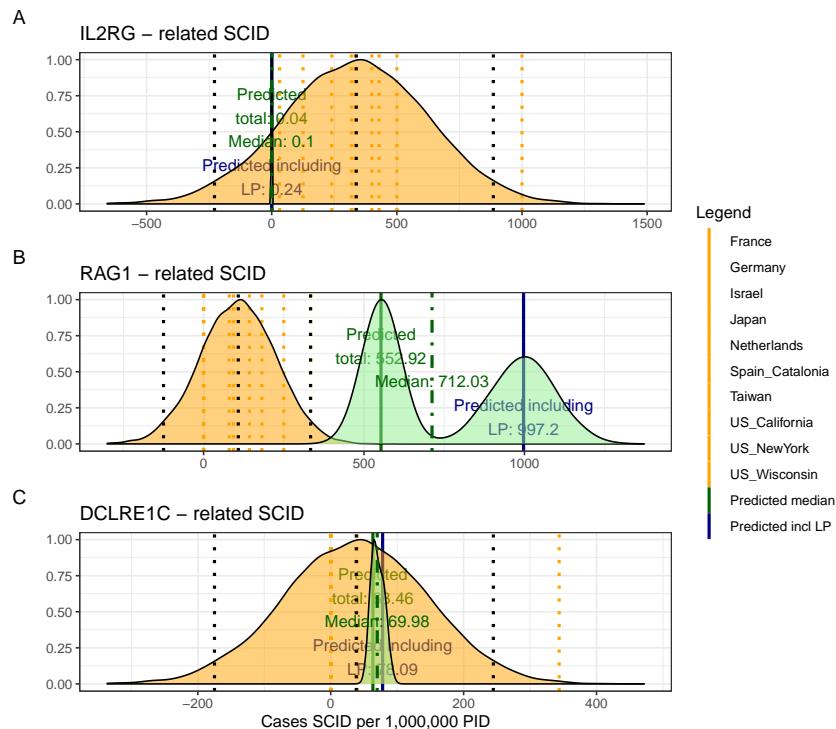


Figure S4: Combined SCID-specific Predictions and Observed Rates per 1,000,000 PID. The figure presents density distributions for the predicted SCID case counts (per 1,000,000 PID) for three genes: *IL2RG*, *RAG1*, and *DCLRE1C*. Country-specific rates (displayed as dotted vertical lines) are overlaid with the overall predicted distributions for pathogenic and likely pathogenic variants (solid lines with annotated medians). For *IL2RG*, the low predicted value is consistent with the high deleteriousness of loss-of-function variants in this X-linked gene, while *RAG1* exhibits considerably higher predicted counts, reflecting its lower penetrance in an autosomal recessive context.

874 **6.2 Hierarchical Clustering of Enrichment Scores for Major**
 875 **Disease Categories**

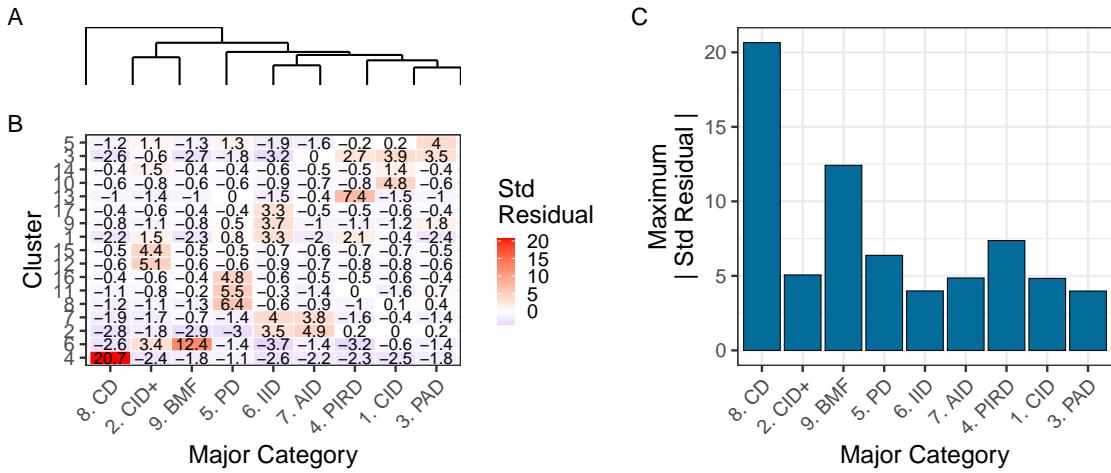


Figure S5: Hierarchical clustering of enrichment scores. The heatmap displays standardised residuals for major disease categories (x-axis) across network clusters (y-axis). A dendrogram groups similar disease categories, and the bar plot shows the maximum absolute residual per category. (8) CD and (9)BMF show the highest values, indicating significant enrichment or depletion ($\text{residuals} > |2|$). Definitions in **Box 2.1**.

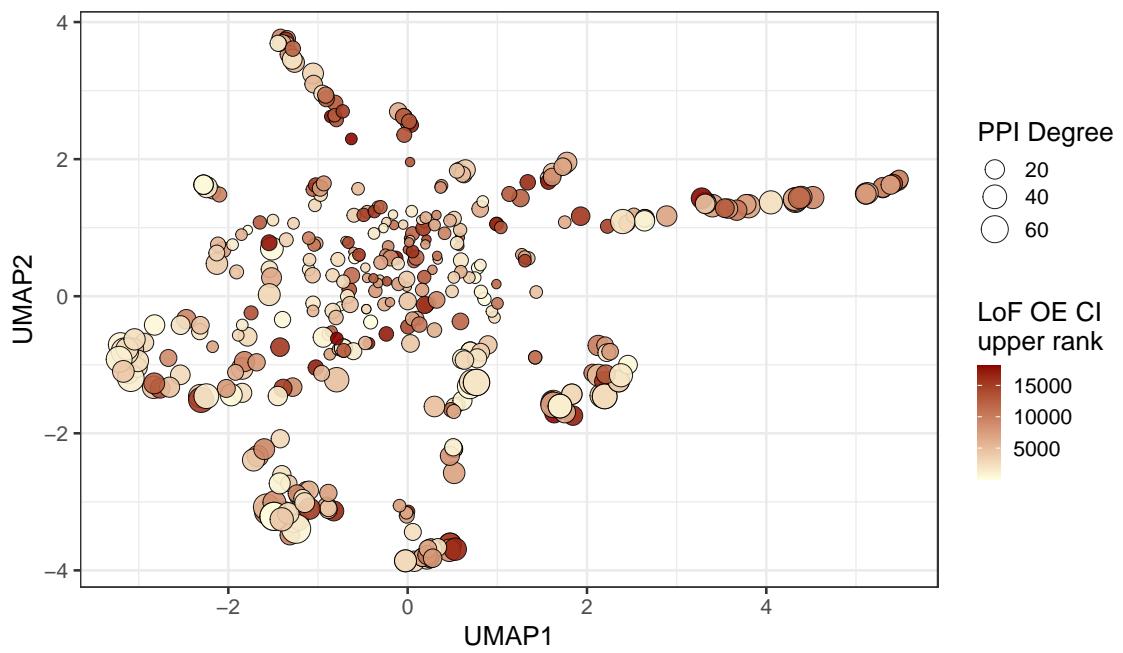


Figure S6: **Analysis of PPI degree versus LOEUF upper rank with UMAP embedding of the PPI network.** The relationship between PPI degree (size) and LOEUF upper rank (color) across gene clusters. No clear patterns are evident.

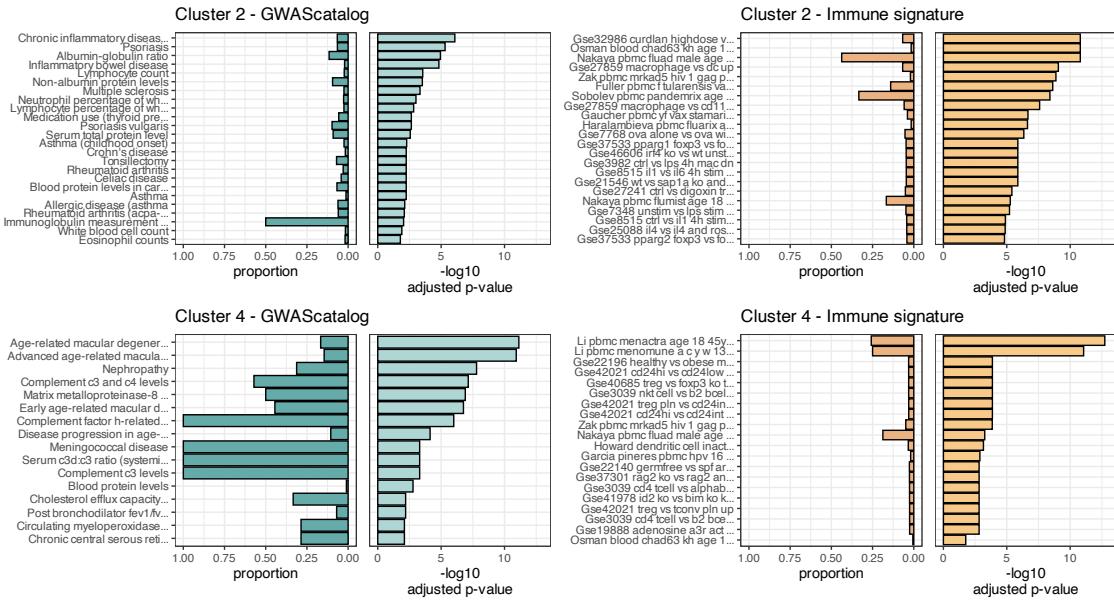


Figure S7: Composite Enrichment Profiles for IEI Gene Sets. We selected the top two enriched clusters (as per **Figure 4**) and performed functional enrichment analysis derived from known disease associations. For each gene set, the left panel displays the proportion of input genes overlapping with a 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.

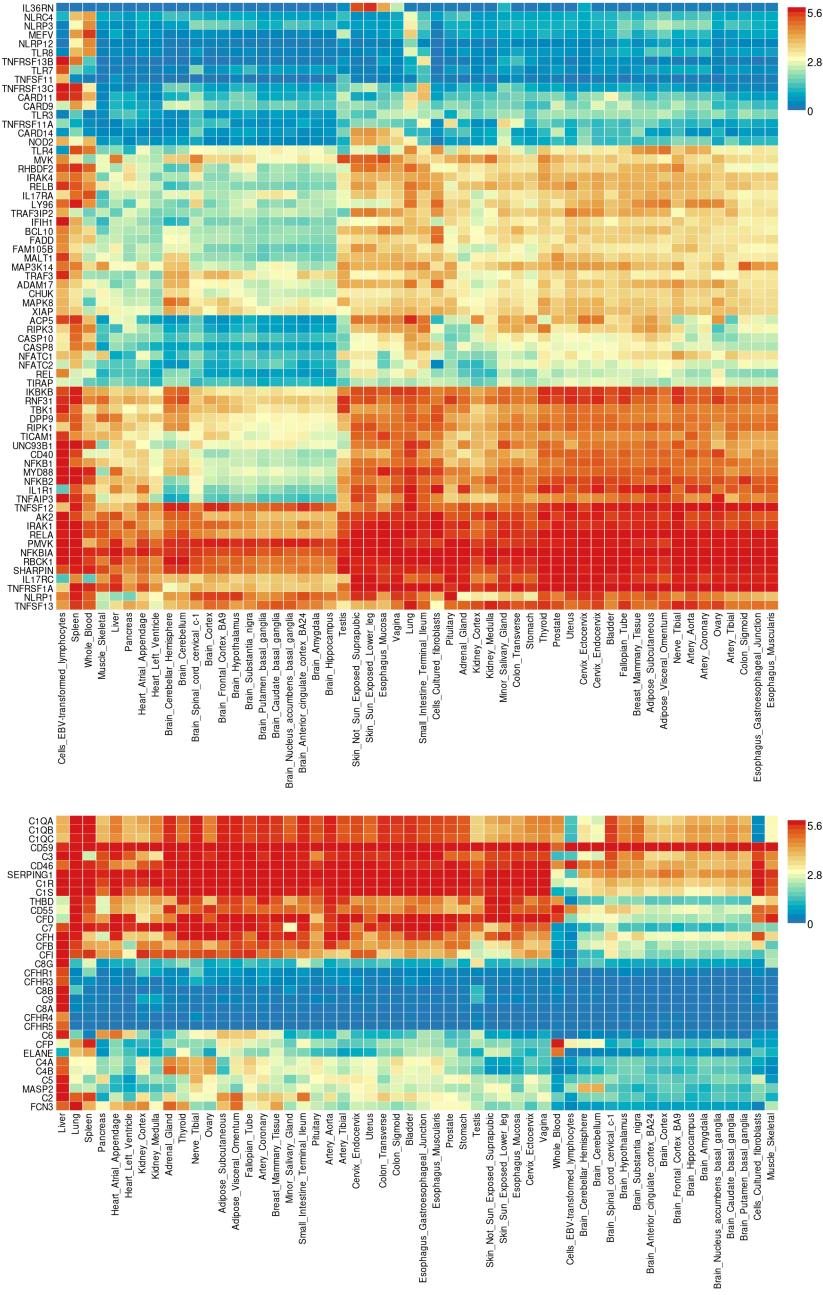


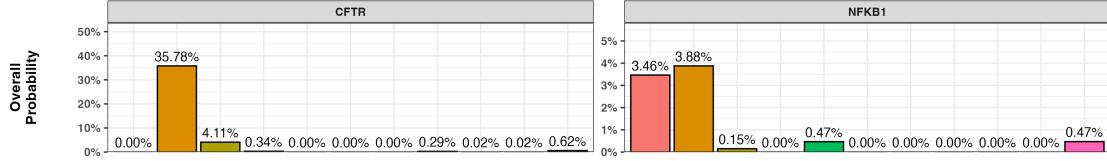
Figure S8: **Gene Expression Heatmaps for IEI Genes.** GTEx v8 data from 54 tissue types display the average expression per tissue label (log₂ transformed) for the IEI gene panels. Top: Cluster 2; Bottom: Cluster 4.

6.3 Interpretation of ClinVar Variant Observations

Recessive and Dominant Disease Genes

A

Overall Probability of an Affected Birth by ClinVar Category



B

Example, Total Expected Cases in UK population size (~ 69.4M)

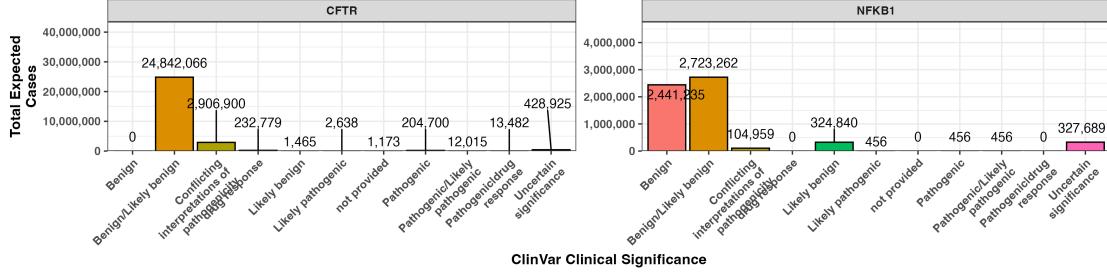


Figure S9: Combined bar charts summarizing the genome-wide analysis of ClinVar clinical significance for the PID gene panel. Panel (A) shows the overall probability of an affected birth by variant classification, and (B) displays the total expected number of cases per classification, both stratified by gene. These integrated results illustrate the variability in variant observations across genes and underpin our validation of the probability estimation framework.

6.4 Novel PID classifications

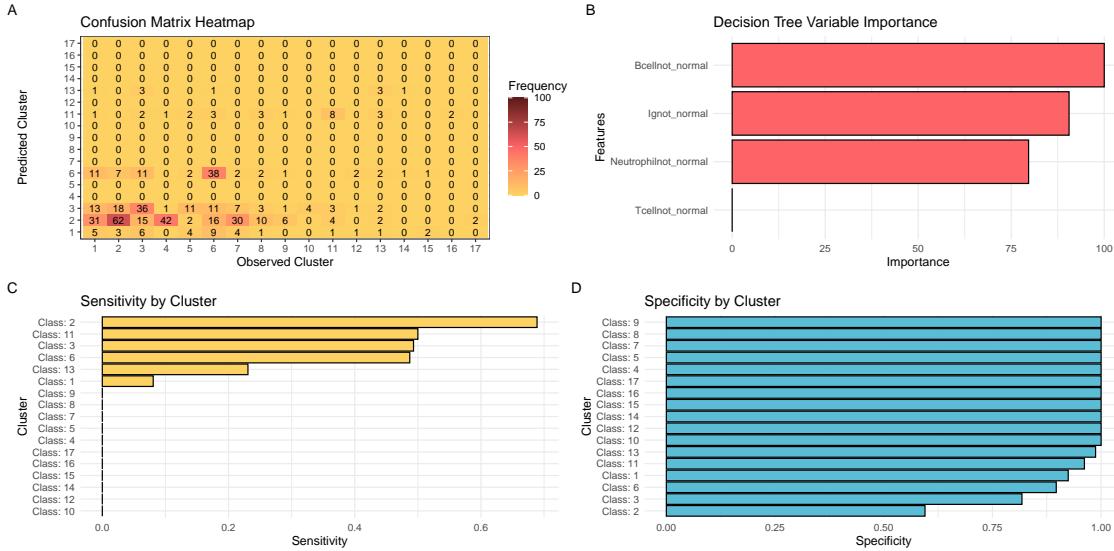


Figure S10: Model performance for fine-tuned PID classification. (A) Confusion matrix heatmap comparing observed and predicted PPI clusters. (B) Variable importance plot ranking immunophenotypic features contributing to the classifier. (C) Per-class sensitivity and (D) per-class specificity bar plots. These panels collectively demonstrate the performance of the decision tree classifier in stratifying PID genes based on immunophenotypic and PPI features.

6.5 Probability of observing AlphaMissense pathogenicity

878

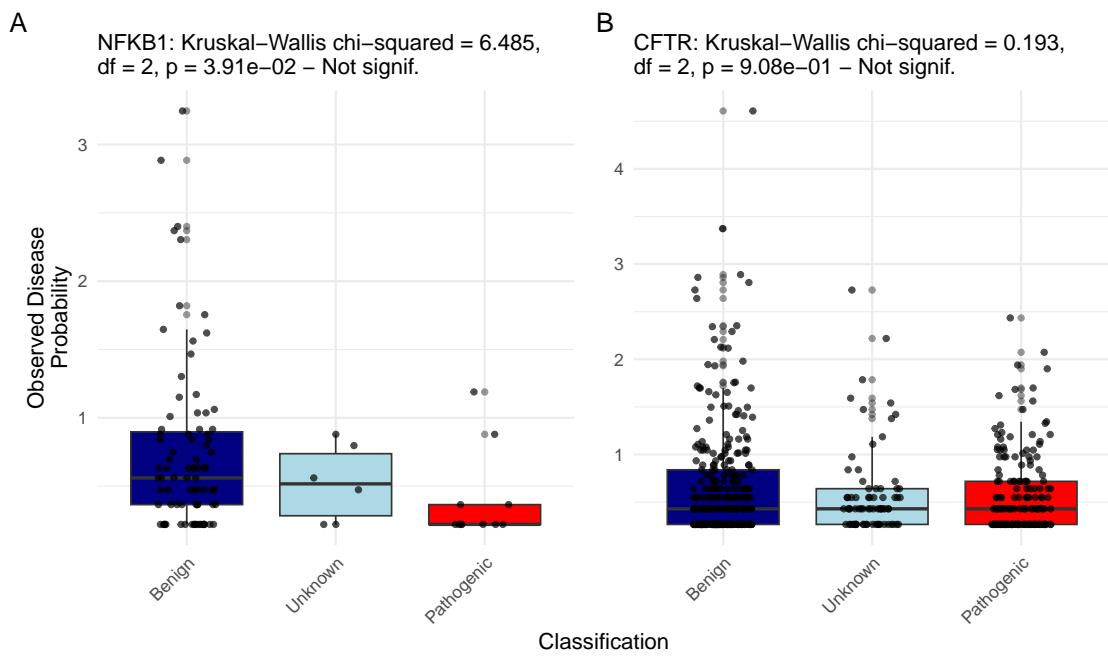


Figure S11: Observed Disease Probability by Clinical Classification with AlphaMissense. The figure displays the Kruskal-Wallis test results for NFKB1 and CFTR, showing no significant differences.

879 7 Clinical Genetics Application

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

890 7.1 Methods

891 7.1.1 Quality Control:

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

900 7.1.2 Prior Probability Calculation:

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

1. **Convert the ACMG Score:** The evidence score S is normalised to a frac-
tional support level:

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

905 This value reflects the strength of the pathogenic support.

2. **Assign a Minimal Risk (ϵ):** 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:

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

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

- 908 3. **Adjust the Allele Frequency:** The observed allele frequency p_i is then in-
909 creased by ϵ to yield an adjusted frequency:

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

908 This adjusted frequency reflects both the empirical observation and the ACMG
909 evidence.

- 910 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

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

911 **7.1.3 Deriving the Confidence Interval (CI)**

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

- 914 • The observed allele frequency and its adjustment via ϵ .
915 • The potential misclassification of variants (e.g. a VUS might be miscalled,
916 contributing to FP or FN counts).
917 • Missing sequence data at known pathogenic sites.

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

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

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

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

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

925 where $z = 1.96$ for a 95% confidence level. This interval integrates uncertainty from
 926 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 α and β 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:

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

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

929 Both methods integrate the uncertainty from the observed data, the adjustment
 930 via ϵ from the ACMG evidence score, and the potential misclassification or missing
 931 sequence data. In our analysis, the resulting 95% CI for pathogenicity is derived from
 932 such propagation of uncertainty, ensuring that all outcomes (TP, FP, TN, FN) are
 933 reflected in the final confidence bounds.

934 7.2 Results

935 We illustrate our method with two examples:

936 **Example 1: Missing Sequence Information** In one case, a known pathogenic
937 nucleotide position in *GENE_XYZ* exhibited low quality reads and did not pass QC.
938 This missing information prevents confirmation of the absence of the known variant (a
939 potential false negative), thereby widening the uncertainty in our probability estimate.
940 In such cases, the adjusted allele frequency is calculated with additional variance,
941 leading to a broader CI. For instance, if the observed allele frequency is 1.0×10^{-5}
942 and after adjusting with the ACMG score the estimated occurrence probability is
943 1.0×10^{-5} , the propagated uncertainty might yield a 95% CI of [0.70, 0.85]. This
944 broader interval reflects the impact of missing sequence data on our confidence.

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

956 Table S1 shows the final variant results for a male patient carrying an X-linked
957 loss-of-function (LOF) variant in *GENE_XYZ* where all known pathogenic positions
958 were confirmed as reference alleles. For the variant c. 1234del (p.Glu412Argfs*5), the
959 observed allele frequency is 1.2×10^{-5} . After applying the ACMG evidence score
960 adjustment (for a VUS lacking a ClinVar classification), the adjusted allele frequency
961 remains consistent with the observed data. The resulting occurrence probability is
962 1.2×10^{-5} , and by propagating the uncertainty from the allele frequency, evidence
963 score adjustment, and the full range of possible outcomes (TP, FP, TN, FN), we
964 derive a 95% CI for causality of [0.92, 0.97]. This confirms the variant as the top
965 causal variant in this patient, with no evidence of additional alternative pathogenic
966 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	1.2×10^{-5}
Occurrence Probability	1.2×10^{-5}
95% CI for Causality	[0.92, 0.97]
Clinical Interpretation	Top causal variant confirmed; no evidence of additional alternative pathogenic variants