

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.

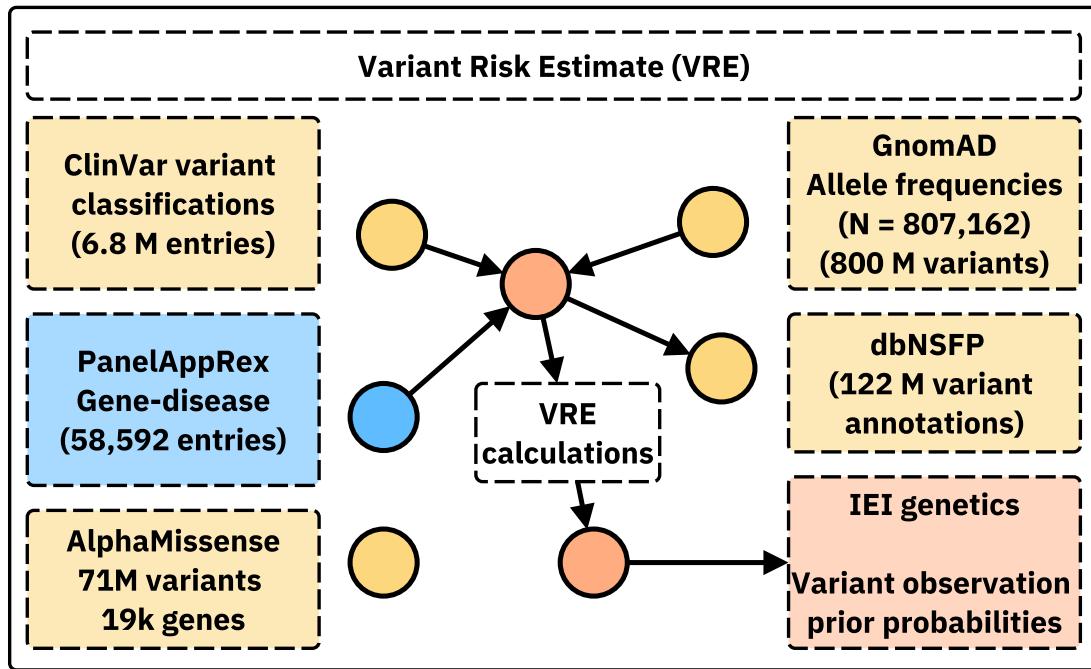
April 16, 2025

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.....	²⁸
²¹ ACAT Aggregated Cauchy Association Test	²⁸
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	²⁰
²⁷ CI Confidence Interval.....	¹⁶
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 InDel Insertion/Deletion	⁵
³⁹ IUIS International Union of Immunological Societies	⁵
⁴⁰ LD Linkage Disequilibrium	²²
⁴¹ LOEUF Loss-Of-function Observed/Expected Upper bound Fraction	¹²
⁴² LOF Loss-of-Function	¹⁹
⁴³ MOI Mode of Inheritance	⁴
⁴⁴ NFKB1 Nuclear Factor Kappa B Subunit 1	⁵
⁴⁵ OMIM Online Mendelian Inheritance in Man	²⁵
⁴⁶ PID Primary Immunodeficiency	⁴
⁴⁷ PPI Protein-Protein Interaction	⁵
⁴⁸ SNV Single Nucleotide Variant	⁴
⁴⁹ SKAT Sequence Kernel Association Test.....	²⁸
⁵⁰ STRINGdb Search Tool for the Retrieval of Interacting Genes/Proteins.....	⁵
⁵¹ HSD Honestly Significant Difference	¹²
⁵² UMAP Uniform Manifold Approximation and Projection	¹⁹
⁵³ UniProt Universal Protein Resource	⁵
⁵⁴ VEP Variant Effect Predictor.....	⁵
⁵⁵ XL X-Linked	⁴

92 1 Introduction

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

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

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

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

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

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

2 Methods

2.1 Dataset

Data from Genome Aggregation Database (gnomAD) v4 comprised 807,162 individuals, including 730,947 exomes and 76,215 genomes (7). This dataset provided 786,500,648 SNVs and 122,583,462 InDels, with variant type counts of 9,643,254 synonymous, 16,412,219 missense, 726,924 nonsense, 1,186,588 frameshift and 542,514 canonical splice site variants. ClinVar data were obtained from the variant summary dataset (as of: 16 March 2025) available from the NCBI FTP site, and included 6,845,091 entries, which were processed into 91,319 gene classification groups and a total of 38,983 gene classifications; for example, the gene *A1BG* contained four variants classified as likely benign and 102 total entries (13). For our analysis phase we also used dbNSFP which consisted of a number of annotations for 121,832,908 SNVs (15). The PanelAppRex core model contained 58,592 entries consisting of 52 sets of annotations, including the gene name, disease-gene panel ID, diseases-related features, confidence measurements. (2) A Protein-Protein Interaction (PPI) network data was provided by Search Tool for the Retrieval of Interacting Genes/Proteins (STRINGdb), consisting of 19,566 proteins and 505,968 interactions (16). The Human Genome Variation Society (HGVS) nomenclature is used with Variant Effect Predictor (VEP)-based codes for variant IDs. We carried out validations for disease cohorts with Nuclear Factor Kappa B Subunit 1 (*NFKB1*) (17–20) and Cystic Fibrosis Transmembrane Conductance Regulator (*CFTR*) (21–23) to demonstrate applications in AD and AR disease genes, respectively. AlphaMissense includes pathogenicity prediction classifications for 71 million variants in 19 thousand human genes (12; 26). We used these scores to compared against the probability of observing the same given variants. **Box 2.1** list the definitions from the International Union of Immunological

¹⁷⁴ 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	

¹⁷⁵

¹⁷⁶ 2.2 Variant Class Observation Probability

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

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

¹⁷⁷ where p is the allele frequency, $q = 1 - p$, p^2 represents the homozygous dominant,
¹⁷⁸ $2pq$ the heterozygous, and q^2 the homozygous recessive genotype frequencies. For dis-
¹⁷⁹ ease phenotypes, particularly under AR MOI, the risk is traditionally linked to the
¹⁸⁰ homozygous state (p^2); however, to account for compound heterozygosity across mul-
¹⁸¹ tiple variants, we extend this by incorporating the contribution from other pathogenic
¹⁸² alleles.

¹⁸³ Our computational pipeline estimated the probability of observing a disease-associated
¹⁸⁴ genotype for each variant and aggregated these probabilities by gene and ClinVar
¹⁸⁵ classification. This approach included all variant classifications, not limited solely to
¹⁸⁶ those deemed “pathogenic”, and explicitly conditioned the classification on the given
¹⁸⁷ phenotype, recognising that a variant could only be considered pathogenic relative to
¹⁸⁸ 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:

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

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

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

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

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

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

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

214 **2.3 Validation of Autosomal Dominant Estimates Using *NFKB1***

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

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

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

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

226 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

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

²²⁷ with a median value of approximately 118 and a 95% confidence interval of 70 to 181
²²⁸ cases (derived from posterior sampling).

²²⁹ **4. Bayesian Adjustment.** Recognising that the clinical cohort likely represents
²³⁰ nearly all CVID cases (besides first-second degree relatives), two Bayesian adjust-
²³¹ ments were performed:

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

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

²³² 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

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

²³³ which yielded a Bayesian mixture adjusted median estimate of 67 cases with a
²³⁴ 95% credible interval of approximately 43 to 99 cases.

²³⁵ **5. Predicted Total Genotype Counts.** The predicted total synthetic genotype
²³⁶ count (before adjustment) was 456, whereas the predicted total genotypes adjusted
²³⁷ for *synth_flag* was 0. This higher synthetic count was set based on a minimal risk
²³⁸ threshold, ensuring that at least one genotype is assumed to exist (e.g. accounting for
²³⁹ a potential unknown de novo variant) even when no variant is observed in gnomAD
²⁴⁰ (as per **section 2.2**).

²⁴¹ **6. Validation Test.** Thus, the expected number of *NFKB1*-related CVID cases
²⁴² derived from our genome-wide probability estimates was compared with the observed
²⁴³ counts from the UK-based PID cohort. This comparison validates our framework for
²⁴⁴ estimating disease incidence in AD disorders.

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

²⁴⁶ To validate our framework for AR diseases, we focused on Cystic Fibrosis (CF).
²⁴⁷ For comparability sizes between the validation studies, we analysed the most com-
²⁴⁸ mon SNV in the *CFTR* gene, typically reported as “p.Arg117His” (GRCh38 Chr
²⁴⁹ 7:117530975 G/A, MANE Select HGVS p.ENST00000003084.11: p.Arg117His). Our
²⁵⁰ goal was to validate our genome-wide probability estimates by comparing the ex-
²⁵¹ pected number of CF cases attributable to the p.Arg117His variant in *CFTR* with
²⁵² 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.

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

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

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

296 **2.6 Protein Network and Genetic Constraint Interpretation**

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

311 **2.7 Gene Set Enrichment Test**

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

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

319 **2.8 Deriving novel PID classifications by genetic PPI and**
320 **clinical features**

321 Immunophenotypic data were extracted from the IUIS IEI dataset, systematically
322 cleaned (see database link) and simplified by binarising clinical features (T cell, B
323 cell, Ig, Neutrophil) into normal and not normal categories. The PPI network em-
324 beddings from STRINGdb were reused for assigning each gene to its PPI cluster.
325 Associations between PPI clusters and the binarised clinical features were quantified
326 using chi-square tests. A decision tree model was subsequently constructed to predict
327 PPI clusters based on the clinical features. Hyperparameters (complexity parameter,
328 minimum split and bucket sizes, and maximum tree depth) were optimised by 5-fold
329 cross validation using `caret` with `rpart`. Terminal node assignments from the result-
330 ing tree defined novel PID classification groups, which were relabelled according to
331 the predominant abnormal feature within each group.

332 **2.9 Probability of observing AlphaMissense pathogenicity**

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

341 **3 Results**

342 **3.1 Observation Probability Across Disease Genes**

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

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

358 3.2 Validation studies

359 3.2.1 Validation of Dominant Disease Occurrence with *NFKB1*

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

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

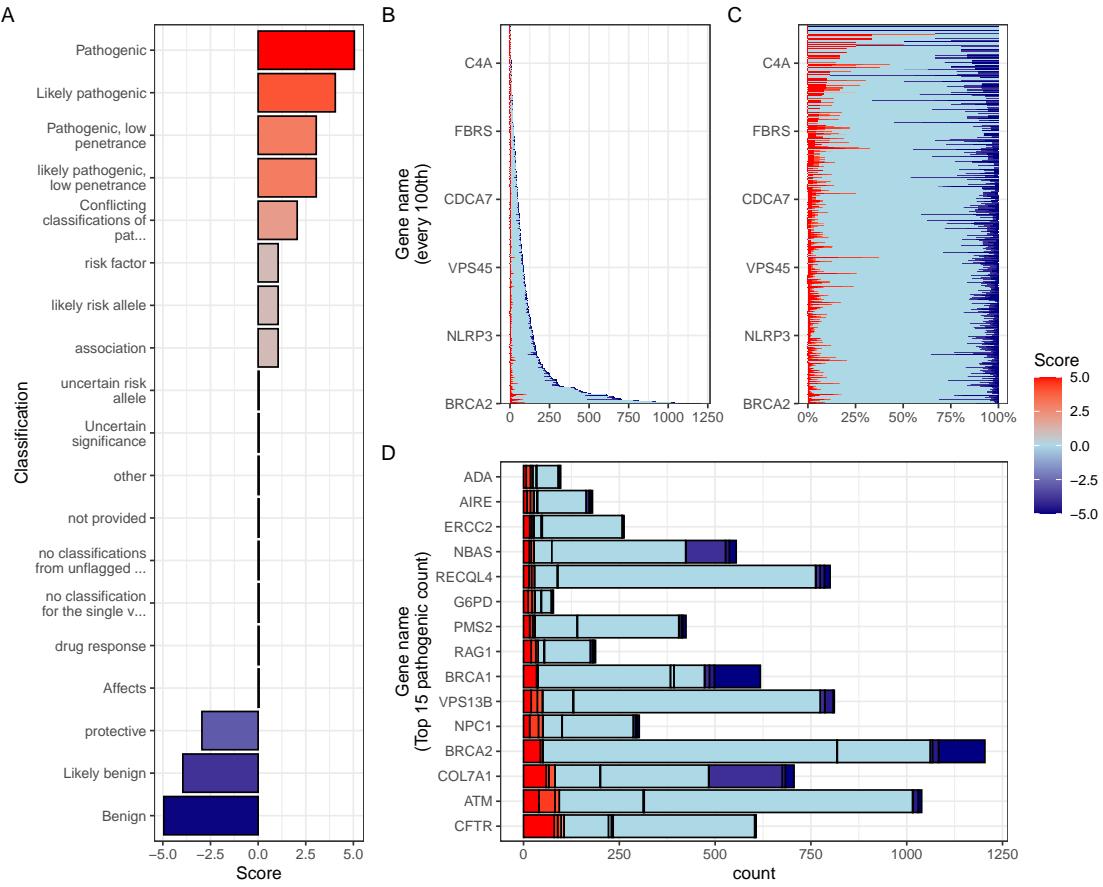


Figure 1: Summary of ClinVar clinical significance classifications in the PID gene panel. (A) Shows the numeric score coding for each classification. Panels (B) and (C) display the tally of classifications per gene as absolute counts and as percentages, respectively. (D) Highlights the top 15 genes with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

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

378 3.2.2 Validation of Recessive Disease Occurrence with *CFTR*

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

384 In contrast, the nationally reported number of CF cases was 714, as recorded in the
 385 UK Cystic Fibrosis Registry 2023 Annual Data Report (21). To account for factors

such as reduced penetrance and the mortality-adjusted expected genotype, we derived a Bayesian-adjusted estimate via posterior simulation. Our Bayesian approach yielded a median estimate of 740 cases (95% Confidence Interval (CI): 696, 786) and a mixture-based estimate of 727 cases (95% CI: 705, 750). **Figure S1 (B)** illustrates the close concordance between the predicted values, the Bayesian-adjusted estimates, and the national report supports the validity of our approach for estimating disease.

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

3.2.3 Interpretation of ClinVar Variant Observations

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

3.2.4 Validation of SCID-specific Disease Occurrence

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

424 **3.3 Genetic constraint in high-impact protein networks**

425 We next examined genetic constraint in high-impact protein networks across the whole
426 IEI gene set of over 500 known disease-gene phenotypes (1). By integrating ClinVar
427 variant classification scores with PPI data, we quantified the pathogenic burden per
428 gene and assessed its relationship with network connectivity and genetic constraint
429 (7; 16).

430 **3.3.1 Score-Positive-Total within IEI PPI network**

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

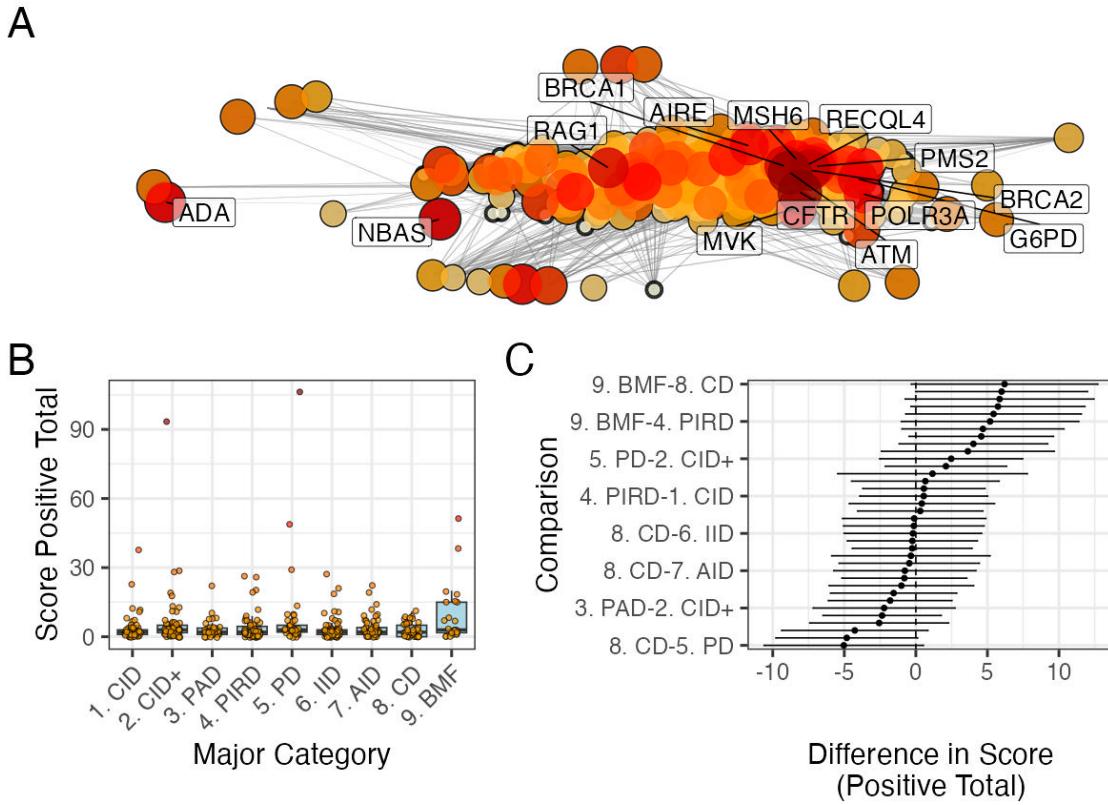


Figure 2: PPI network and score positive total ClinVar significance variants. (A) PPI network of disease-associated genes. Node size and colour represent the log-transformed score positive total, the top 15 genes/proteins with the highest probability of being observed in disease are labelled. (B) Distribution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

438 3.3.2 Association Analysis of Score-Positive-Total across IEI Categories

439 We checked for any statistical enrichment in score positive totals, which represents
 440 the expected observation of pathogenicity, between the IEI categories. The one-way
 441 ANOVA revealed an effect of major disease category on score positive total ($F(8, 500) =$
 442 2.82, $p = 0.0046$), indicating that group means were not identical, which we observed
 443 in **Figure 2 (B)**. However, despite some apparent differences in median scores across
 444 categories (i.e. 9. Bone Marrow Failure (BMF)), the Tukey HSD post hoc compar-
 445 isons **Figure 2 (C)** showed that all pairwise differences had 95% confidence intervals
 446 overlapping zero, suggesting that individual group differences were not significant.

447 3.3.3 UMAP Embedding of the PPI Network

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

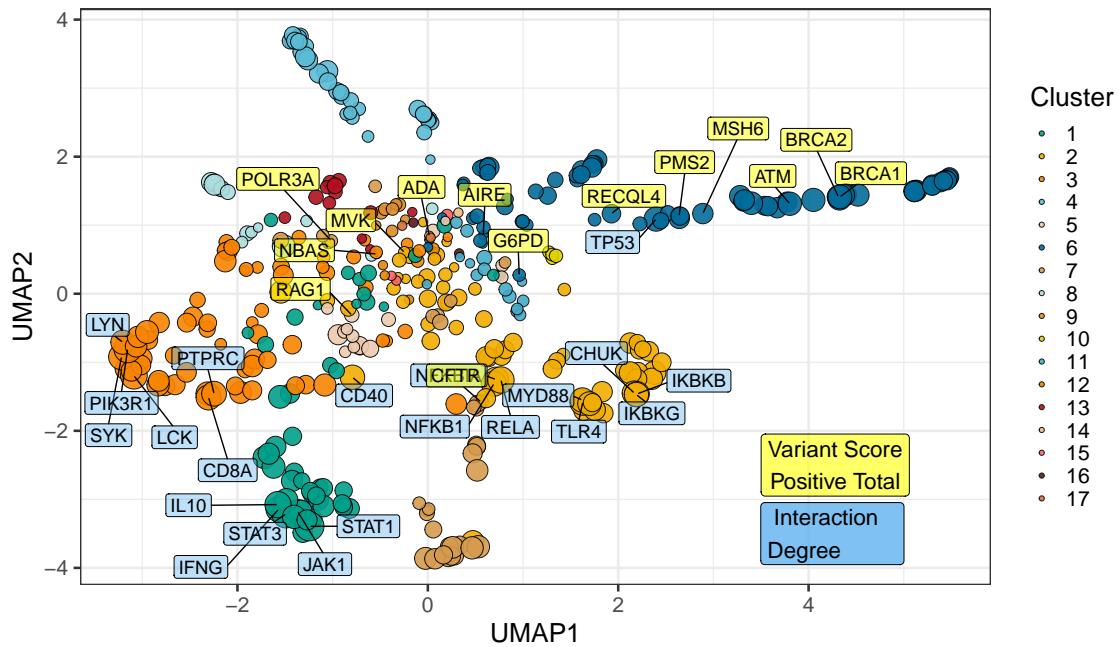


Figure 3: **UMAP embedding of the PPI network (p_umap)**. The plot projects the high-dimensional protein-protein interaction network into two dimensions, with nodes coloured by cluster and sized by interaction degree. Blue labels indicate hub genes (degree above the 95th percentile) and yellow labels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

458 3.3.4 Hierarchical Clustering of Enrichment Scores for Major Disease Cat- 459 egories

460 **Figure S5** presents a heatmap of standardised residuals for major disease categories
 461 across network clusters, as per **Figure 3**. A dendrogram clusters similar disease cate-

462 gories, while the accompanying bar plot displays the maximum absolute standardised
463 residual for each category. Notably, (8) Complement Deficiencies (CD) shows the
464 highest maximum enrichment, followed by (9) BMF. While all maximum values
465 exceed 2, the threshold for significance, this likely reflects the presence of protein
466 clusters with strong damaging variant scores rather than uniform significance across
467 all categories (i.e. genes from cluster 4 in 8 CD).

468 **3.3.5 PPI Connectivity, LOEUF Constraint and Enriched Network Clus-
469 ter Analysis**

470 Based on the preliminary insight from **Figure S5**, we evaluated the relationship
471 between network connectivity (PPI degree) and LOEUF constraint (LOEUF upper rank)
472 Karczewski et al. (7) using Spearman's rank correlation. Overall, there was a weak
473 but significant negative correlation ($\rho = -0.181, p = 0.00024$) at the global scale,
474 indicating that highly connected genes tend to be more constrained. A supplementary
475 analysis (**Figure S6**) did not reveal distinct visual associations between network
476 clusters and constraint metrics, likely due to the high network density. However
477 once stratified by gene clusters, the natural biological scenario based on quantitative
478 PPI evidence (16), some groups showed strong correlations; for instance, cluster 2
479 ($\rho = -0.375, p = 0.000994$) and cluster 4 ($\rho = -0.800, p < 0.000001$), while others did
480 not. This indicated that shared mechanisms within pathway clusters may underpin
481 genetic constraints, particularly for LOF intolerance. We observe that the score
482 positive total metric effectively summarises the aggregate pathogenic burden across
483 IEI genes, serving as a robust indicator of genetic constraint and highlighting those
484 with elevated disease relevance.

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

489 **3.4 New Insight from Functional Enrichment**

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

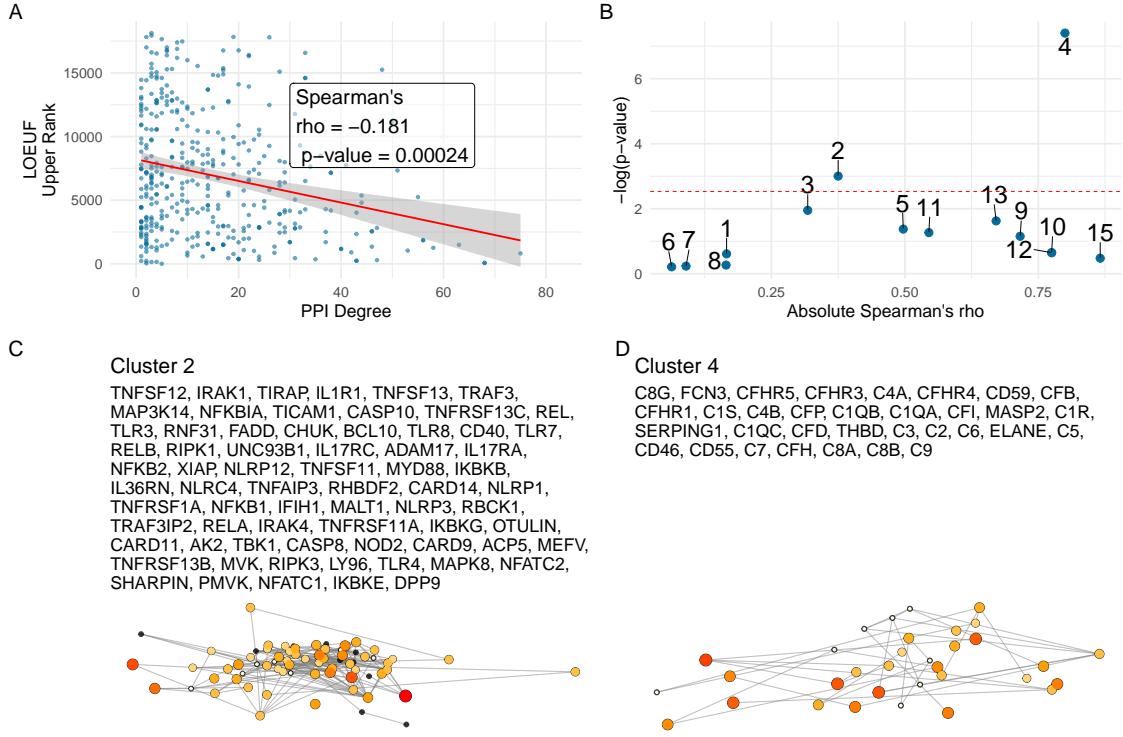


Figure 4: **Correlation between PPI degree and LOEUF upper rank.** (A) Ananlysis across all genes revealed a weak, significant negative correlation between PPI degree and LOEUF upper rank. (B) The cluster-wise analysis showed that clusters 2 and 4 exhibited moderate to strong correlations, while other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

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

3.5 Genome-wide Gene Distribution and Locus-specific Variant Occurrence

Figure 5 (A) shows a genome-wide karyoplot of all IEI panel genes across GRCh38, with colour-coding based on MOI. Figures (B) and (C) display zoomed-in locus plots for *NFKB1* and *CFTR*, respectively. In **Figure 5 (B)**, the probability of observing

variants with known classifications is high only for variants such as p.Ala475Gly, which are considered benign in the AD *NFKB1* gene that is intolerant to LOF. In **Figure 5 (C)**, high probabilities of observing patients with pathogenic variants in *CFTR* are evident, reproducing this well-established phenomenon. Furthermore, the analysis of Linkage Disequilibrium (LD) using R^2 shows that high LD regions can be modelled effectively, allowing independent variant signals to be distinguished.

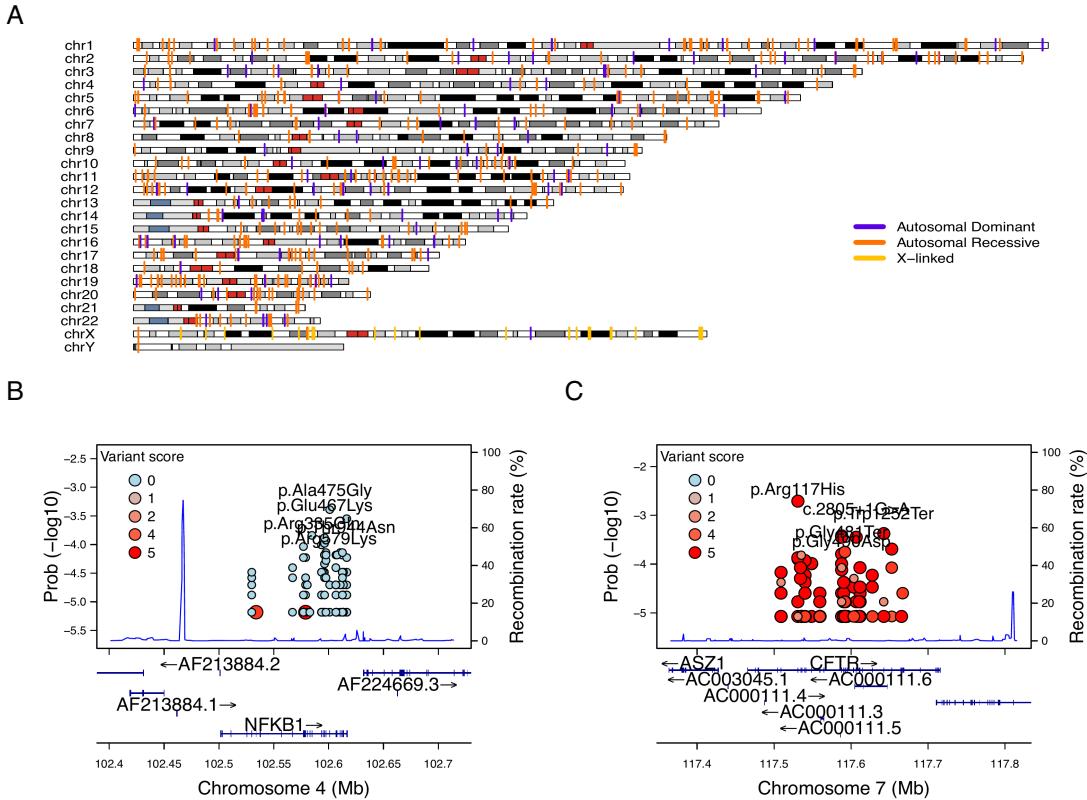


Figure 5: Genome-wide IEI, variant occurrence probability and LD by R^2 .
(A) Genome-wide karyoplot of all IEI panel genes mapped to GRCh38, with colours indicating MOI. (B) Zoomed-in locus plot for *NFKB1* showing variant observation probabilities; only benign variants such exhibit high probabilities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

520 3.6 Novel PID classifications derived from genetic PPI and 521 clinical features

We recategorised 315 immunophenotypic features from the original IUIS IEI annotations, reducing detailed descriptions (e.g. “decreased cd8, normal or decreased cd4”),

524 first to minimal labels (e.g.“low”), and second to binary outcomes (normal vs. not-
 525 normal) for T cells, B cells, neutrophils, and immunoglobulins (**Figure 6**). These
 526 simplified profiles were integrated with PPI network clustering from STRINGdb to
 527 refine PID gene groupings. Chi-square analyses confirmed significant associations be-
 528 tween specific clinical abnormalities and PPI clusters (**Figure S10**). A decision tree
 529 classifier, with hyperparameters optimised via 5-fold cross validation, demonstrated
 530 high sensitivity and specificity, as shown in the confusion matrices and variable impor-
 531 tance metrics (**Figure S11**). The resulting novel PID classifications, illustrated by
 532 the decision tree and gene group distributions (**Figure 7**), provide a more coherent
 533 and data-driven framework for categorising PID genes.

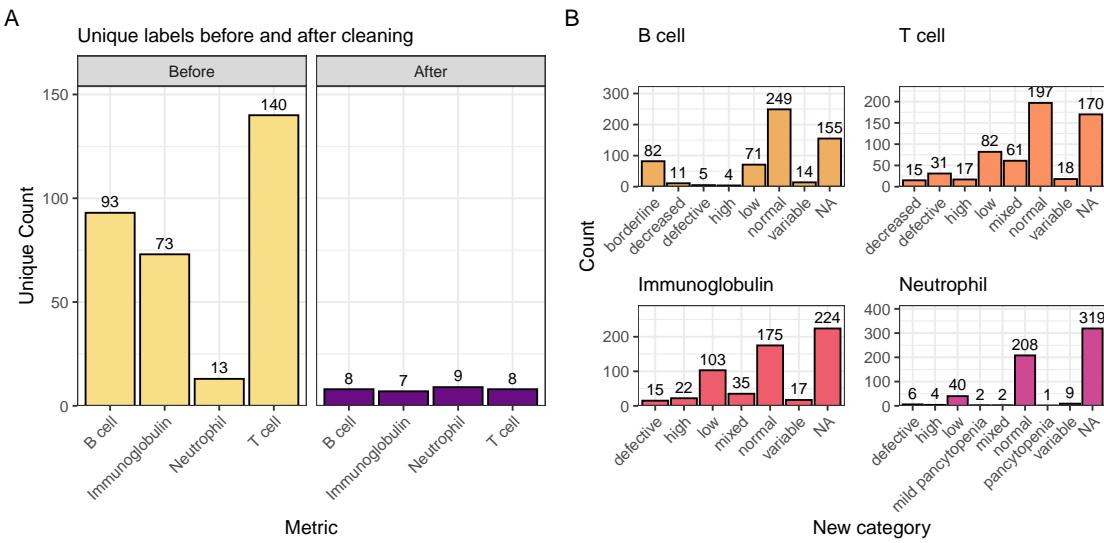


Figure 6: Distribution of immunophenotypic features before and after recategorisation. The original IUIS IEI descriptions contain information such as T cell-related “decreased cd8, normal or decreased cd4 cells” which we recategorise as “low”. The bar plot shows the count of unique labels for each status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

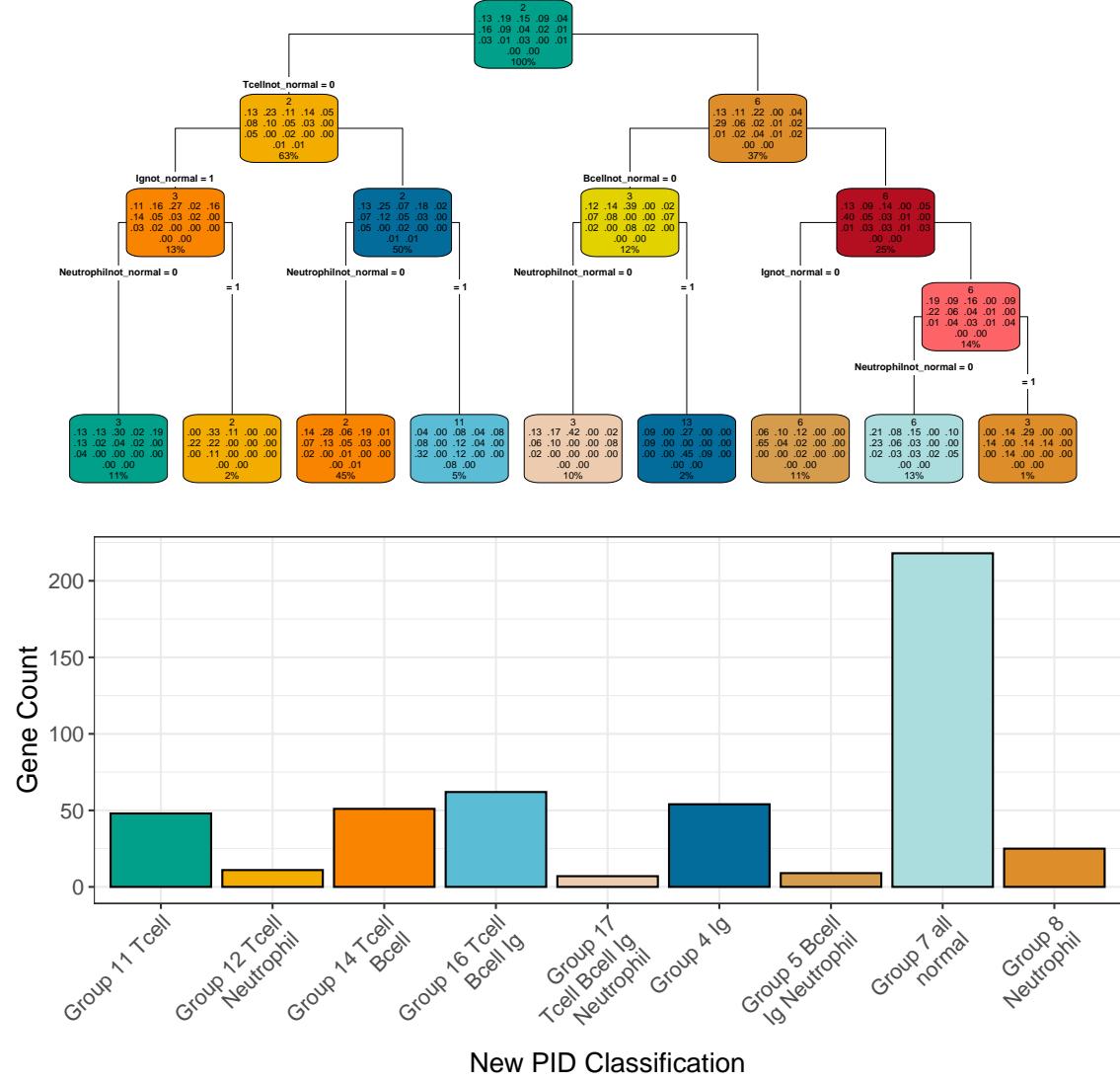


Figure 7: **Fine-tuned model for PID classification.** (Top) In each terminal node, the top block indicates the number of genes in the node; the middle block shows the fitted class probabilities (which sum to 1); and the bottom block displays the percentage of the total sample in that node. These metrics summarise the model’s assignment based on immunophenotypic and PPI features. (Bottom) 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.

534 **3.6.1 Integration of Variant Probabilities into IEI Genetics Data**

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

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

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	HPO term
All				All								
1. CID	1. T-B+ SCID	CD3z deficiency	CD247	AR	2	1/0/33/1	15/0/139/218	186780	P20963	HP:0002715; HP:0005403	Abnormalit	
1. CID	1. T-B+ SCID	CD3d deficiency	CD3D	AR	3	1/0/34/0	20/1/101/162/234	186790	P04234	HP:0002715; HP:0005403	Abnormalit	
1. CID	1. T-B+ SCID	CD3e deficiency	CD3E	AR	3	1/2/29/2	26/1/14/173/346	186830	P07766	HP:0002715; HP:0005403	Abnormalit	
1. CID	1. T-B+ SCID	Coronin-1A deficiency	CORO1A	AR	2	1/1/43/2	19/1/4/226/378	605000	P31146	HP:0002715; HP:0005403	Abnormalit	
1. CID	1. T-B+ SCID	γc-deficiency (common gamma chain SCID, CD132 deficiency)	IL2RG	XL	3	1/1/16/28	184/104/244/414	308380	P31185	HP:0002715; HP:0005403	Abnormalit	
1. CID	1. T-B+ SCID	IL7Ra deficiency	IL7R	AR	12	1/2/81/14	81/2/26/438/598	146661	P16871	HP:0002715; HP:0005403	Abnormalit	
1. CID	1. T-B+ SCID	ITPKB deficiency	ITPKB	AR	3	1/2/15/9	0/2/1/30/40	query	P27987	HP:0002715; HP:0005403	Abnormalit	
1. CID	1. T-B+ SCID	JAK3 deficiency	JAK3	AR	12	1/5/131/13	182/19/197/1528	600173	P52333	HP:0002715; HP:0005403	Abnormalit	
1. CID	1. T-B+ SCID	LAT deficiency	LAT	AR	1	1/0/38/3	14/1/2/138/242	602354	O43561	HP:0002715; HP:0005403	Abnormalit	

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

548 **3.7 Probability of observing AlphaMissense pathogenicity**

549 AlphaMissense provides pathogenicity scores for all possible amino acid substitutions;
 550 however, our results in **Figure 9** show that the most probable observations in pa-
 551 tients occur predominantly for benign or unknown variants. This finding places the
 552 likelihood of disease-associated substitutions into perspective and offers a data-driven
 553 foundation for future improvements in variant prediction. The values in **Figure 9 (A)**
 554 can be directly compared to **Figure 1 (D)** to view the distribution of classifications.
 555 A Kruskal-Wallis test was used to compare the observed disease probability across

556 clinical classification groups and no significant differences were detected. In general,
 557 most variants in patients are classified as benign or unknown, indicating limited dis-
 558 criminative power in the current classification, such that pathogenicity prediction
 559 does not infer observation prediction ([Figure S12](#)).

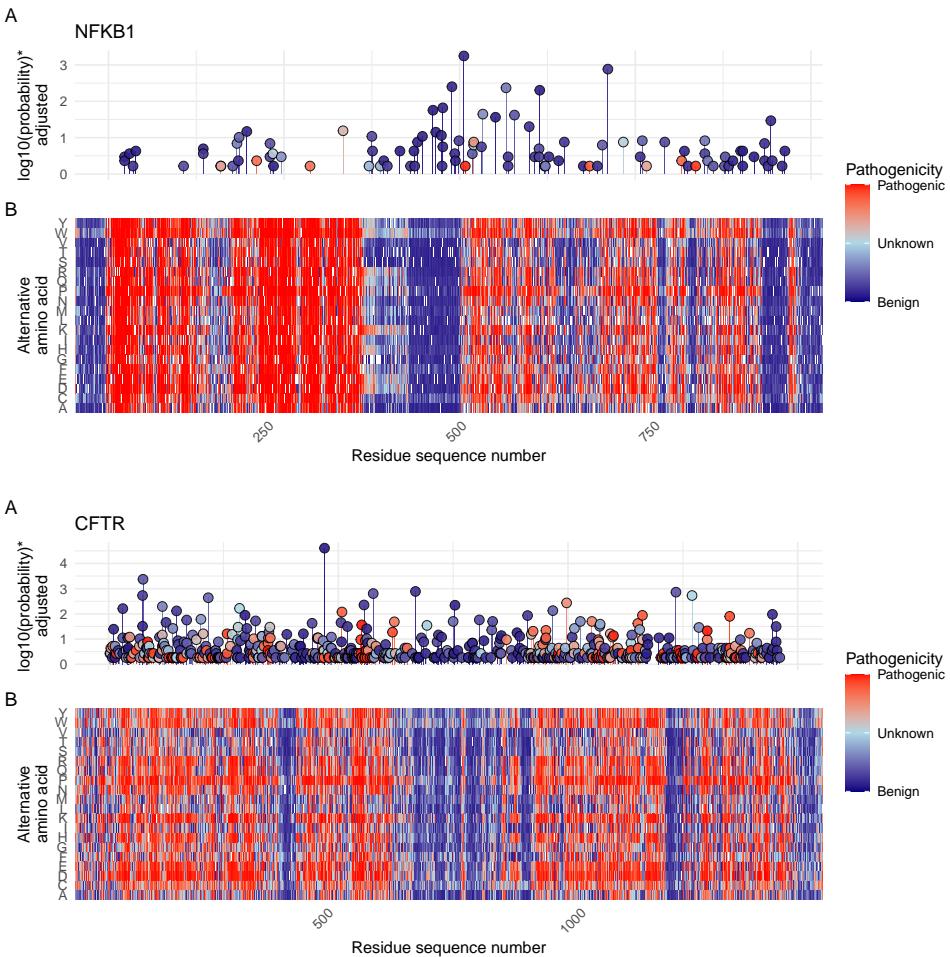


Figure 9: **(A) Probabilities of observing a patient with (B) AlphaMissense-derived pathogenicity scores.** Although AlphaMissense provides scores for every possible amino acid substitution, the most frequently observed variants in patients tend to be classified as benign or of unknown significance. 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.

560 **4 Discussion**

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

568 Our approach yielded two key results. First, our detailed, per-variant pre-calculated
569 results provide prior probabilities of observing disease-associated variants across all
570 MOI for any gene-disease combination. Second, the score positive total metric effec-
571 tively summarises the aggregate pathogenic burden across genes, serving as a robust
572 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:

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

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

575 Our framework addresses this challenge by combining variant classifications, pop-
576 ulation allele frequencies, and curated gene-disease associations. While imperfect on
577 an individual level, these sources exhibit predictable aggregate behaviour, supported
578 by James-Stein estimation principles (28). Curated gene-disease associations help
579 identify genes that explainable for most disease cases, allowing us to approximate
580 $P(G | D)$ close to one. In this way, we obtain robust estimates of $P(G)$ (the fre-
581 quency of disease-associated genotypes), even when exact values of penetrance and
582 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.

619 5 Conclusion

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

623 Acknowledgements

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

638 Competing interest

639 We declare no competing interest.

640 References

- 641 [1] Stuart G. Tangye, Waleed Al-Herz, Aziz Bousfiha, Charlotte Cunningham-
642 Rundles, Jose Luis Franco, Steven M. Holland, Christoph Klein, Tomohiro Morio,
643 Eric Oksenhendler, Capucine Picard, Anne Puel, Jennifer Puck, Mikko R. J.
644 Seppänen, Raz Somech, Helen C. Su, Kathleen E. Sullivan, Troy R. Torgerson,
645 and Isabelle Meyts. Human Inborn Errors of Immunity: 2022 Update
646 on the Classification from the International Union of Immunological Societies
647 Expert Committee. *Journal of Clinical Immunology*, 42(7):1473–1507, October
648 2022. ISSN 0271-9142, 1573-2592. doi: 10.1007/s10875-022-01289-3. URL
649 <https://link.springer.com/10.1007/s10875-022-01289-3>.
- 650 [2] Dylan Lawless. PanelAppRex aggregates disease gene panels and facilitates
651 sophisticated search. March 2025. doi: 10.1101/2025.03.20.25324319. URL
652 <http://medrxiv.org/lookup/doi/10.1101/2025.03.20.25324319>.

- 653 [3] Antonio Rueda Martin, Eleanor Williams, Rebecca E. Foulger, Sarah Leigh,
654 Louise C. Daugherty, Olivia Niblock, Ivone U. S. Leong, Katherine R. Smith,
655 Oleg Gerasimenko, Eik Haraldsdottir, Ellen Thomas, Richard H. Scott, Emma
656 Baple, Arianna Tucci, Helen Brittain, Anna De Burca, Kristina Ibañez, Dalia
657 Kasperaviciute, Damian Smedley, Mark Caulfield, Augusto Rendon, and Ellen M.
658 McDonagh. PanelApp crowdsources expert knowledge to establish consensus
659 diagnostic gene panels. *Nature Genetics*, 51(11):1560–1565, November 2019.
660 ISSN 1061-4036, 1546-1718. doi: 10.1038/s41588-019-0528-2. URL <https://www.nature.com/articles/s41588-019-0528-2>.
- 662 [4] Oliver Mayo. A Century of Hardy–Weinberg Equilibrium. *Twin Research
and Human Genetics*, 11(3):249–256, June 2008. ISSN 1832-4274, 1839-
663 2628. doi: 10.1375/twin.11.3.249. URL https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal_article.
- 664 [5] Nikita Abramovs, Andrew Brass, and May Tassabehji. Hardy-Weinberg Equi-
665 librium in the Large Scale Genomic Sequencing Era. *Frontiers in Genetics*,
666 11:210, March 2020. ISSN 1664-8021. doi: 10.3389/fgene.2020.00210. URL
667 <https://www.frontiersin.org/article/10.3389/fgene.2020.00210/full>.
- 668 [6] Johannes Zschocke, Peter H. Byers, and Andrew O. M. Wilkie. Mendelian
669 inheritance revisited: dominance and recessiveness in medical genetics. *Nature
Reviews Genetics*, 24(7):442–463, July 2023. ISSN 1471-0056, 1471-0064.
doi: 10.1038/s41576-023-00574-0. URL <https://www.nature.com/articles/s41576-023-00574-0>.
- 670 [7] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings,
671 Jessica Alfoldi, Qingbo Wang, Ryan L Collins, Kristen M Laricchia, Andrea
672 Ganna, Daniel P Birnbaum, et al. The mutational constraint spectrum quantified
673 from variation in 141,456 humans. *Nature*, 581(7809):434–443, 2020.
- 674 [8] Sarah L. Bick, Aparna Nathan, Hannah Park, Robert C. Green, Monica H. Wo-
675 jcik, and Nina B. Gold. Estimating the sensitivity of genomic newborn screen-
676 ing for treatable inherited metabolic disorders. *Genetics in Medicine*, 27(1):
677 101284, January 2025. ISSN 10983600. doi: 10.1016/j.gim.2024.101284. URL
678 <https://linkinghub.elsevier.com/retrieve/pii/S1098360024002181>.
- 679 [9] Benjamin D. Evans, Piotr Słowiński, Andrew T. Hattersley, Samuel E. Jones,
680 Seth Sharp, Robert A. Kimmitt, Michael N. Weedon, Richard A. Oram,
681 Krasimira Tsaneva-Atanasova, and Nicholas J. Thomas. Estimating disease
682 prevalence in large datasets using genetic risk scores. *Nature Communications*,
683 12(1):6441, November 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-26501-7.
684 URL <https://www.nature.com/articles/s41467-021-26501-7>.
- 685 [10] William B. Hannah, Mitchell L. Drumm, Keith Nykamp, Tiziano Prampano,
686 Robert D. Steiner, and Steven J. Schrodi. Using genomic databases to de-
687 termine the frequency and population-based heterogeneity of autosomal reces-
688 sive conditions. *Genetics in Medicine Open*, 2:101881, 2024. ISSN 29497744.

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

- 695 [11] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov,
696 Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek,
697 Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J.
698 Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh
700 Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy,
701 Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer,
702 Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray
703 Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein
704 structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August
705 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL
706 <https://www.nature.com/articles/s41586-021-03819-2>.
- 707 [12] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Taylor
708 Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias
709 Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hassabis,
710 Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense variant
711 effect prediction with AlphaMissense. *Science*, 381(6664):eadg7492, September
712 2023. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.adg7492. URL
713 <https://www.science.org/doi/10.1126/science.adg7492>.
- 714 [13] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen Chao,
715 Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Wonhee
716 Jang, Karen Karapetyan, Kenneth Katz, Chunlei Liu, Zenith Maddipatla, Adriana
717 Malheiro, Kurt McDaniel, Michael Ovetsky, George Riley, George Zhou,
718 J Bradley Holmes, Brandi L Kattman, and Donna R Maglott. ClinVar: improving
719 access to variant interpretations and supporting evidence. *Nucleic Acids Research*,
720 46(D1):D1062–D1067, January 2018. ISSN 0305-1048, 1362-4962. doi:
721 10.1093/nar/gkx1153. URL <http://academic.oup.com/nar/article/46/D1/D1062/4641904>.
- 722 [14] The UniProt Consortium, Alex Bateman, Maria-Jesus Martin, Sandra Orchard,
723 Michele Magrane, Aduragbemi Adesina, Shadab Ahmad, Bowler-Barnett, and
724 Others. UniProt: the Universal Protein Knowledgebase in 2025. *Nucleic Acids Research*,
725 53(D1):D609–D617, January 2025. ISSN 0305-1048, 1362-4962.
726 doi: 10.1093/nar/gkae1010. URL <https://academic.oup.com/nar/article/53/D1/D609/7902999>.
- 727 [15] Xiaoming Liu, Chang Li, Chengcheng Mou, Yibo Dong, and Yicheng Tu.
728 dbNSFP v4: a comprehensive database of transcript-specific functional predictions
729 and annotations for human nonsynonymous and splice-site SNVs. *Genome Medicine*,
730 12(1):103, December 2020. ISSN 1756-994X. doi: 10.1186/s13073-020-00803-9.
731 URL <https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-020-00803-9>.

- 735 [16] Damian Szklarczyk, Katerina Nastou, Mikaela Koutrouli, Rebecca Kirsch, Far-
736 rokh Mehryary, Radja Hachilif, Dewei Hu, Matteo E Peluso, Qingyao Huang,
737 Tao Fang, et al. The string database in 2025: protein networks with directional-
738 ity of regulation. *Nucleic Acids Research*, 53(D1):D730–D737, 2025.
- 739 [17] Paul Tuijnenburg, Hana Lango Allen, Siobhan O. Burns, Daniel Greene,
740 Machiel H. Jansen, and Others. Loss-of-function nuclear factor B subunit
741 1 (NFKB1) variants are the most common monogenic cause of common vari-
742 able immunodeficiency in Europeans. *Journal of Allergy and Clinical Im-*
743 *munology*, 142(4):1285–1296, October 2018. ISSN 00916749. doi: 10.1016/
744 j.jaci.2018.01.039. URL <https://linkinghub.elsevier.com/retrieve/pii/S0091674918302860>.
- 745 [18] WHO Scientific Group et al. Primary immunodeficiency diseases: report of a
746 who scientific group. *Clin. Exp. Immunol.*, 109(1):1–28, 1997.
- 747 [19] Charlotte Cunningham-Rundles and Carol Bodian. Common variable immunod-
748 eficiency: clinical and immunological features of 248 patients. *Clinical immunol-*
749 *ogy*, 92(1):34–48, 1999.
- 750 [20] Eric Oksenhendler, Laurence Gérard, Claire Fieschi, Marion Malphettes, Gael
751 Mouillot, Roland Jaussaud, Jean-François Viallard, Martine Gardembas, Lionel
752 Galicier, Nicolas Schleinitz, et al. Infections in 252 patients with common variable
753 immunodeficiency. *Clinical Infectious Diseases*, 46(10):1547–1554, 2008.
- 754 [21] Y Naito, F Adams, S Charman, J Duckers, G Davies, and S Clarke. Uk cystic
755 fibrosis registry 2023 annual data report. *London: Cystic Fibrosis Trust*, 2023.
- 756 [22] Carlo Castellani, CFTR2 team, et al. Cftr2: how will it help care? *Paediatric*
757 *respiratory reviews*, 14:2–5, 2013.
- 758 [23] Hartmut Grasemann and Felix Ratjen. Cystic fibrosis. *New England Journal*
759 *of Medicine*, 389(18):1693–1707, 2023. doi: 10.1056/NEJMra2216474. URL
760 <https://www.nejm.org/doi/full/10.1056/NEJMra2216474>.
- 761 [24] Kyoko Watanabe, Erdogan Taskesen, Arjen Van Bochoven, and Danielle
762 Posthuma. Functional mapping and annotation of genetic associations with
763 FUMA. *Nature Communications*, 8(1):1826, November 2017. ISSN 2041-1723.
764 doi: 10.1038/s41467-017-01261-5. URL <https://www.nature.com/articles/s41467-017-01261-5>.
- 765 [25] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir,
766 Pablo Tamayo, and Jill P. Mesirov. Molecular signatures database (MSigDB)
767 3.0. *Bioinformatics*, 27(12):1739–1740, June 2011. ISSN 1367-4811, 1367-
768 4803. doi: 10.1093/bioinformatics/btr260. URL <https://academic.oup.com/bioinformatics/article/27/12/1739/257711>.
- 769
770
771

- 772 [26] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-
773 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias
774 Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-
775 abis, Pushmeet Kohli, and Žiga Avsec. Predictions for alphanonsense, September
776 2023. URL <https://doi.org/10.5281/zenodo.8208688>.
- 777 [27] Dylan Lawless. Variant risk estimate probabilities for iei genes. March 2025. doi:
778 10.5281/zenodo.15111584. URL <https://doi.org/10.5281/zenodo.15111584>.
- 779 [28] Bradley Efron and Carl Morris. Stein’s Estimation Rule and Its Competitors—
780 An Empirical Bayes Approach. *Journal of the American Statistical Association*,
781 68(341):117, March 1973. ISSN 01621459. doi: 10.2307/2284155. URL <https://www.jstor.org/stable/2284155?origin=crossref>.
- 783 [29] Yaowu Liu, Sixing Chen, Zilin Li, Alanna C Morrison, Eric Boerwinkle, and
784 Xihong Lin. Acat: a fast and powerful p value combination method for rare-
785 variant analysis in sequencing studies. *The American Journal of Human Genetics*,
786 104(3):410–421, 2019.
- 787 [30] Xihao Li, Zilin Li, Hufeng Zhou, Sheila M Gaynor, Yaowu Liu, Han Chen, Ryan
788 Sun, Rounak Dey, Donna K Arnett, Stella Aslibekyan, et al. Dynamic incorpora-
789 tion of multiple in silico functional annotations empowers rare variant association
790 analysis of large whole-genome sequencing studies at scale. *Nature genetics*, 52
791 (9):969–983, 2020.
- 792 [31] Michael C Wu, Seunggeun Lee, Tianxi Cai, Yun Li, Michael Boehnke, and Xi-
793 hong Lin. Rare-variant association testing for sequencing data with the sequence
794 kernel association test. *The American Journal of Human Genetics*, 89(1):82–93,
795 2011.
- 796 [32] Seunggeun Lee, Mary J Emond, Michael J Bamshad, Kathleen C Barnes, Mark J
797 Rieder, Deborah A Nickerson, David C Christiani, Mark M Wurfel, and Xihong
798 Lin. Optimal unified approach for rare-variant association testing with applica-
799 tion to small-sample case-control whole-exome sequencing studies. *The American
800 Journal of Human Genetics*, 91(2):224–237, 2012.
- 801 [33] Augustine Kong, Gudmar Thorleifsson, Michael L Frigge, Bjarni J Vilhjalmsson,
802 Alexander I Young, Thorgeir E Thorgeirsson, Stefania Benonisdottir, Asmundur
803 Oddsson, Bjarni V Halldorsson, Gisli Masson, et al. The nature of nurture:
804 Effects of parental genotypes. *Science*, 359(6374):424–428, 2018.
- 805 [34] Laurence J Howe, Michel G Nivard, Tim T Morris, Ailin F Hansen, Humaira
806 Rasheed, Yoonsu Cho, Geetha Chittoor, Penelope A Lind, Teemu Palviainen,
807 Matthijs D van der Zee, et al. Within-sibship gwas improve estimates of direct
808 genetic effects. *BioRxiv*, pages 2021–03, 2021.
- 809 [35] Sue Richards, Nazneen Aziz, Sherri Bale, David Bick, Soma Das, Julie Gastier-
810 Foster, Wayne W Grody, Madhuri Hegde, Elaine Lyon, Elaine Spector, et al.

- 811 Standards and guidelines for the interpretation of sequence variants: a joint
812 consensus recommendation of the american college of medical genetics and ge-
813 nomics and the association for molecular pathology. *Genetics in medicine*, 17
814 (5):405–423, 2015.
- 815 [36] Sean V Tavtigian, Steven M Harrison, Kenneth M Boucher, and Leslie G
816 Biesecker. Fitting a naturally scaled point system to the acmng/amp variant
817 classification guidelines. *Human mutation*, 41(10):1734–1737, 2020.
- 818 [37] Quan Li and Kai Wang. Intervar: clinical interpretation of genetic variants by
819 the 2015 acmng-amp guidelines. *The American Journal of Human Genetics*, 100
820 (2):267–280, 2017.
- 821 [38] Brent S Pedersen, Joe M Brown, Harriet Dashnow, Amelia D Wallace, Matt
822 Velinder, Martin Tristani-Firouzi, Joshua D Schiffman, Tatiana Tvrzik, Rong
823 Mao, D Hunter Best, et al. Effective variant filtering and expected candidate
824 variant yield in studies of rare human disease. *NPJ Genomic Medicine*, 6(1):1–8,
825 2021.
- 826 [39] Carl A Anderson, Fredrik H Pettersson, Geraldine M Clarke, Lon R Cardon,
827 Andrew P Morris, and Krina T Zondervan. Data quality control in genetic
828 case-control association studies. *Nature protocols*, 5(9):1564–1573, 2010. URL
829 <https://doi.org/10.1038/nprot.2010.116>.
- 830 [40] David T Miller, Kristy Lee, Noura S Abul-Husn, Laura M Amendola, Kyle Broth-
831 ers, Wendy K Chung, Michael H Gollob, Adam S Gordon, Steven M Harrison,
832 Ray E Hershberger, et al. Acmng sf v3. 2 list for reporting of secondary findings
833 in clinical exome and genome sequencing: a policy statement of the american
834 college of medical genetics and genomics (acmng). *Genetics in Medicine*, 25(8):
835 100866, 2023.

836 **6 Supplemental**

837 **6.1 Validation studies**

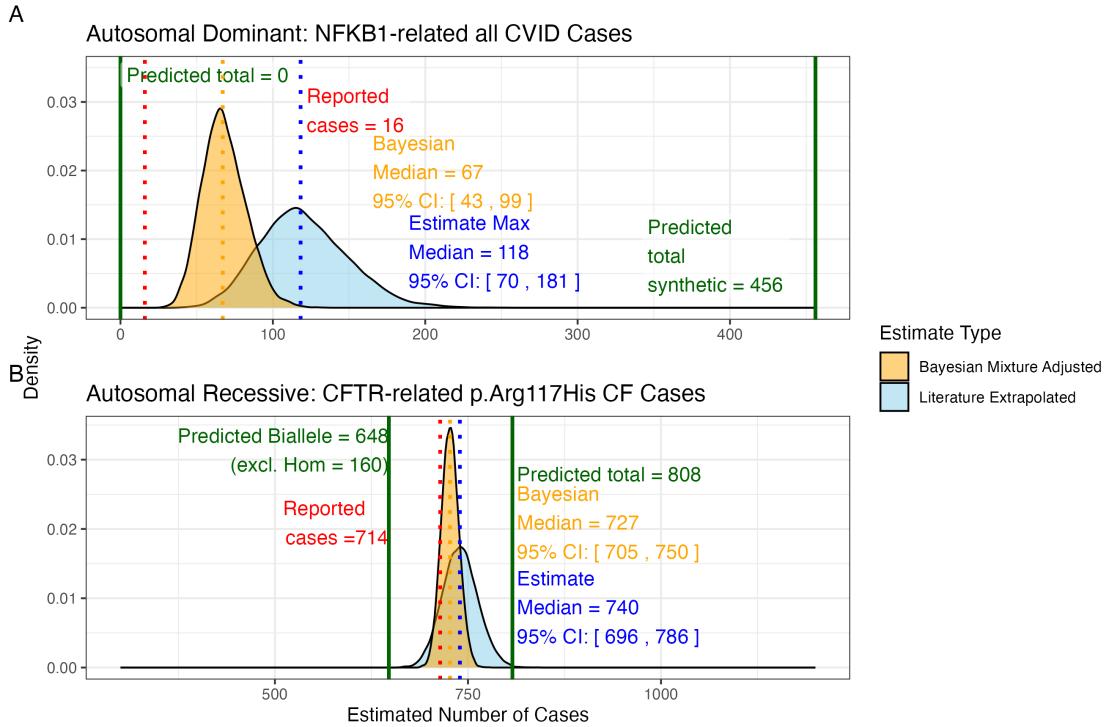


Figure S1: **Prior probabilities compared to validation disease cohort metrics.** (A) Density distributions for the number of *NFKB1*-related CVID cases in the UK. Our model (green) predicted 456 cases, which falls between the observed cohort count (red) of 390 and the upper extrapolated values. The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

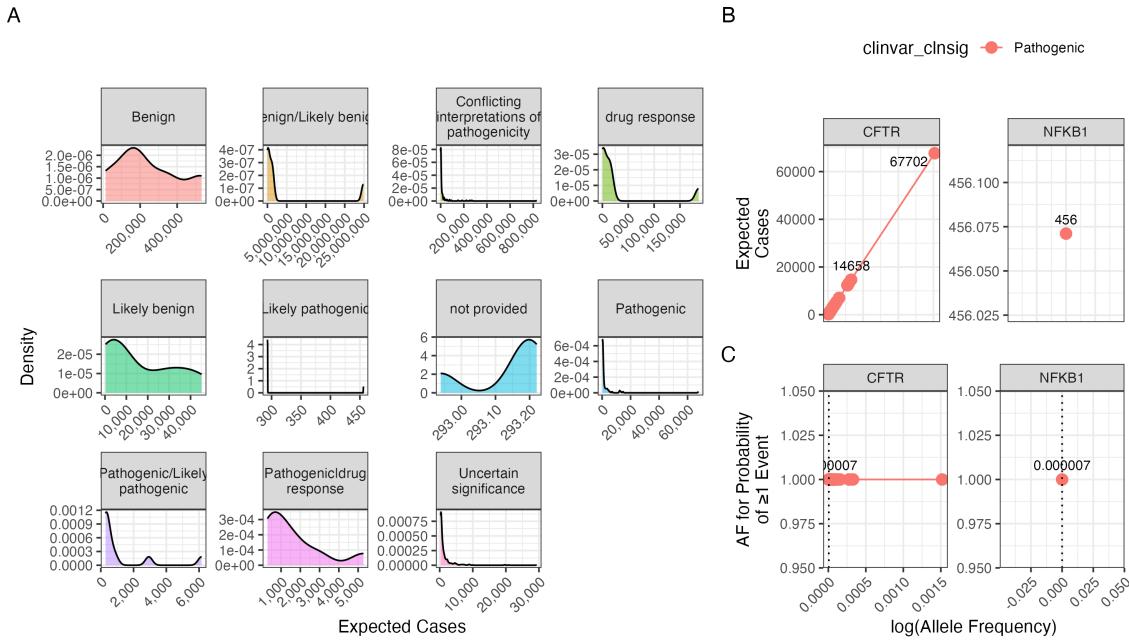


Figure S2: Interpretation of probability of observing a variant classification.
The result from the chosen validation genes *CFTR* and *NFKB1* are shown. Case counts are dependant on the population size and phenotype. (A) The density plots of expected observations by ClinVar clinical significance. We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

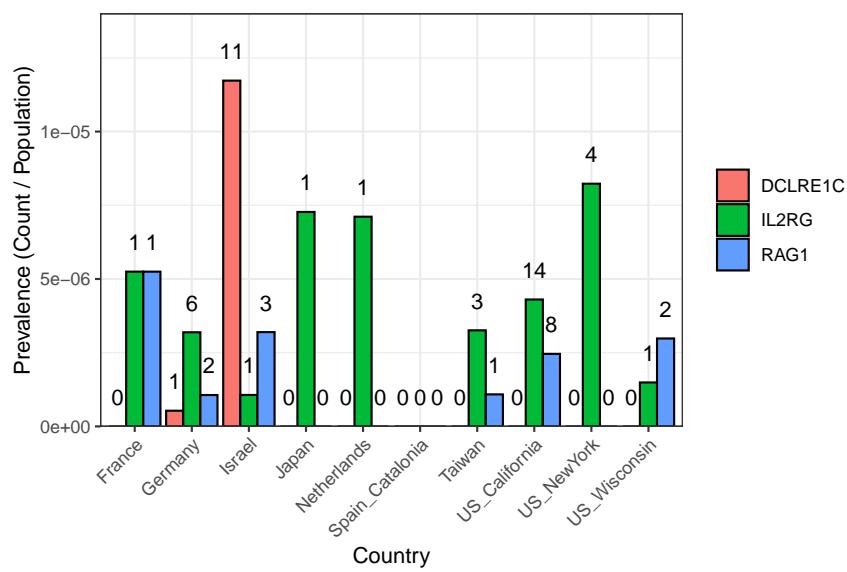


Figure S3: SCID-specific gene comparison across regions. The bar plot shows the prevalence of SCID-related cases (count divided by population) for each gene and country (or region), with numbers printed above the bars representing the actual counts in the original cohort (ranging from 0 to 11 per region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

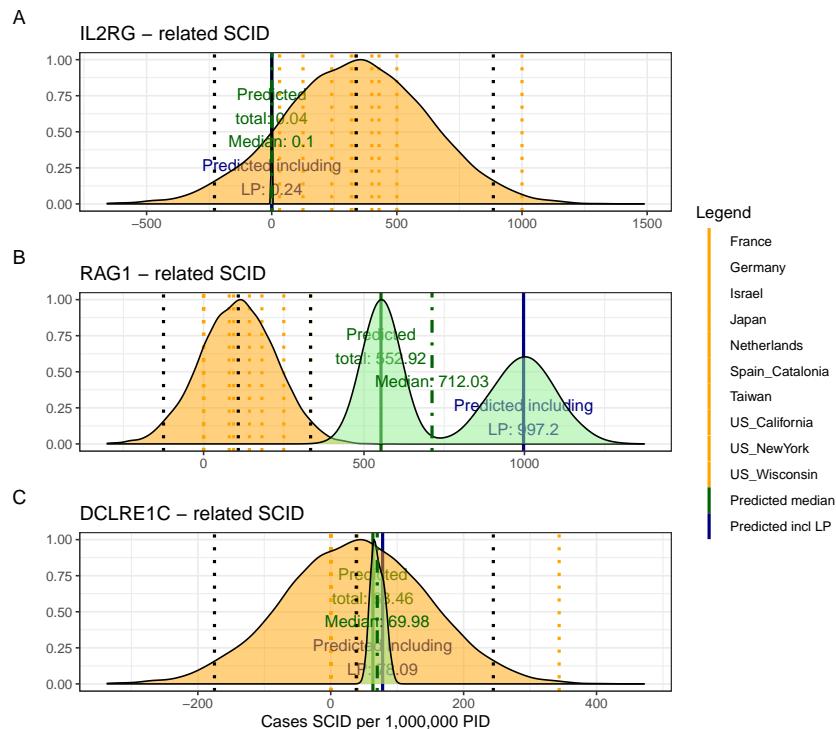


Figure S4: Combined SCID-specific Predictions and Observed Rates per 1,000,000 PID. The figure presents density distributions for the predicted SCID case counts (per 1,000,000 PID) for three genes: *IL2RG*, *RAG1*, and *DCLRE1C*. Country-specific rates (displayed as dotted vertical lines) are overlaid 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.

838 **6.2 Hierarchical Clustering of Enrichment Scores for Major**
 839 **Disease Categories**

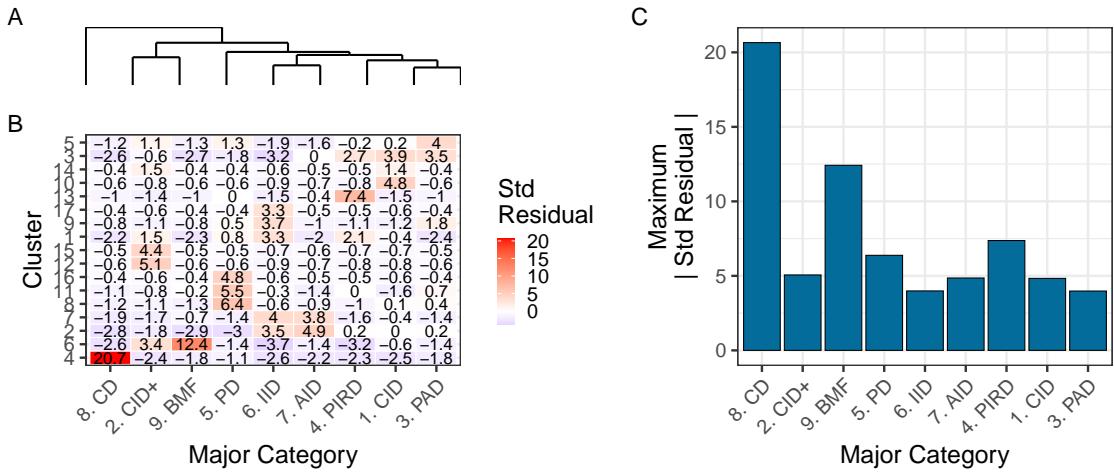


Figure S5: Hierarchical clustering of enrichment scores. The heatmap displays standardised residuals for major disease categories (x-axis) across network clusters (y-axis). A dendrogram groups similar disease categories, and the bar plot shows the maximum absolute residual per category. (8) CD and (9)BMF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

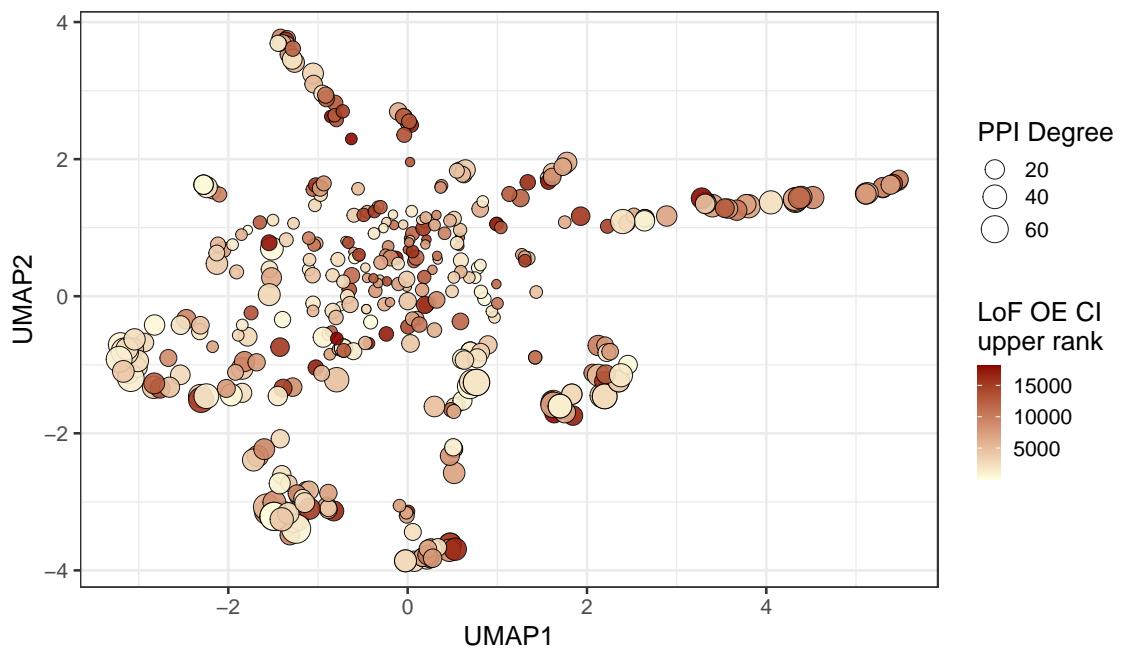


Figure S6: **Analysis of PPI degree versus LOEUF upper rank with UMAP embedding of the PPI network.** The relationship between PPI degree (size) and LOEUF upper rank (color) across gene clusters. No clear patterns are evident.

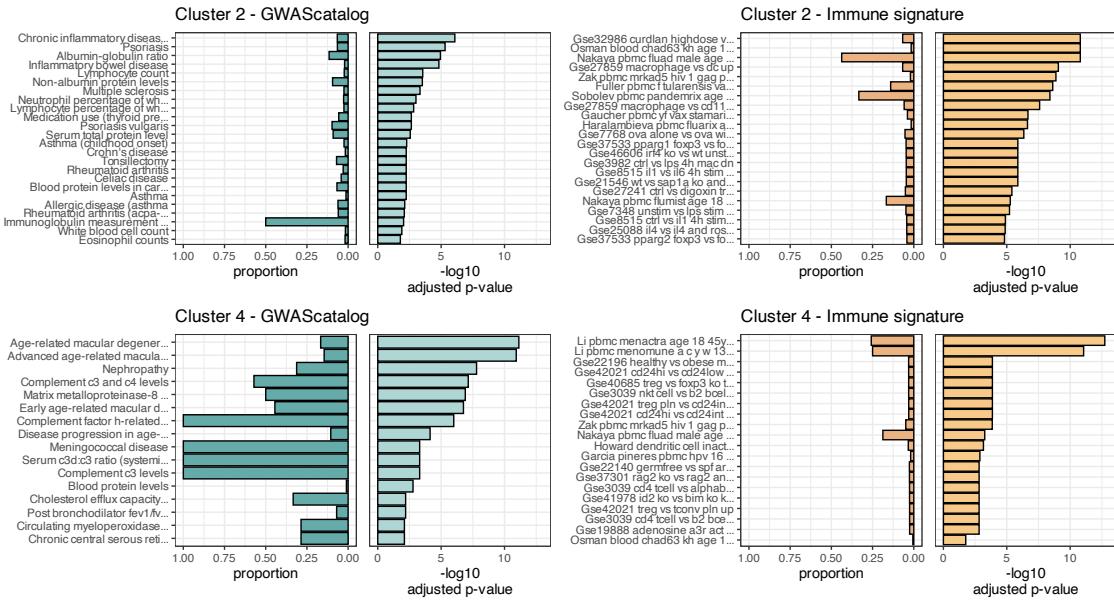


Figure S7: Composite Enrichment Profiles for IEI Gene Sets. We selected the top two enriched clusters (as per **Figure 4**) and performed functional enrichment analysis derived from known disease associations. For each gene set, the left panel displays the proportion of input genes overlapping with a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

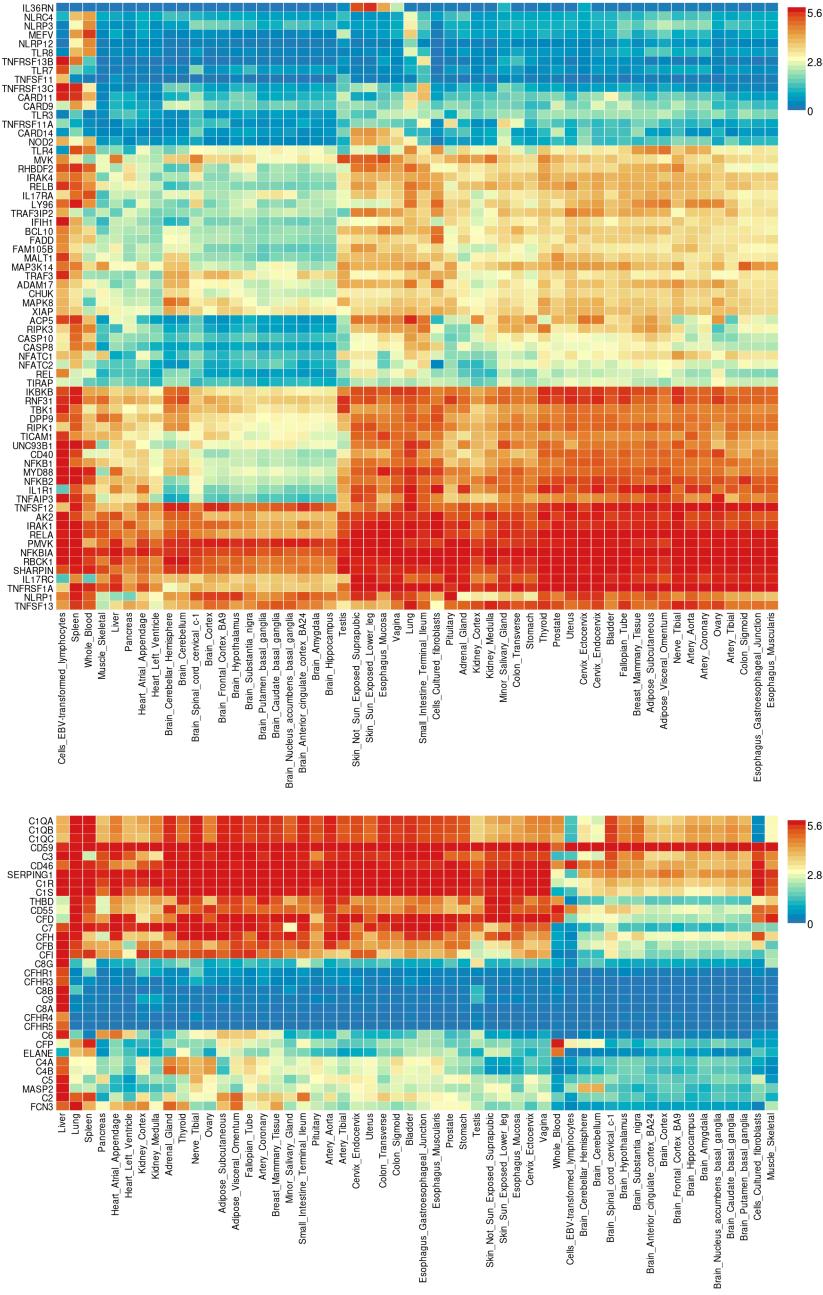
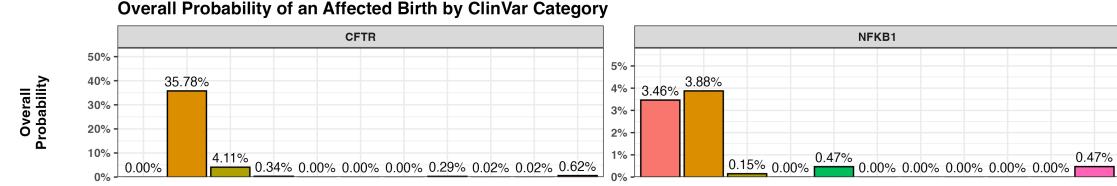


Figure S8: **Gene Expression Heatmaps for IEI Genes.** GTEx v8 data from 54 tissue types display the average expression per tissue label (log₂ transformed) for the IEI gene panels. Top: Cluster 2; Bottom: Cluster 4.

6.3 Interpretation of ClinVar Variant Observations

Recessive and Dominant Disease Genes

A



B

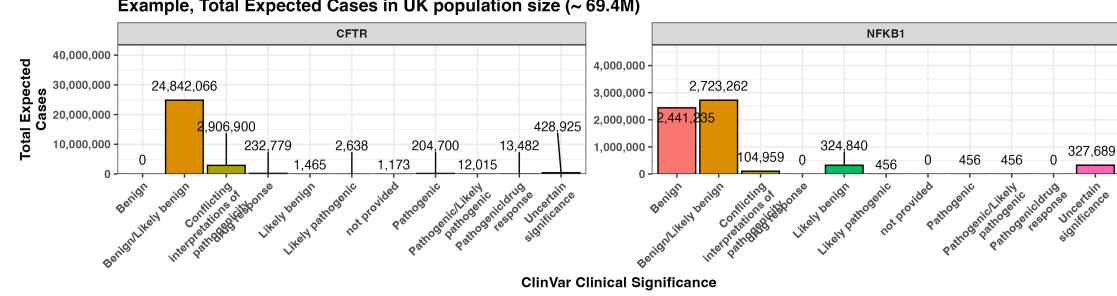


Figure S9: Combined bar charts summarizing the genome-wide analysis of ClinVar clinical significance for the PID gene panel. Panel (A) shows the overall probability of an affected birth by variant classification, and (B) displays the total expected number of cases per classification, both stratified 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

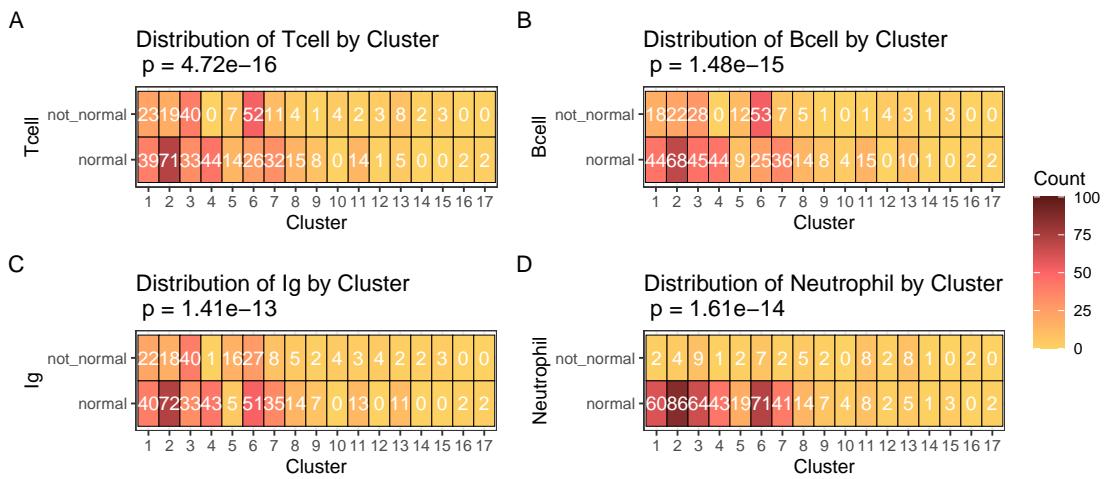


Figure S10: Heatmaps of clinical feature distributions by PPI cluster. The heatmaps display the count of observations for each clinical feature (T cell, B cell, Ig, Neutrophil) in relation to the PPI clusters, with p-values from chi-square tests annotated in the titles.

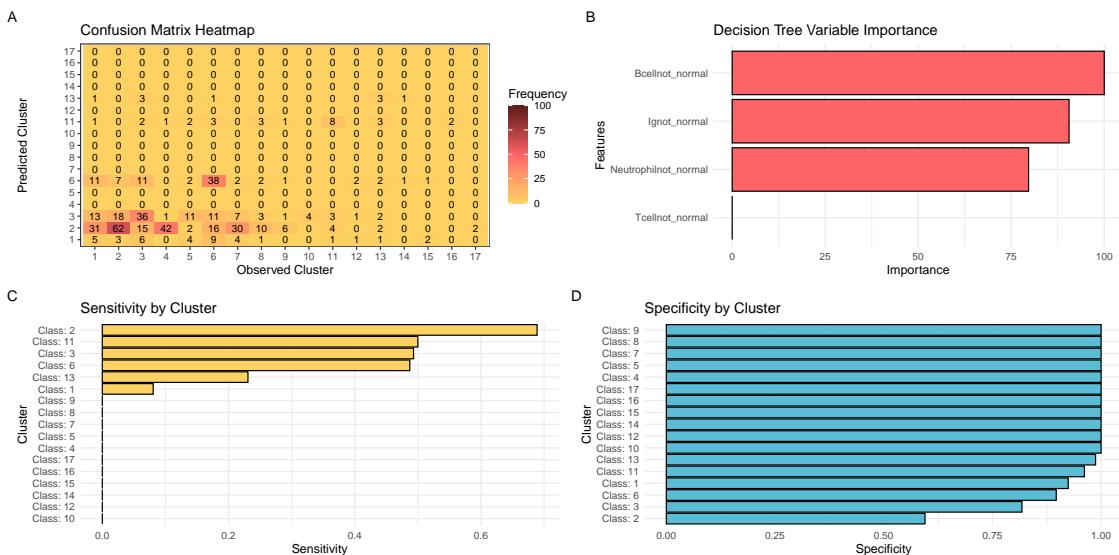


Figure S11: Model performance for fine-tuned PID classification. (A) Confusion matrix heatmap comparing observed and predicted PPI clusters. (B) Variable importance plot ranking immunophenotypic features contributing to the classifier. (C) Per-class sensitivity and (D) per-class specificity bar plots. These panels collectively demonstrate the performance of the decision tree classifier in stratifying PID genes based on immunophenotypic and PPI features.

842 6.5 Probability of observing AlphaMissense pathogenicity

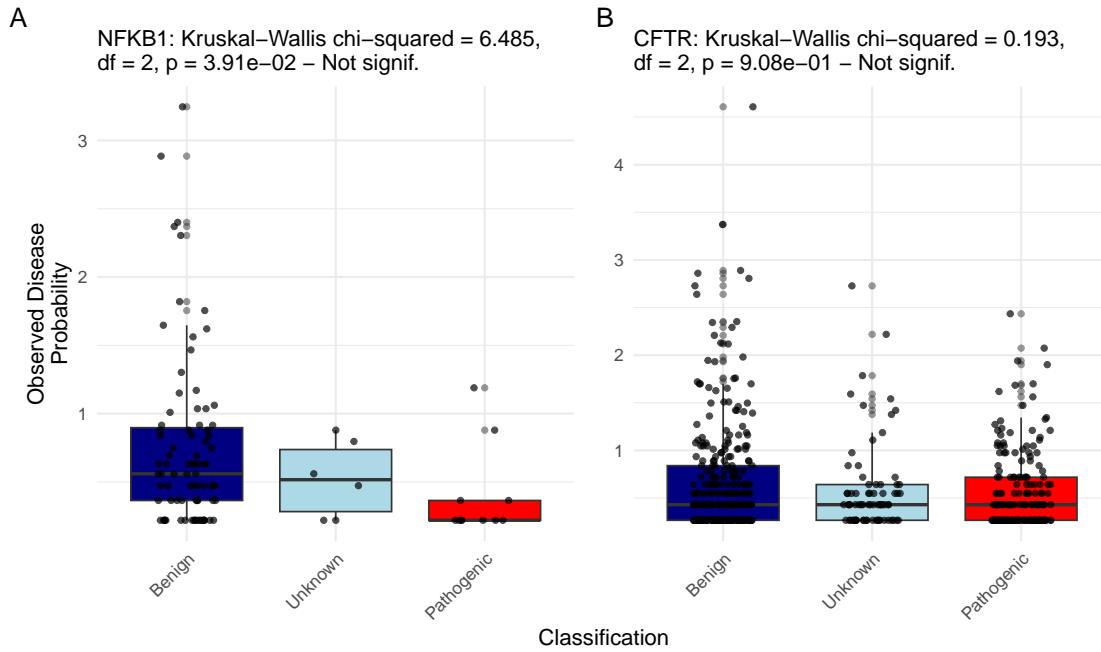


Figure S12: **Observed Disease Probability by Clinical Classification with AlphaMissense.** The figure displays the Kruskal–Wallis test results for NFKB1 and CFTR, showing no significant differences.

843 7 Clinical Genetics Application

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

854 7.1 Methods

855 7.1.1 Quality Control:

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

864 7.1.2 Prior Probability Calculation:

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

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

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

This value reflects the strength of the pathogenic support.

- 869 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:

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

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

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

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

872 This adjusted frequency reflects both the empirical observation and the ACMG
873 evidence.

- 874 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

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

875 7.1.3 Deriving the Confidence Interval (CI)

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

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

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

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

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

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

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

889 where $z = 1.96$ for a 95% confidence level. This interval integrates uncertainty from
 890 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:

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

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

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

898 7.2 Results

899 We illustrate our method with two examples:

900 **Example 1: Missing Sequence Information** In one case, a known pathogenic
901 nucleotide position in *GENE_XYZ* exhibited low quality reads and did not pass QC.
902 This missing information prevents confirmation of the absence of the known variant (a
903 potential false negative), thereby widening the uncertainty in our probability estimate.
904 In such cases, the adjusted allele frequency is calculated with additional variance,
905 leading to a broader CI. For instance, if the observed allele frequency is 1.0×10^{-5}
906 and after adjusting with the ACMG score the estimated occurrence probability is
907 1.0×10^{-5} , the propagated uncertainty might yield a 95% CI of [0.70, 0.85]. This
908 broader interval reflects the impact of missing sequence data on our confidence.

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

920 Table S1 shows the final variant results for a male patient carrying an X-linked
921 loss-of-function (LOF) variant in *GENE_XYZ* where all known pathogenic positions
922 were confirmed as reference alleles. For the variant c. 1234del (p.Glu412Argfs*5), the
923 observed allele frequency is 1.2×10^{-5} . After applying the ACMG evidence score
924 adjustment (for a VUS lacking a ClinVar classification), the adjusted allele frequency
925 remains consistent with the observed data. The resulting occurrence probability is
926 1.2×10^{-5} , and by propagating the uncertainty from the allele frequency, evidence
927 score adjustment, and the full range of possible outcomes (TP, FP, TN, FN), we
928 derive a 95% CI for causality of [0.92, 0.97]. This confirms the variant as the top
929 causal variant in this patient, with no evidence of additional alternative pathogenic
930 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