

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 19, 2025

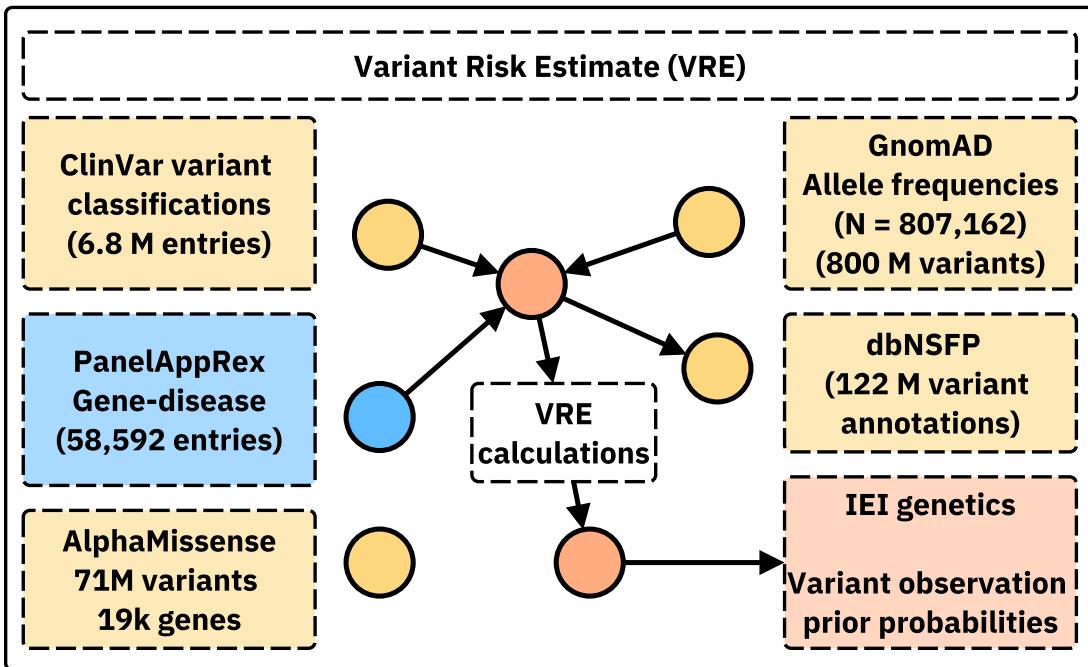
1

Abstract

We present a framework to quantify the prior probability of observing disease-causing variants across all genes and inheritance modes. First, we computed genome-wide occurrence probabilities by integrating population allele frequencies, variant classifications, and Hardy-Weinberg expectations under autosomal dominant, recessive, and X-linked inheritance. Second, we used these priors to derive posterior probabilities for pathogenicity in individual patients, integrating both observed variants (true positives) and unsequenced (false negative) positions into a single, calibrated conclusion. Third, we summarised variant probabilities for 557 genes responsible for inborn errors of immunity (IEI), now integrated into a public database. Fourth, we derived new data-driven IEI classifications using protein-protein interactions and curated clinical features, aligned to immunophenotypes. Finally, we validated the framework in national-scale cohorts of autosomal dominant, recessive, and X-linked disorders, showing close concordance with observed case numbers. The resulting datasets support Bayesian variant interpretation and evidence-weighted decision-making in clinical genetics.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 and IEI-genetics project at <https://github.com/iei-genetics/iei-genetics.github.io>. 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.....	³⁵
²¹ ACAT	Aggregated Cauchy Association Test	³⁵
²² AD	Autosomal Dominant.....	⁴
²³ ANOVA	Analysis of Variance	¹⁴
²⁴ AR	Autosomal Recessive	⁴
²⁵ BMF	Bone Marrow Failure.....	²³
²⁶ CD	Complement Deficiencies	²⁴
²⁷ CI	Confidence Interval.....	⁹
²⁸ CF	Cystic Fibrosis	¹²
²⁹ CFTR	Cystic Fibrosis Transmembrane Conductance Regulator.....	⁶
³⁰ CVID	Common Variable Immunodeficiency.....	¹⁰
³¹ dbNSFP	database for Non-Synonymous Functional Predictions	⁵
³² GE	Genomics England	⁵
³³ gnomAD	Genome Aggregation Database	⁵
³⁴ gVCF	genomic variant call format	⁹
³⁵ HGVS	Human Genome Variation Society.....	⁶
³⁶ HPC	High-Performance Computing.....	⁸
³⁷ HSD	Honestly Significant Difference	¹⁴
³⁸ HWE	Hardy-Weinberg Equilibrium	⁴
³⁹ IEI	Inborn Errors of Immunity	⁴
⁴⁰ Ig	Immunoglobulin	²⁸
⁴¹ InDel	Insertion/Deletion	⁵
⁴² IUIS	International Union of Immunological Societies	⁶
⁴³ LD	Linkage Disequilibrium	²⁶
⁴⁴ LOEUF	Loss-Of-function Observed/Expected Upper bound Fraction	¹⁴
⁴⁵ LOF	Loss-of-Function	²³
⁴⁶ MOI	Mode of Inheritance	⁴
⁴⁷ NFKB1	Nuclear Factor Kappa B Subunit 1	⁶
⁴⁸ OMIM	Online Mendelian Inheritance in Man	³²
⁴⁹ PID	Primary Immunodeficiency	⁴
⁵⁰ PPI	Protein-Protein Interaction	⁶
⁵¹ QC	Quality Control	¹⁷
⁵² SNV	Single Nucleotide Variant	⁴
⁵³ SKAT	Sequence Kernel Association Test.....	³⁵
⁵⁴ STRINGdb	Search Tool for the Retrieval of Interacting Genes/Proteins.....	⁶
⁵⁵ TP	true positive.....	⁹
⁵⁶ FP	false positive.....	⁹

⁹³	TN true negative	9
⁹⁵	FN false negative	9
⁹⁷	UMAP Uniform Manifold Approximation and Projection	23
¹⁰⁰	UniProt Universal Protein Resource	5
¹⁰²	VCF variant call format	9
¹⁰⁴	VEP Variant Effect Predictor	6
¹⁰⁶	XL X-Linked	4
¹⁰⁷		

¹⁰⁸ 1 Introduction

¹⁰⁹ In this study, we focused on reporting the probability of disease observation through
¹¹⁰ genome-wide assessments of gene-disease combinations. Our central hypothesis was
¹¹¹ that by using highly curated annotation data including population allele frequen-
¹¹² cies, disease phenotypes, Mode of Inheritance (MOI) patterns, and variant classi-
¹¹³ fications and by applying rigorous calculations based on Hardy-Weinberg Equilib-
¹¹⁴ rium (HWE), we could accurately estimate the expected probabilities of observing
¹¹⁵ disease-associated variants. Among other benefits, this knowledge can be used to
¹¹⁶ derive genetic diagnosis confidence by incorporating these new priors.

¹¹⁷ In this report, we focused on known Inborn Errors of Immunity (IEI) genes, also re-
¹¹⁸ ferred to as the Primary Immunodeficiency (PID) or Monogenic Inflammatory Bowel
¹¹⁹ Disease genes ([1–3](#)) to validate our approach and demonstrate its clinical relevance.
¹²⁰ This application to a well-established genotype-phenotype set, comprising over 500
¹²¹ gene-disease associations, underscores its utility ([1](#)).

¹²² Quantifying the risk that a newborn inherits a disease-causing variant is a fun-
¹²³ damental challenge in genomics. Classical statistical approaches grounded in HWE
¹²⁴ ([4; 5](#)) have long been used to calculate genetic MOI probabilities for Single Nucleotide
¹²⁵ Variant (SNV)s. However, applying these methods becomes more complex when ac-
¹²⁶ counting for different MOI, such as Autosomal Recessive (AR) versus Autosomal
¹²⁷ Dominant (AD) or X-Linked (XL) disorders. In AR conditions, for example, the
¹²⁸ occurrence probability must incorporate both the homozygous state and compound
¹²⁹ heterozygosity, whereas for AD and XL disorders, a single pathogenic allele is suffi-
¹³⁰ cient to cause disease. Advances in genetic research have revealed that MOI can be
¹³¹ even more complex ([6](#)). Mechanisms such as dominant negative effects, haploinsuffi-
¹³² ciency, mosaicism, and digenic or epistatic interactions can further modulate disease
¹³³ risk and clinical presentation, underscoring the need for nuanced approaches in risk
¹³⁴ estimation. Karczewski et al. ([7](#)) made significant advances; however, the remain-
¹³⁵ ing challenge lay in applying the necessary statistical genomics data across all MOI
¹³⁶ for any gene-disease combination. Similar approaches have been reported for disease
¹³⁷ such Wilson disease, Mucopolysaccharidoses, Primary ciliary dyskinesia, and treat-
¹³⁸ able metabolic diseasesse, ([8; 9](#)), as reviewed by Hannah et al. ([10](#)).

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

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

165 2 Methods

166 2.1 Dataset

167 Data from Genome Aggregation Database (gnomAD) v4 comprised 807,162 indi-
168 viduals, including 730,947 exomes and 76,215 genomes (7). This dataset provided
169 786,500,648 SNVs and 122,583,462 InDels, with variant type counts of 9,643,254 syn-
170 onymous, 16,412,219 missense, 726,924 nonsense, 1,186,588 frameshift and 542,514
171 canonical splice site variants. ClinVar data were obtained from the variant summary
172 dataset (as of: 16 March 2025) available from the NCBI FTP site, and included
173 6,845,091 entries, which were processed into 91,319 gene classification groups and a
174 total of 38,983 gene classifications; for example, the gene *A1BG* contained four vari-
175 ants classified as likely benign and 102 total entries (13). For our analysis phase
176 we also used dbNSFP which consisted of a number of annotations for 121,832,908
177 SNVs (15). The PanelAppRex core model contained 58,592 entries consisting of

178 52 sets of annotations, including the gene name, disease-gene panel ID, diseases-
 179 related features, confidence measurements. (2) A Protein-Protein Interaction (PPI)
 180 network data was provided by Search Tool for the Retrieval of Interacting Genes/Pro-
 181 teins (STRINGdb), consisting of 19,566 proteins and 505,968 interactions (16). The
 182 Human Genome Variation Society (HGVS) nomenclature is used with Variant Effect
 183 Predictor (VEP)-based codes for variant IDs. We carried out validations for disease
 184 cohorts with Nuclear Factor Kappa B Subunit 1 (*NFKB1*) (17–20) and Cystic Fibrosis
 185 Transmembrane Conductance Regulator (*CFTR*) (21–23) to demonstrate applications
 186 in AD and AR disease genes, respectively. AlphaMissense includes pathogenicity pre-
 187 diction classifications for 71 million variants in 19 thousand human genes (12; 26).
 188 We used these scores to compared against the probability of observing the same given
 189 variants. **Box 2.1** list the definitions from the International Union of Immunological
 190 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

191

192 2.2 Variant Class Observation Probability

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

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

193 where p is the allele frequency, $q = 1 - p$, p^2 represents the homozygous dominant,
 194 $2pq$ the heterozygous, and q^2 the homozygous recessive genotype frequencies. For
 195 disease phenotypes, particularly under AR MOI, the risk is traditionally linked to
 196 the homozygous state (p^2); however, to account for compound heterozygosity across
 197 multiple variants, we allocated the overall gene-level risk proportionally among vari-
 198 ants.

199 Our computational pipeline estimated the probability of observing a disease-associated
 200 genotype for each variant and aggregated these probabilities by gene and ClinVar
 201 classification. This approach included all variant classifications, not limited solely to
 202 those deemed “pathogenic”, and explicitly conditioned the classification on the given

203 phenotype, recognising that a variant could only be considered pathogenic relative to
 204 a defined clinical context. The core calculations proceeded as follows:

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

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

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

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

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

213 **2. Occurrence Probability Based on MOI.** The probability that an individual
 214 was affected by a variant depended on the MOI relative to a specific phenotype.
 215 Specifically, we calculated the occurrence probability $p_{\text{disease},i}$ for each variant as follows:

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

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

- For **AR** variants, disease is expected to manifest when two pathogenic alleles were present. In this case, we accounted for both the homozygous state and the possibility of compound heterozygosity. We allocated the overall gene-level risk (P_{tot}^2) proportionally by variant allele frequency:

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

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

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

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

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

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

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

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

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

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

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

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

246 **2.3 Integrating observed true positives and unobserved false**
247 **negatives into a single, actionable conclusion**

248 In this section, we detail our approach to integrating sequencing data with prior clas-
249 sification evidence (e.g. pathogenic on ClinVar) to calculate the posterior probability
250 of a complete successful genetic diagnosis. Our method is designed to account for
251 possible outcomes of true positive (TP), true negative (TN), and false negative (FN),
252 by first ensuring that all nucleotides corresponding to known variant classifications
253 (benign, pathogenic, etc.) have been accurately sequenced. This implies the use of ge-
254 nomic variant call format (gVCF)-style data which refer to variant call format (VCF)s
255 that contain a record for every position in the genome (or interval of interest) regard-
256 less of whether a variant was detected at that site or not. Only after confirming that
257 these positions match the reference alleles (or novel unaccounted variants are classi-
258 fied) do we calculate the probability that additional, alternative pathogenic variants
259 (those not observed in the sequencing data) could be present. Our Confidence In-
260 terval (CI) for pathogenicity thus incorporates uncertainty from the entire process,
261 including the tally of TP, TN, and FN outcomes. We ignore the contribution of false
262 positive (FP)s as a separate task to be tackled in the future.

263 We estimated, for every query (e.g. gene or disease-panel), the posterior proba-
264 bility that at least one constituent allele is both damaging and causal in the proband.
265 The workflow comprises four consecutive stages.

266 **(i) Data pre-processing** All coding and canonical splice-region variants for *NFKB1*
267 were extracted from the gVCF. Sites corresponding to previously reported pathogenic
268 alleles were checked for read depth ≥ 10 and genotype quality ≥ 20 . Positions that
269 failed this check were labelled *missing*, thus separating true reference calls from un-
270 informative sequence.

271 **(ii) Evidence mapping and occurrence probability** PanelAppRex variants
272 were annotated with ClinVar clinical significance. Each label was converted to an
273 ordinal evidence score $S_i \in [-5, 5]$ (Table S2) and rescaled to a pathogenic weight
274 $W_i = \text{rescale}(S_i; -5, 5 \rightarrow 0, 1)$. The HWE-based pipeline of Section 2.2 supplied a
275 per-variant occurrence probability p_i . The adjusted prior was

$$p_i^* = W_i p_i, \quad \text{and} \quad \text{flag}_i \in \{\text{present}, \text{missing}\}.$$

276 **(iii) Prior specification** In a hypothetical cohort of $n = 200$ diploid individuals
277 the count of allele i follows a Beta–Binomial model. Marginalising the Binomial yields
278 the Beta prior

$$\pi_i \sim \text{Beta}(\alpha_i, \beta_i), \quad \alpha_i = \text{round}(2np_i^*) + \tilde{w}_i, \quad \beta_i = 2n - \text{round}(2np_i^*) + 1,$$

279 where $\tilde{w}_i = \max(1, S_i + 1)$ contributes an additional pseudo-count whenever $S_i >$
 280 0.

281 **(iv) Posterior simulation and aggregation** For each variant i we drew $M =$
 282 10 000 realisations $\pi_i^{(m)}$ and normalised within each iteration,

$$\tilde{\pi}_i^{(m)} = \frac{\pi_i^{(m)}}{\sum_j \pi_j^{(m)}}.$$

283 Variants with $S_i > 3$ were deemed *causal*. Their mean posterior share $\bar{\pi}_i =$
 284 $M^{-1} \sum_m \tilde{\pi}_i^{(m)}$ and 95% credible interval were retained. The probability that a dam-
 285 aging causal allele is physically present was obtained by a second layer:

$$P^{(m)} = \sum_{i: S_i > 3} \tilde{\pi}_i^{(m)} G_i^{(m)}, \quad G_i^{(m)} \sim \text{Bernoulli}(g_i),$$

286 with $g_i = 1$ for present variants, $g_i = 0$ for reference calls, and $g_i = p_i$ for missing
 287 variants. The gene-level estimate is the median of $\{P^{(m)}\}_{m=1}^M$ and its 2.5th/97.5th
 288 percentiles.

289 **Scenario analysis** Two scenarios were explored: (1) *observed variants only* includ-
 290 ing only one known TP pathogenic variant, p.Ser237Ter, and (2) *inclusion of the*
 291 *additional plausible yet unsequenced splice-donor allele c.159+1G>A* as a FN. All
 292 subsequent steps were identical.

293 2.4 Validation of Autosomal Dominant Estimates Using *NFKB1*

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

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

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

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

305 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

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

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

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

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

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

311 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

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

312 which yielded a Bayesian mixture adjusted median estimate of 67 cases with a
313 95% credible interval of approximately 43 to 99 cases.

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

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

324 **2.5 Validation Study for Autosomal Recessive CF Using CFTR**

325 To validate our framework for AR diseases, we focused on Cystic Fibrosis (CF).
 326 For comparability sizes between the validation studies, we analysed the most com-
 327 mon SNV in the *CFTR* gene, typically reported as “p.Arg117His” (GRCh38 Chr
 328 7:117530975 G/A, MANE Select HGVSp ENST00000003084.11: p.Arg117His). Our
 329 goal was to validate our genome-wide probability estimates by comparing the ex-
 330 pected number of CF cases attributable to the p.Arg117His variant in *CFTR* with
 331 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:

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

332 using a large effective allele count (AN_{eff}) for illustration. By generating 10,000 posterior
333 samples of p , we obtained a distribution of the literature-based adjusted expected
334 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:

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

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

339 **2.6 Validation of SCID-specific Estimates Using PID–SCID 340 Genes**

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

346 **1. Incidence Conversion.** Based on literature, PID occurs in approximately 1 in
347 1,000 births, whereas SCID occurs in approximately 1 in 100,000 births. Consequently,
348 in a population of 1,000,000 births there are about 1,000 PID cases and 10 SCID cases.
349 To express SCID-related variant counts on a per-million PID scale, the observed SCID
350 counts were multiplied by 100. For example, if a gene is expected to cause SCID in
351 10 cases within the total PID population, then on a per-million PID basis the count
352 is $10 \times 100 = 1,000$ cases (across all relevant genes).

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

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

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

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

375 **2.7 Protein Network and Genetic Constraint Interpretation**

376 A PPI network was constructed using protein interaction data from STRINGdb (16).
377 We previously prepared and reported on this dataset consisting of 19,566 proteins and
378 505,968 interactions (<https://github.com/DylanLawless/ProteoMCLustR>). Node
379 attributes were derived from log-transformed score-positive-total values, which in-
380 formed both node size and colour. Top-scoring nodes (top 15 based on score) were
381 labelled to highlight prominent interactions. To evaluate group differences in score-
382 positive-total across major disease categories, one-way Analysis of Variance (ANOVA)
383 was performed followed by Tukey Honestly Significant Difference (HSD) post hoc tests
384 (and non-parametric Dunn's test for confirmation). GnomAD v4.1 constraint metrics
385 data was used for the PPI analysis and was sourced from Karczewski et al. (7). This
386 provided transcript-level metrics, such as observed/expected ratios, Loss-Of-function
387 Observed/Expected Upper bound Fraction (LOEUF), pLI, and Z-scores, quantifying
388 loss-of-function and missense intolerance, along with confidence intervals and related
389 annotations for 211,523 observations.

390 **2.8 Gene Set Enrichment Test**

391 To test for overrepresentation of biological functions, the prioritised genes were com-
392 pared against gene sets from MsigDB (including hallmark, positional, curated, motif,
393 computational, GO, oncogenic, and immunologic signatures) and WikiPathways using
394 hypergeometric tests with FUMA (24; 25). The background set consisted of 24,304
395 genes. Multiple testing correction was applied per data source using the Benjamini-
396 Hochberg method, and gene sets with an adjusted P-value ≤ 0.05 and more than one
397 overlapping gene are reported.

398 **2.9 Deriving novel PID classifications by genetic PPI and
399 clinical features**

400 We recategorised 315 immunophenotypic features from the original IUIS IEI anno-
401 tations, reducing the original multi-level descriptors (e.g. “decreased cd8, normal or
402 decreased cd4”) first to minimal labels (e.g.“low”) and second to binary outcomes (nor-
403 mal vs. not-normal) for T cells, B cells, neutrophils, and immunoglobulins Each gene
404 was mapped to its PPI cluster derived from STRINGdb and UMAP embeddings from
405 previous steps. We first tested for non-random associations between these four binary
406 immunophenotypes and PPI clusters using χ^2 tests. To generate a data-driven PID
407 classification, we trained a decision tree (rpart) to predict PPI cluster membership
408 from the four immunophenotypic features plus the traditional IUIS Major and Subcat-
409 egory labels. Hyperparameters (complexity parameter = 0.001, minimum split = 10,
410 minimum bucket = 5, maximum depth = 30) were optimised via five-fold cross vali-
411 dation using the caret framework. Terminal node assignments were then relabelled
412 according to each group’s predominant abnormal feature profile.

413 **2.10 Probability of observing AlphaMissense pathogenicity**

414 We obtained the subset pathogenicity predictions from AlphaMissense via the Al-
415 phaFold database and whole genome data from the studies data repository(12; 26).
416 The AlphaMissense data (genome-aligned and amino acid substitutions) were merged
417 with the panel variants based on genomic coordinate and HGVS annotation. Occur-
418 rence probabilities were log-transformed and adjusted (y-axis displaying $\log_{10}(\text{occurrence}$
419 $\text{prob} + 1\text{e-}5) + 5$), to visualise the distribution of pathogenicity scores across the
420 residue sequence. A Kruskal-Wallis test was used to compare the observed disease
421 probability across clinical classification groups.

422 3 Results

423 3.1 Observation probability across disease genes

424 Our study integrated large-scale annotation databases with gene panels from PanelAppRex to systematically assess disease genes by MOI. By combining population
425 allele frequencies with ClinVar clinical classifications, we computed an expected observation
426 probability for each SNV, representing the likelihood of encountering a variant
427 of a specific pathogenicity for a given phenotype. We report these probabilities for
428 54,814 ClinVar variant classifications across 557 genes (linked dataset (27)).

430 In practice, our approach computed a simple observation probability for every
431 SNV across the genome and was applicable to any disease-gene panel. Here, we focused
432 on panels related to Primary Immunodeficiency or Monogenic Inflammatory
433 Bowel Disease, using PanelAppRex panel ID 398 as a case study. **Figure 1** displays all reported ClinVar variant classifications for this panel. The resulting natural
434 scaling system (-5 to +5) accounts for the frequently encountered combinations of
435 classification labels (e.g. benign to pathogenic). The resulting data set (27) is briefly
436 shown in **Table 1** to illustrate that our method yielded estimations of the probability
437 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

439 3.2 Integrating observed true positives and unobserved false 440 negatives into a single, actionable conclusion

441 Having previously established a probabilistic framework for estimating the prior probability
442 of observing disease-associated variants under different inheritance modes, we
443 then applied this model to a specific patient to demonstrate it’s potential for clinical
444 genetics. For each gene, we used the computed prior probabilities for all variant
445 classifications, (e.g. benign, uncertain, and pathogenic). We verified that all known
446 pathogenic positions have been sequenced and observed as reference (true negatives),

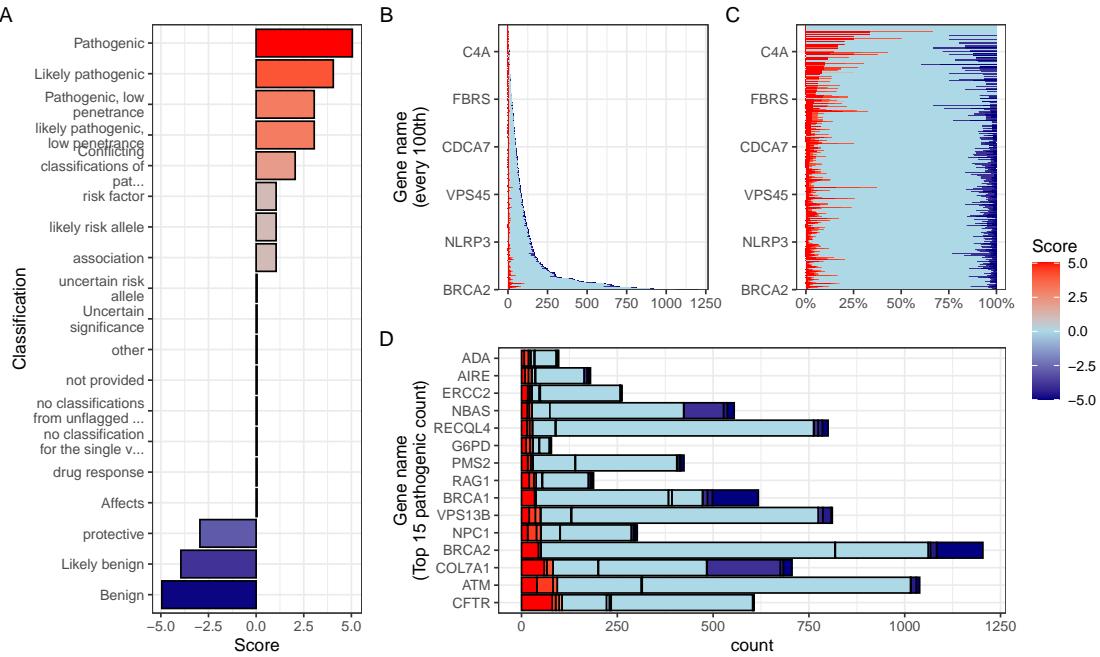


Figure 1: Summary of ClinVar clinical significance classifications in the PID gene panel. (A) Shows the numeric score coding for each classification. Panels (B) and (C) display the tally of classifications per gene as absolute counts and as percentages, respectively. (D) Highlights the top 15 genes with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

and identify any positions that were either observed as variant (true positives) or not assessable due to missing sequence data or failed Quality Control (QC). These missing sites represent potential false negatives. By jointly modelling the observed and unobserved space, the method yielded a calibrated, evidence-weighted probability that at least one damaging causal variant could be present in the gene.

Figure S1 shows the results of the first scenario, in which only one known pathogenic variant, *NFKB1* p.Ser237Ter, was observed and all previously reported pathogenic positions were successfully sequenced and confirmed as reference. In this setting, the model assigned the full posterior probability to the observed allele, yielding 100 % confidence that all present evidence supported a single, true positive causal explanation.

Figure 2 shows the second scenario, where the same pathogenic variant *NFKB1* p.Ser237Ter was present, but coverage was incomplete at three additional sites known to harbour potentially pathogenic variants. Among these was the likely-pathogenic splice-site variant c.159+1G>A, which was not captured in the sequencing data. The panels of **Figure 2 (A-F)** illustrate the stepwise integration of observed and missing evidence, culminating in a posterior probability that reflects both confirmed findings and residual uncertainty. **Table 2** lists the metrics used to reach the conclusion for reporting the clinical genetics results.

Bayesian integration of every annotated *NFKB1* allele yielded the first quantitative credible intervals for gene-level pathogenic attribution that (i) preserve Hardy-Weinberg expectations, (ii) accommodate AD, AR, XL inheritance, and (iii) carry explicit uncertainty for non-sequenced (or failed QC) genomic positions. **Figure 2** (A) depicts the prior landscape where occurrence probabilities are partitioned by observed or missing status and by causal or non-causal evidence, with colour reflecting the underlying ClinVar score. **Figure 2** (B) shows posterior normalisation which concentrates probability density on two high-confidence (high evidence score) alleles since the benign variants are, by definition, non-causal. **Figure 2** (C) shows the resulting per-variant probability of being simultaneously damaging and causal; only the confirmed present (true positive) nonsense variant p.Ser237Ter and the (false negative) hypothetical splice-donor c.159+1G>A retain substantial support. Restricting the view to causal candidates in **Figure 2** (D) confirms that posterior mass is distributed almost equally between these two variants. **Figure 2** (E) decomposes the total damaging probability into observed (52 %) and missing (48 %) sources, whereas **Figure 2** (F) summarises the gene-level posterior: inclusion of the splice-site allele (scenario 2) produces a median probability of 0.75 with a 95 % credible interval of 0.50–0.92, compared with 0.38 (0.24–0.54) when the analysis is limited to sequenced alleles (scenario 1). Numerically, the present variant p.Ser237Ter accounts for 0.38 (95 % CrI 0.21–0.55) of the posterior share, the the potentially causal but missing splice-donor allele contributes 0.37 (0.24–0.53), and the remaining four alleles together explain a negligible % (**Table 2**).

Table 2: Result of clinical genetics diagnosis scenario 2. The proband carried three observed variants, including the known pathogenic p.Ser237Ter (true positive), and lacked coverage at three additional sites, including likely-pathogenic splice-donor c.159+1G>A (false negative). The damaging-only posterior probabilities for these two variants were 0.382 and 0.351, resulting in total probability (prob causal) of causal diagnosis given the existing evidence of 0.521 (95% CI: 0.248–0.787).

Variant	Flag	Class	Evidence Score	Occurrence Prob	Adj Occ Prob	Alpha	Beta	Lower	Median	Upper	Posterior Share	Prob Causal
p.Ser237Ter	present	causal	5.0	0.000	0	6.0	369	0.003	0.091	0.561	0.382	0.382
c.159+1G>A	missing	causal	4.5	0.000	0	5.5	367	NA	NA	NA	0.351	0.351
p.Thr567Ile	missing	other	-5.0	0.002	0	1.0	331	NA	NA	NA	0.000	0.000
p.Arg231His	present	other	0.0	0.000	0	1.0	375	0.003	0.091	0.561	0.000	0.000
p.Gly650Arg	present	other	0.0	0.000	0	1.0	355	0.003	0.091	0.561	0.000	0.000
p.Val236Ile	missing	other	0.0	0.000	0	1.0	357	NA	NA	NA	0.000	0.000
Total	NA	NA	NA	NA	NA	NA	NA	0.248	0.521	0.787	NA	0.521

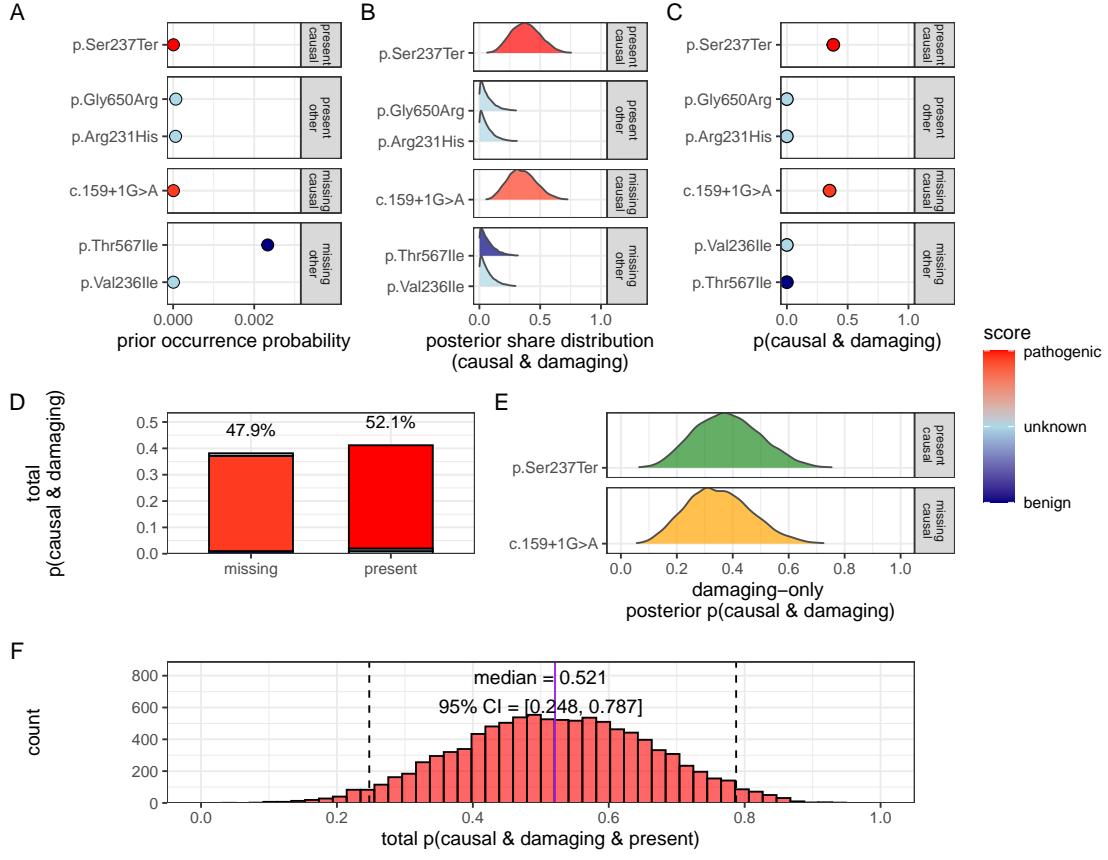


Figure 2: Quantification of gene-level pathogenic attribution for *NFKB1* (scenario 2). The proband carried three known heterozygous variants, including pathogenic p.Ser237Ter, and had incomplete coverage at three additional loci, including likely-pathogenic splice-site variant c.159+1G>A. The sequential steps towards the posterior probability of complete genetic diagnosis are shown: (A) Prior occurrence probabilities, stratified by observed/missing and causal/non-causal status. (B) Posterior distributions of normalised variant weights $\tilde{\pi}_i$. (C) Per-variant posterior probability of being both damaging and causal. (D) Posterior distributions for causal variants only. (E) Decomposition of total pathogenic probability into observed (green) and missing (orange) sources. (F) Gene-level posterior showing the probability that at least one damaging causal allele is present; median 0.75, 95 % credible interval 0.50–0.92.

488 **3.3 Validation studies**

489 **3.3.1 Validation of dominant disease occurrence with *NFKB1***

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

497 Our model calculated 0 known pathogenic variant *NFKB1*-related CVID cases
498 in the UK with a minimal risk of 456 unknown de novo variants. In the reference
499 cohort, 16 *NFKB1* CVID cases were reported. We additionally wanted to account for
500 potential under-reporting in the reference study. We used an extrapolated national
501 CVID prevalence which yielded a median estimate of 118 cases (95% CI: 70–181),
502 while a Bayesian-adjusted mixture estimate produced a median of 67 cases (95% CI:
503 43–99). **Figure S2 (A)** illustrates that our predicted values reflect these ranges and
504 are closer to the observed count. This case supports the validity of our integrated
505 probability estimation framework for AD disorders, and represents a challenging ex-
506 ample where pathogenic SNV are not reported in the reference population of gnomAD.
507 Our min-max values successfully contained the true reported values.

508 **3.3.2 Validation of recessive disease occurrence with *CFTR***

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

514 In contrast, the nationally reported number of CF cases was 714, as recorded
515 in the UK Cystic Fibrosis Registry 2023 Annual Data Report (21). To account for
516 factors such as reduced penetrance and the mortality-adjusted expected genotype,
517 we derived a Bayesian-adjusted estimate via posterior simulation. Our Bayesian ap-
518 proach yielded a median estimate of 740 cases (95% CI: 696, 786) and a mixture-
519 based estimate of 727 cases (95% CI: 705, 750). **Figure S2 (B)** illustrates the close
520 concordance between the predicted values, the Bayesian-adjusted estimates, and the
521 national report supports the validity of our approach for estimating disease.

522 **Figure S3** shows the final values for these genes of interest in a given population
523 size and phenotype. It reveals that an allele frequency threshold of approximately
524 0.000007 is required to observe a single heterozygous disease-causing variant carrier in
525 the UK population for both genes. However, owing to the AR MOI pattern of *CFTR*,
526 this threshold translates into more than 100,000 heterozygous carriers, compared to

527 only 456 carriers for the AD gene *NFKB1*. Note that this allele frequency threshold,
528 being derived from the current reference population, represents a lower bound that
529 can become more precise as public datasets continue to grow. This marked difference
530 underscores the significant impact of MOI patterns on population carrier frequencies
531 and the observed disease prevalence.

532 3.3.3 Interpretation of ClinVar variant observations

533 **Figure S10** shows the two validation study PID genes, representing AR and domi-
534 nant MOI. **Figure S10 (A)** illustrates the overall probability of an affected birth by
535 ClinVar variant classification, whereas **Figure S10 (B)** depicts the total expected
536 number of cases per classification for an example population, here the UK, of approx-
537 imately 69.4 million.

538 3.3.4 Validation of SCID-specific disease occurrence

539 Given that SCID is a subset of PID, our probability estimates reflect the likelihood of
540 observing a genetic variant as a diagnosis when the phenotype is PID. However, we
541 additionally tested our results against SCID cohorts in **Figure S5**. The summarised
542 raw cohort data for SCID-specific gene counts are summarised and compared across
543 countries in **Figure S4**. True counts for *IL2RG* and *DCLRE1C* from ten distinct
544 locations yielded 95% confidence intervals surrounding our predicted values. For
545 *IL2RG*, the prediction was low (approximately 1 case per 1,000,000 PID), as expected
546 since loss-of-function variants in this X-linked gene are highly deleterious and rarely
547 observed in gnomAD. In contrast, the predicted value for *RAG1* was substantially
548 higher (553 cases per 1,000,000 PID) than the observed counts (ranging from 0 to
549 200). We attributed this discrepancy to the lower penetrance and higher background
550 frequency of *RAG1* variants in recessive inheritance, whereby reference studies may
551 underreport the true national incidence. Overall, we argued that agreement within
552 an order of magnitude was tolerable given the inherent uncertainties from reference
553 studies arising from variable penetrance and allele frequencies.

554 3.4 Genetic constraint in high-impact protein networks

555 We next examined genetic constraint in high-impact protein networks across the whole
556 IEI gene set of over 500 known disease-gene phenotypes (1). By integrating ClinVar
557 variant classification scores with PPI data, we quantified the pathogenic burden per
558 gene and assessed its relationship with network connectivity and genetic constraint
559 (7; 16).

560 **3.4.1 Score-positive-total within IEI PPI network**

561 The ClinVar classifications reported in **Figure 1** were scaled -5 to +5 based on their
 562 pathogenicity. We were interested in positive (potentially damaging) but not negative
 563 (benign) scoring variants, which are statistically incidental in this analysis. We tallied
 564 gene-level positive scores to give the score positive total metric. **Figure 3 (A)** shows
 565 the PPI network of disease-associated genes, where node size and colour encode the
 566 score positive total (log-transformed). The top 15 genes with the highest total prior
 567 probabilities of being observed with disease are labelled (as per **Figure 1**).

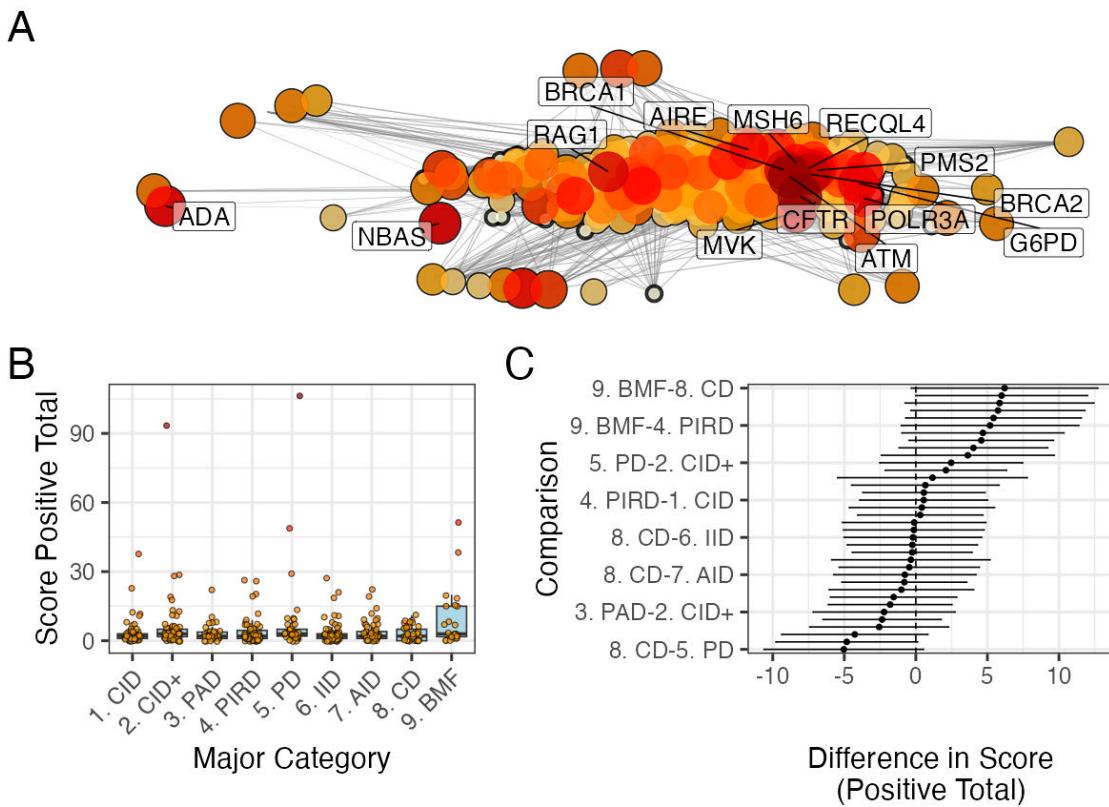


Figure 3: **PPI network and score positive total ClinVar significance variants.** (A) PPI network of disease-associated genes. Node size and colour represent the log-transformed score positive total, the top 15 genes/proteins with the highest probability of being observed in disease are labelled. (B) 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.

568 **3.4.2 Association analysis of score-positive-total across IEI categories**

569 We checked for any statistical enrichment in score positive totals, which represents
 570 the expected observation of pathogenicity, between the IEI categories. The one-way

571 ANOVA revealed an effect of major disease category on score positive total ($F(8, 500) =$
572 2.82, $p = 0.0046$), indicating that group means were not identical, which we observed
573 in **Figure 3 (B)**. However, despite some apparent differences in median scores across
574 categories (i.e. 9. Bone Marrow Failure (BMF)), the Tukey HSD post hoc compar-
575 isons **Figure 3 (C)** showed that all pairwise differences had 95% confidence intervals
576 overlapping zero, suggesting that individual group differences were not significant.

577 **3.4.3 UMAP embedding of the PPI network**

578 To address the density of the PPI network for the IEI gene panel, we applied Uniform
579 Manifold Approximation and Projection (UMAP) (**Figure 4**). Node sizes reflect
580 interaction degree, a measure of evidence-supported connectivity ([16](#)). We tested
581 for a correlation between interaction degree and score positive total. In **Figure**
582 **4**, gene names with degrees above the 95th percentile are labelled in blue, while
583 the top 15 genes by score positive total are labelled in yellow (as per **Figure 1**).
584 Notably, genes with high pathogenic variant loads segregated from highly connected
585 nodes, suggesting that Loss-of-Function (LOF) in hub genes is selectively constrained,
586 whereas damaging variants in lower-degree genes yield more specific effects. This
587 observation was subsequently tested empirically.

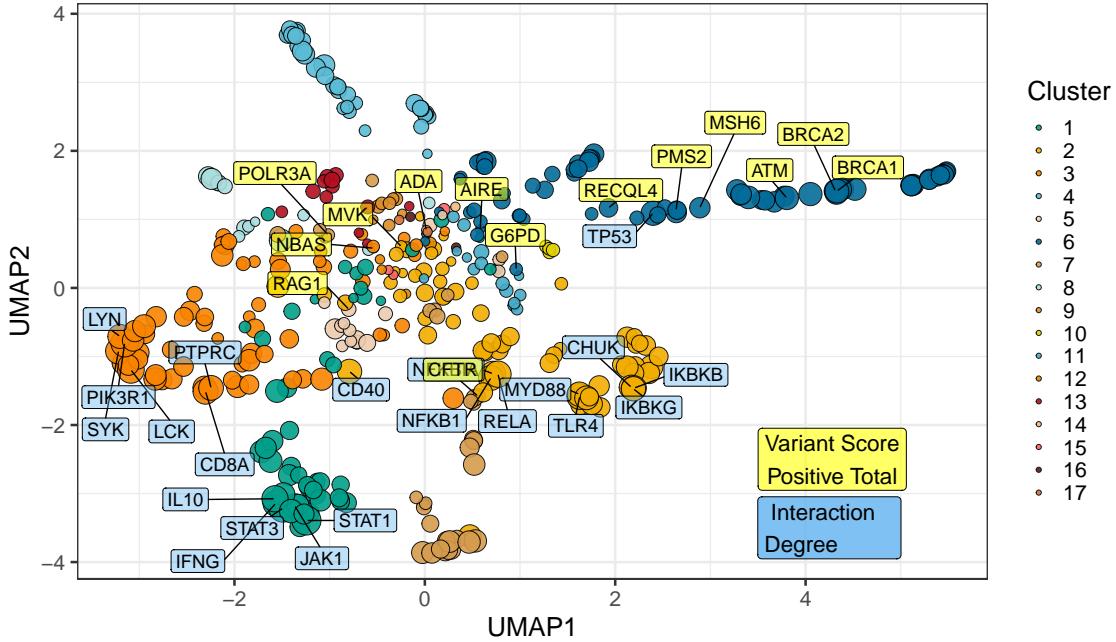


Figure 4: **UMAP embedding of the PPI network (p_umap).** The plot projects the high-dimensional protein-protein interaction network into two dimensions, with nodes coloured by cluster and sized by interaction degree. Blue labels indicate hub genes (degree above the 95th percentile) and yellow labels 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.

588 **3.4.4 Hierarchical clustering of enrichment scores for major disease cate-**
 589 **gories**

590 **Figure S6** presents a heatmap of standardised residuals for major disease categories
 591 across network clusters, as per **Figure 4**. A dendrogram clusters similar disease cate-
 592 gories, while the accompanying bar plot displays the maximum absolute standardised
 593 residual for each category. Notably, (8) Complement Deficiencies (CD) shows the
 594 highest maximum enrichment, followed by (9) BMF. While all maximum values
 595 exceed 2, the threshold for significance, this likely reflects the presence of protein
 596 clusters with strong damaging variant scores rather than uniform significance across
 597 all categories (i.e. genes from cluster 4 in 8 CD).

598 **3.4.5 PPI connectivity, LOEUF constraint and enriched network cluster**
 599 **analysis**

600 Based on the preliminary insight from **Figure S6**, we evaluated the relationship
 601 between network connectivity (PPI degree) and LOEUF constraint (LOEUF upper rank)
 602 Karczewski et al. (7) using Spearman's rank correlation. Overall, there was a weak

603 but significant negative correlation ($\rho = -0.181$, $p = 0.00024$) at the global scale,
 604 indicating that highly connected genes tend to be more constrained. A supplementary
 605 analysis (**Figure S7**) did not reveal distinct visual associations between network
 606 clusters and constraint metrics, likely due to the high network density. However
 607 once stratified by gene clusters, the natural biological scenario based on quantitative
 608 PPI evidence (16), some groups showed strong correlations; for instance, cluster 2
 609 ($\rho = -0.375$, $p = 0.000994$) and cluster 4 ($\rho = -0.800$, $p < 0.000001$), while others did
 610 not. This indicated that shared mechanisms within pathway clusters may underpin
 611 genetic constraints, particularly for LOF intolerance. We observe that the score
 612 positive total metric effectively summarises the aggregate pathogenic burden across
 613 IEI genes, serving as a robust indicator of genetic constraint and highlighting those
 614 with elevated disease relevance.

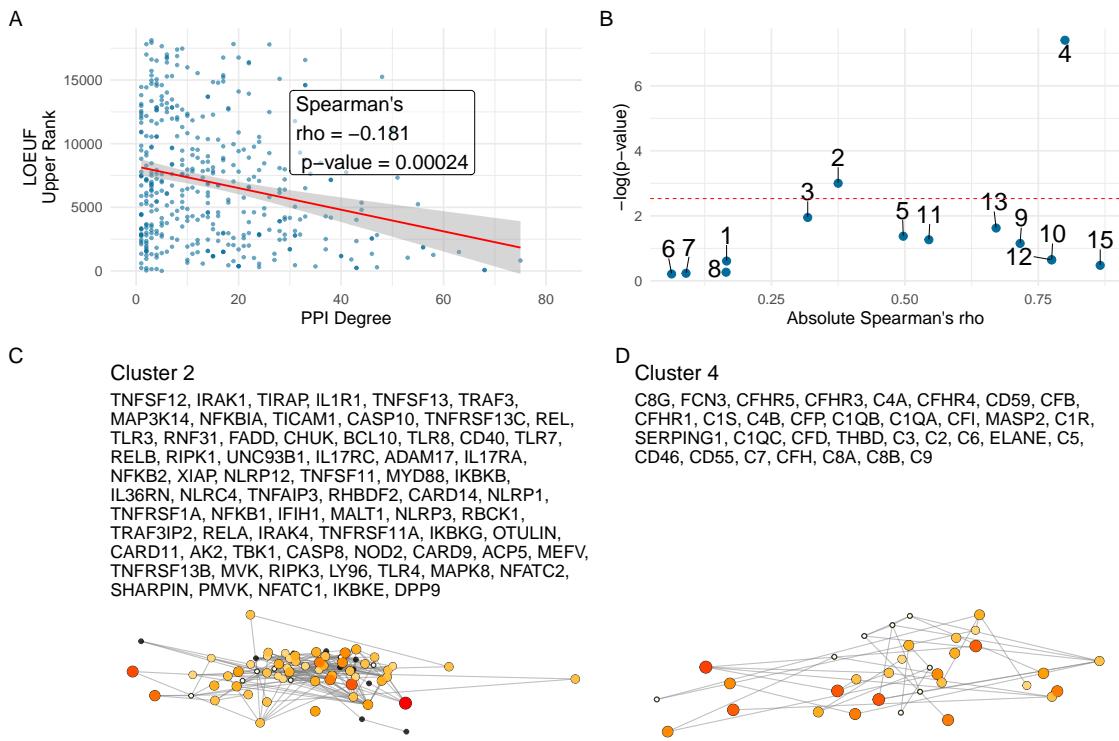


Figure 5: **Correlation between PPI degree and LOEUF upper rank.** (A) Ananlysis across all genes revealed a weak, significant negative correlation between PPI degree and LOEUF upper rank. (B) The cluster-wise analysis showed that clusters 2 and 4 exhibited moderate to strong correlations, while other 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.

615 **Figure 5 (C, D)** shows the re-plotted PPI networks for clusters with significant
 616 correlations between PPI degree and LOEUF upper rank. In these networks, node
 617 size is scaled by a normalised variant score, while node colour reflects the variant
 618 score according to a predefined palette.

619 **3.5 New insight from functional enrichment**

620 To interpret the functional relevance of our prioritised IEI gene sets with the highest
621 load of damaging variants (i.e. clusters 2 and 4 in **Figure 5**), we performed func-
622 tional enrichment analysis for known disease associations using MsigDB with FUMA
623 (i.e. GWAScatalog and Immunologic Signatures) (24). Composite enrichment pro-
624 files (**Figure S8**) reveal that our enriched PPI clusters were associated with distinct
625 disease-related phenotypes, providing functional insights beyond traditional IUIS IEI
626 groupings (1). The gene expression profiles shown in **Figure S9** (GTEx v8 54 tissue
627 types) offer the tissue-specific context for these associations. Together, these results
628 enable the annotation of IEI gene sets with established disease phenotypes, supporting
629 a data-driven classification of IEI.

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

639 **3.6 Genome-wide gene distribution and locus-specific variant
640 occurrence**

641 **Figure 6 (A)** shows a genome-wide karyoplot of all IEI panel genes across GRCh38,
642 with colour-coding based on MOI. Figures (B) and (C) display zoomed-in locus plots
643 for *NFKB1* and *CFTR*, respectively. In **Figure 6 (B)**, the probability of observing
644 variants with known classifications is high only for variants such as p.Ala475Gly,
645 which are considered benign in the AD *NFKB1* gene that is intolerant to LOF. In
646 **Figure 6 (C)**, high probabilities of observing patients with pathogenic variants in
647 *CFTR* are evident, reproducing this well-established phenomenon. Furthermore, the
648 analysis of Linkage Disequilibrium (LD) using R^2 shows that high LD regions can be
649 modelled effectively, allowing independent variant signals to be distinguished.

650 **3.7 Novel PID classifications derived from genetic PPI and
651 clinical features**

652 We recategorised 315 immunophenotypic features from the original IUIS IEI annota-
653 tions, reducing detailed descriptions (e.g. “decreased cd8, normal or decreased cd4”),
654 first to minimal labels (e.g.“low”), and second to binary outcomes (normal vs. not-
655 normal) for T cells, B cells, neutrophils, and immunoglobulins (**Figure 7**). These

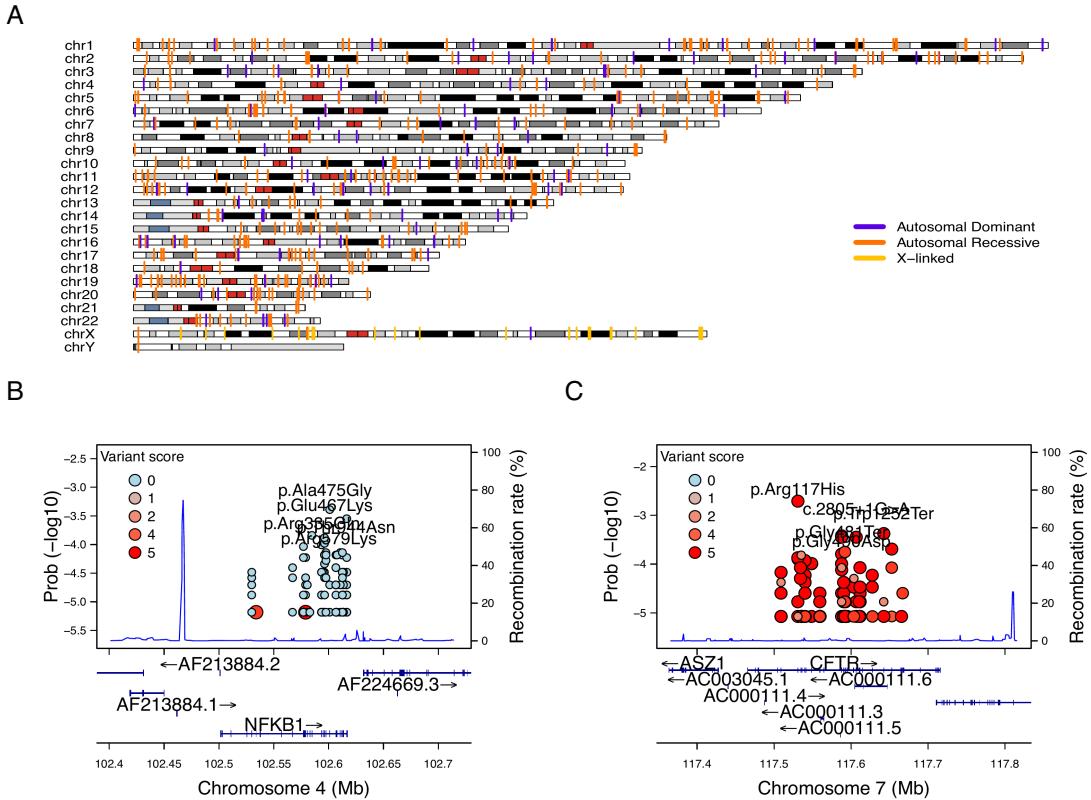


Figure 6: Genome-wide IEI, variant occurrence probability and LD by R^2 .
(A) Genome-wide karyoplot of all IEI panel genes mapped to GRCh38, with colours indicating MOI. (B) Zoomed-in locus plot for *NFKB1* showing variant observation probabilities; only benign variants such exhibit high probabilities 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.

simplified profiles were integrated with PPI network clustering from STRINGdb to refine PID gene groupings. Chi-square analyses confirmed significant associations between specific clinical abnormalities and PPI clusters (**Figure ??**). A decision tree classifier, with hyperparameters optimised via 5-fold cross validation, demonstrated high sensitivity and specificity, as shown in the confusion matrices and variable importance metrics (**Figure S11**). The resulting novel PID classifications, illustrated by the decision tree and gene group distributions (**Figure 10**), provide a more coherent and data-driven framework for categorising PID genes.

664 **3.8 Novel PID classifications derived from genetic PPI and**
665 **clinical features**

666 We recategorised 315 immunophenotypic features from the original IUIS IEI annotations,
667 reducing detailed descriptions (e.g. “decreased CD8, normal or decreased CD4”) to minimal labels (e.g. “low”) and then binarising them (normal vs. not-normal) for T
668 cells, B cells, Immunoglobulin (Ig) and neutrophils (**Figure 7**). These simplified profiles
669 were mapped onto STRINGdb PPI clusters, revealing non-random distributions
670 ($\chi^2 < 1e-13$; **Figure 8**), indicating that network context captures key immunopheno-
671 notypic variation.

673 We next compared four classifiers under 5-fold cross-validation to determine which
674 features predicted PPI clustering. As shown in **Figure 9**, the fully combined model
675 achieved the highest accuracy among the four: (i) phenotypes only (33 %) (i.e. T
676 cell, B cell, Ig, Neutrophil); (ii) phenotypes + IUIS major category (50 %) (e.g. CID.
677 See **Box 2.1** for more); (iii) IUIS major + subcategory only (59 %) (e.g. CID, T-B+
678 SCID); and (iv) phenotypes + IUIS major + subcategory (61 %). This demonstrated
679 that incorporating both traditional IUIS classifications and core immunophenotypic
680 markers into the PPI-based framework yielded the most robust discrimination of PID
681 gene clusters. Variable importance analysis highlighted abnormality status for Ig and
682 T cells were among the top ten features in addition to the other IUIS major and sub
683 categories. Per-class specificity remained uniform across the classes while sensitivity
684 dropped.

685 The PPI and immunophenotype model yielded 17 data-driven PID groups, whereas
686 incorporating the full complement of IUIS categories expanded this to 33 groups. For
687 clarity, we only demonstrate the decision tree from the smaller 17-group model in
688 **Figure 10**. Each terminal node is annotated by its predominant immunophenotypic
689 signature (for example, “group 65 with abnormal T cell and B cell features”), and the
690 full resulting gene counts per 33 class are plotted in **Figure 10**. Although, less user-
691 friendly, this data-driven taxonomy both aligns with and refines traditional IUIS IEI
692 classifications to provide a scaffold for large-scale computational analyses. Because
693 this framework is fully reproducible, alternative PPI embeddings or incorporate addi-
694 tional molecular annotations can readily swapped to continue building on these PID
695 classification schemes.

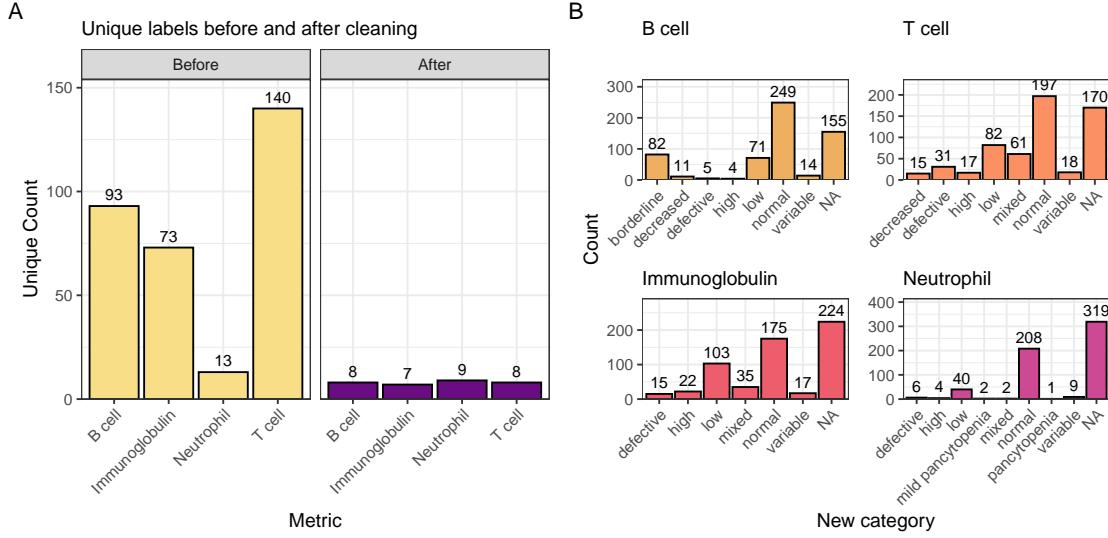


Figure 7: Distribution of immunophenotypic features before and after recategorisation. The original IUIS IEI descriptions contain information such as T cell-related “decreased cd8, normal or decreased cd4 cells” which we recategorise as “low”. The bar plot shows the count of unique labels for each status (normal, not_normal) across the T cell, B cell, Ig, and Neutrophil features.

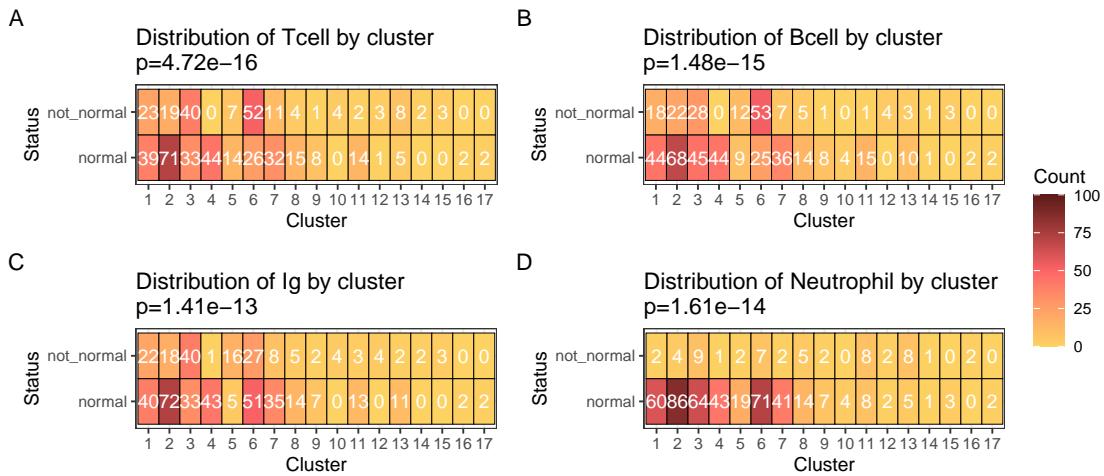


Figure 8: Heatmaps of clinical feature distributions by PPI cluster. The heatmaps display the count of observations for abnormality of each clinical feature (A) T cell, (B) B cell, (C) Immunoglobulin, (D) Neutrophil, in relation to the PPI clusters, with p-values from chi-square tests annotated in the titles.

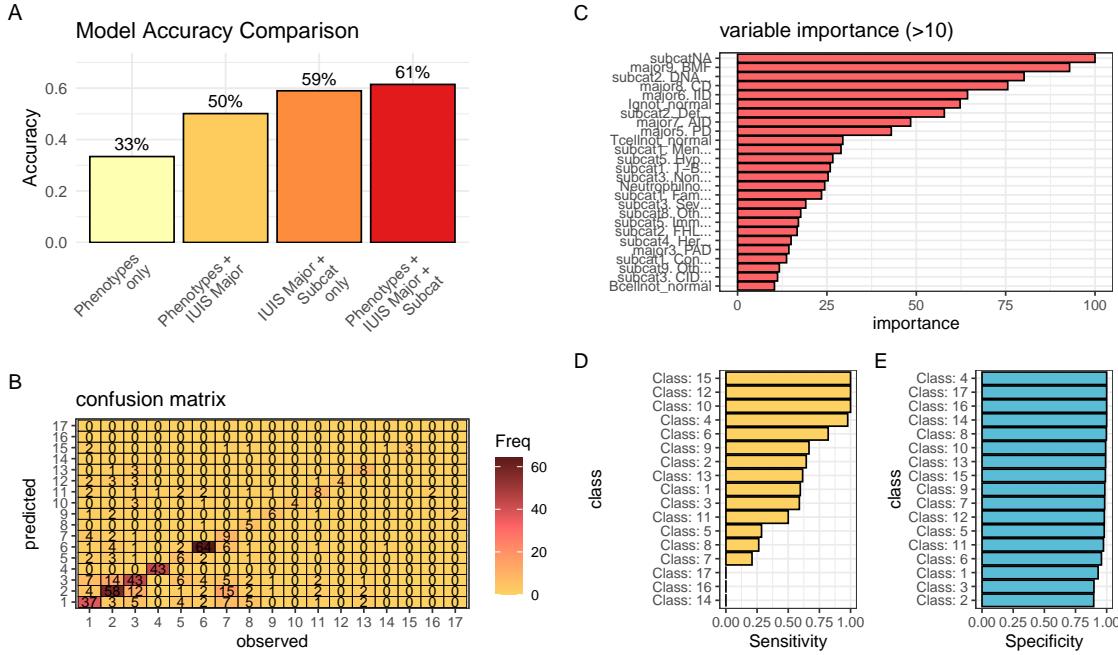


Figure 9: Performance comparison of PID classifiers. Classification predicting PPI cluster membership from IUIS major category, subcategory, and immunological features. (A) Overall accuracy for four rpart models used to predict PPI clustering. The combined model achieves 61.4 % accuracy, exceeding all simpler approaches. Nodes were split to minimize Gini impurity, pruned by cost-complexity ($cp = 0.001$), and validated via 5-fold cross-validation. (B-E) The summary statistics from the top model are detailed.

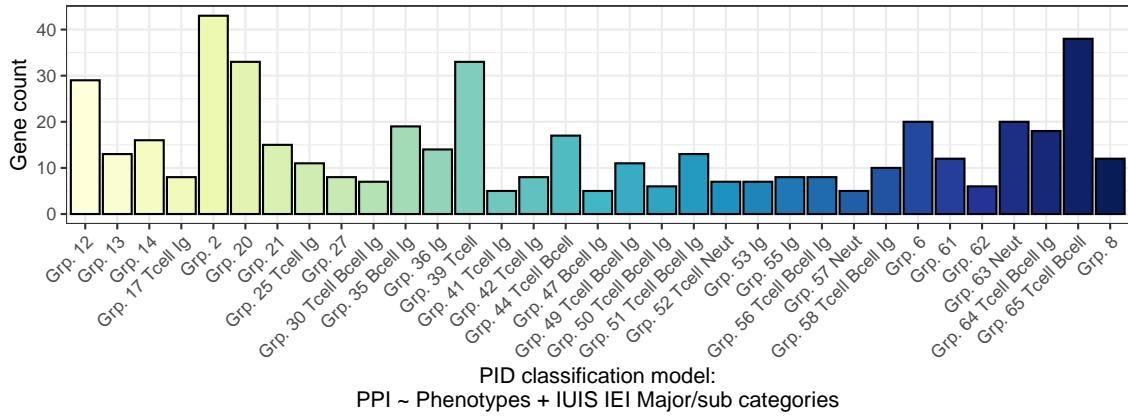
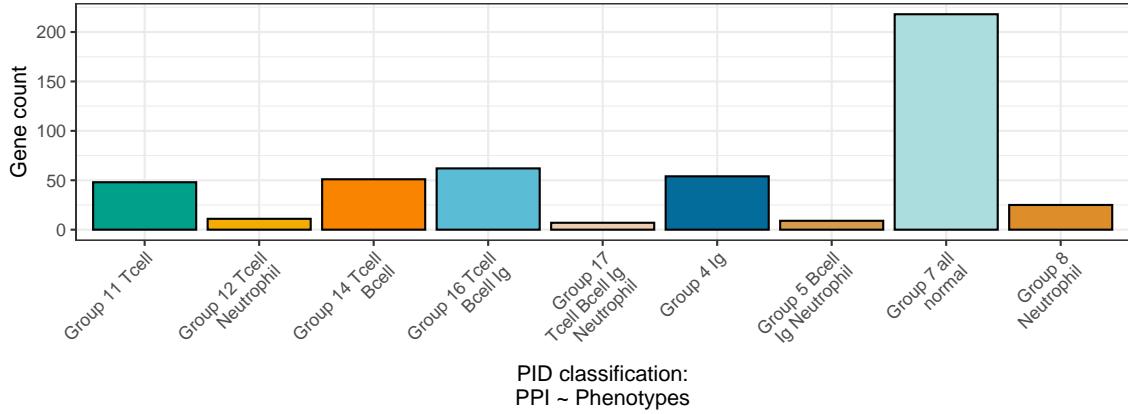
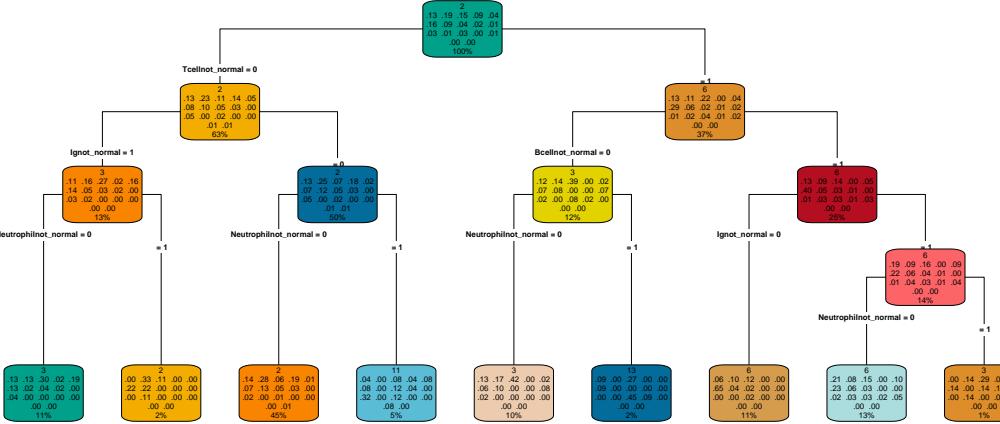


Figure 10: **Fine-tuned model for PID classification.** (Top) In each terminal node, the top block indicates the number of genes in the node; the middle block shows the fitted class probabilities (which sum to 1); and the bottom block displays the percentage of the total sample in that node. These metrics summarise the model’s assignment based on immunophenotypic and PPI features. (Middle) Bar plot presenting the distribution of novel PID classifications, where group labels denote the predominant abnormal clinical feature(s) (e.g. T cell, B cell, Ig, Neutrophil) characterising each group. (Bottom) The complete model including the traditional IUIS IEI categories.

696 3.8.1 Integration of Variant Probabilities into IEI Genetics Data

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

704 **Figure 11** shows the interface summarising integrated variant data. Server-side
 705 pre-calculation of summary statistics minimises browser load, while clinical signifi-
 706 cance is converted to numerical metrics. Key quantiles (min, Q1, median, Q3, max)
 707 for each gene are rendered as sparkline box plots, and dynamic URLs link table entries
 708 to external databases (e.g. ClinVar, Online Mendelian Inheritance in Man (OMIM),
 709 AlphaFold).

The screenshot displays a table titled "Viewer Zoom" with a search bar at the top. The table has 13 columns: Major category, Subcategory, Disease, Genetic defect, Inheritance, Gene score, Prior prob of pathogenicity, ClinVar SNV classification, ClinVar all variant reports, OMIM, Alpha Missense / Uniprot ID, HPO combined, and HPO term. The rows represent different genetic conditions, such as SCID, IL2RG deficiency, IL7Ra deficiency, ITPKB deficiency, JAK3 deficiency, and LAT deficiency. Each row contains a sparkline box plot for the "Prior prob of pathogenicity" column, which visualizes the distribution of observed pathogenicity scores. The "ClinVar all variant reports" column contains dynamic URLs. The "OMIM" column lists OMIM IDs (e.g., P20963, P04234, P07766, P31146, P31185). The "Alpha Missense / Uniprot ID" column lists UniProt IDs (e.g., 186780, 186790, 186830, 605000, 308380, 600173). The "HPO combined" and "HPO term" columns provide information about the clinical phenotype, including terms like "Abnormalit" and "T lymphoc". A footer at the bottom indicates "1-10 of 591 rows" and "Show 10".

Figure 11: **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.

710 3.9 Probability of observing AlphaMissense pathogenicity

711 AlphaMissense provides pathogenicity scores for all possible amino acid substitutions;
 712 however, our results in **Figure 12** show that the most probable observations in pa-
 713 tients occur predominantly for benign or unknown variants. This finding places the
 714 likelihood of disease-associated substitutions into perspective and offers a data-driven
 715 foundation for future improvements in variant prediction. The values in **Figure 12**
 716 (**A**) can be directly compared to **Figure 1 (D)** to view the distribution of classifi-
 717 cations. A Kruskal-Wallis test was used to compare the observed disease probability

718 across clinical classification groups and no significant differences were detected. In
 719 general, most variants in patients are classified as benign or unknown, indicating
 720 limited discriminative power in the current classification, such that pathogenicity
 721 prediction does not infer observation prediction (**Figure S12**).

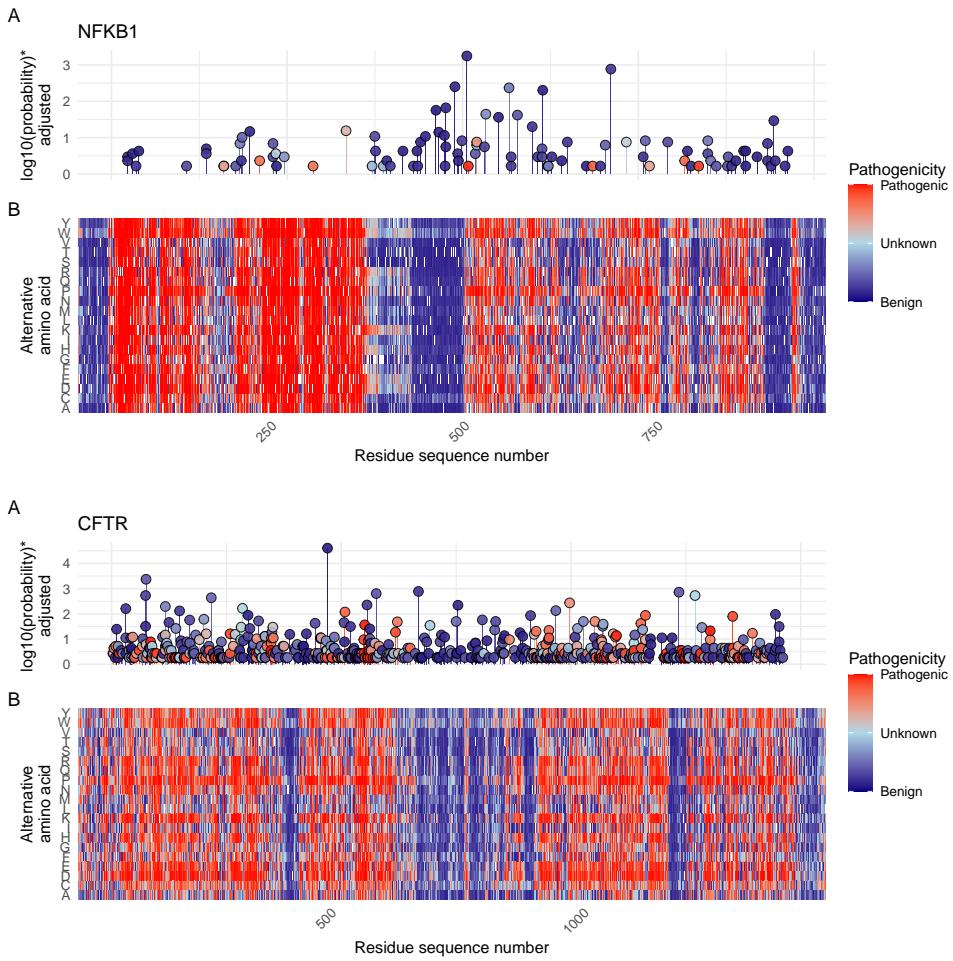


Figure 12: **(A) Probabilities of observing a patient with (B) AlphaMissense-derived pathogenicity scores.** Although AlphaMissense provides scores for every possible amino acid substitution, the most frequently observed variants in patients tend to be classified as benign or of unknown significance. 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.

722 4 Discussion

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

730 Our approach yielded two key results. First, our detailed, per-variant pre-calculated
731 results provide prior probabilities of observing disease-associated variants across all
732 MOI for any gene-disease combination. Second, the score positive total metric effec-
733 tively summarises the aggregate pathogenic burden across genes, serving as a robust
734 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:

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

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

737 Our framework addresses this challenge by combining variant classifications, pop-
738 ulation allele frequencies, and curated gene-disease associations. While imperfect on
739 an individual level, these sources exhibit predictable aggregate behaviour, supported
740 by James-Stein estimation principles (28). Curated gene-disease associations help
741 identify genes that explainable for most disease cases, allowing us to approximate
742 $P(G | D)$ close to one. In this way, we obtain robust estimates of $P(G)$ (the fre-
743 quency of disease-associated genotypes), even when exact values of penetrance and
744 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.

The results presented in the *NFKB1* case study demonstrate how our Bayesian framework quantifies the final posterior probability that a damaging causal variant is

783 present, by integrating both observed variant calls and unobserved, but potentially
784 pathogenic, sites into a single conclusion. This result applies across all disease-genes.
785 The model partitions posterior probability between a true positive (`p.Ser237Ter`) and
786 a likely causal splice-site variant that was not captured by sequencing. Despite the
787 latter being unobserved, the analysis attributes 48 % of total pathogenic probability
788 to this missing position—an insight that would be unavailable under conventional
789 approaches limited to called variants.

790 This framework enables calibrated probabilistic conclusions at two levels: per-
791 variant and per-gene. Rather than issuing a binary classification, the method produces
792 interpretable posterior intervals that reflect both sequencing completeness and prior
793 pathogenic evidence. For clinicians, this means that residual uncertainty—arising
794 from false negatives or incomplete genomic coverage—is made visible, quantifiable,
795 and actionable.

796 A key output is the ability to generate structured variant reports (e.g. Table ??) that
797 transparently communicate which alleles are supported, which are excluded, and
798 which remain unassessed but potentially disease-causing. These summaries facilitate
799 informed consent, prioritise follow-up sequencing (e.g. long-read or targeted assays),
800 and provide a reproducible, evidence-weighted foundation for diagnostic decision-
801 making.

802 More broadly, the same approach generalises to genome-wide applications. The
803 precomputed priors for over 50,000 variants across 557 immune-related genes now
804 enable quantitative confidence statements for any patient–phenotype combination.
805 Because the model explicitly incorporates allele frequency, clinical classification,
806 inheritance mode, and sequence quality, it provides a scalable template for delivering
807 uncertainty-aware interpretation in clinical genomics.

808 5 Conclusion

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

812 Acknowledgements

813 We acknowledge Genomics England for providing public access to the PanelApp data.
814 The use of data from Genomics England panelapp was licensed under the Apache
815 License 2.0. The use of data from UniProt was licensed under Creative Commons
816 Attribution 4.0 International (CC BY 4.0). ClinVar asks its users who distribute or
817 copy data to provide attribution to them as a data source in publications and websites
818 (13). dbNSFP version 4.4a is licensed under the Creative Commons Attribution-
819 NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0); while we cite

820 this dataset as used our research publication, it is not used for the final version which
821 instead used ClinVar and gnomAD directly. GnomAD is licensed under Creative
822 Commons Zero Public Domain Dedication (CC0 1.0 Universal). GnomAD request
823 that usages cites the gnomAD flagship paper (7) and any online resources that include
824 the data set provide a link to the browser, and note that tool includes data from the
825 gnomAD v4.1 release. AlphaMissense asks to cite Cheng et al. (12) for usage in
826 research, with data available from Cheng et al. (26).

827 Competing interest

828 We declare no competing interest.

829 References

- 830 [1] Stuart G. Tangye, Waleed Al-Herz, Aziz Bousfiha, Charlotte Cunningham-
831 Rundles, Jose Luis Franco, Steven M. Holland, Christoph Klein, Tomohiro Morio,
832 Eric Oksenhendler, Capucine Picard, Anne Puel, Jennifer Puck, Mikko R. J.
833 Seppänen, Raz Somech, Helen C. Su, Kathleen E. Sullivan, Troy R. Torger-
834 son, and Isabelle Meyts. Human Inborn Errors of Immunity: 2022 Update
835 on the Classification from the International Union of Immunological Societies
836 Expert Committee. *Journal of Clinical Immunology*, 42(7):1473–1507, October
837 2022. ISSN 0271-9142, 1573-2592. doi: 10.1007/s10875-022-01289-3. URL
838 <https://link.springer.com/10.1007/s10875-022-01289-3>.
- 839 [2] Dylan Lawless. PanelAppRex aggregates disease gene panels and facilitates
840 sophisticated search. March 2025. doi: 10.1101/2025.03.20.25324319. URL
841 <http://medrxiv.org/lookup/doi/10.1101/2025.03.20.25324319>.
- 842 [3] Antonio Rueda Martin, Eleanor Williams, Rebecca E. Foulger, Sarah Leigh,
843 Louise C. Daugherty, Olivia Niblock, Ivone U. S. Leong, Katherine R. Smith,
844 Oleg Gerasimenko, Eik Haraldsdottir, Ellen Thomas, Richard H. Scott, Emma
845 Baple, Arianna Tucci, Helen Brittain, Anna De Burca, Kristina Ibañez, Dalia
846 Kasperaviciute, Damian Smedley, Mark Caulfield, Augusto Rendon, and Ellen M.
847 McDonagh. PanelApp crowdsources expert knowledge to establish consensus
848 diagnostic gene panels. *Nature Genetics*, 51(11):1560–1565, November 2019.
849 ISSN 1061-4036, 1546-1718. doi: 10.1038/s41588-019-0528-2. URL <https://www.nature.com/articles/s41588-019-0528-2>.
- 850 [4] Oliver Mayo. A Century of Hardy–Weinberg Equilibrium. *Twin Research
851 and Human Genetics*, 11(3):249–256, June 2008. ISSN 1832-4274, 1839-
852 2628. doi: 10.1375/twin.11.3.249. URL https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal_article.

- 855 [5] Nikita Abramovs, Andrew Brass, and May Tassabehji. Hardy-Weinberg Equi-
856 librium in the Large Scale Genomic Sequencing Era. *Frontiers in Genetics*,
857 11:210, March 2020. ISSN 1664-8021. doi: 10.3389/fgene.2020.00210. URL
858 <https://www.frontiersin.org/article/10.3389/fgene.2020.00210/full>.
- 859 [6] Johannes Zschocke, Peter H. Byers, and Andrew O. M. Wilkie. Mendelian
860 inheritance revisited: dominance and recessiveness in medical genetics. *Nature
861 Reviews Genetics*, 24(7):442–463, July 2023. ISSN 1471-0056, 1471-0064.
862 doi: 10.1038/s41576-023-00574-0. URL <https://www.nature.com/articles/s41576-023-00574-0>.
- 864 [7] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings,
865 Jessica Alföldi, Qingbo Wang, Ryan L Collins, Kristen M Laricchia, Andrea
866 Ganna, Daniel P Birnbaum, et al. The mutational constraint spectrum quantified
867 from variation in 141,456 humans. *Nature*, 581(7809):434–443, 2020.
- 868 [8] Sarah L. Bick, Aparna Nathan, Hannah Park, Robert C. Green, Monica H. Wo-
869 jcik, and Nina B. Gold. Estimating the sensitivity of genomic newborn screen-
870 ing for treatable inherited metabolic disorders. *Genetics in Medicine*, 27(1):
871 101284, January 2025. ISSN 10983600. doi: 10.1016/j.gim.2024.101284. URL
872 <https://linkinghub.elsevier.com/retrieve/pii/S1098360024002181>.
- 873 [9] Benjamin D. Evans, Piotr Słowiński, Andrew T. Hattersley, Samuel E. Jones,
874 Seth Sharp, Robert A. Kimmitt, Michael N. Weedon, Richard A. Oram,
875 Krasimira Tsaneva-Atanasova, and Nicholas J. Thomas. Estimating disease
876 prevalence in large datasets using genetic risk scores. *Nature Communications*,
877 12(1):6441, November 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-26501-7.
878 URL <https://www.nature.com/articles/s41467-021-26501-7>.
- 879 [10] William B. Hannah, Mitchell L. Drumm, Keith Nykamp, Tiziano Prampano,
880 Robert D. Steiner, and Steven J. Schrödi. Using genomic databases to de-
881 termine the frequency and population-based heterogeneity of autosomal reces-
882 sive conditions. *Genetics in Medicine Open*, 2:101881, 2024. ISSN 29497744.
883 doi: 10.1016/j.gimo.2024.101881. URL <https://linkinghub.elsevier.com/retrieve/pii/S2949774424010276>.
- 885 [11] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov,
886 Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek,
887 Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J.
888 Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh
889 Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy,
890 Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamás Berghammer,
891 Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray
892 Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate pro-
893 tein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August
894 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL
895 <https://www.nature.com/articles/s41586-021-03819-2>.

- 896 [12] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-
897 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias
898 Sergeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-
899 abis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense vari-
900 ant effect prediction with AlphaMissense. *Science*, 381(6664):eadg7492, Septem-
901 ber 2023. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.adg7492. URL
902 <https://www.science.org/doi/10.1126/science.adg7492>.
- 903 [13] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen Chao,
904 Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Wonhee
905 Jang, Karen Karapetyan, Kenneth Katz, Chunlei Liu, Zenith Maddipatla, Adri-
906 ana Malheiro, Kurt McDaniel, Michael Ovetsky, George Riley, George Zhou,
907 J Bradley Holmes, Brandi L Kattman, and Donna R Maglott. ClinVar: im-
908 proving access to variant interpretations and supporting evidence. *Nucleic Acids
909 Research*, 46(D1):D1062–D1067, January 2018. ISSN 0305-1048, 1362-4962. doi:
910 10.1093/nar/gkx1153. URL [http://academic.oup.com/nar/article/46/D1/
911 D1062/4641904](http://academic.oup.com/nar/article/46/D1/D1062/4641904).
- 912 [14] The UniProt Consortium, Alex Bateman, Maria-Jesus Martin, Sandra Orchard,
913 Michele Magrane, Aduragbemi Adesina, Shadab Ahmad, Bowler-Barnett, and
914 Others. UniProt: the Universal Protein Knowledgebase in 2025. *Nucleic
915 Acids Research*, 53(D1):D609–D617, January 2025. ISSN 0305-1048, 1362-4962.
916 doi: 10.1093/nar/gkae1010. URL [https://academic.oup.com/nar/article/
917 53/D1/D609/7902999](https://academic.oup.com/nar/article/53/D1/D609/7902999).
- 918 [15] Xiaoming Liu, Chang Li, Chengcheng Mou, Yibo Dong, and Yicheng Tu.
919 dbNSFP v4: a comprehensive database of transcript-specific functional pre-
920 dictions and annotations for human nonsynonymous and splice-site SNVs.
921 *Genome Medicine*, 12(1):103, December 2020. ISSN 1756-994X. doi: 10.
922 1186/s13073-020-00803-9. URL [https://genomemedicine.biomedcentral.
923 com/articles/10.1186/s13073-020-00803-9](https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-020-00803-9).
- 924 [16] Damian Szklarczyk, Katerina Nastou, Mikaela Koutrouli, Rebecca Kirsch, Far-
925 rokh Mehryary, Radja Hachilif, Dewei Hu, Matteo E Peluso, Qingyao Huang,
926 Tao Fang, et al. The string database in 2025: protein networks with directional-
927 ity of regulation. *Nucleic Acids Research*, 53(D1):D730–D737, 2025.
- 928 [17] Paul Tijnenburg, Hana Lango Allen, Siobhan O. Burns, Daniel Greene,
929 Machiel H. Jansen, and Others. Loss-of-function nuclear factor B subunit
930 1 (NFKB1) variants are the most common monogenic cause of common vari-
931 able immunodeficiency in Europeans. *Journal of Allergy and Clinical Im-
932 munology*, 142(4):1285–1296, October 2018. ISSN 00916749. doi: 10.1016/
933 j.jaci.2018.01.039. URL [https://linkinghub.elsevier.com/retrieve/pii/
934 S0091674918302860](https://linkinghub.elsevier.com/retrieve/pii/S0091674918302860).
- 935 [18] WHO Scientific Group et al. Primary immunodeficiency diseases: report of a
936 who scientific group. *Clin. Exp. Immunol.*, 109(1):1–28, 1997.

- 937 [19] Charlotte Cunningham-Rundles and Carol Bodian. Common variable immunodeficiency: clinical and immunological features of 248 patients. *Clinical immunology*, 92(1):34–48, 1999.
- 940 [20] Eric Oksenhendler, Laurence Gérard, Claire Fieschi, Marion Malphettes, Gael
941 Mouillot, Roland Jaussaud, Jean-François Viallard, Martine Gardembas, Lionel
942 Galicier, Nicolas Schleinitz, et al. Infections in 252 patients with common variable
943 immunodeficiency. *Clinical Infectious Diseases*, 46(10):1547–1554, 2008.
- 944 [21] Y Naito, F Adams, S Charman, J Duckers, G Davies, and S Clarke. Uk cystic
945 fibrosis registry 2023 annual data report. *London: Cystic Fibrosis Trust*, 2023.
- 946 [22] Carlo Castellani, CFTR2 team, et al. Cftr2: how will it help care? *Paediatric
947 respiratory reviews*, 14:2–5, 2013.
- 948 [23] Hartmut Grasemann and Felix Ratjen. Cystic fibrosis. *New England Journal
949 of Medicine*, 389(18):1693–1707, 2023. doi: 10.1056/NEJMra2216474. URL
950 <https://www.nejm.org/doi/full/10.1056/NEJMra2216474>.
- 951 [24] Kyoko Watanabe, Erdogan Taskesen, Arjen Van Bochoven, and Danielle
952 Posthuma. Functional mapping and annotation of genetic associations with
953 FUMA. *Nature Communications*, 8(1):1826, November 2017. ISSN 2041-1723.
954 doi: 10.1038/s41467-017-01261-5. URL <https://www.nature.com/articles/s41467-017-01261-5>.
- 955 [25] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir,
956 Pablo Tamayo, and Jill P. Mesirov. Molecular signatures database (MSigDB)
957 3.0. *Bioinformatics*, 27(12):1739–1740, June 2011. ISSN 1367-4811, 1367-
958 4803. doi: 10.1093/bioinformatics/btr260. URL <https://academic.oup.com/bioinformatics/article/27/12/1739/257711>.
- 959 [26] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-
960 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias
961 Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-
962 abis, Pushmeet Kohli, and Žiga Avsec. Predictions for alphanonsense, September
963 2023. URL <https://doi.org/10.5281/zenodo.8208688>.
- 964 [27] Dylan Lawless. Variant risk estimate probabilities for ie genes. March 2025. doi:
965 10.5281/zenodo.15111584. URL <https://doi.org/10.5281/zenodo.15111584>.
- 966 [28] Bradley Efron and Carl Morris. Stein’s Estimation Rule and Its Competitors—
967 An Empirical Bayes Approach. *Journal of the American Statistical Association*,
968 68(341):117, March 1973. ISSN 01621459. doi: 10.2307/2284155. URL <https://www.jstor.org/stable/2284155?origin=crossref>.
- 969 [29] Yaowu Liu, Sixing Chen, Zilin Li, Alanna C Morrison, Eric Boerwinkle, and
970 Xihong Lin. Acat: a fast and powerful p value combination method for rare-
971 variant analysis in sequencing studies. *The American Journal of Human Genetics*,
972 104(3):410–421, 2019.

- 976 [30] Xihao Li, Zilin Li, Hufeng Zhou, Sheila M Gaynor, Yaowu Liu, Han Chen, Ryan
977 Sun, Rounak Dey, Donna K Arnett, Stella Aslibekyan, et al. Dynamic incorpora-
978 tion of multiple *in silico* functional annotations empowers rare variant association
979 analysis of large whole-genome sequencing studies at scale. *Nature genetics*, 52
980 (9):969–983, 2020.
- 981 [31] Michael C Wu, Seunggeun Lee, Tianxi Cai, Yun Li, Michael Boehnke, and Xi-
982 hong Lin. Rare-variant association testing for sequencing data with the sequence
983 kernel association test. *The American Journal of Human Genetics*, 89(1):82–93,
984 2011.
- 985 [32] Seunggeun Lee, Mary J Emond, Michael J Bamshad, Kathleen C Barnes, Mark J
986 Rieder, Deborah A Nickerson, David C Christiani, Mark M Wurfel, and Xihong
987 Lin. Optimal unified approach for rare-variant association testing with applica-
988 tion to small-sample case-control whole-exome sequencing studies. *The American
989 Journal of Human Genetics*, 91(2):224–237, 2012.
- 990 [33] Augustine Kong, Gudmar Thorleifsson, Michael L Frigge, Bjarni J Vilhjalmsson,
991 Alexander I Young, Thorgeir E Thorgeirsson, Stefania Benonisdottir, Asmundur
992 Oddsson, Bjarni V Halldorsson, Gisli Masson, et al. The nature of nurture:
993 Effects of parental genotypes. *Science*, 359(6374):424–428, 2018.
- 994 [34] Laurence J Howe, Michel G Nivard, Tim T Morris, Ailin F Hansen, Humaira
995 Rasheed, Yoonsu Cho, Geetha Chittoor, Penelope A Lind, Teemu Palviainen,
996 Matthijs D van der Zee, et al. Within-sibship gwas improve estimates of direct
997 genetic effects. *BioRxiv*, pages 2021–03, 2021.
- 998 [35] Sue Richards, Nazneen Aziz, Sherri Bale, David Bick, Soma Das, Julie Gastier-
999 Foster, Wayne W Grody, Madhuri Hegde, Elaine Lyon, Elaine Spector, et al.
1000 Standards and guidelines for the interpretation of sequence variants: a joint
1001 consensus recommendation of the american college of medical genetics and ge-
1002 nomics and the association for molecular pathology. *Genetics in medicine*, 17
1003 (5):405–423, 2015.
- 1004 [36] Sean V Tavtigian, Steven M Harrison, Kenneth M Boucher, and Leslie G
1005 Biesecker. Fitting a naturally scaled point system to the acmg/amp variant
1006 classification guidelines. *Human mutation*, 41(10):1734–1737, 2020.
- 1007 [37] Quan Li and Kai Wang. Intervar: clinical interpretation of genetic variants by
1008 the 2015 acmg-amp guidelines. *The American Journal of Human Genetics*, 100
1009 (2):267–280, 2017.
- 1010 [38] Brent S Pedersen, Joe M Brown, Harriet Dashnow, Amelia D Wallace, Matt
1011 Velinder, Martin Tristani-Firouzi, Joshua D Schiffman, Tatiana Tvrzik, Rong
1012 Mao, D Hunter Best, et al. Effective variant filtering and expected candidate
1013 variant yield in studies of rare human disease. *NPJ Genomic Medicine*, 6(1):1–8,
1014 2021.

- 1015 [39] Carl A Anderson, Fredrik H Pettersson, Geraldine M Clarke, Lon R Cardon,
1016 Andrew P Morris, and Krina T Zondervan. Data quality control in genetic
1017 case-control association studies. *Nature protocols*, 5(9):1564–1573, 2010. URL
1018 <https://doi.org/10.1038/nprot.2010.116>.
- 1019 [40] David T Miller, Kristy Lee, Noura S Abul-Husn, Laura M Amendola, Kyle Broth-
1020 ers, Wendy K Chung, Michael H Gollob, Adam S Gordon, Steven M Harrison,
1021 Ray E Hershberger, et al. Acmg sf v3. 2 list for reporting of secondary findings
1022 in clinical exome and genome sequencing: a policy statement of the american
1023 college of medical genetics and genomics (acmg). *Genetics in Medicine*, 25(8):
1024 100866, 2023.

1025 **6 Supplemental**

1026 **6.1 Integrating observed true positives and unobserved false**
 1027 **negatives into a single, actionable conclusion**

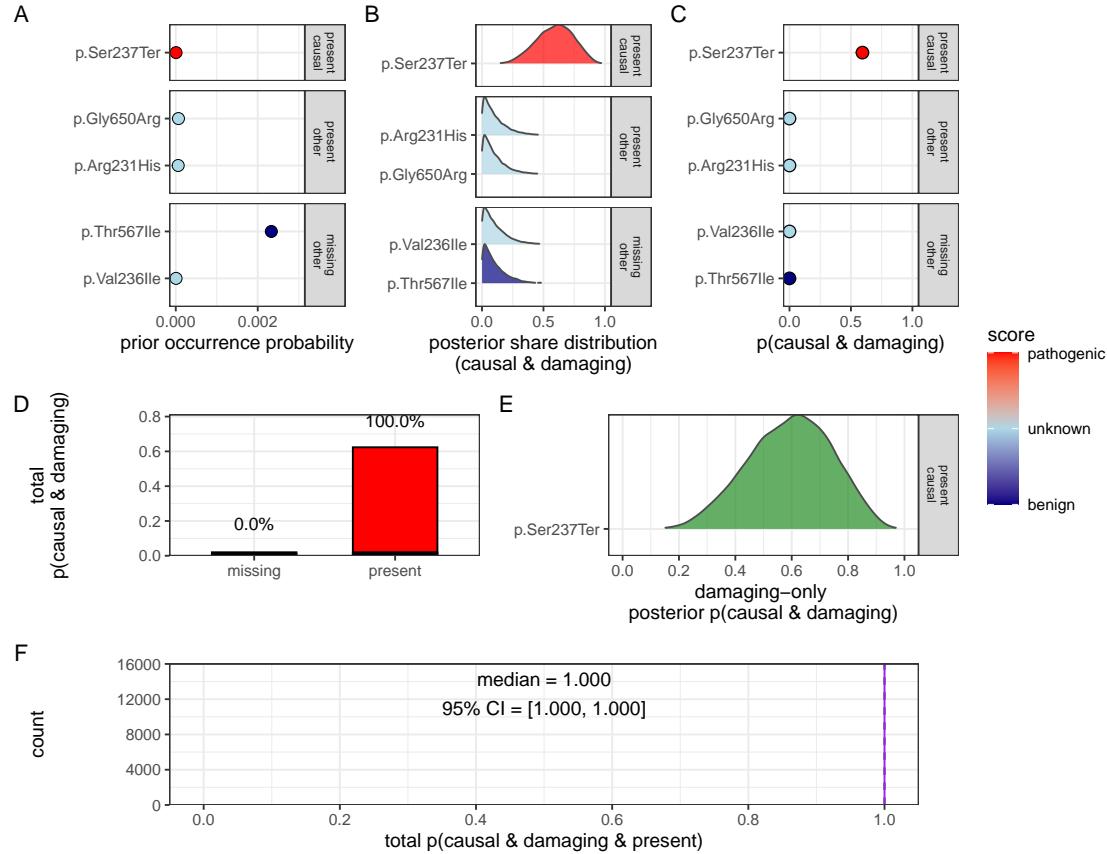


Figure S1: Quantification of gene-level pathogenic attribution for *NFKB1* (scenario 1). Only one known pathogenic variant, p.Ser237Ter, was observed and all previously reported pathogenic positions were successfully sequenced and confirmed as reference (true negatives). Panels (A–F) follow the same structure as scenario 2 described in **Figure 2**, culminating in a gene-level posterior probability of 1 (95 % credible interval: 0.99–1.00), with full support assigned to the observed allele given the available evidence.

1028 **6.2 Validation studies**

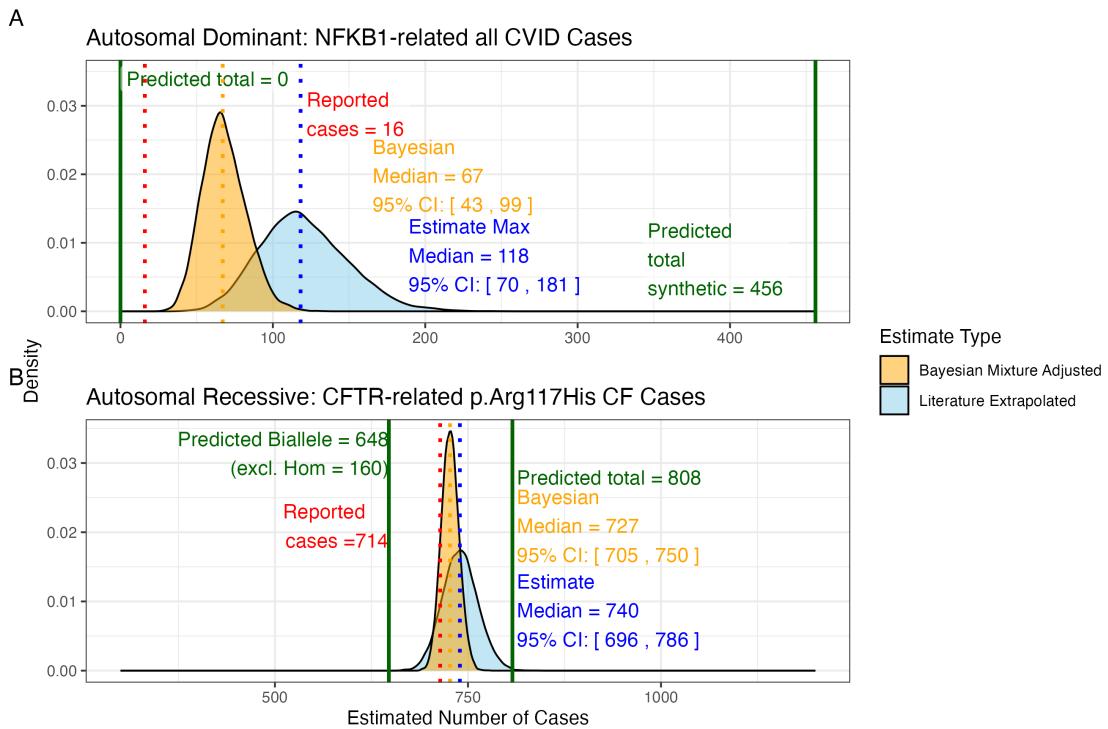


Figure S2: Prior probabilities compared to validation disease cohort metrics.

(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.

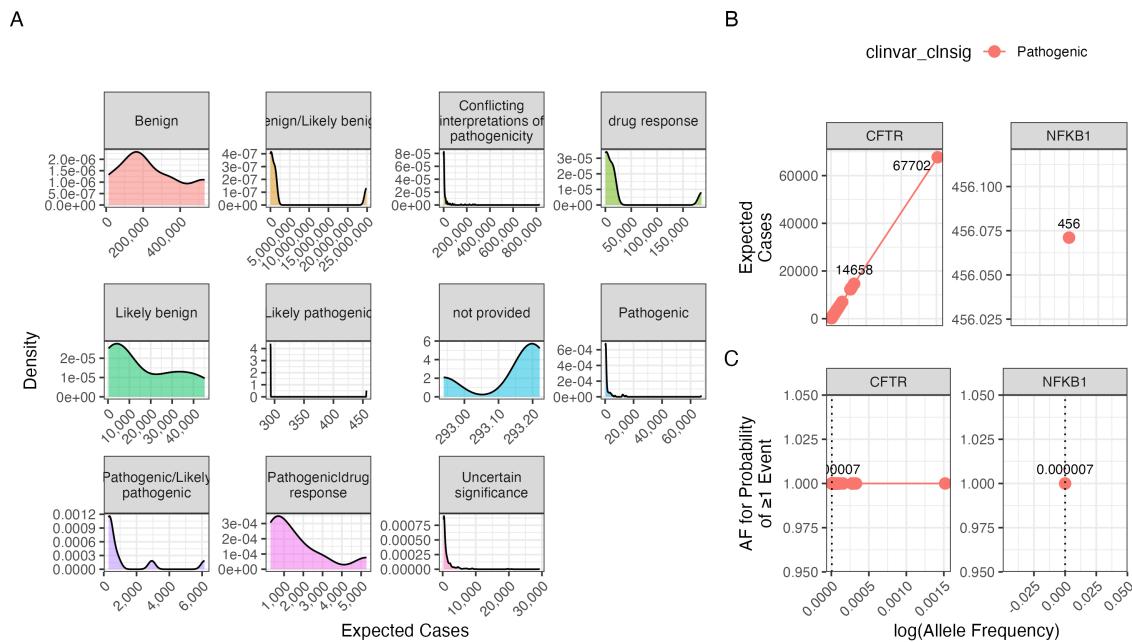


Figure S3: Interpretation of probability of observing a variant classification. The result from the chosen validation genes *CFTR* and *NFKB1* are shown. Case counts are dependant on the population size and phenotype. (A) The density plots of expected observations by ClinVar clinical significance. We 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.

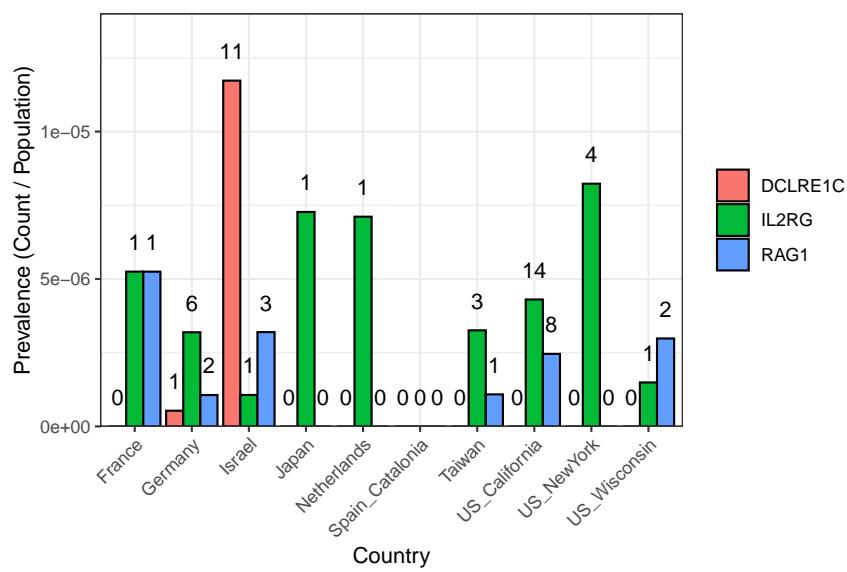


Figure S4: **SCID-specific gene comparison across regions.** The bar plot shows the prevalence of SCID-related cases (count divided by population) for each gene and country (or region), with numbers printed above the bars representing the actual counts in the original cohort (ranging from 0 to 11 per region and gene).

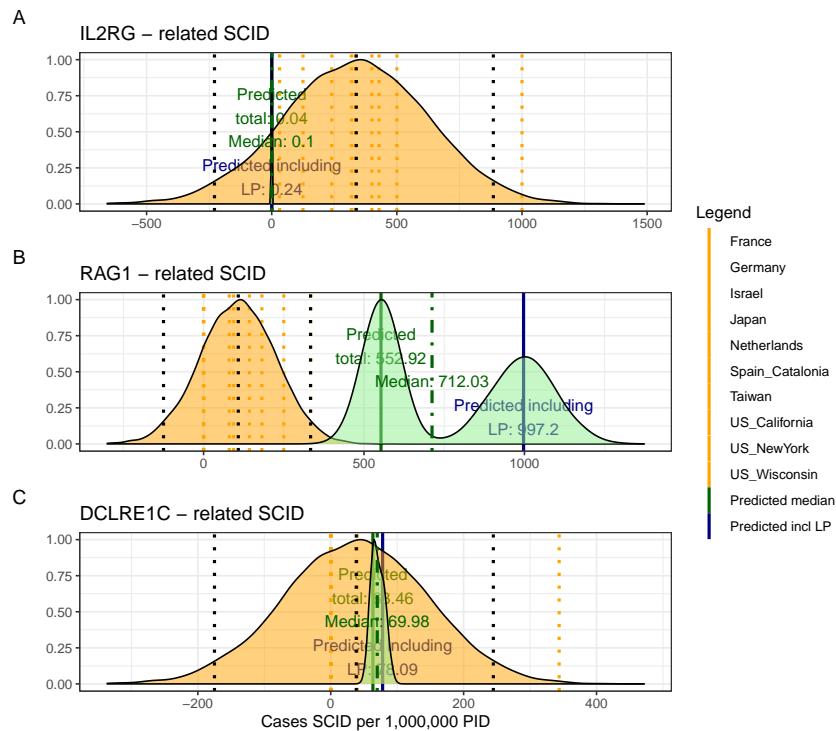


Figure S5: Combined SCID-specific Predictions and Observed Rates per 1,000,000 PID. The figure presents density distributions for the predicted SCID case counts (per 1,000,000 PID) for three genes: *IL2RG*, *RAG1*, and *DCLRE1C*. Country-specific rates (displayed as dotted vertical lines) are overlaid 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.

1029 **6.3 Hierarchical Clustering of Enrichment Scores for Major**
 1030 **Disease Categories**

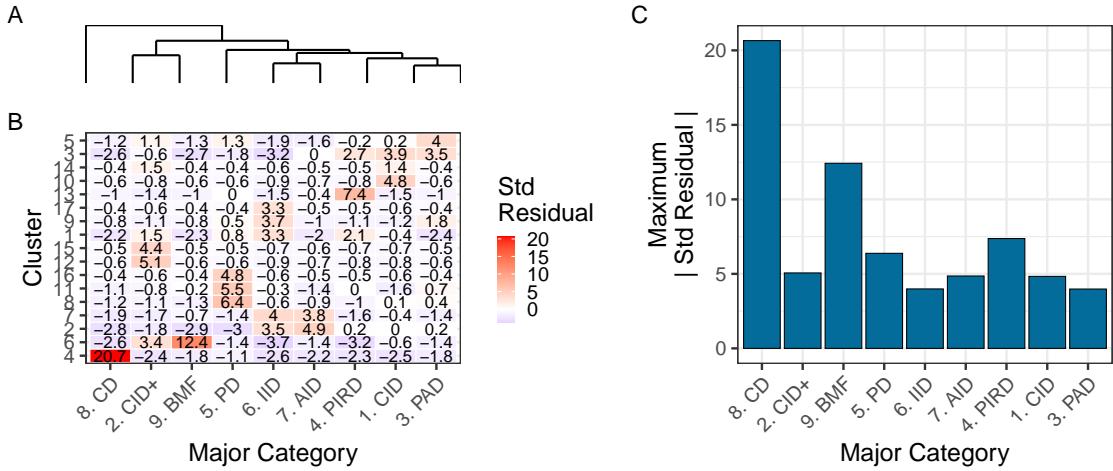


Figure S6: **Hierarchical clustering of enrichment scores.** The heatmap displays standardised residuals for major disease categories (x-axis) across network clusters (y-axis). A dendrogram groups similar disease categories, and the bar plot shows the maximum absolute residual per category. (8) CD and (9)BMF show the highest values, indicating significant enrichment or depletion ($\text{residuals} > |2|$). Definitions in **Box 2.1**.

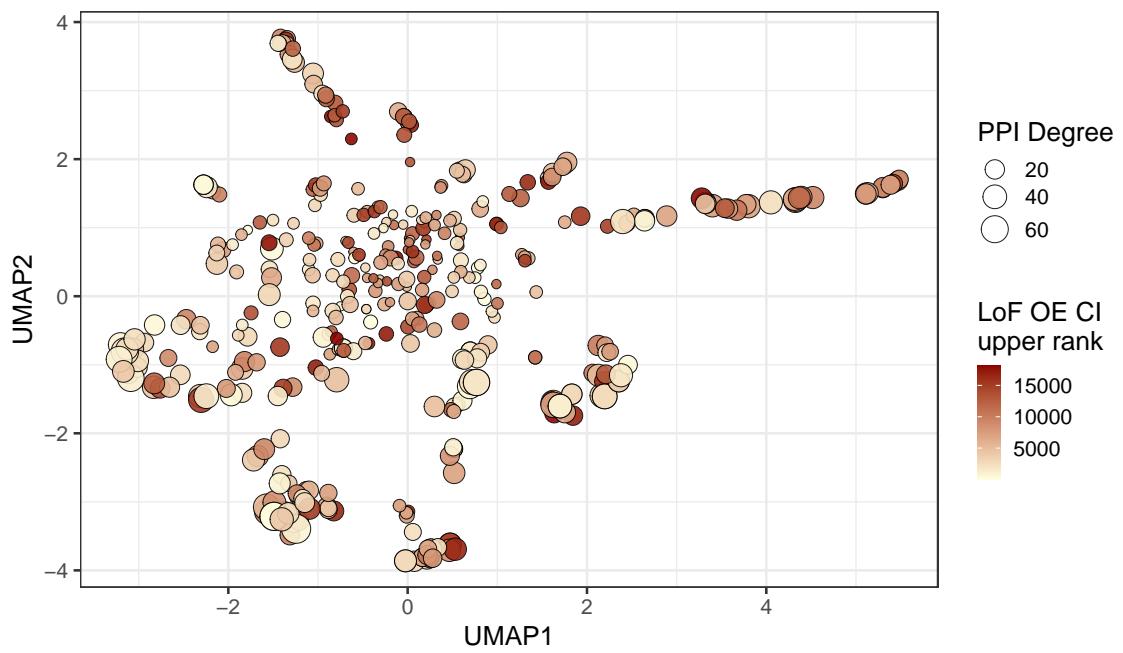


Figure S7: **Analysis of PPI degree versus LOEUF upper rank with UMAP embedding of the PPI network.** The relationship between PPI degree (size) and LOEUF upper rank (color) across gene clusters. No clear patterns are evident.

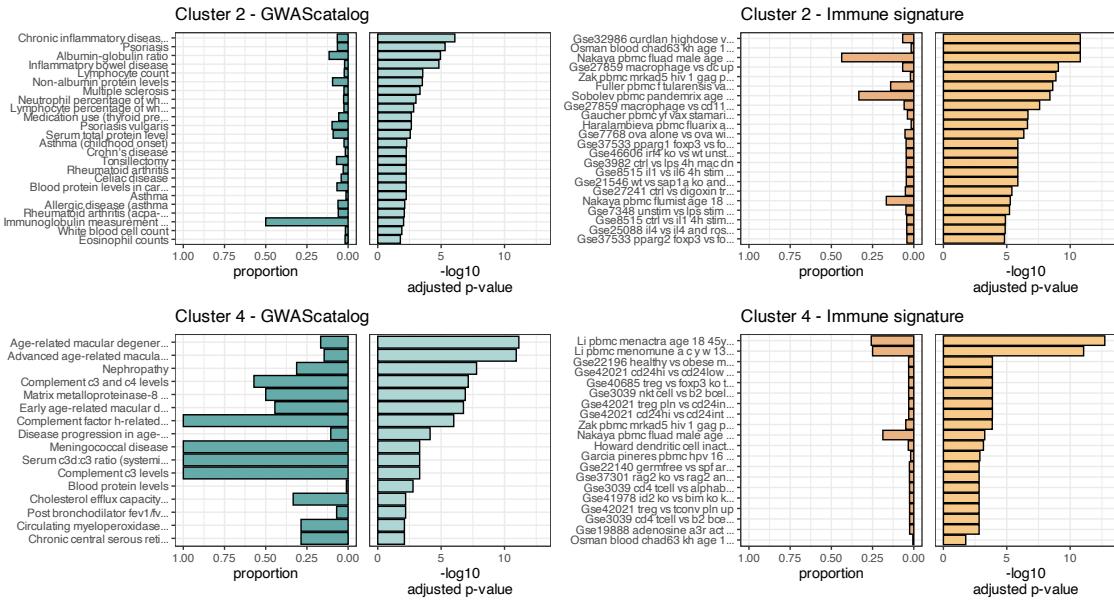


Figure S8: Composite Enrichment Profiles for IEI Gene Sets. We selected the top two enriched clusters (as per [Figure 5](#)) and performed functional enrichment analysis derived from known disease associations. For each gene set, the left panel displays the proportion of input genes overlapping with 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.

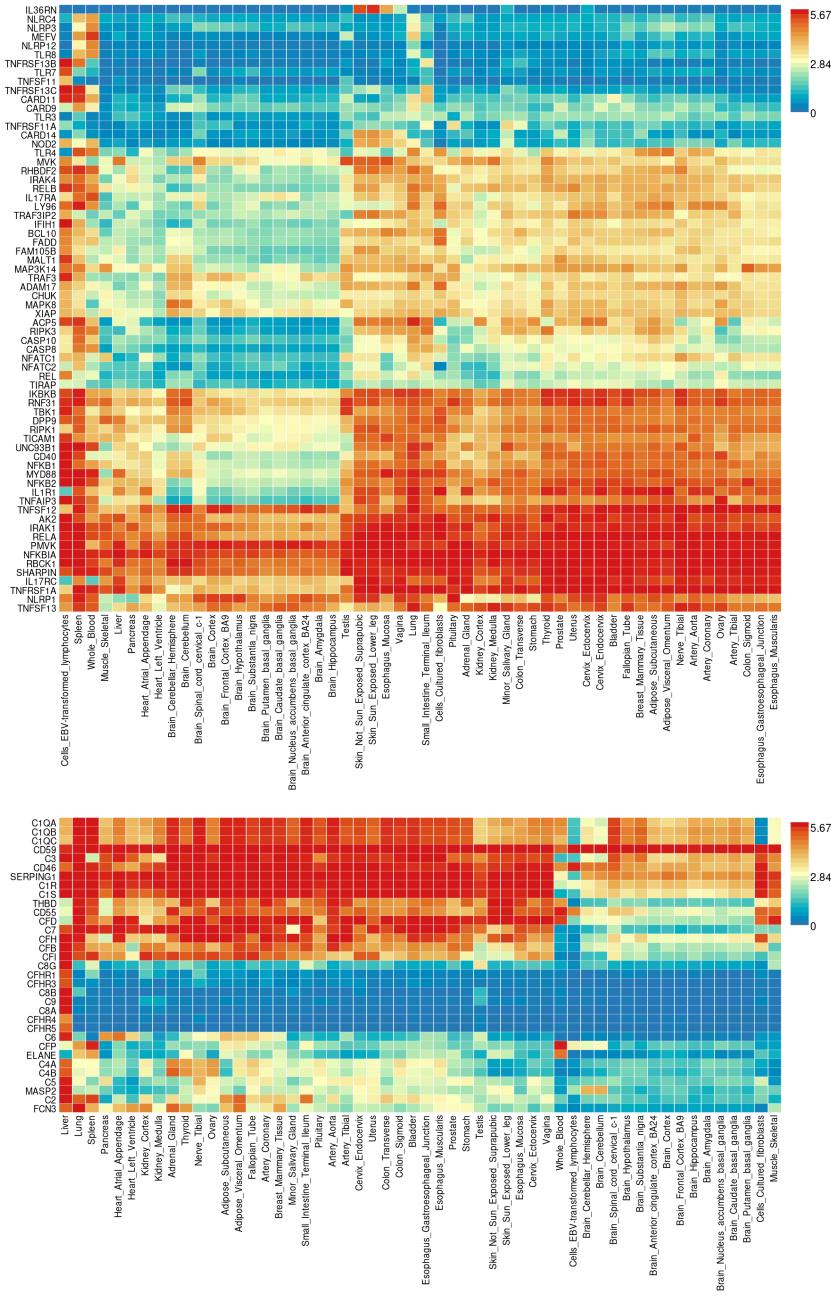


Figure S9: Gene Expression Heatmaps for IEI Genes. GTEx v8 data from 54 tissue types display the average expression per tissue label (\log_2 transformed) for the IEI gene panels. Top: Cluster 2; Bottom: Cluster 4.

6.4 Interpretation of ClinVar Variant Observations

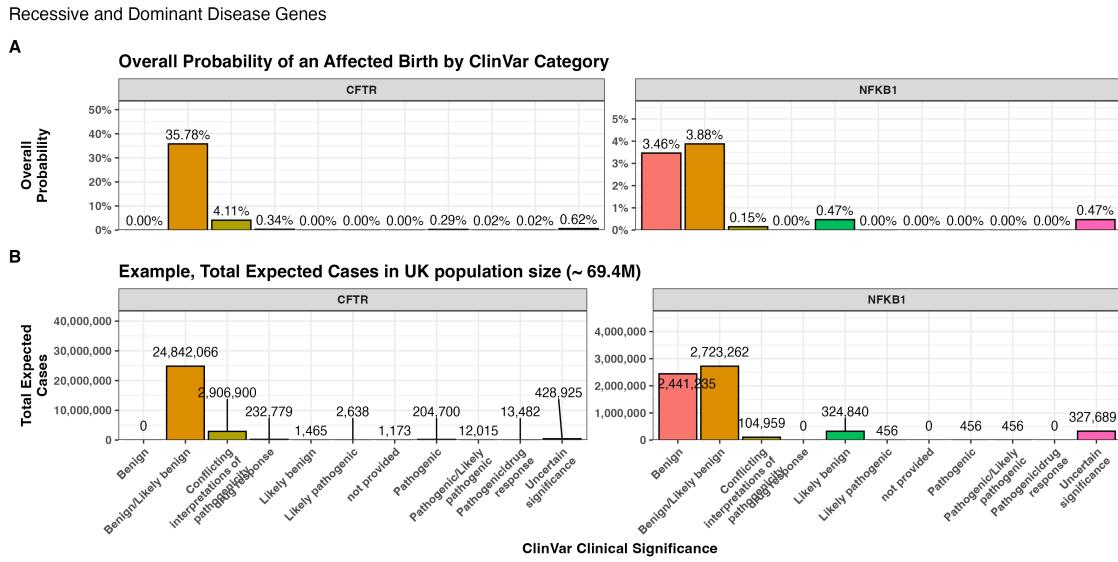


Figure S10: Combined bar charts summarizing the genome-wide analysis of ClinVar clinical significance for the PID gene panel. Panel (A) shows the overall probability of an affected birth by variant classification, and (B) displays the total expected number of cases per classification, both stratified by gene. These integrated results illustrate the variability in variant observations across genes and underpin our validation of the probability estimation framework.

6.5 Novel PID classifications

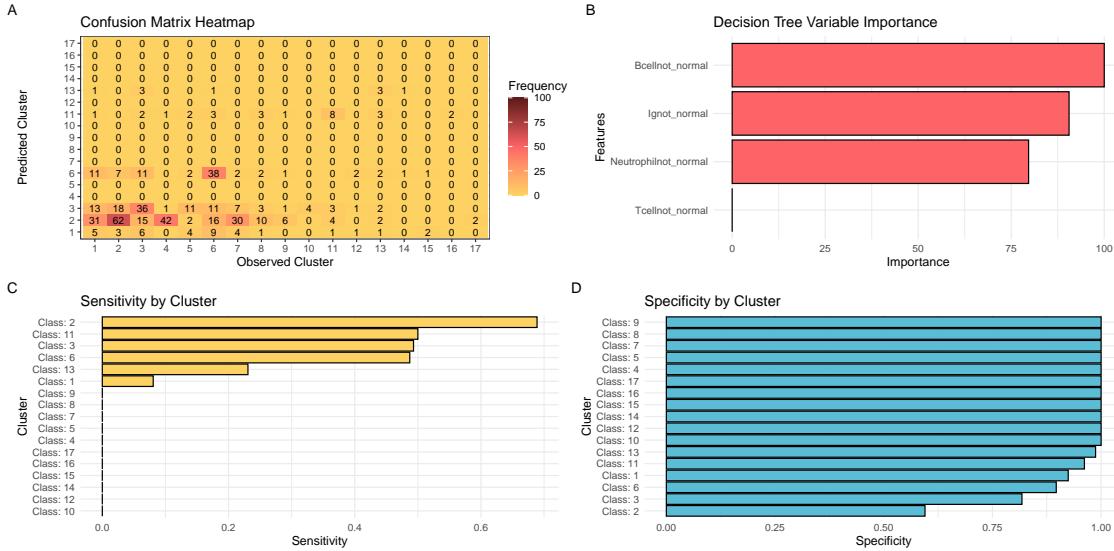


Figure S11: Model performance for fine-tuned PID classification. (A) Confusion matrix heatmap comparing observed and predicted PPI clusters. (B) Variable importance plot ranking immunophenotypic features contributing to the classifier. (C) Per-class sensitivity and (D) per-class specificity bar plots.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

1033 6.6 Probability of observing AlphaMissense pathogenicity

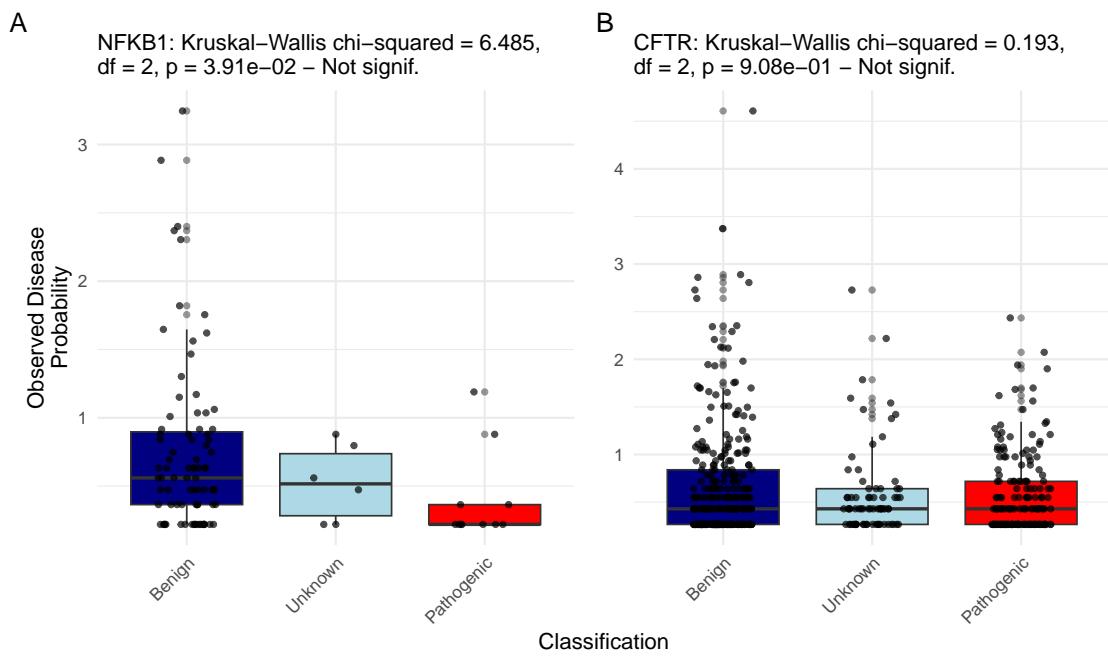


Figure S12: Observed Disease Probability by Clinical Classification with AlphaMissense. The figure displays the Kruskal-Wallis test results for NFKB1 and CFTR, showing no significant differences.