

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

Dylan Lawless^{*1}, Simon Boutry², Ali Saadat², and Jacques Fellay²

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

April 22, 2025

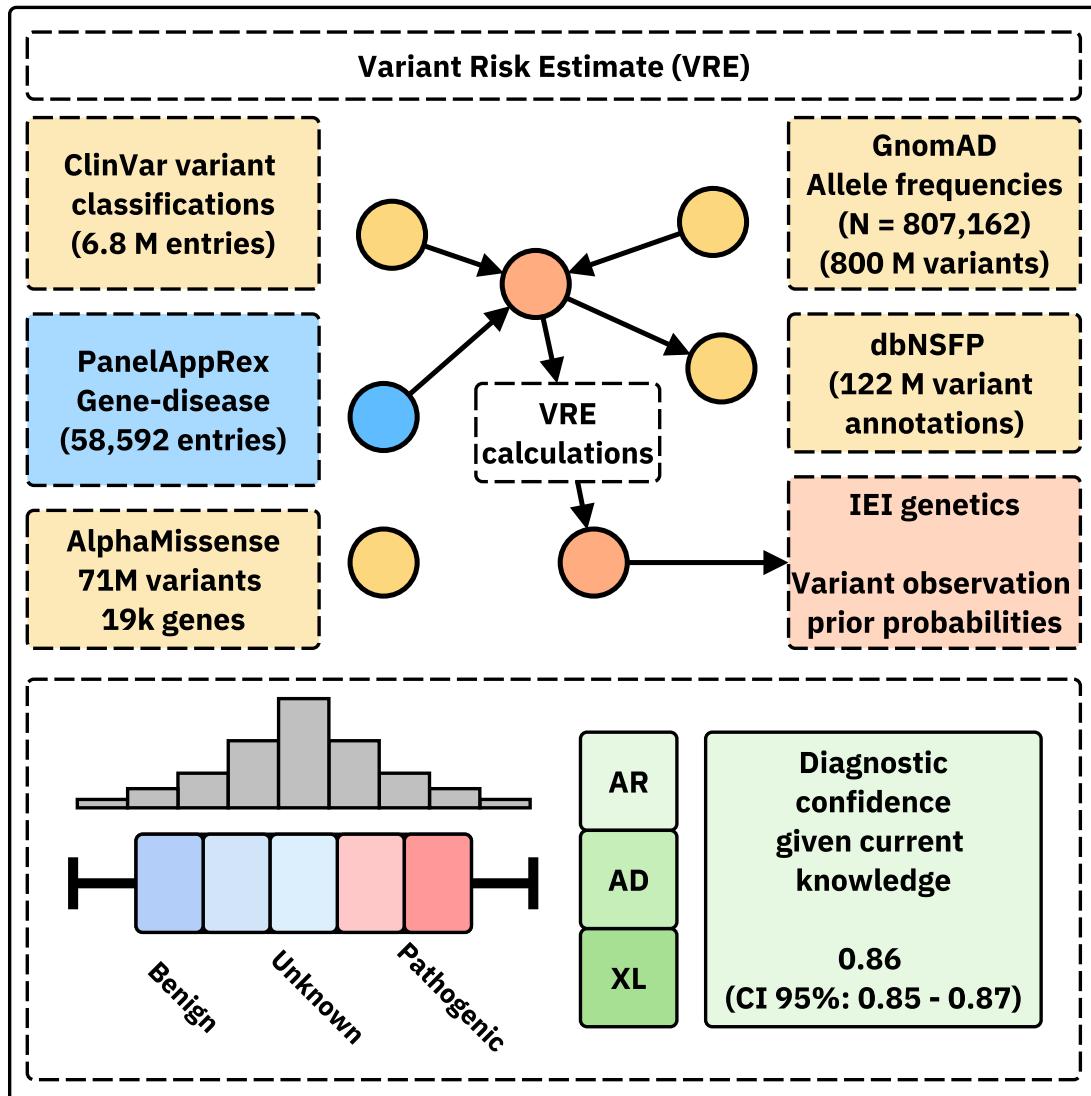
¹

Abstract

We present a framework to quantify the prior probability of observing known 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, both pathogenic variants and missing causal candidates were tested to identify the most likely genetic disease determinant and provide a clear confidence range for the overall diagnosis. 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, showing close concordance with observed case numbers. The resulting datasets supported causal variant interpretation and evidence-aware 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 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 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 ical abstract.

Graph-

²⁰ **Acronyms**

²¹ ACMG	American College of Medical Genetics and Genomics.....	³⁷
²² ACAT	Aggregated Cauchy Association Test	³⁷
²³ AD	Autosomal Dominant.....	⁵
²⁶ AF	Allele Frequency.....	⁴
²⁸ ANOVA	Analysis of Variance.....	¹⁷
³⁰ AR	Autosomal Recessive	⁵
³³ BMF	Bone Marrow Failure.....	²⁶
³⁴ CD	Complement Deficiencies	²⁷
³⁶ CI	Confidence Interval.....	²⁴
³⁸ CrI	Credible Interval	¹¹
⁴¹ CF	Cystic Fibrosis	¹⁵
⁴³ CFTR	Cystic Fibrosis Transmembrane Conductance Regulator.....	⁷
⁴⁴ CVID	Common Variable Immunodeficiency.....	¹³
⁴⁶ DCLRE1C	DNA Cross-Link Repair 1C.....	⁷
⁴⁸ dbNSFP	database for Non-Synonymous Functional Predictions	⁶
⁵¹ GE	Genomics England	⁶
⁵³ gnomAD	Genome Aggregation Database	⁶
⁵⁵ gVCF	genomic variant call format	¹¹
⁵⁶ HGVS	Human Genome Variation Society.....	⁶
⁵⁸ HPC	High-Performance Computing.....	¹⁰
⁶¹ HSD	Honestly Significant Difference	¹⁷
⁶³ HWE	Hardy-Weinberg Equilibrium	⁵
⁶⁶ IEI	Inborn Errors of Immunity	⁵
⁶⁷ Ig	Immunoglobulin	³⁰
⁶⁸ IL2RG	Interleukin 2 Receptor Subunit Gamma.....	⁷
⁷¹ InDel	Insertion/Deletion	⁶
⁷³ IUIS	International Union of Immunological Societies	⁶
⁷⁴ LD	Linkage Disequilibrium	²⁹
⁷⁶ LOEUF	Loss-Of-function Observed/Expected Upper bound Fraction	¹⁷
⁷⁸ LOF	Loss-of-Function	¹⁷
⁸¹ MOI	Mode of Inheritance	⁵
⁸³ NFKB1	Nuclear Factor Kappa B Subunit 1	⁷
⁸⁴ OMIM	Online Mendelian Inheritance in Man	³⁵
⁸⁶ PID	Primary Immunodeficiency	⁵
⁸⁸ PPI	Protein-Protein Interaction	⁶
⁹¹ pLI	Probability of being Loss-of-function Intolerant	¹⁷
⁹³ QC	Quality Control	¹¹

94	RAG1 Recombination activating gene 1	7
96	SCID Severe Combined Immunodeficiency	7
98	SNV Single Nucleotide Variant	5
100	SKAT Sequence Kernel Association Test.....	37
103	STRINGdb Search Tool for the Retrieval of Interacting Genes/Proteins.....	6
104	TP true positive.....	4
106	FP false positive.....	4
108	TN true negative.....	4
110	FN false negative.....	4
113	TNFAIP3 Tumor necrosis factor, alpha-induced protein 3	7
114	UMAP Uniform Manifold Approximation and Projection	18
116	UniProt Universal Protein Resource.....	6
118	VCF variant call format.....	11
120	VEP Variant Effect Predictor.....	6
123	VRE variant risk estimate.....	7
124	XL X-Linked	5
126		

127 1 Introduction

128 Accurately determining the probability that a patient harbours a disease-causing
 129 genetic variant remains a foundational challenge in clinical and statistical genetics.
 130 For over a century, the primary focus has been on identifying true positive (TP)s,
 131 pathogenic causal variants observed in affected individuals. Peer review and classifica-
 132 tion frameworks also work to suppress false positive (FP)s. However, two critical com-
 133 ponents of the genetic landscape have received far less attention: false negative (FN),
 134 where pathogenic variants are missed due to technical or interpretive limitations, and
 135 true negative (TN)s, which represent the vast majority of benign or non-causal vari-
 136 ants. TNs are more commonly used in contexts such as cancer screening, where a
 137 negative result can provide reassurance that a panel of known actionable variants has
 138 been checked. Yet outside these specific uses, their broader statistical and clinical
 139 value is rarely leveraged. From a statistical perspective, FNs and TNs are an un-
 140 tapped goldmine. They hold essential information about what is not observed, what
 141 should be expected under baseline assumptions, and how confident one can be in
 142 the absence of a pathogenic finding. Yet these dimensions are rarely quantified, leav-
 143 ing current variant interpretation frameworks biased toward known TPs and lacking
 144 principled priors for genome-wide disease probability estimation.

145 In this study, we focused on reporting the probability of disease observation
 146 through genome-wide assessments of gene-disease combinations. Our central hypoth-
 147 esis was that by using highly curated annotation data including population Allele

148 Frequency (AF)s, disease phenotypes, Mode of Inheritance (MOI) patterns, and variant
149 classifications and by applying rigorous calculations based on Hardy-Weinberg
150 Equilibrium (HWE), we could accurately estimate the expected probabilities of ob-
151 serving disease-associated variants. Among other benefits, this knowledge can be used
152 to derive genetic diagnosis confidence by incorporating these new priors.

153 To demonstrate, we focused on known Inborn Errors of Immunity (IEI) genes,
154 also referred to as the Primary Immunodeficiency (PID) or Monogenic Inflammatory
155 Bowel Disease genes (1–3) to validate our approach and demonstrate its clinical rel-
156 evance. This application to a well-established genotype-phenotype set, comprising
157 over 500 gene-disease associations, underscores its utility (1).

158 Quantifying the risk that a newborn inherits a disease-causing variant is a fun-
159 damental challenge in genomics. Classical statistical approaches grounded in HWE
160 (4; 5) have long been used to calculate genetic MOI probabilities for Single Nucleotide
161 Variant (SNV)s. However, applying these methods becomes more complex when ac-
162 counting for different MOI, such as Autosomal Recessive (AR) versus Autosomal
163 Dominant (AD) or X-Linked (XL) disorders. In AR conditions, for example, the
164 occurrence probability must incorporate both the homozygous state and compound
165 heterozygosity, whereas for AD and XL disorders, a single pathogenic allele is suffi-
166 cient to cause disease. Advances in genetic research have revealed that MOI can be
167 even more complex (6). Mechanisms such as dominant negative effects, haploinsuffi-
168 ciency, mosaicism, and digenic or epistatic interactions can further modulate disease
169 risk and clinical presentation, underscoring the need for nuanced approaches in risk
170 estimation. Karczewski et al. (7) made significant advances; however, the remaining
171 challenge lay in applying the necessary statistical genomics data across all MOI for
172 any gene-disease combination, which our current work aims to address. Similar ap-
173 proaches have been reported for disease such Wilson disease, Mucopolysaccharidoses,
174 Primary ciliary dyskinesia, and treatable metabolic disease, (8; 9), as reviewed by
175 Hannah et al. (10).

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

186 This gap is not only due conceptual limitations, but to the historical absence
187 of large, harmonised reference datasets. Only recently have resources become avail-
188 able to support rigorous genome-wide probability estimation. These include high-
189 resolution population allele frequencies (e.g. gnomAD v4 (7)), curated variant clas-
190 sifications (e.g. ClinVar (13)), functional annotations (e.g. UniProt (14)), and

pathogenicity prediction models (e.g. AlphaMissense (12)). We previously introduced PanelAppRex to aggregate gene panel data from multiple sources, including Genomics England (GE) PanelApp, ClinVar, and Universal Protein Resource (UniProt), thereby enabling advanced natural searches for clinical and research applications (2; 3; 13; 14). It automatically retrieves expert-curated panels, such as those from the NHS National Genomic Test Directory and the 100,000 Genomes Project, and converts them into machine-readable formats for rapid variant discovery and interpretation. Together, these resources now make it possible to model the expected distribution of variant types, frequencies, and classifications across the genome.

By reframing variant interpretation as a problem of calibrated expectation rather than solely reactive confirmation, our framework empowers clinicians and researchers to anticipate both observed and unobserved pathogenic burdens. This scalable, genome-wide approach promises to streamline diagnostic workflows, reduce uncertainty in inconclusive cases, inform statistical models and genetic epidemiology studies, and accelerate the integration of genetic insights into patient care.

2 Methods

2.1 Dataset

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

230 The following genes were used for disease cohort validations and examples. We
231 used the two most commonly reported genes from the IEI panel Nuclear Factor
232 Kappa B Subunit 1 (*NFKB1*) (18–21) and Cystic Fibrosis Transmembrane Conductance
233 Regulator (*CFTR*) (22–24) to demonstrate applications in AD and AR disease
234 genes, respectively. We used Severe Combined Immunodeficiency (SCID)-specific
235 genes AR DNA Cross-Link Repair 1C (*DCLRE1C*), AR Recombination activating
236 gene 1 (*RAG1*), XL Interleukin 2 Receptor Subunit Gamma (*IL2RG*) to demonstrate
237 a IEI subset disease phenotype of SCID. We also used AD Tumor necrosis factor,
238 alpha-induced protein 3 (*TNFAIP3*) for other examples comparable to *NFKB1* since
239 it is also causes AD pro-inflammatory disease but has more known ClinVar classifica-
240 tions at higher AF than *NFKB1*.

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

241

2.2 Variant class observation probability

To quantify the likelihood that an individual harbours a variant with a given disease classification, we compute the variant-level occurrence probability (variant risk estimate (VRE)) for each variant. As a starting point, we considered the classical HWE for a biallelic locus:

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

243 where p is the allele frequency, $q = 1 - p$, p^2 represents the homozygous dominant,
244 $2pq$ the heterozygous, and q^2 the homozygous recessive genotype frequencies. For
245 disease phenotypes, particularly under AR MOI, the risk is traditionally linked to
246 the homozygous state (p^2); however, to account for compound heterozygosity across
247 multiple variants, we allocated the overall gene-level risk proportionally among vari-
248 ants.

249 Our computational pipeline estimated the probability of observing a disease-associated
250 genotype for each variant and aggregated these probabilities by gene and ClinVar
251 classification. This approach included all variant classifications, not limited solely to

252 those deemed “pathogenic”, and explicitly conditioned the classification on the given
253 phenotype, recognising that a variant could only be considered pathogenic relative to
254 a defined clinical context. The core calculations proceeded as follows:

255 **1. Allele frequency and total variant frequency.** For each variant i in a gene,
256 the allele frequency was denoted as p_i . For each gene (any genomic region or set),
257 we defined the total variant frequency (summing across all reported variants in that
258 gene) as:

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

259 Note that, because each calculation is confined to one gene, no additional scaling
260 was required for our primary analyses (P_{tot}). However, if this same unscaled
261 summation is applied across regions or variant sets of differing size or dosage sensi-
262 tivity, it can bias burden estimates. In such cases, normalisation by region length or
263 incorporation of gene- or region-specific dosage constraints is recommended.

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

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

266 where $\max(AN)$ was the maximum allele number observed for that gene. This
267 adjustment ensured that a nonzero risk was incorporated even in the absence of
268 observed variants in the reference database.

269 **2. Occurrence probability based on MOI.** The probability that an individual
270 is affected by a variant depends on the MOI. For **AD** and **XL** variants, a single
271 pathogenic allele suffices:

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

272 For **AR** variants, disease manifests when two pathogenic alleles are present, either
273 as homozygotes or as compound heterozygotes. We use:

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

274 Under HWE, the overall gene-level probability of an AR genotype is

$$P_{\text{AR}} = P_{\text{tot}}^2 = \sum_i p_i^2 + 2 \sum_{i < j} p_i p_j,$$

275 where $P_{\text{tot}} = \sum_i p_i$. A naïve per-variant assignment

$$p_i^2 + 2 p_i (P_{\text{tot}} - p_i)$$

276 would, when summed over all i , double-count the compound heterozygous terms.
 277 To partition P_{AR} among variants without double counting, we allocate risk in propor-
 278 tion to each variant's allele frequency:

$$p_{\text{disease},i} = \frac{p_i}{P_{\text{tot}}} \times P_{\text{tot}}^2 = p_i P_{\text{tot}}.$$

279 This ensures

$$\sum_i p_{\text{disease},i} = \sum_i p_i P_{\text{tot}} = P_{\text{tot}}^2,$$

280 recovering the correct AR risk while attributing each variant its fair share of
 281 homozygous and compound-heterozygous contributions.

282 More simply, for AD or XL conditions a single pathogenic allele suffices, so the
 283 classification risk (e.g. benign, pathogenic) equals its population frequency. For AR
 284 conditions two pathogenic alleles are required - either two copies of the same variant
 285 or one copy each of two different variants, so we divide the overall recessive risk among
 286 variants according to each variant's share of the total classification frequency in that
 287 gene.

288 **3. Expected case numbers and case detection probability.** Given a popu-
 289 lation with N births (e.g. as seen in our validation studies, $N = 69\,433\,632$), the
 290 expected number of cases attributable to variant i was calculated as:

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

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

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

293 **4. Aggregation by gene and ClinVar classification.** For each gene and for
294 each ClinVar classification (e.g. “Pathogenic”, “Likely pathogenic”, “Uncertain sig-
295 nificance”, etc.), we aggregated the results across all variants. The classification
296 grouping can be substituted by any alternative score system. The total expected
297 cases for a given group was:

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

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

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

300 **5. Data processing and implementation.** We implemented the calculations
301 within a High-Performance Computing (HPC) pipeline and provided an example
302 for a single dominant disease gene, *TNFAIP3*, in the source code to enhance repro-
303 ducibility. Variant data were imported in chunks from the annotation database for
304 all chromosomes (1-22, X, Y, M).

305 For each data chunk, the relevant fields were gene name, position, allele number,
306 allele frequency, ClinVar classification, and HGVS annotations. Missing classifications
307 (denoted by “.”) were replaced with zeros and allele frequencies were converted to
308 numeric values. Subsequently, the variant data were merged with gene panel data
309 from PanelAppRex to obtain the disease-related MOI mode for each gene. For each
310 gene, if no variant was observed for a given ClinVar classification (i.e. $p_i = 0$), a
311 minimal risk was assigned as described above. Finally, we computed the occurrence
312 probability, expected cases, and the probability of observing at least one case of
313 disease using the equations presented.

314 The final results were aggregated by gene and ClinVar classification and used to
315 generate summary statistics that reviewed the predicted disease observation proba-
316 bilities. We define the *VRE* as the prior probability of observing a variant classified
317 as the cause of disease

318 **6. Score-positive-total.** For use as a simple summary statistic on the resulting
319 user-interface, we defined the *score-positive-total* as the total number of positively
320 scored variant classifications within a given region (gene, locus, or variant set). Using
321 the ClinVar classification scale from -5 (benign) to +5 (pathogenic), we included only
322 scores > 0 , corresponding to some evidence of pathogenicity. This ClinVar scoring
323 approach represents the simplest method for replication without adding additional
324 databases to our dataset. It is modular and can be substituted with any similar

325 evidence-based classification system. Variants with scores ≤ 0 were excluded from
326 the tally, as benign classifications do not inform the likelihood of disease. The score-
327 positive-total can be normalised by region size to provide a relative estimate of prior
328 likelihood that a phenotype is due to known pathogenic variants.

329 **2.3 Integrating observed true positives and unobserved false
330 negatives into a single, actionable conclusion**

331 In this section, we detail our approach to integrating sequencing data with prior classi-
332 fication evidence (e.g. pathogenic on ClinVar) to calculate the posterior probability of
333 a complete successful genetic diagnosis. Our method is designed to account for possi-
334 ble outcomes of TP, TN, and FN, by first ensuring that all nucleotides corresponding
335 to known variant classifications (benign, pathogenic, etc.) have been accurately se-
336 quenced. This implies the use of genomic variant call format (gVCF)-style data which
337 refer to variant call format (VCF)s that contain a record for every position in the
338 genome (or interval of interest) regardless of whether a variant was detected at that
339 site or not. Only after confirming that these positions match the reference alleles (or
340 novel unaccounted variants are classified) do we calculate the probability that addi-
341 tional, alternative pathogenic variants (those not observed in the sequencing data)
342 could be present. Our Credible Interval (CrI) for pathogenicity thus incorporates
343 uncertainty from the entire process, including the tally of TP, TN, and FN outcomes.
344 We ignore the contribution of FPs as a separate task to be tackled in the future.

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

348 **(i) Data pre-processing.** We synthesized an example patient in a disease cohort
349 of 200 cases. We made several scenarios where a causal genetic diagnosis based on
350 the available data is either simple, difficult, or impossible. Our example focused
351 on a proband two representative genes for AD IEI: *NFKB1* and *TNFAIP3*. All
352 coding and canonical splice-region variants for *NFKB1* were extracted from the gVCF.
353 We assumed a typical Quality Control (QC) scenario, where sites corresponding to
354 previously reported pathogenic alleles were checked for read depth ≥ 10 and genotype
355 quality ≥ 20 . Positions that failed this check were labelled *missing*, thus separating
356 true reference calls from non-sequenced or uninformative sequence.

357 **(ii) Evidence mapping and occurrence probability.** PanelAppRex variants
358 were annotated with ClinVar clinical significance. Each label was converted to an ordi-
359 nial evidence score $S_i \in [-5, 5]$ and rescaled to a pathogenic weight $W_i = \text{rescale}(S_i; -5, 5 \rightarrow$
360 $0, 1)$. This scoring system can be replaced with any comparable alternative. The
361 HWE-based pipeline of Section 2.2 supplied a per-variant occurrence probability p_i .
362 The adjusted prior was

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

363 **(iii) Prior specification.** In a hypothetical cohort of $n = 200$ diploid individuals
364 the count of allele i follows a Beta-Binomial model. Marginalising the Binomial yields
365 the Beta prior

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

366 where $\tilde{w}_i = \max(1, S_i + 1)$ contributes an additional pseudo-count whenever $S_i >$
367 0.

368 **(iv) Posterior simulation and aggregation.** For each variant i we drew $M =$
369 10 000 realisations $\pi_i^{(m)}$ and normalised within each iteration,

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

370 Variants with $S_i > 4$ were deemed to have evidence as *causal* (pathogenic or likely
371 pathogenic). We note that an alternative evidence score or conditional threshold can
372 be substituted for this step. Their mean posterior share $\bar{\pi}_i = M^{-1} \sum_m \tilde{\pi}_i^{(m)}$ and 95%
373 CrI were retained. The probability that a damaging causal allele is physically present
374 was obtained by a second layer:

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

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

378 **(v) Scenario analysis.** The three scenarios were explored for a causal genetic di-
379 agnosis that is either simple, difficult, or impossible given the existing data. The
380 proband spiked data had either: (1) known classified variants only, including only
381 one known TP pathogenic variant, *NFKB1* p.Ser237Ter, (2) inclusion of an ad-
382 dditional plausible yet non-sequenced splice-donor allele *NFKB1* c.159+1G>A (likely
383 pathogenic) as a FN, and (3) where no known causal variants were present for a pa-
384 tient, one representative variant from each distinct ClinVar classification was selected

385 and marked as unsequenced to emulate a range of putative FNs. The selected vari-
386 ants were: *TNFAIP3* p.Cys243Arg (pathogenic), p.Tyr246Ter (likely pathogenic),
387 p.His646Pro (conflicting interpretations of pathogenicity), p.Thr635Ile (uncertain
388 significance), p.Arg162Trp (not provided), p.Arg280Trp (likely benign), p.Ile207Leu
389 (benign/likely benign), and p.Lys304Glu (benign). All subsequent steps were identi-
390 cal.

391 2.4 Validation of autosomal dominant estimates using *NFKB1*

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

396 **1. Reference dataset.** We used a reference dataset reported by Tuijnenburg
397 et al. (18) to build a validation model in an AD disease gene. This study performed
398 whole-genome sequencing of 846 predominantly sporadic, unrelated PID cases from
399 the NIHR BioResource-Rare Diseases cohort. There were 390 CVID cases in the co-
400 hort. The study identified *NFKB1* as one of the genes most strongly associated with
401 PID. Sixteen novel heterozygous variants including truncating, missense, and gene
402 deletion variants, were found in *NFKB1* among the CVID cases.

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

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

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

406 **3. National estimate based on literature.** Based on literature (18; 19; 21), the
407 prevalence of CVID in the general population was estimated as

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

408 For a UK population of $N_{\text{UK}} \approx 69\,433\,632$, the expected total number of CVID
409 cases was

$$410 E_{\text{CVID}} \approx \frac{69\,433\,632}{25\,000} \approx 2777.$$

411 Assuming that the proportion of CVID cases attributable to *NFKB1* is equivalent
412 to the cohort estimate, the literature extrapolated estimate is Estimated *NFKB1* cases \approx

413 $2777 \times 0.041 \approx 114$, with a median value of approximately 118 and a 95% confidence
414 interval of 70 to 181 cases (derived from posterior sampling).

415 **4. Bayesian adjustment.** Recognising that the sequenced cohort cases likely
416 captures the majority of *NFKB1*-related patients (apart from close relatives), but may
417 still miss rare or geographically dispersed variants, we combined the cohort-based and
418 literature-based estimates using two complementary Bayesian approaches:

419 1. **Weighted adjustment (emphasising the cohort, $w = 0.9$):** We assigned
420 90% weight to the directly observed cohort count (16) and 10% to the extrapolated
421 population estimate (114), thereby accounting, illustratively, for a small fraction of unobserved cases while retaining confidence in our well-characterised cohort:
423

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

424 yielding a 95% CrI of roughly 21 to 33 cases.

425 2. **Mixture adjustment (equal weighting, $w = 0.5$):** To reflect greater uncertainty about how representative the cohort is, we combined cohort and population prevalences equally. We sampled from the posterior distribution of the cohort prevalence,
428

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

429 and mixed this with the literature-based rate at 50% each (18; 19; 21). This
430 yields a median estimate of 67 cases and a wider 95% CrI of approximately 43 to
431 99 cases, capturing uncertainty in both under-ascertainment and over-generalisation.

432 **5. Predicted total genotype counts.** The predicted total synthetic genotype
433 count (before adjustment) was 456, whereas the predicted total genotypes adjusted
434 for `synth_flag` was 0. This higher synthetic count was set based on a minimal risk
435 threshold, ensuring that at least one genotype is assumed to exist (e.g. accounting for
436 a potential unknown de novo variant) even when no variant is observed in gnomAD
437 (as per **section 2.2**).

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

⁴⁴² **2.5 Validation study for autosomal recessive CF using *CFTR***

⁴⁴³ To validate our framework for AR diseases, we focused on Cystic Fibrosis (CF).
⁴⁴⁴ For comparability sizes between the validation studies, we analysed the most com-
⁴⁴⁵ mon SNV in the *CFTR* gene, typically reported as p.Arg117His (GRCh38 Chr
⁴⁴⁶ 7:117530975 G/A, MANE Select HGVSp ENST00000003084.11: p.Arg117His). Our
⁴⁴⁷ goal was to validate our genome-wide probability estimates by comparing the ex-
⁴⁴⁸ pected number of CF cases attributable to the p.Arg117His variant in *CFTR* with
⁴⁴⁹ the nationally reported case count in a well-characterised disease cohort (22–24).

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:

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

⁴⁵⁰ **2. Mortality adjustment.** Since CF patients experience increased mortality, we
⁴⁵¹ adjusted the expected genotype counts using an exponential survival model (22–24).
⁴⁵² With an annual mortality rate $\lambda \approx 0.004$ and a median age of 22 years, the survival
⁴⁵³ factor was computed as:

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

⁴⁵⁴ Thus, the mortality-adjusted expected genotype count became:

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

455 **3. Bayesian uncertainty simulation.** To incorporate uncertainty in the allele
456 frequency p , we modelled p as a beta-distributed random variable:

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

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

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

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

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

468 **2.6 Validation of SCID-specific estimates using PID–SCID
469 genes**

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

474 **1. Incidence conversion.** Based on literature, PID occurs in approximately 1 in
475 1,000 births, whereas SCID occurs in approximately 1 in 100,000 births. Consequently,
476 in a population of 1,000,000 births there are about 1,000 PID cases and 10 SCID cases.
477 To express SCID-related variant counts on a per-million PID scale, the observed SCID
478 counts were multiplied by 100. For example, if a gene is expected to cause SCID in
479 10 cases within the total PID population, then on a per-million PID basis the count
480 is $10 \times 100 = 1,000$ cases (across all relevant genes).

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

487 **3. Integration with prior probability estimates.** The predicted genotype oc-
488 currence probabilities were derived from our framework across the PID gene panel.
489 These probabilities were then converted to expected case counts per million PID
490 cases by multiplying by 1,000,000. For instance, if the probability of observing a
491 pathogenic variant in *IL2RG* is p , the expected SCID-related count becomes $p \times 10^6$.
492 Similar conversions are applied for all relevant SCID genes.

493 **4. Bayesian Uncertainty and Comparison with Observed Data.** To address
494 uncertainty in the SCID-specific estimates, a Bayesian uncertainty simulation was
495 performed for each gene to generate a distribution of predicted case counts on a per-
496 million PID scale. The resulting median estimates and 95% CIs were then compared
497 against known national SCID counts compiled from independent registries. This
498 comparison permuted a direct evaluation of our framework's accuracy in predicting
499 the occurrence of SCID-associated variants within a PID cohort.

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

503 **2.7 Protein network and genetic constraint interpretation**

504 A PPI network was constructed using protein interaction data from STRINGdb (16).
505 We previously prepared and reported on this dataset consisting of 19,566 proteins and
506 505,968 interactions (<https://github.com/DylanLawless/ProteoMCLustR>). Node
507 attributes were derived from log-transformed score-positive-total values, which in-
508 formed both node size and colour. Top-scoring nodes (top 15 based on score) were
509 labelled to highlight prominent interactions. To evaluate group differences in score-
510 positive-total across major disease categories, one-way Analysis of Variance (ANOVA)
511 was performed followed by Tukey Honestly Significant Difference (HSD) post hoc
512 tests (and non-parametric Dunn's test for confirmation). GnomAD v4.1 constraint
513 metrics data was used for the PPI analysis and was sourced from Karczewski et al.
514 (7). This provided transcript-level metrics, such as observed/expected ratios, Loss-Of-
515 function Observed/Expected Upper bound Fraction (LOEUF), Probability of being
516 Loss-of-function Intolerant (pLI), and Z-scores, quantifying Loss-of-Function (LOF)
517 and missense intolerance, along with confidence intervals and related annotations for
518 211,523 observations.

519 **2.8 Gene set enrichment test**

520 To test for overrepresentation of biological functions, the prioritised genes were com-
521 pared against gene sets from MsigDB (including hallmark, positional, curated, motif,
522 computational, GO, oncogenic, and immunologic signatures) and WikiPathways using
523 hypergeometric tests with FUMA (25; 26). The background set consisted of 24,304
524 genes. Multiple testing correction was applied per data source using the Benjamini-
525 Hochberg method, and gene sets with an adjusted P-value ≤ 0.05 and more than one
526 overlapping gene are reported.

527 **2.9 Deriving novel PID classifications by genetic PPI and**
528 **clinical features**

529 We recategorised 315 immunophenotypic features from the original IUIS IEI annota-
530 tions, reducing the original multi-level descriptors (e.g. “decreased CD8, normal or
531 decreased CD4”) first to minimal labels (e.g.“low”) and second to binary outcomes
532 (normal vs. not-normal) for T cells, B cells, neutrophils, and immunoglobulins. Each
533 gene was mapped to its PPI cluster derived from STRINGdb and Uniform Manifold
534 Approximation and Projection (UMAP) embeddings from previous steps. We first
535 tested for non-random associations between these four binary immunophenotypes and
536 PPI clusters using χ^2 tests. To generate a data-driven PID classification, we trained
537 a decision tree (rpart) to predict PPI cluster membership from the four immunophe-
538 notypic features plus the traditional IUIS Major and Subcategory labels. Hyperpa-
539 rameters (complexity parameter = 0.001, minimum split = 10, minimum bucket = 5,
540 maximum depth = 30) were optimised via five-fold cross validation using the caret
541 framework. Terminal node assignments were then relabelled according to each group’s
542 predominant abnormal feature profile.

543 **2.10 Probability of observing AlphaMissense pathogenicity**

544 We obtained the subset pathogenicity predictions from AlphaMissense via the Al-
545 phaFold database and whole genome data from the studies data repository(12; 17).
546 The AlphaMissense data (genome-aligned and amino acid substitutions) were merged
547 with the panel variants based on genomic coordinate and HGVSc annotation. Occur-
548 rence probabilities were log-transformed and adjusted (y-axis displaying $\log_{10}(\text{occurrence}$
549 $\text{prob} + 1\text{e}-5) + 5$), to visualise the distribution of pathogenicity scores across the
550 residue sequence. A Kruskal-Wallis test was used to compare the observed disease
551 probability across clinical classification groups.

552

3 Results

553

3.1 Observation probability across disease genes

554 Our study integrated large-scale annotation databases with gene panels from PanelAppRex to systematically assess disease genes by MOI. By combining population
555 allele frequencies with ClinVar clinical classifications, we computed an expected obser-
556 vation probability for each SNV, representing the likelihood of encountering a variant
557 of a specific pathogenicity for a given phenotype. We report these probabilities for
558 54,814 ClinVar variant classifications across 557 genes (linked dataset (27)).

560 We focused on panels related to Primary Immunodeficiency or Monogenic Inflam-
561 matory Bowel Disease, using PanelAppRex panel ID 398. **Figure 1** displays all
562 reported ClinVar variant classifications for this panel. The resulting natural scaling
563 system (-5 to +5) accounts for the frequently encountered combinations of classifica-
564 tion labels (e.g. benign to pathogenic). The resulting data set (27) is briefly shown
565 in **Table 1** to illustrate that our method yielded estimations of the probability of
566 observing a variant with a particular ClinVar classification.

557 **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 “Occurrence Prob” represents
the calculated probability of observing the corresponding variant class for a given
phenotype. MOI shows the gene-disease-specific mode of inheritance. Additional
columns, such as population allele frequency, are not shown. (27)

Gene	Panel ID	ClinVar Clinical Significance	GRCh38 Pos	HGVSc	HGVSp	MOI	Occurrence Prob
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
...

567

3.2 Integrating observed true positives and unobserved false 568 negatives into a single, actionable conclusion

569 Having previously established a probabilistic framework for estimating the prior prob-
570 ability of observing disease-associated variants under different inheritance modes, we
571 then applied this model to an example patient to demonstrate it’s potential for clin-
572 ical genetics. The algorithm first verified that all known pathogenic positions have
573 been sequenced and observed as reference (true negatives), and identified any posi-
574 tions that were either observed as variant (true positives) or not assessable due to

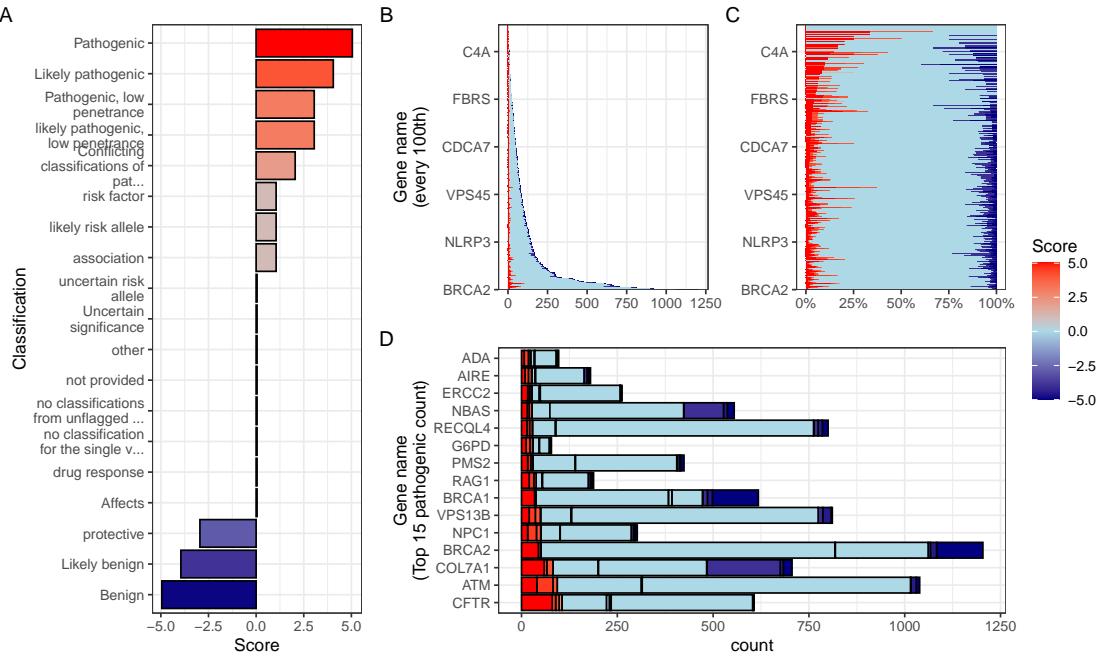


Figure 1: Summary of ClinVar clinical significance classifications in the PID gene panel. (A) Shows the numeric score coding for each classification. Panels (B) and (C) display the tally of classifications per gene as absolute counts and as percentages, respectively. (D) Highlights the top 15 genes with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

missing sequence data of failed 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 a gene.

3.3 Scenario one - simple diagnosis

We present the results from three scenarios for an example single-case patient being investigated for the genetic diagnosis of IEI. **Figure S1** shows the results of the first simple scenario, in which only one known pathogenic variant, *NFKB1* p.Ser237Ter, was observed and all other 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. The most strongly supported observed variant was p.Ser237Ter (posterior: 0.594). The strongest (probability of observing) non-sequenced variant was a benign variant p.Thr567Ile (posterior: 0). The total probability of a causal diagnosis given the available evidence was 1 (95% CI: 1–1) (**Table S1**).

591 **3.4 Scenario two - complex diagnosis**

592 **Figure 2** shows the second more complex scenario, where the same pathogenic variant
593 *NFKB1* p.Ser237Ter was present, but coverage was incomplete at three additional
594 sites of known classified variants. Among these was the likely-pathogenic splice-site
595 variant *NFKB1* c.159+1G>A, which was not captured in the sequencing data. The
596 panels of **Figure 2 (A–F)** illustrate the stepwise integration of observed and missing
597 evidence, culminating in a posterior probability that reflects both confirmed findings
598 and residual uncertainty. **Table 2** lists the final conclusion for reporting the clinical
599 genetics results. **Table S2** lists the main metrics used to reach the conclusion.

600 Bayesian integration of every annotated allele yielded the quantitative CrIs for
601 pathogenic attribution that (i) preserve Hardy-Weinberg expectations, (ii) accommo-
602 date AD, AR, XL inheritance, and (iii) carry explicit uncertainty for non-sequenced
603 (or failed QC) genomic positions. **Figure 2 (A)** depicts the prior landscape where
604 occurrence probabilities are partitioned by observed or missing status and by causal
605 or non-causal evidence, with colour reflecting the underlying ClinVar score. **Fig-**
606 **ure 2 (B)** shows posterior normalisation which concentrates probability density on
607 two high-confidence (high evidence score) alleles since the benign variants are, by
608 definition, non-causal. **Figure 2 (C)** shows the resulting per-variant probability
609 of being simultaneously damaging and causal; only the confirmed present (true posi-
610 tive) nonsense variant p.Ser237Ter and the (false negative) hypothetical splice-donor
611 c.159+1G>A retain substantial support. Restricting the view to causal candidates in
612 **Figure 2 (D)** confirms that posterior mass is distributed across these two variants.
613 **Figure 2 (E)** decomposes the total damaging probability into observed (approxi-
614 mately 40 %) and missing (approximately 34 %) sources, whereas **Figure 2 (F)** sum-
615 marises the gene-level posterior: inclusion of the splice-site allele (scenario 2) produces
616 a median probability of 0.542 with a 95 % CrI of 0.264–0.8.

617 Numerically, the present variant p.Ser237Ter accounts for 0.399 of the posterior
618 share, and the potentially causal but missing splice-donor allele c.159+1G>A con-
619 tributes 0.339. The remaining alleles together explain a negligible share (**Table S2**).
620 Thus, we can report that in this patients' scenario the probability of correct genetic
621 diagnosis due to *NFKB1* p.Ser237Ter is 0.542 (95 % CrI of 0.264–0.8) given that a
622 likely alternative remains to be confirmed for this patient. Upon confirmation that
623 the second variant is not present, the confidence will rise to 1 (95 % CrI of 1–1) as
624 shown in scenario one.

Table 2: Final variant report for clinical genetics scenario 2. Reported causal: p.Ser237Ter (posterior 0.377). Undetected causal: c.159+1G>A (posterior 0.364). The total probability of a causal diagnosis given the available evidence was 0.511 (95% CI: 0.237–0.774).

Parameter	present	missing
Gene	NFKB1	NFKB1
HGVSc	c.710C>G	c.159+1G>A
HGVSp	p.Ser237Ter	.
Inheritance	AD	AD
Patient sex	Male	Male
gnomAD frequency	6.57e-06	6.57e-06
95% CI lower	0.003	NA
p(median)	0.090	NA
95% CI upper	0.551	NA
Posterior p(causal)	0.377	0.364
Interpretation	Reported causal; variant observed	Reported causal; variant not detected — consider follow-up
Summary	Overall probability of correct causal diagnosis due to SNV/INDEL given the currently available evidence: 0.511 (95% CI 0.237–0.774).	

625 3.5 Scenario three - currently impossible diagnosis

626 **Figure S2** shows the third scenario, in which no observed variants were detected in
 627 the proband for *NFKB1*. Instead, a broad range of plausible FN we detected as miss-
 628 ing for the gene *TNFAIP3*. The strongest (probability of observing and pathogenic)
 629 of these non-sequenced variants was p.Cys243Arg (posterior: 0.347). However, the
 630 total probability of a causal diagnosis for the patient *given the available evidence* was
 631 0 (95% CI: 0–0) since these missing variants must be accounted for (**Table S3**). Upon
 632 confirmation, these probabilities can update, as per scenario two.

Gene: NFKB1

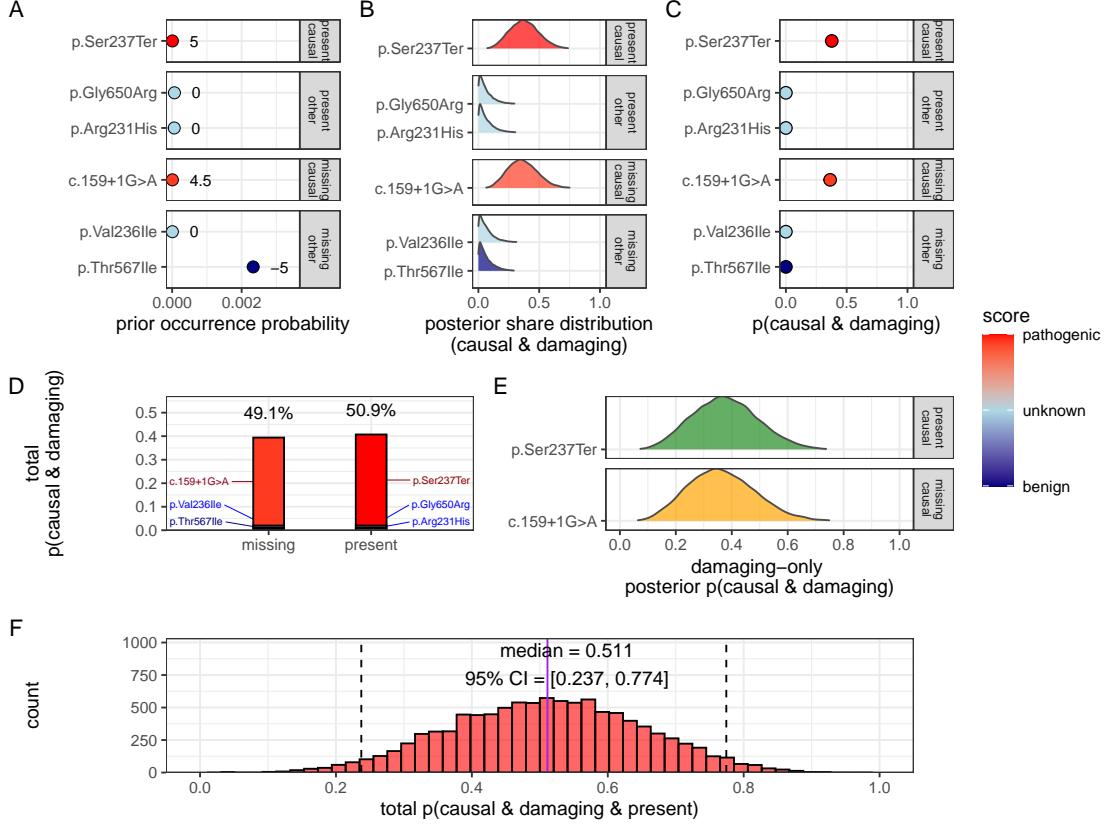


Figure 2: **Quantification of present (TP) and missing (FN) causal genetic variants for disease in *NFKB1* (scenario 2).** The example proband carried three known heterozygous variants, including pathogenic p.Ser237Ter, and had incomplete coverage at three additional loci, including likely-pathogenic 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. Pathogenicity scores (-5 to +5) are annotated. (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 0.54, 95 % CrI 0.26-0.80. This result can be compared to scenarios one and three in Figures S1 and S2, respectively.

633 **3.6 Validation studies**

634 **3.6.1 Validation of dominant disease occurrence with *NFKB1***

635 To validate our genome-wide probability estimates for AD disorders, we focused
636 on *NFKB1*. We used a reference dataset from Tuijnenburg et al. (18), in which
637 whole-genome sequencing of 846 PID patients identified *NFKB1* as one of the genes
638 most strongly associated with the disease, with 16 *NFKB1*-related CVID cases at-
639 tributed to AD heterozygous variants. Our goal was to compare the predicted num-
640 ber of *NFKB1*-related CVID cases with the reported count in this well-characterised
641 national-scale cohort.

642 Our model calculated 0 known pathogenic variant *NFKB1*-related CVID cases
643 in the UK with a minimal risk of 456 unknown de novo variants. In the reference
644 cohort, 16 *NFKB1* CVID cases were reported. We additionally wanted to account for
645 potential under-reporting in the reference study. We used an extrapolated national
646 CVID prevalence which yielded a median estimate of 118 cases (95% CI: 70–181),
647 while a Bayesian-adjusted mixture estimate produced a median of 67 cases (95% CI:
648 43–99). **Figure S3 (A)** illustrates that our predicted values reflect these ranges and
649 are closer to the observed count. This case supports the validity of our integrated
650 probability estimation framework for AD disorders, and represents a challenging ex-
651 ample where pathogenic SNV are not reported in the reference population of gnomAD.
652 Our min-max values successfully contained the true reported values.

653 **3.6.2 Validation of recessive disease occurrence with *CFTR***

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

659 In contrast, the nationally reported number of CF cases was 714, as recorded in the
660 UK Cystic Fibrosis Registry 2023 Annual Data Report (22). To account for factors
661 such as reduced penetrance and the mortality-adjusted expected genotype, we derived
662 a Bayesian-adjusted estimate via posterior simulation. Our Bayesian approach yielded
663 a median estimate of 740 cases (95% Confidence Interval (CI): 696, 786) and a
664 mixture-based estimate of 727 cases (95% CI: 705, 750). **Figure S3 (B)** illustrates
665 the close concordance between the predicted values, the Bayesian-adjusted estimates,
666 and the national report supports the validity of our approach for estimating disease.

667 **Figure S4** shows the final values for these genes of interest in a given population
668 size and phenotype. It reveals that an allele frequency threshold of approximately
669 0.000007 is required to observe a single heterozygous disease-causing variant carrier in
670 the UK population for both genes. However, owing to the AR MOI pattern of *CFTR*,
671 this threshold translates into more than 100,000 heterozygous carriers, compared to

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

677 3.6.3 Interpretation of ClinVar variant observations

678 **Figure S13** shows the two validation study PID genes, representing AR and domi-
679 nant MOI. **Figure S13 (A)** illustrates the overall probability of an affected birth by
680 ClinVar variant classification, whereas **Figure S13 (B)** depicts the total expected
681 number of cases per classification for an example population, here the UK, of approx-
682 imately 69.4 million.

683 3.6.4 Validation of SCID-specific disease occurrence

684 Given that SCID is a subset of PID, our probability estimates reflect the likelihood of
685 observing a genetic variant as a diagnosis when the phenotype is PID. However, we
686 additionally tested our results against SCID cohorts in **Figure S6**. The summarised
687 raw cohort data for SCID-specific gene counts are summarised and compared across
688 countries in **Figure S5**. True counts for *IL2RG* and *DCLRE1C* from ten distinct lo-
689 cations yielded 95% CI surrounding our predicted values. For *IL2RG*, the prediction
690 was low (approximately 1 case per 1,000,000 PID), as expected since loss-of-function
691 variants in this XL gene are highly deleterious and rarely observed in gnomAD. In con-
692 trast, the predicted value for *RAG1* was substantially higher (553 cases per 1,000,000
693 PID) than the observed counts (ranging from 0 to 200). We attributed this discrep-
694 ency to the lower penetrance and higher background frequency of *RAG1* variants in
695 recessive inheritance, whereby reference studies may underreport the true national
696 incidence. Overall, we report that agreement within an order of magnitude was tol-
697 erable given the inherent uncertainties from reference studies arising from variable
698 penetrance and allele frequencies.

699 3.7 Genetic constraint in high-impact protein networks

700 We next examined genetic constraint in high-impact protein networks across the whole
701 IEI gene set of over 500 known disease-gene phenotypes (1). By integrating ClinVar
702 variant classification scores with PPI data, we quantified the pathogenic burden per
703 gene and assessed its relationship with network connectivity and genetic constraint
704 (7; 16).

705 **3.7.1 Score-positive-total within IEI PPI network**

706 The ClinVar classifications reported in **Figure 1** were scaled -5 to +5 based on
707 their pathogenicity. We were interested in positive (potentially damaging) but not
708 negative (benign) scoring variants, which are statistically incidental in this analysis.
709 We tallied gene-level positive scores to give the score-positive-total metric. **Figure S7**
710 (**A**) shows the PPI network of disease-associated genes, where node size and colour
711 encode the score-positive-total (log-transformed). The top 15 genes with the highest
712 total prior probabilities of being observed with disease are labelled (as per **Figure**
713 **1**).

714 **3.7.2 Association analysis of score-positive-total across IEI categories**

715 We checked for any statistical enrichment in score-positive-totals, which represents the
716 expected observation of pathogenicity, between the IEI categories. One-way ANOVA
717 revealed an effect of major disease category on score-positive-total ($F(8, 500) =$
718 2.82, $p = 0.0046$), indicating that group means were not identical, which we ob-
719 served in **Figure S7 (B)**. However, despite some apparent differences in median
720 scores across categories (i.e. 9. Bone Marrow Failure (BMF)), the Tukey HSD post
721 hoc comparisons **Figure S7 (C)** showed that all pairwise differences had 95% CIs
722 overlapping zero, suggesting that individual group differences were not significant.

723 **3.7.3 UMAP embedding of the PPI network**

724 To address the density of the PPI network for the IEI gene panel, we applied UMAP
725 (**Figure 3**). Node sizes reflect interaction degree, a measure of evidence-supported
726 connectivity (**16**). We tested for a correlation between interaction degree and score-
727 positive-total. In **Figure 3**, gene names with degrees above the 95th percentile are
728 labelled in blue, while the top 15 genes by score-positive-total are labelled in yellow
729 (as per **Figure 1**). Notably, genes with high pathogenic variant loads segregated from
730 highly connected nodes, suggesting that LOF in hub genes is selectively constrained,
731 whereas damaging variants in lower-degree genes yield more specific effects. This
732 observation was subsequently tested empirically.

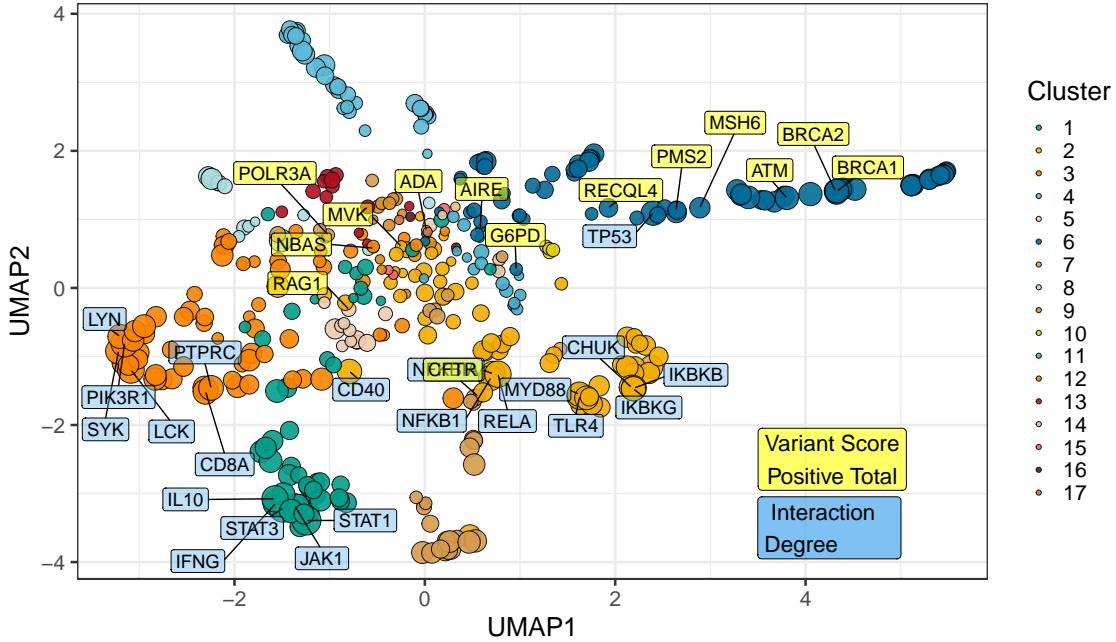


Figure 3: **UMAP embedding of the PPI network.** The plot projects the high-dimensional protein-protein interaction network into two dimensions, with nodes coloured by cluster and sized by interaction degree. Blue labels indicate hub genes (degree above the 95th percentile) and yellow labels mark the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

733 **3.7.4 Hierarchical clustering of enrichment scores for major disease cate-**
 734 **gories**

735 **Figure S8** presents a heatmap of standardised residuals for major disease categories
 736 across network clusters, as per **Figure 3**. A dendrogram clusters similar disease cate-
 737 gories, while the accompanying bar plot displays the maximum absolute standardised
 738 residual for each category. Notably, (8) Complement Deficiencies (CD) shows the
 739 highest maximum enrichment, followed by (9) BMF. While all maximum values
 740 exceed 2, the threshold for significance, this likely reflects the presence of protein
 741 clusters with strong damaging variant scores rather than uniform significance across
 742 all categories (i.e. genes from cluster 4 in 8 CD).

743 **3.7.5 PPI connectivity, LOEUF constraint and enriched network cluster**
 744 **analysis**

745 Based on the preliminary insight from **Figure S8**, we evaluated the relationship
 746 between network connectivity (PPI degree) and LOEUF constraint (LOEUF upper rank)
 747 Karczewski et al. (7) using Spearman's rank correlation. Overall, there was a weak

748 but significant negative correlation ($\rho = -0.181$, $p = 0.00024$) at the global scale,
749 indicating that highly connected genes tend to be more constrained. A supplementary
750 analysis (**Figure S9**) did not reveal distinct visual associations between network
751 clusters and constraint metrics, likely due to the high network density. However
752 once stratified by gene clusters, the natural biological scenario based on quantitative
753 PPI evidence (16), some groups showed strong correlations; for instance, cluster 2
754 ($\rho = -0.375$, $p = 0.000994$) and cluster 4 ($\rho = -0.800$, $p < 0.000001$), while others did
755 not. This indicated that shared mechanisms within pathway clusters may underpin
756 genetic constraints, particularly for LOF intolerance. We observe that the score-
757 positive-total metric effectively summarises the aggregate pathogenic burden across
758 IEI genes, serving as a robust indicator of genetic constraint and highlighting those
759 with elevated disease relevance.

760 **Figure S12 (C, D)** shows the re-plotted PPI networks for clusters with significant
761 correlations between PPI degree and LOEUF upper rank. In these networks, node
762 size is scaled by a normalised variant score, while node colour reflects the variant
763 score according to a predefined palette.

764 3.8 New insight from functional enrichment

765 To interpret the functional relevance of our prioritised IEI gene sets with the highest
766 load of damaging variants (i.e. clusters 2 and 4 in **Figure S12**), we performed
767 functional enrichment analysis for known disease associations using MsigDB with
768 FUMA (i.e. GWAScatalog and Immunologic Signatures) (25). Composite enrichment
769 profiles (**Figure S10**) reveal that our enriched PPI clusters were associated with
770 distinct disease-related phenotypes, providing functional insights beyond traditional
771 IUIS IEI groupings (1). The gene expression profiles shown in **Figure S11** (GTEx v8
772 54 tissue types) offer the tissue-specific context for these associations. Together, these
773 results enable the annotation of IEI gene sets with established disease phenotypes,
774 supporting a data-driven classification of IEI.

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

784 **3.9 Genome-wide gene distribution and linkage disequilibrium**
 785

786 **Figure 4 (A)** shows a genome-wide karyoplot of all IEI panel genes across GRCh38,
 787 with colour-coding based on MOI. Figures (B) and (C) display zoomed-in locus plots
 788 for *NFKB1* and *CFTR*, respectively. In **Figure 4 (B)**, the probability of observing
 789 variants with known classifications is high only for variants such as p.Ala475Gly,
 790 which are considered benign in the AD *NFKB1* gene that is intolerant to LOF. In
 791 **Figure 4 (C)**, high probabilities of observing patients with pathogenic variants in
 792 *CFTR* are evident, reproducing this well-established phenomenon. Furthermore, the
 793 analysis of Linkage Disequilibrium (LD) using R^2 shows that high LD regions can be
 794 modelled effectively, allowing independent variant signals to be distinguished.

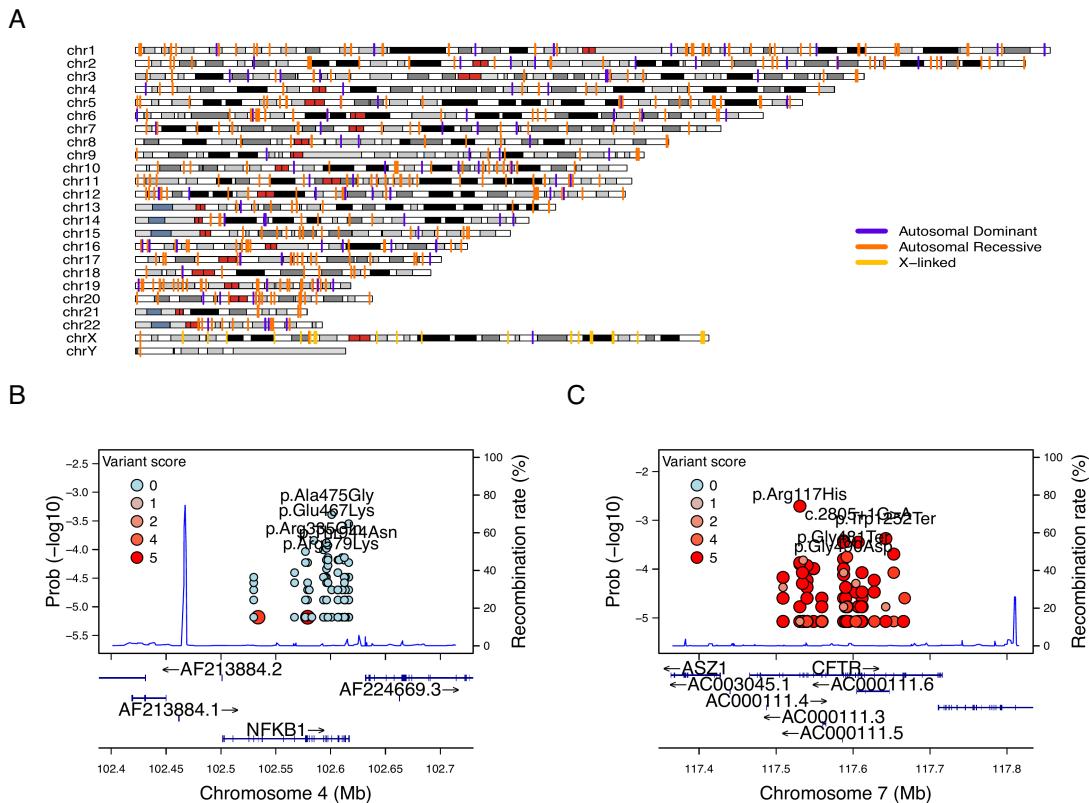


Figure 4: **Genome-wide IEI, variant occurrence probability and LD by R^2 .**
 (A) Genome-wide karyoplot of all IEI panel genes mapped to GRCh38, with colours indicating MOI. (B) Zoomed-in locus plot example for *NFKB1* showing variant observation probabilities; only benign variants such exhibit high probabilities in this AD gene intolerant to LOF. (C) Locus plot example 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.

795 **3.10 Novel PID classifications derived from genetic PPI and**
796 **clinical features**

797 We recategorised 315 immunophenotypic features from the original IUIS IEI annotations,
798 reducing detailed descriptions (e.g. “decreased CD8, normal or decreased CD4”) to minimal labels
799 (e.g. “low”) and then binarising them (normal vs. not-normal) for
800 T cells, B cells, Immunoglobulin (Ig) and neutrophils (**Figure 5**). These simplified
801 profiles were mapped onto STRINGdb PPI clusters, revealing non-random distributions
802 ($\chi^2 < 1e-13$; **Figure S14**), indicating that network context captures key
803 immunophenotypic variation.

804 We next compared four classifiers under 5-fold cross-validation to determine which
805 features predicted PPI clustering. As shown in **Figure S15**, the fully combined model
806 achieved the highest accuracy among the four: (i) phenotypes only (33 %) (i.e. T
807 cell, B cell, Ig, Neutrophil); (ii) phenotypes + IUIS major category (50 %) (e.g. CID.
808 See **Box 2.1** for more); (iii) IUIS major + subcategory only (59 %) (e.g. CID, T-B+
809 SCID); and (iv) phenotypes + IUIS major + subcategory (61 %). This demonstrated
810 that incorporating both traditional IUIS IEI classifications and core immunopheno-
811 typic markers into the PPI-based framework yielded the most robust discrimination
812 of PID gene clusters. Variable importance analysis highlighted abnormality status for
813 Ig and T cells were among the top ten features in addition to the other IUIS major
814 and sub categories. Per-class specificity remained uniform across the classes while
815 sensitivity dropped.

816 The PPI and immunophenotype model yielded 17 data-driven PID groups, whereas
817 incorporating the full complement of IUIS categories expanded this to 33 groups. For
818 clarity, we only demonstrate the decision tree from the smaller 17-group model in
819 **Figure 6**. Each terminal node is annotated by its predominant immunophenotypic
820 signature (for example, “group 65 with abnormal T cell and B cell features”), and
821 the full resulting gene counts per 33 class are plotted in **Figure 6**. Although, less
822 user-friendly, this data-driven taxonomy both aligns with and refines traditional IUIS
823 IEI classifications to provide a scaffold for large-scale computational analyses. Be-
824 cause this framework is fully reproducible, alternative PPI embeddings incorporate
825 additional molecular annotations can readily swapped to continue building on these
826 IEI classification schemes.

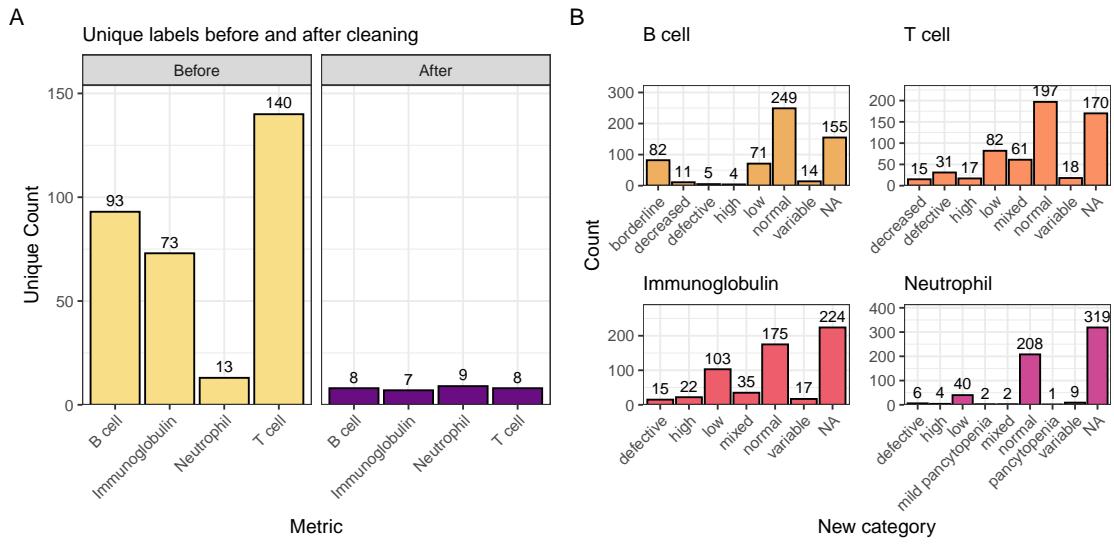


Figure 5: Distribution of immunophenotypic features before and after recategorisation. The original IUIS IEI descriptions contain information such as T cell-related “decreased CD8, normal or decreased CD4 cells” which we recategorise as “low”. The bar plot shows the count of unique labels for each status (normal, not_normal) across the T cell, B cell, Ig, and Neutrophil features.

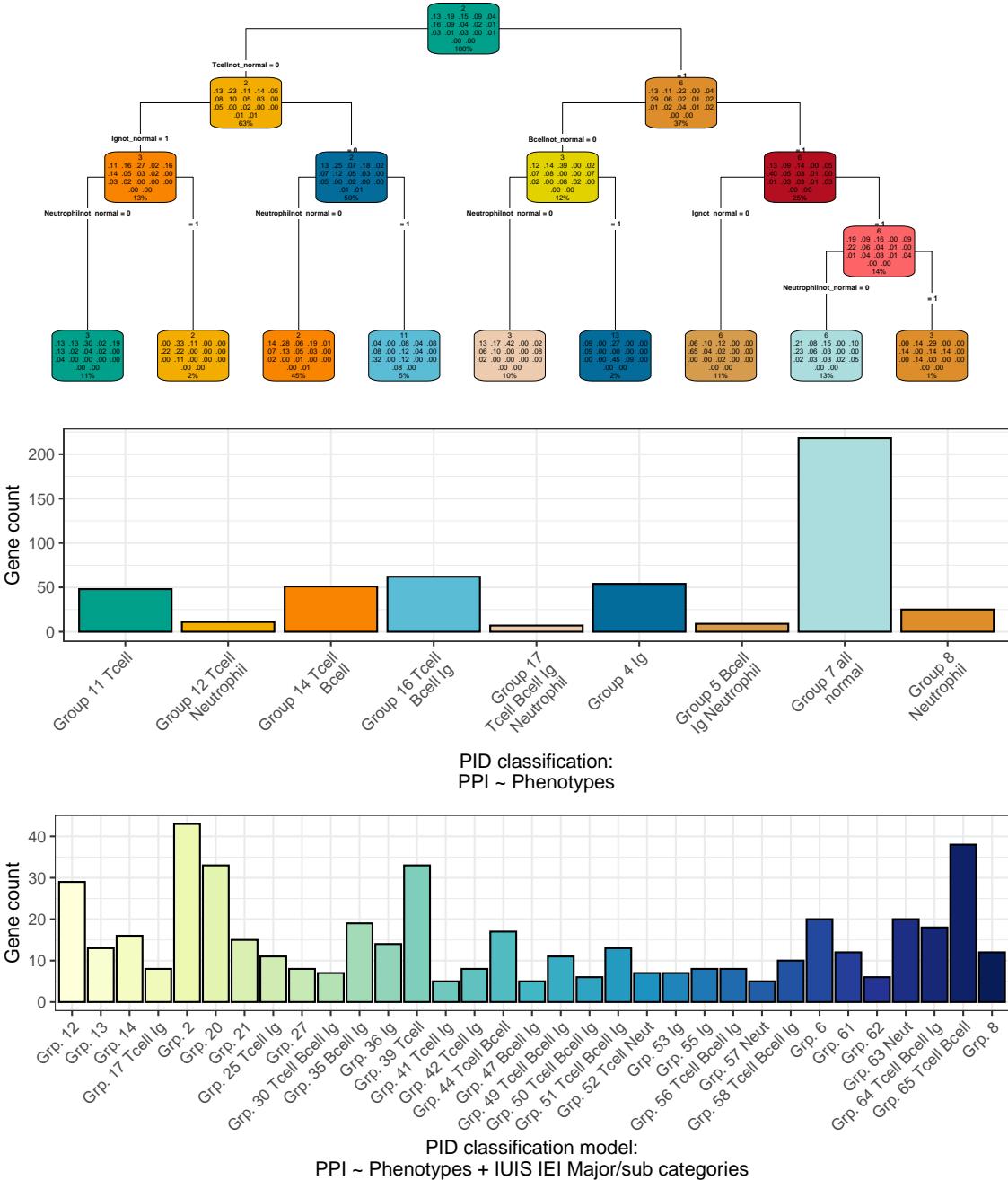


Figure 6: Fine-tuned model for PID classification. (Top) In each terminal node, the top block indicates the number of genes in the node; the middle block shows the fitted class probabilities (which sum to 1); and the bottom block displays the percentage of the total sample in that node. These metrics 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.

827 **3.11 Probability of observing AlphaMissense pathogenicity**

828 AlphaMissense provides pathogenicity scores for all possible amino acid substitutions;
829 however, our results in **Figure 7** show that the most probable observations in pa-
830 tients occur predominantly for benign or unknown variants. This finding places the
831 likelihood of disease-associated substitutions into perspective and offers a data-driven
832 foundation for future improvements in variant prediction. The values in **Figure 7**
833 (**A**) can be directly compared to **Figure 1 (D)** to view the distribution of classifi-
834 cations. A Kruskal-Wallis test was used to compare the observed disease probability
835 across clinical classification groups and no significant differences were detected. In
836 general, most variants in patients are classified as benign or unknown, indicating
837 limited discriminative power in the current classification, such that pathogenicity pre-
838 diction does not infer observation prediction (**Figure S17**). Inverse correlation likely
839 depends on factors like MOI and intolerance to LOF.

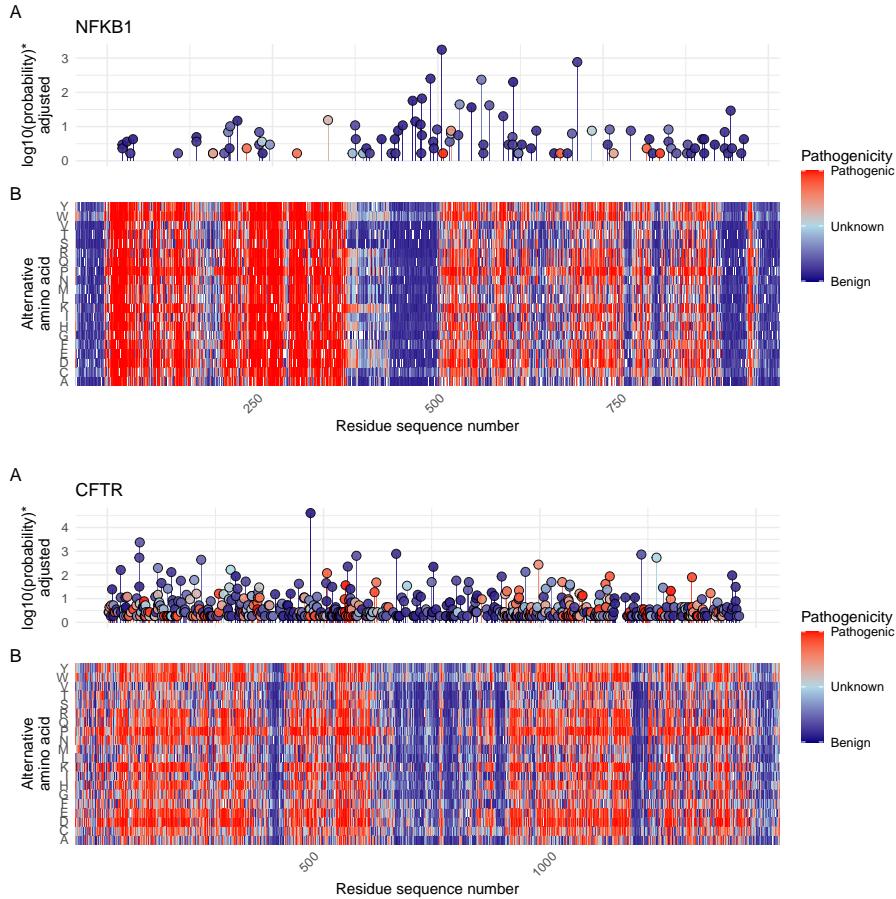


Figure 7: (A) Probabilities of observing a patient with (B) AlphaMissense-derived pathogenicity scores. Although AlphaMissense provides scores for every possible amino acid substitution, the most frequently observed variants in patients tend to be classified as benign or of unknown significance. This juxtaposition contextualises the likelihood of disease-associated substitutions and underlines prospects for refining predictive models. *Axis scaling for visibility near zero. Higher point indicates higher probability.

3.12 Integration of variant probabilities into IEI genetics data

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

Figure 8 shows the interface summarising integrated variant data. We include pre-calculated summary statistics and clinical significance as numerical metrics. Key quantiles (min, Q1, median, Q3, max) for each gene are rendered as sparkline box plots, and dynamic URLs link table entries to external databases (e.g. ClinVar, Online Mendelian Inheritance in Man (OMIM), AlphaFold) as per **Section 3.1**. The prepared data are available for bioinformatic application (27) as per **Section 3.2**.

Major category	Subcategory	Disease	Genetic defect	Inheritance	Gene score	Prior prob of observing pathogenic	ClinVar SNV classification	ClinVar all variant reports	OMIM	Alpha Missense / Uniprot ID	HPO combined	HPO term
All				All								
1. CID	1. T-B+ SCID	CD3z deficiency	CD247	AR	2			110/331/1	19,0/133/218	186780	P20963	HP-0002715; Abnormality; HP-0005403
1. CID	1. T-B+ SCID	CD3d deficiency	CD3D	AR	3			2,1/34,0	20,1/102/234	186790	P04234	HP-0002715; Abnormality; HP-0005403
1. CID	1. T-B+ SCID	CD3e deficiency	CD3E	AR	3			1,1/29,2	26,1/111/246	186830	P07766	HP-0002715; Abnormality; HP-0005403
1. CID	1. T-B+ SCID	Coronin-1A deficiency	CORO1A	AR	2			1,1/143,2	19,1/141/236/376	605000	P31146	HP-0002715; Abnormality; HP-0005403
1. CID	1. T-B+ SCID	gc deficiency (common gamma chain SCID, CD132 deficiency)	IL2RG	XL	3			1,1/16/28	194/104/244/414	303380	P31785	HP-0002715; Abnormality; HP-0005403
1. CID	1. T-B+ SCID	IL7Ra deficiency	IL7R	AR	12			8,2/181/14	8,1/29/139/196	146661	P16871	HP-0002715; Abnormality; HP-0005403
1. CID	1. T-B+ SCID	ITPKB deficiency	ITPKB	AR	3			1,1/15,9	0,1/2/130/40	Q8093	P27987	HP-0002715; Abnormality; HP-0005403
1. CID	1. T-B+ SCID	JAK3 deficiency	JAK3	AR	12			9,5/131/13	152/89/627/158	600173	P52333	HP-0002715; Abnormality; HP-0005403
1. CID	1. T-B+ SCID	LAT deficiency	LAT	AR	1			1,1/39,3	54,1/139/242	602354	O43561	HP-0002715; Abnormality; HP-0005403

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

4 Discussion

Our study presents, to our knowledge, the first comprehensive framework for calculating prior probabilities of observing disease-associated variants and the first to demonstrate the method for an evidence-aware genetic diagnosis CrI. By integrating large-scale genomic annotations, including population allele frequencies from gnomAD (7), variant classifications from ClinVar (13), and functional annotations from resources such as dbNSFP, with classical Hardy-Weinberg-based calculations, we derived robust

861 estimates for 54,814 ClinVar variant classifications across 557 IEI genes implicated in
862 PID and monogenic inflammatory bowel disease (1; 2).

863 A major deficit in current clinical genetics is the focus on confirming only the
864 presence of TP variants. Our approach yielded three key results to overcome this hurdle.
865 First, our detailed, per-variant pre-calculated results provide prior probabilities
866 of observing disease-associated variants across all MOI for any gene-disease combi-
867 nation. Second, the score-positive-total metric effectively summarises the aggregate
868 pathogenic burden across genes (before observing a patient), serving as a robust in-
869 dicator of genetic constraint and highlighting those with elevated disease relevance.
870 Building on this foundation, our third key result is a clinically applicable method to
871 estimate the probability that a patient carries a damaging causal variant, combining
872 observed and potentially unobserved variants into a single, interpretable result.

873 In the example scenarios, this enabled high-confidence attribution to a known
874 pathogenic variant while simultaneously capturing the contribution of a likely-pathogenic
875 splice-site variant missed by sequencing. This insight not is achievable using conven-
876 tional approaches which focus on detecting TP. The quantification of residual uncer-
877 tainty enables structured reporting that highlights supported, excluded, and plausible-
878 but-unseen variants, making the results actionable for clinical decision-making. These
879 outputs are suitable for diagnostic reports, support reanalysis and follow-up testing,
880 and generalise to any phenotype using the accompanying genome-wide priors. By
881 combining variant classification, allele frequency, MOI, and sequencing quality met-
882 rics, our method offers a scalable foundation for uncertainty-aware diagnostics in
883 clinical genomics.

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:

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

884 When penetrance is incomplete, defined as $P(D | G)$, the risk becomes: $P(D) =$
885 $P(G) P(D | G)$. In more realistic scenarios where multiple variants contribute to
886 disease, $P(G | D)$ denotes the fraction of cases attributable to a given variant. This
887 leads to:

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

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

Our framework addresses this challenge by combining variant classifications, population allele frequencies, and curated gene-disease associations. While imperfect on an individual level, these sources exhibit predictable aggregate behaviour, supported by James-Stein estimation principles (28). Curated gene-disease associations help identify genes that explainable for most disease cases, allowing us to approximate $P(G | D)$ close to one. In this way, we obtain robust estimates of $P(G)$ (the frequency of disease-associated genotypes), even when exact values of penetrance and 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(s) 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 CIs 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 same 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 due the the complication of requiring multiple large-scale reference databases which have only become available recently (2; 7; 12; 13).

Our results refine 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 QC 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 framework is currently focused (but not restricted) on SNVs and does not incorporate numerous other complexities of genetic disease, such

928 as structural variants, de novo variants, hypomorphic alleles, overdominance, variable
929 penetrance, tissue-specific expression, the Wahlund effect, pleiotropy, and others (6).
930 In certain applications, more refined estimates would benefit from including factors
931 such as embryonic lethality, condition-specific penetrance, and age of onset (10). Our
932 analysis also relies on simplifying assumptions of random mating, an effectively infinite
933 population, and the absence of migration, novel mutations, or natural selection. A
934 challenging remaining feature is the accurate implementation of LD which depends
935 on the whole genome population-based pairwise genotype matrix which is (based on
936 gnomAD data) approximately 80TB in size. We note that we used the reference global
937 population AFs, which is more generalisable but less accurate than population-specific
938 AF values.

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

943 5 Conclusion

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

947 Acknowledgements

948 We acknowledge Genomics England for providing public access to the PanelApp data.
949 The use of data from Genomics England panelapp was licensed under the Apache
950 License 2.0. The use of data from UniProt was licensed under Creative Commons
951 Attribution 4.0 International (CC BY 4.0). ClinVar asks its users who distribute or
952 copy data to provide attribution to them as a data source in publications and websites
953 (13). dbNSFP version 4.4a is licensed under the Creative Commons Attribution-
954 NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0); while we cite
955 this dataset as used our research publication, it is not used for the final version which
956 instead used ClinVar and gnomAD directly. GnomAD is licensed under Creative
957 Commons Zero Public Domain Dedication (CC0 1.0 Universal). GnomAD request
958 that usages cites the gnomAD flagship paper (7) and any online resources that include
959 the data set provide a link to the browser, and note that tool includes data from the
960 gnomAD v4.1 release. AlphaMissense asks to cite Cheng et al. (12) for usage in
961 research, with data available from Cheng et al. (17).

962 Competing interest

963 We declare no competing interest.

964 References

- 965 [1] Stuart G. Tangye, Waleed Al-Herz, Aziz Bousfiha, Charlotte Cunningham-
966 Rundles, Jose Luis Franco, Steven M. Holland, Christoph Klein, Tomohiro Morio,
967 Eric Oksenhendler, Capucine Picard, Anne Puel, Jennifer Puck, Mikko R. J.
968 Seppänen, Raz Somech, Helen C. Su, Kathleen E. Sullivan, Troy R. Torger-
969 son, and Isabelle Meyts. Human Inborn Errors of Immunity: 2022 Update
970 on the Classification from the International Union of Immunological Societies
971 Expert Committee. *Journal of Clinical Immunology*, 42(7):1473–1507, October
972 2022. ISSN 0271-9142, 1573-2592. doi: 10.1007/s10875-022-01289-3. URL
973 <https://link.springer.com/10.1007/s10875-022-01289-3>.
- 974 [2] Dylan Lawless. PanelAppRex aggregates disease gene panels and facilitates
975 sophisticated search. March 2025. doi: 10.1101/2025.03.20.25324319. URL
976 <http://medrxiv.org/lookup/doi/10.1101/2025.03.20.25324319>.
- 977 [3] Antonio Rueda Martin, Eleanor Williams, Rebecca E. Foulger, Sarah Leigh,
978 Louise C. Daugherty, Olivia Niblock, Ivone U. S. Leong, Katherine R. Smith,
979 Oleg Gerasimenko, Eik Haraldsdottir, Ellen Thomas, Richard H. Scott, Emma
980 Baple, Arianna Tucci, Helen Brittain, Anna De Burca, Kristina Ibañez, Dalia
981 Kasperaviciute, Damian Smedley, Mark Caulfield, Augusto Rendon, and Ellen M.
982 McDonagh. PanelApp crowdsources expert knowledge to establish consensus
983 diagnostic gene panels. *Nature Genetics*, 51(11):1560–1565, November 2019.
984 ISSN 1061-4036, 1546-1718. doi: 10.1038/s41588-019-0528-2. URL <https://www.nature.com/articles/s41588-019-0528-2>.
- 985 [4] Oliver Mayo. A Century of Hardy–Weinberg Equilibrium. *Twin Research
986 and Human Genetics*, 11(3):249–256, June 2008. ISSN 1832-4274, 1839-
987 2628. doi: 10.1375/twin.11.3.249. URL https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal_article.
- 988 [5] Nikita Abramovs, Andrew Brass, and May Tassabehji. Hardy-Weinberg Equi-
989 librium in the Large Scale Genomic Sequencing Era. *Frontiers in Genetics*,
990 11:210, March 2020. ISSN 1664-8021. doi: 10.3389/fgene.2020.00210. URL
991 <https://www.frontiersin.org/article/10.3389/fgene.2020.00210/full>.
- 992 [6] Johannes Zschocke, Peter H. Byers, and Andrew O. M. Wilkie. Mendelian
993 inheritance revisited: dominance and recessiveness in medical genetics. *Nature
994 Reviews Genetics*, 24(7):442–463, July 2023. ISSN 1471-0056, 1471-0064.
995 doi: 10.1038/s41576-023-00574-0. URL <https://www.nature.com/articles/s41576-023-00574-0>.

- 999 [7] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings,
1000 Jessica Alföldi, Qingbo Wang, Ryan L Collins, Kristen M Laricchia, Andrea
1001 Ganna, Daniel P Birnbaum, et al. The mutational constraint spectrum quantified
1002 from variation in 141,456 humans. *Nature*, 581(7809):434–443, 2020.
- 1003 [8] Sarah L. Bick, Aparna Nathan, Hannah Park, Robert C. Green, Monica H. Wo-
1004 jcik, and Nina B. Gold. Estimating the sensitivity of genomic newborn screen-
1005 ing for treatable inherited metabolic disorders. *Genetics in Medicine*, 27(1):
1006 101284, January 2025. ISSN 10983600. doi: 10.1016/j.gim.2024.101284. URL
1007 <https://linkinghub.elsevier.com/retrieve/pii/S1098360024002181>.
- 1008 [9] Benjamin D. Evans, Piotr Słowiński, Andrew T. Hattersley, Samuel E. Jones,
1009 Seth Sharp, Robert A. Kimmitt, Michael N. Weedon, Richard A. Oram,
1010 Krasimira Tsaneva-Atanasova, and Nicholas J. Thomas. Estimating disease
1011 prevalence in large datasets using genetic risk scores. *Nature Communications*,
1012 12(1):6441, November 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-26501-7.
1013 URL <https://www.nature.com/articles/s41467-021-26501-7>.
- 1014 [10] William B. Hannah, Mitchell L. Drumm, Keith Nykamp, Tiziano Prampano,
1015 Robert D. Steiner, and Steven J. Schrödi. Using genomic databases to de-
1016 termine the frequency and population-based heterogeneity of autosomal reces-
1017 sive conditions. *Genetics in Medicine Open*, 2:101881, 2024. ISSN 29497744.
1018 doi: 10.1016/j.gimo.2024.101881. URL <https://linkinghub.elsevier.com/retrieve/pii/S2949774424010276>.
- 1019 [11] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov,
1020 Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek,
1021 Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J.
1022 Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh
1023 Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy,
1024 Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamás Berghammer,
1025 Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray
1026 Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate pro-
1027 tein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August
1028 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL
1029 <https://www.nature.com/articles/s41586-021-03819-2>.
- 1030 [12] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-
1031 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias
1032 Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-
1033 abis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense vari-
1034 ant effect prediction with AlphaMissense. *Science*, 381(6664):eadg7492, Septem-
1035 ber 2023. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.adg7492. URL
1036 <https://www.science.org/doi/10.1126/science.adg7492>.
- 1037 [13] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen Chao,
1038 Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Wonhee

- 1040 Jang, Karen Karapetyan, Kenneth Katz, Chunlei Liu, Zenith Maddipatla, Adri-
1041 ana Malheiro, Kurt McDaniel, Michael Ovetsky, George Riley, George Zhou,
1042 J Bradley Holmes, Brandi L Kattman, and Donna R Maglott. ClinVar: im-
1043 proving access to variant interpretations and supporting evidence. *Nucleic Acids*
1044 *Research*, 46(D1):D1062–D1067, January 2018. ISSN 0305-1048, 1362-4962. doi:
1045 10.1093/nar/gkx1153. URL <http://academic.oup.com/nar/article/46/D1/D1062/4641904>.
- 1046
- 1047 [14] The UniProt Consortium, Alex Bateman, Maria-Jesus Martin, Sandra Orchard,
1048 Michele Magrane, Aduragbemi Adesina, Shadab Ahmad, Bowler-Barnett, and
1049 Others. UniProt: the Universal Protein Knowledgebase in 2025. *Nucleic*
1050 *Acids Research*, 53(D1):D609–D617, January 2025. ISSN 0305-1048, 1362-4962.
1051 doi: 10.1093/nar/gkae1010. URL <https://academic.oup.com/nar/article/53/D1/D609/7902999>.
- 1052
- 1053 [15] Xiaoming Liu, Chang Li, Chengcheng Mou, Yibo Dong, and Yicheng Tu.
1054 dbNSFP v4: a comprehensive database of transcript-specific functional pre-
1055 dictions and annotations for human nonsynonymous and splice-site SNVs.
1056 *Genome Medicine*, 12(1):103, December 2020. ISSN 1756-994X. doi: 10.
1057 1186/s13073-020-00803-9. URL <https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-020-00803-9>.
- 1058
- 1059 [16] Damian Szkłarczyk, Katerina Nastou, Mikaela Koutrouli, Rebecca Kirsch, Far-
1060 rokh Mehryary, Radja Hachilif, Dewei Hu, Matteo E Peluso, Qingyao Huang,
1061 Tao Fang, et al. The string database in 2025: protein networks with directional-
1062 ity of regulation. *Nucleic Acids Research*, 53(D1):D730–D737, 2025.
- 1063
- 1064 [17] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-
1065 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias
1066 Sergeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-
1067 abis, Pushmeet Kohli, and Žiga Avsec. Predictions for alphanonsense, September
1068 2023. URL <https://doi.org/10.5281/zenodo.8208688>.
- 1069
- 1070 [18] Paul Tuijnenburg, Hana Lango Allen, Siobhan O. Burns, Daniel Greene,
1071 Machiel H. Jansen, and Others. Loss-of-function nuclear factor B subunit
1 (NFKB1) variants are the most common monogenic cause of common vari-
1072 able immunodeficiency in Europeans. *Journal of Allergy and Clinical Im-*
1073 *munology*, 142(4):1285–1296, October 2018. ISSN 00916749. doi: 10.1016/
1074 j.jaci.2018.01.039. URL <https://linkinghub.elsevier.com/retrieve/pii/S0091674918302860>.
- 1075
- 1076 [19] WHO Scientific Group et al. Primary immunodeficiency diseases: report of a
1077 who scientific group. *Clin. Exp. Immunol.*, 109(1):1–28, 1997.
- 1078
- 1079 [20] Charlotte Cunningham-Rundles and Carol Bodian. Common variable immunod-
1080 eficiency: clinical and immunological features of 248 patients. *Clinical immunol-*
1081 *ogy*, 92(1):34–48, 1999.

- 1080 [21] Eric Oksenhendler, Laurence Gérard, Claire Fieschi, Marion Malphettes, Gael
1081 Mouillot, Roland Jaussaud, Jean-François Viallard, Martine Gardembas, Lionel
1082 Galicier, Nicolas Schleinitz, et al. Infections in 252 patients with common variable
1083 immunodeficiency. *Clinical Infectious Diseases*, 46(10):1547–1554, 2008.
- 1084 [22] Y Naito, F Adams, S Charman, J Duckers, G Davies, and S Clarke. Uk cystic
1085 fibrosis registry 2023 annual data report. *London: Cystic Fibrosis Trust*, 2023.
- 1086 [23] Carlo Castellani, CFTR2 team, et al. Cftr2: how will it help care? *Paediatric
1087 respiratory reviews*, 14:2–5, 2013.
- 1088 [24] Hartmut Grasemann and Felix Ratjen. Cystic fibrosis. *New England Journal
1089 of Medicine*, 389(18):1693–1707, 2023. doi: 10.1056/NEJMra2216474. URL
1090 <https://www.nejm.org/doi/full/10.1056/NEJMra2216474>.
- 1091 [25] Kyoko Watanabe, Erdogan Taskesen, Arjen Van Bochoven, and Danielle
1092 Posthuma. Functional mapping and annotation of genetic associations with
1093 FUMA. *Nature Communications*, 8(1):1826, November 2017. ISSN 2041-1723.
1094 doi: 10.1038/s41467-017-01261-5. URL <https://www.nature.com/articles/s41467-017-01261-5>.
- 1095 [26] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir,
1096 Pablo Tamayo, and Jill P. Mesirov. Molecular signatures database (MSigDB)
1097 3.0. *Bioinformatics*, 27(12):1739–1740, June 2011. ISSN 1367-4811, 1367-
1098 4803. doi: 10.1093/bioinformatics/btr260. URL <https://academic.oup.com/bioinformatics/article/27/12/1739/257711>.
- 1099 [27] Dylan Lawless. Variant risk estimate probabilities for iei genes. March 2025. doi:
1100 10.5281/zenodo.15111584. URL <https://doi.org/10.5281/zenodo.15111584>.
- 1101 [28] Bradley Efron and Carl Morris. Stein’s Estimation Rule and Its Competitors—
1102 An Empirical Bayes Approach. *Journal of the American Statistical Association*,
1103 68(341):117, March 1973. ISSN 01621459. doi: 10.2307/2284155. URL <https://www.jstor.org/stable/2284155?origin=crossref>.
- 1104 [29] Yaowu Liu, Sixing Chen, Zilin Li, Alanna C Morrison, Eric Boerwinkle, and
1105 Xihong Lin. Acat: a fast and powerful p value combination method for rare-
1106 variant analysis in sequencing studies. *The American Journal of Human Genetics*,
1107 104(3):410–421, 2019.
- 1108 [30] Xiacao Li, Zilin Li, Hufeng Zhou, Sheila M Gaynor, Yaowu Liu, Han Chen, Ryan
1109 Sun, Rounak Dey, Donna K Arnett, Stella Aslibekyan, et al. Dynamic incorporation
1110 of multiple in silico functional annotations empowers rare variant association
1111 analysis of large whole-genome sequencing studies at scale. *Nature genetics*, 52
1112 (9):969–983, 2020.
- 1113 [31] Michael C Wu, Seunggeun Lee, Tianxi Cai, Yun Li, Michael Boehnke, and Xi-
1114 hong Lin. Rare-variant association testing for sequencing data with the sequence
1115

- 1118 kernel association test. *The American Journal of Human Genetics*, 89(1):82–93,
1119 2011.
- 1120 [32] Seunggeun Lee, Mary J Emond, Michael J Bamshad, Kathleen C Barnes, Mark J
1121 Rieder, Deborah A Nickerson, David C Christiani, Mark M Wurfel, and Xihong
1122 Lin. Optimal unified approach for rare-variant association testing with applica-
1123 tion to small-sample case-control whole-exome sequencing studies. *The American
1124 Journal of Human Genetics*, 91(2):224–237, 2012.
- 1125 [33] Augustine Kong, Gudmar Thorleifsson, Michael L Frigge, Bjarni J Vilhjalmsson,
1126 Alexander I Young, Thorgeir E Thorsteinsson, Stefania Benonisdottir, Asmundur
1127 Oddsson, Bjarni V Halldorsson, Gisli Masson, et al. The nature of nurture:
1128 Effects of parental genotypes. *Science*, 359(6374):424–428, 2018.
- 1129 [34] Laurence J Howe, Michel G Nivard, Tim T Morris, Ailin F Hansen, Humaira
1130 Rasheed, Yoonsoo Cho, Geetha Chittoor, Penelope A Lind, Teemu Palviainen,
1131 Matthijs D van der Zee, et al. Within-sibship gwas improve estimates of direct
1132 genetic effects. *BioRxiv*, pages 2021–03, 2021.
- 1133 [35] Sue Richards, Nazneen Aziz, Sherri Bale, David Bick, Soma Das, Julie Gastier-
1134 Foster, Wayne W Grody, Madhuri Hegde, Elaine Lyon, Elaine Spector, et al.
1135 Standards and guidelines for the interpretation of sequence variants: a joint
1136 consensus recommendation of the american college of medical genetics and ge-
1137 nomics and the association for molecular pathology. *Genetics in medicine*, 17
1138 (5):405–423, 2015.
- 1139 [36] Sean V Tavtigian, Steven M Harrison, Kenneth M Boucher, and Leslie G
1140 Biesecker. Fitting a naturally scaled point system to the acmgs/amp variant
1141 classification guidelines. *Human mutation*, 41(10):1734–1737, 2020.
- 1142 [37] Quan Li and Kai Wang. Intervar: clinical interpretation of genetic variants by
1143 the 2015 acmgs-amp guidelines. *The American Journal of Human Genetics*, 100
1144 (2):267–280, 2017.
- 1145 [38] Brent S Pedersen, Joe M Brown, Harriet Dashnow, Amelia D Wallace, Matt
1146 Velinder, Martin Tristani-Firouzi, Joshua D Schiffman, Tatiana Tvardik, Rong
1147 Mao, D Hunter Best, et al. Effective variant filtering and expected candidate
1148 variant yield in studies of rare human disease. *NPJ Genomic Medicine*, 6(1):1–8,
1149 2021.
- 1150 [39] Carl A Anderson, Fredrik H Pettersson, Geraldine M Clarke, Lon R Cardon,
1151 Andrew P Morris, and Krina T Zondervan. Data quality control in genetic
1152 case-control association studies. *Nature protocols*, 5(9):1564–1573, 2010. URL
1153 <https://doi.org/10.1038/nprot.2010.116>.
- 1154 [40] David T Miller, Kristy Lee, Noura S Abul-Husn, Laura M Amendola, Kyle Broth-
1155 ers, Wendy K Chung, Michael H Gollob, Adam S Gordon, Steven M Harrison,
1156 Ray E Hershberger, et al. Acmg sf v3. 2 list for reporting of secondary findings

1157 in clinical exome and genome sequencing: a policy statement of the american
1158 college of medical genetics and genomics (acmg). *Genetics in Medicine*, 25(8):
1159 100866, 2023.

₁₁₆₀ **6 Supplemental**

₁₁₆₁ **6.1 Integrating observed true positives and unobserved false**
₁₁₆₂ **negatives into a single, actionable conclusion**

Table S1: Result of clinical genetics diagnosis scenario 1 including metadata. The most strongly supported observed variant was p.Ser237Ter (posterior: 0.594). The strongest unsequenced variant was p.Thr567Ile (posterior: 0). The total probability of a causal diagnosis given the available evidence was 1 (95% CI: 1–1).

Variant	Flag	Class	Evidence Score	Occurrence Prob	Adj Occ Prob	Alpha	Beta	Lower	Median	Upper	Posterior Share	Prob Causal
p.Ser237Ter	present	causal	5	0.000	0	6	371	0.004	0.142	0.803	0.594	0.594
p.Thr567Ile	missing	other	-5	0.002	0	1	363	NA	NA	NA	0.000	0.000
p.Arg231His	present	other	0	0.000	0	1	361	0.004	0.142	0.803	0.000	0.000
p.Gly650Arg	present	other	0	0.000	0	1	379	0.004	0.142	0.803	0.000	0.000
p.Val236Ile	missing	other	0	0.000	0	1	351	NA	NA	NA	0.000	0.000
Total	NA	NA	NA	NA	NA	NA	NA	1.000	1.000	1.000	NA	1.000

Table S2: Result of clinical genetics diagnosis scenario 2 including metadata. The most strongly supported observed variant was p.Ser237Ter (posterior: 0.377). The strongest unsequenced variant was c.159+1G>A (posterior: 0.364). The total probability of a causal diagnosis given the available evidence was 0.511 (95% CI: 0.237–0.774).

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	371	0.003	0.090	0.551	0.377	0.377
c.159+1G>A	missing	causal	4.5	0.000	0	5.5	351	NA	NA	NA	0.364	0.364
p.Thr567Ile	missing	other	-5.0	0.002	0	1.0	363	NA	NA	NA	0.000	0.000
p.Arg231His	present	other	0.0	0.000	0	1.0	361	0.003	0.090	0.551	0.000	0.000
p.Gly650Arg	present	other	0.0	0.000	0	1.0	379	0.003	0.090	0.551	0.000	0.000
p.Val236Ile	missing	other	0.0	0.000	0	1.0	351	NA	NA	NA	0.000	0.000
Total	NA	NA	NA	NA	NA	NA	NA	0.237	0.511	0.774	NA	0.511

Table S3: Result of clinical genetics diagnosis scenario 3 including metadata. No observed variants were detected in this scenario. The strongest unsequenced variant was p.Cys243Arg (posterior: 0.366). The total probability of a causal diagnosis given the available evidence was 0 (95% CI: 0–0).

Variant	Flag	Class	Evidence Score	Occurrence Prob	Adj Occ Prob	Alpha	Beta	Lower	Median	Upper	Posterior Share	Prob Causal
p.Cys243Arg	missing	causal	5.0	0.000	0.000	6	341	NA	NA	NA	0.366	0.366
p.Tyr246Ter	missing	causal	4.0	0.000	0.000	5	369	NA	NA	NA	0.284	0.284
p.Lys304Glu	missing	other	-5.0	0.000	0.000	1	353	NA	NA	NA	0.000	0.000
p.Ile207Leu	missing	other	-4.5	0.000	0.000	1	359	NA	NA	NA	0.000	0.000
p.His646Pro	missing	other	0.0	0.002	0.001	1	377	NA	NA	NA	0.000	0.000
p.Arg280Trp	missing	other	-4.0	0.000	0.000	1	357	NA	NA	NA	0.000	0.000
p.Thr635Ile	missing	other	0.0	0.000	0.000	1	349	NA	NA	NA	0.000	0.000
p.Arg162Trp	missing	other	0.0	0.000	0.000	1	369	NA	NA	NA	0.000	0.000
Total	NA	NA	NA	NA	NA	NA	NA	0	0	0	NA	0.000

Gene: *NFKB1*

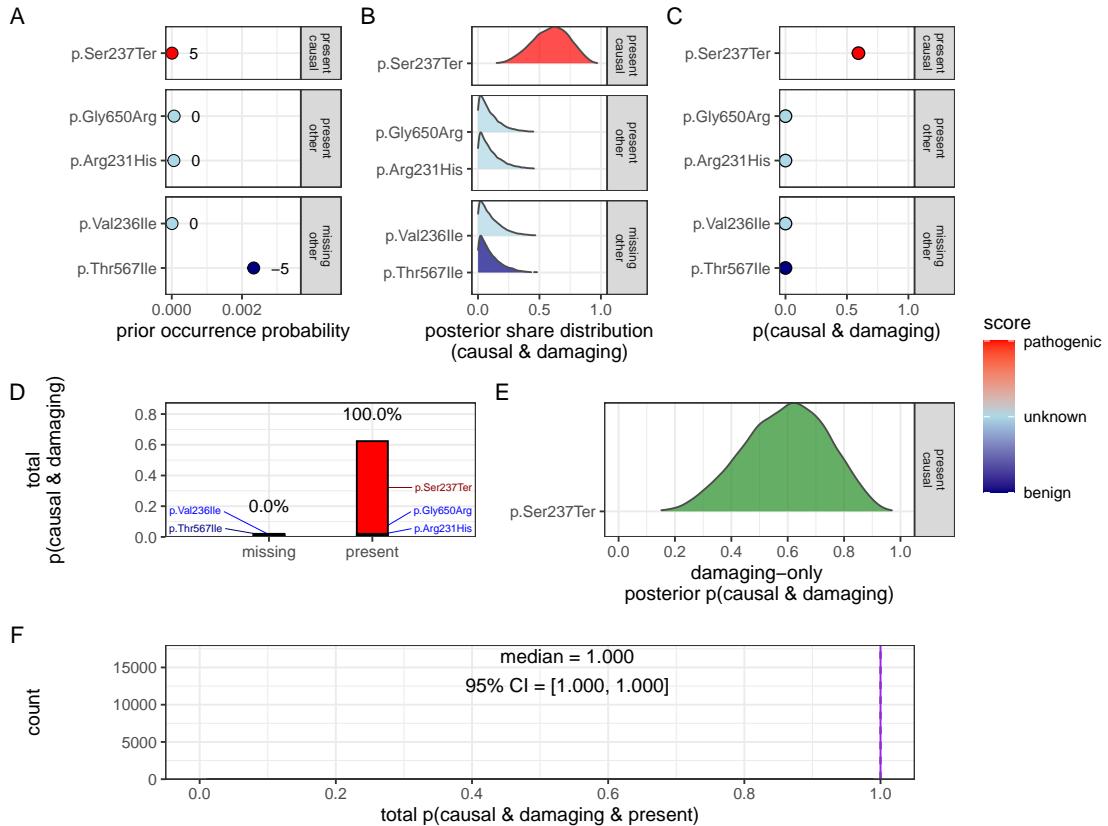


Figure S1: **Quantification of present (TP) and no missing (FN) causal genetic variants for disease in *NFKB1* (scenario 1).** Only one known pathogenic variant, p.Ser237Ter, was observed and all previously reported pathogenic positions were successfully sequenced and confirmed as reference (true negatives). Panels (A–F) follow the same structure as scenario 2 described in **Figure 2**, culminating in a gene-level posterior probability of 1 (95 % CrI: 0.99–1.00), with full support assigned to the observed allele given the available evidence. Pathogenicity scores (-5 to +5) are annotated.

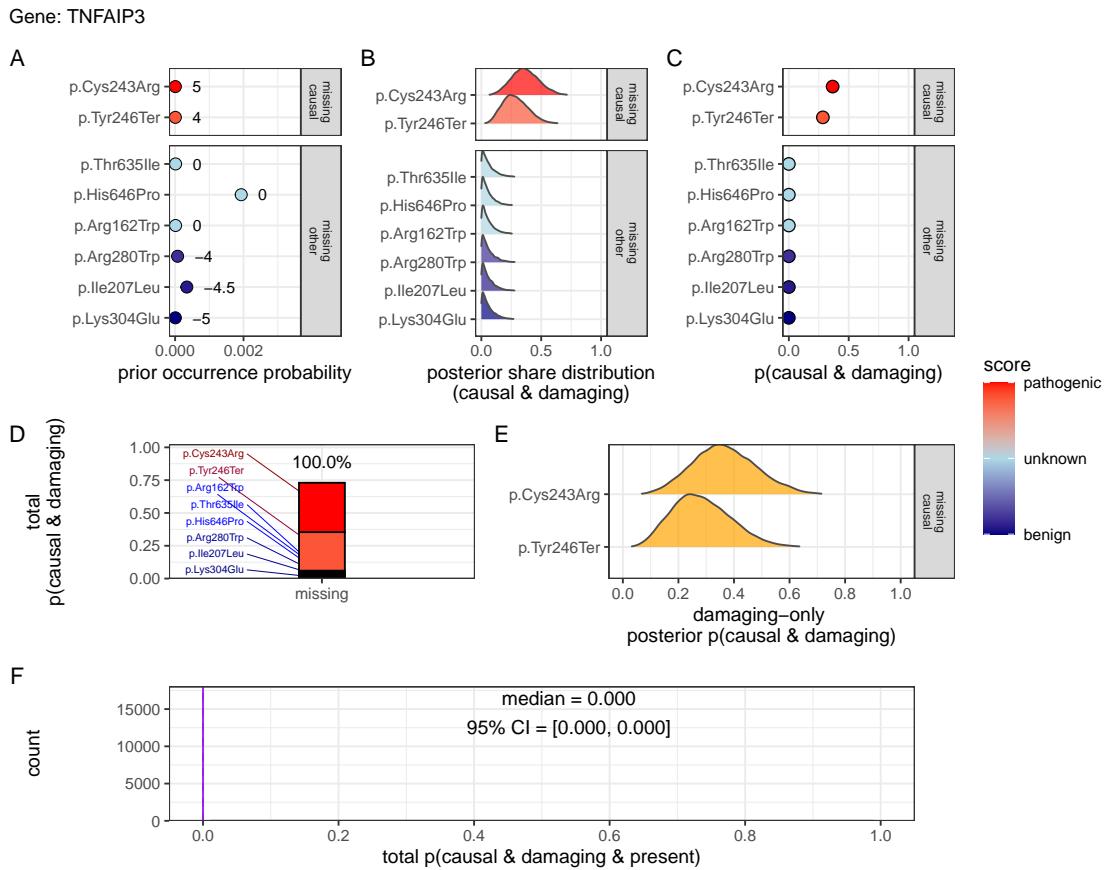


Figure S2: **Quantification of no present (TP) in *NFKB1* and only missing (FN) causal genetic variants for disease in *TNFAIP3* (scenario 3).** No known causal variants were observed in *NFKB1*, but one representative unsequenced allele was selected from each distinct ClinVar classification and treated as a potential false negative. Panels (A–F) follow the same structure as scenario 2 described in **Figure 2**. The posterior reflects uncertainty across multiple plausible but unobserved variants, resulting in low CrI (0–0) and 100% missing overall attribution in contrast to scenarios where known pathogenic variants were observed. For this patient, we have no evidence of a causal variant since the only top candidates are not yet accounted for. Pathogenicity scores (-5 to +5) are annotated in (A).

6.2 Validation studies

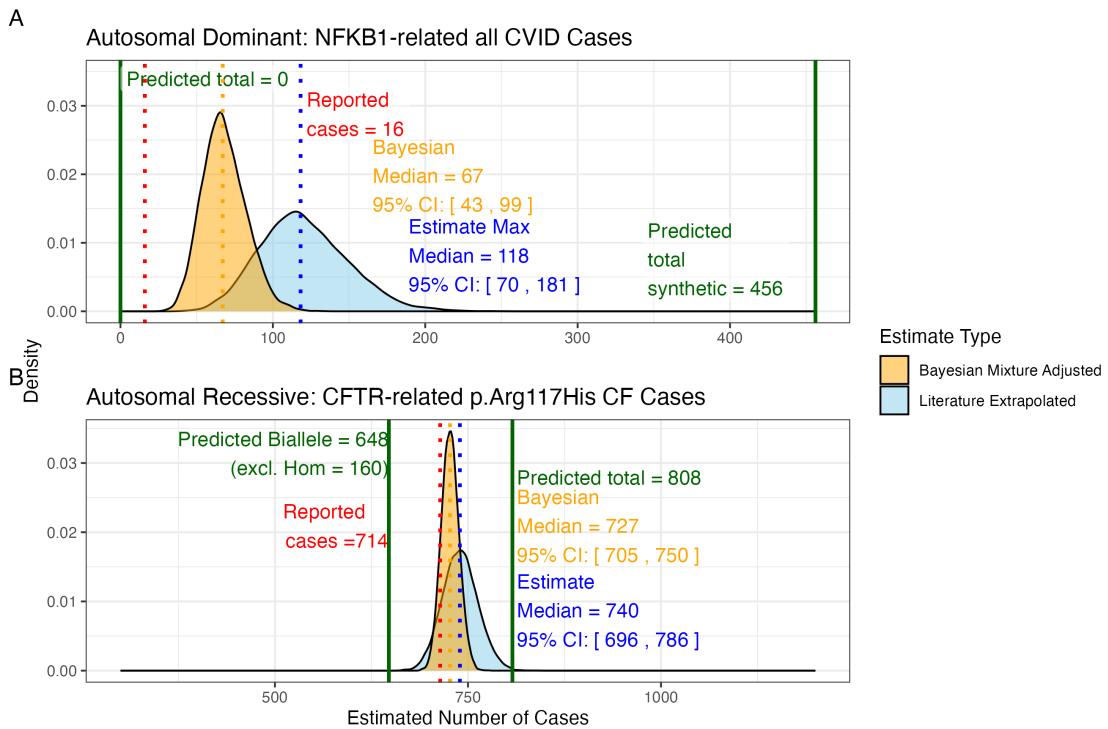


Figure S3: Prior probabilities compared to validation disease cohort metrics.
 (A) Density distributions for the number of *NFKB1*-related CVID cases in the UK. Our model (green) predicted 456 cases, which falls between the observed cohort count (red) of 390 and the upper extrapolated values. The blue 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.

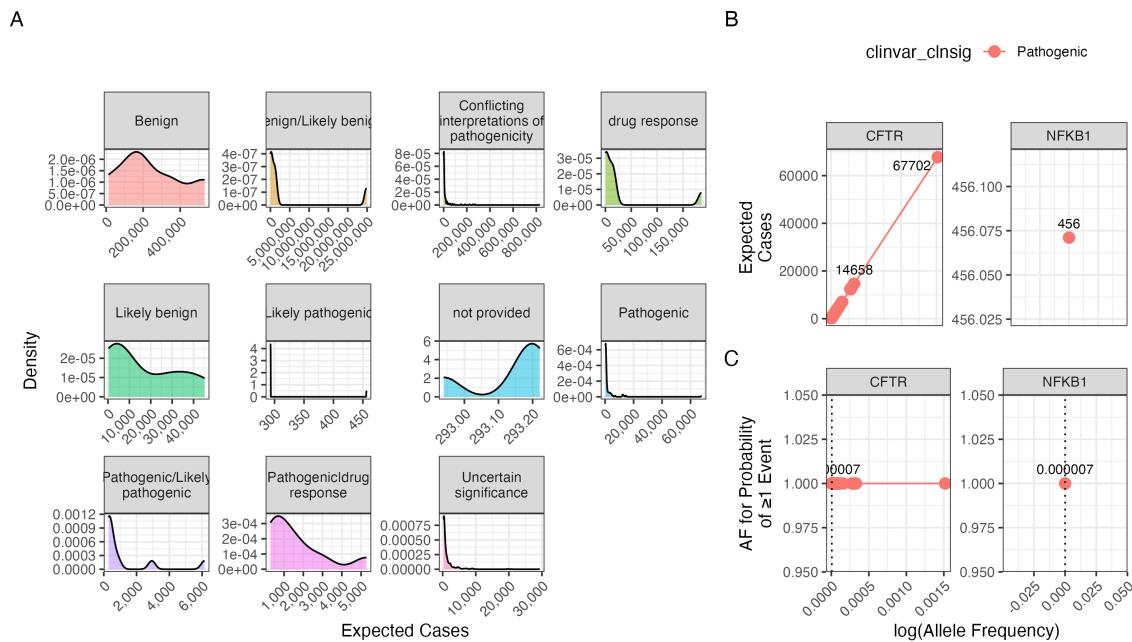


Figure S4: Interpretation of probability of observing a variant classification.
The result from the chosen validation genes *CFTR* and *NFKB1* are shown. Case counts are dependant on the population size and phenotype. (A) The density plots of expected observations by ClinVar clinical significance. We 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.

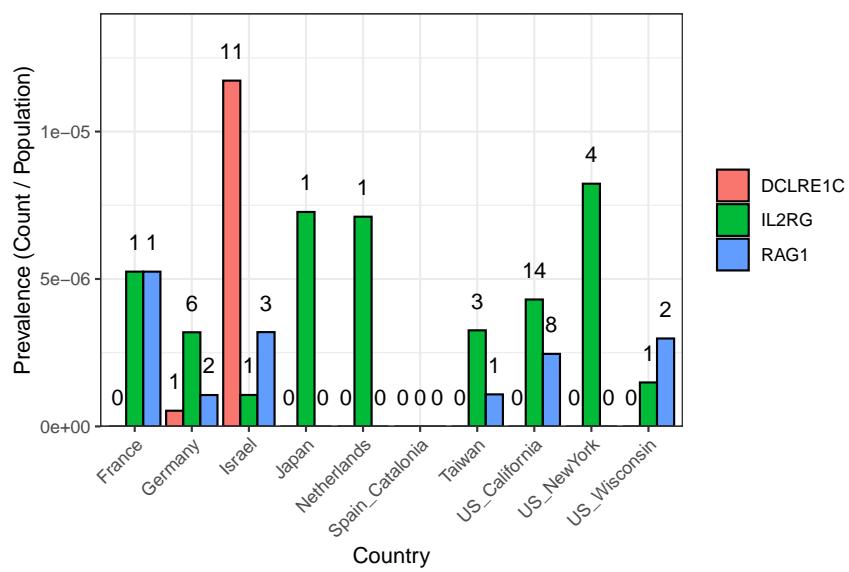


Figure S5: SCID-specific gene comparison across regions. The bar plot shows the prevalence of SCID-related cases (count divided by population) for each gene and country (or region), with numbers printed above the bars representing the actual counts in the original cohort (ranging from 0 to 11 per region and gene).

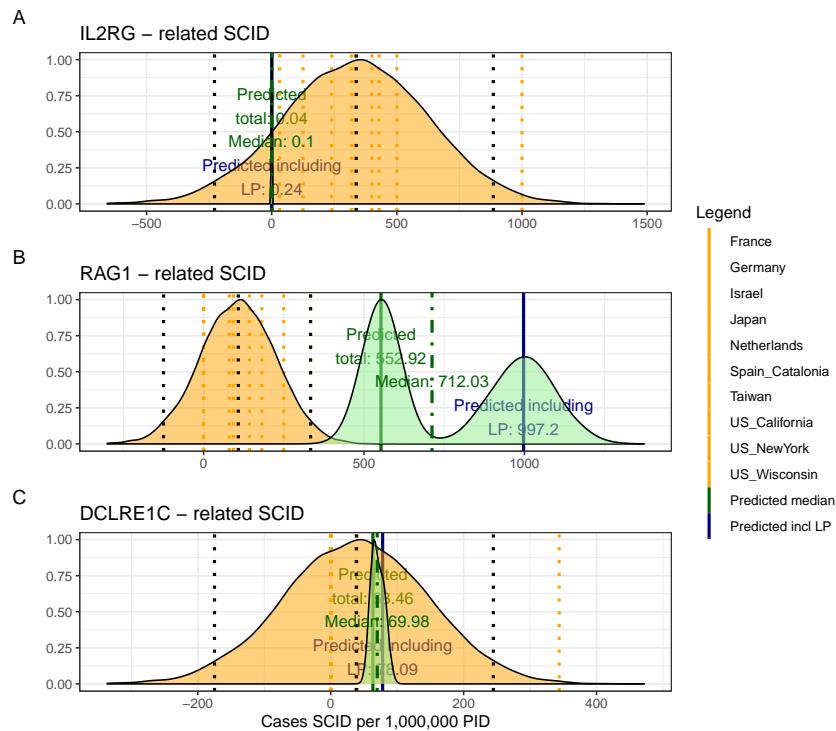


Figure S6: Combined SCID-specific Predictions and Observed Rates per 1,000,000 PID. The figure presents density distributions for the predicted SCID case counts (per 1,000,000 PID) for three genes: *IL2RG*, *RAG1*, and *DCLRE1C*. Country-specific rates (displayed as dotted vertical lines) are overlaid 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.

¹¹⁶⁴ **6.3 Genetic constraint in high-impact protein networks**

¹¹⁶⁵ **6.3.1 Score-positive-total within IEI PPI network**

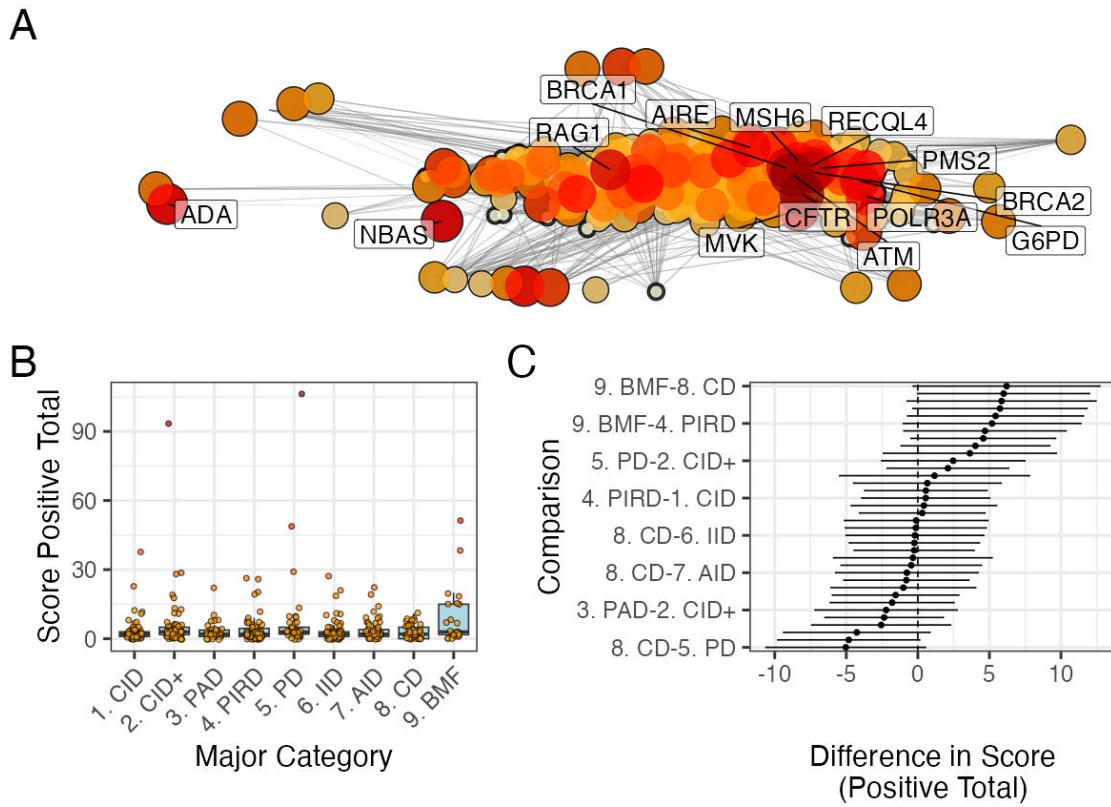


Figure S7: PPI network and score-positive-total ClinVar significance variants. (A) PPI network of disease-associated genes. Node size and colour represent the log-transformed score-positive-total, the top 15 genes/proteins with the highest probability of being observed in disease are labelled. (B) Distribution 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.

¹¹⁶⁶ **6.3.2 Hierarchical Clustering of Enrichment Scores for Major Disease Categories**

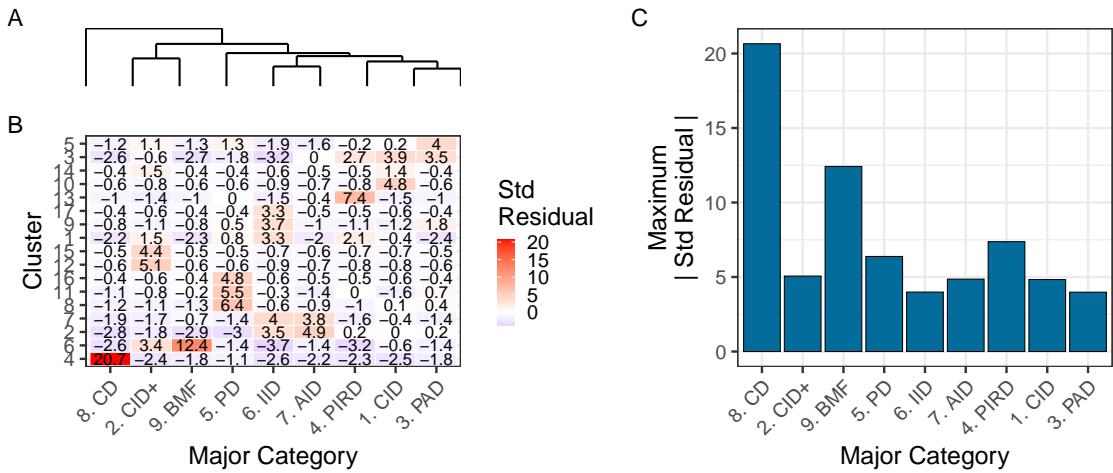


Figure S8: Hierarchical clustering of enrichment scores. The heatmap displays standardised residuals for major disease categories (x-axis) across network clusters (y-axis). A dendrogram groups similar disease categories, and the bar plot shows the maximum absolute residual per category. (8) CD and (9)BMF show the highest values, indicating significant enrichment or depletion (residuals $> |2|$). Definitions in **Box 2.1**.

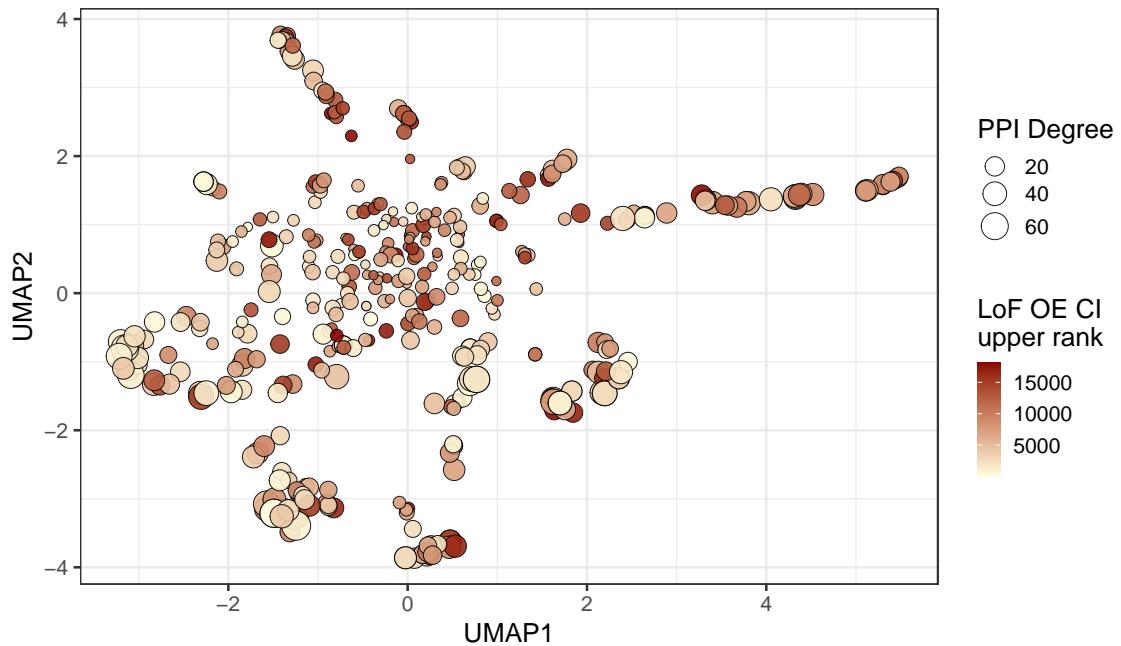


Figure S9: Analysis of PPI degree versus LOEUF upper rank with UMAP embedding of the PPI network. The relationship between PPI degree (size) and LOEUF upper rank (color) across gene clusters. No clear patterns are evident.

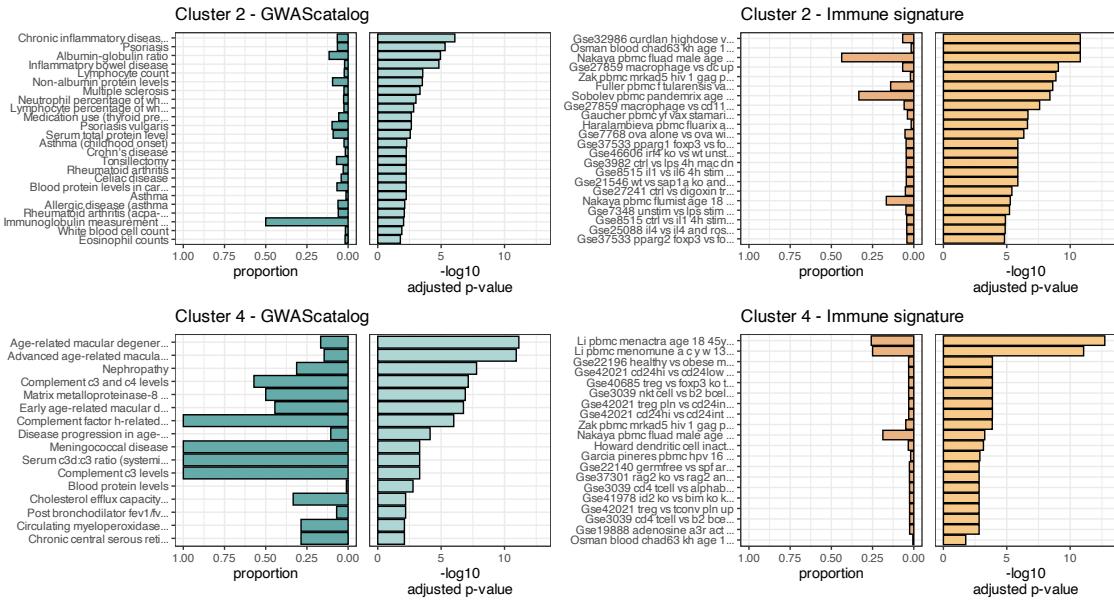


Figure S10: Composite Enrichment Profiles for IEI Gene Sets. We selected the top two enriched clusters (as per [Figure S12](#)) and performed functional enrichment analysis derived from known disease associations. For each gene set, the left panel displays the proportion of input genes overlapping with 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.

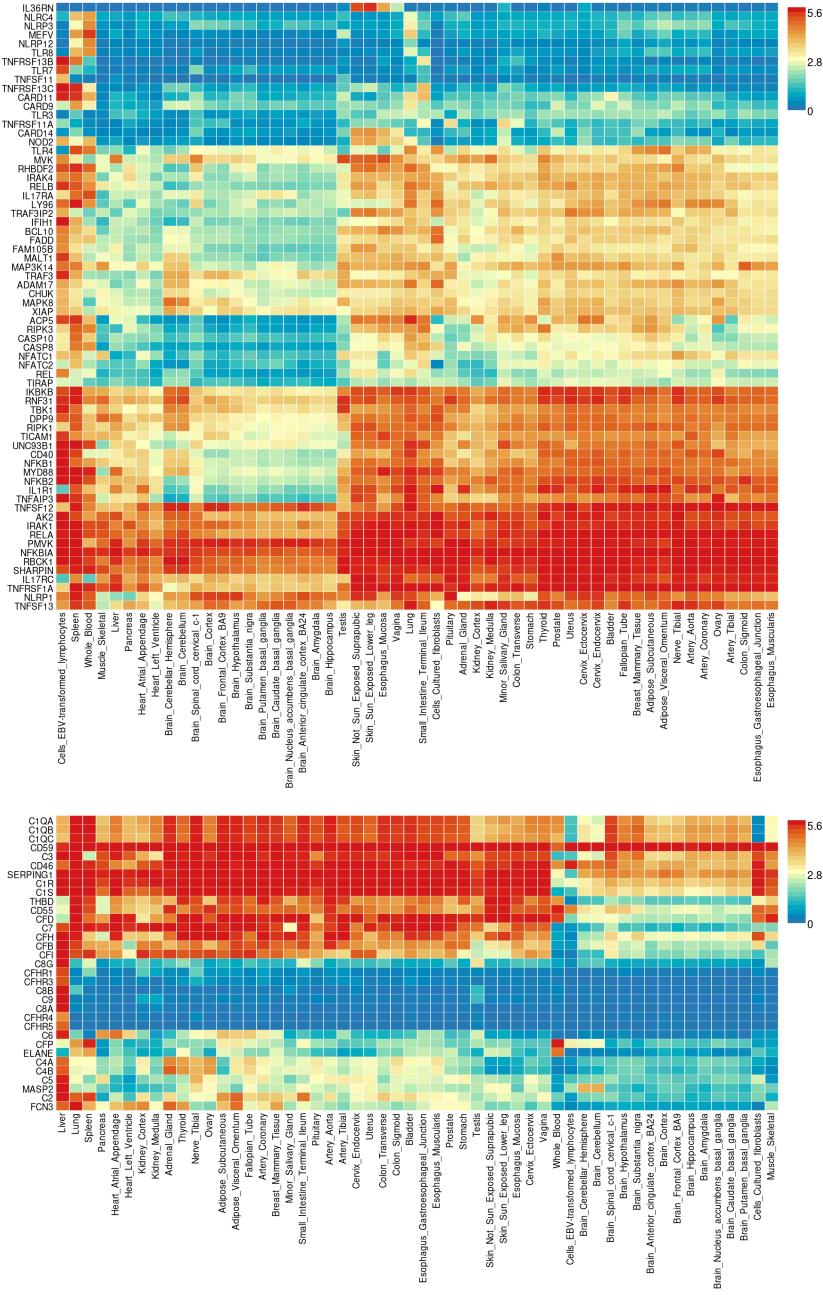


Figure S11: **Gene Expression Heatmaps for IEI Genes.** GTEx v8 data from 54 tissue types display the average expression per tissue label (log₂ transformed) for the IEI gene panels. Top: Cluster 2; Bottom: Cluster 4.

1168 **6.3.3 PPI connectivity, LOEUF constraint and enriched network cluster**
 1169 **analysis**

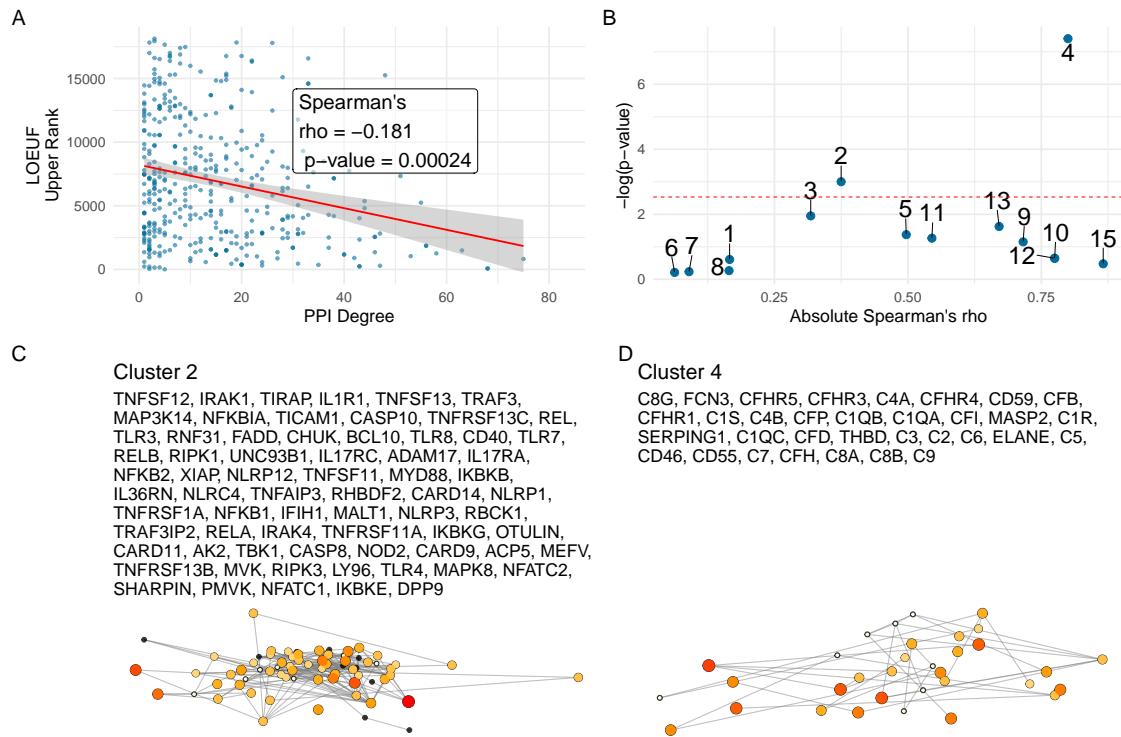


Figure S12: **Correlation between PPI degree and LOEUF upper rank.** (A) Ananlysis across all genes revealed a weak, significant negative correlation between PPI degree and LOEUF upper rank. (B) The cluster-wise analysis showed that clusters 2 and 4 exhibited moderate to strong correlations, while other 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.

1170 **6.4 Interpretation of ClinVar Variant Observations**

Recessive and Dominant Disease Genes

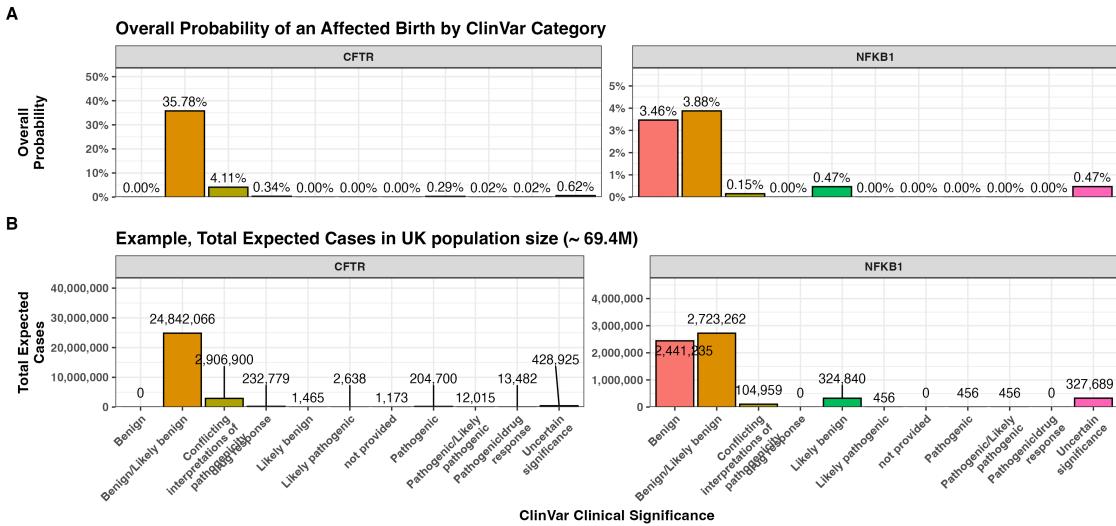


Figure S13: Combined bar charts summarising the genome-wide analysis of ClinVar clinical significance for the PID gene panel. Panel (A) shows the overall probability of an affected birth by variant classification, and (B) displays the total expected number of cases per classification, both stratified by gene. These integrated results illustrate the variability in variant observations across genes and underpin our validation of the probability estimation framework.

¹¹⁷¹ **6.5 Novel PID classifications**

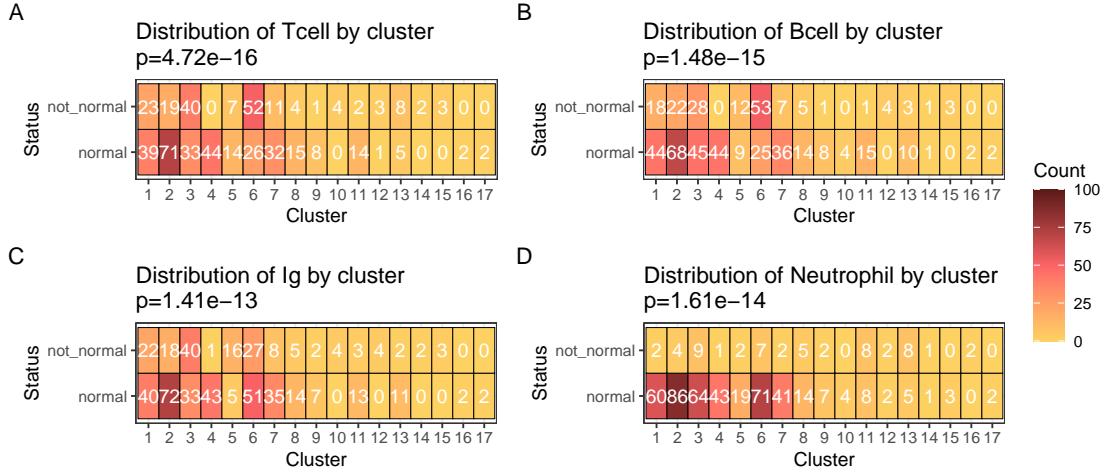


Figure S14: Heatmaps of clinical feature distributions by PPI cluster. The heatmaps display the count of observations for abnormality of each clinical feature (A) T cell, (B) B cell, (C) Immunoglobulin, (D) Neutrophil, in relation to the PPI clusters, with p-values from chi-square tests annotated in the titles.

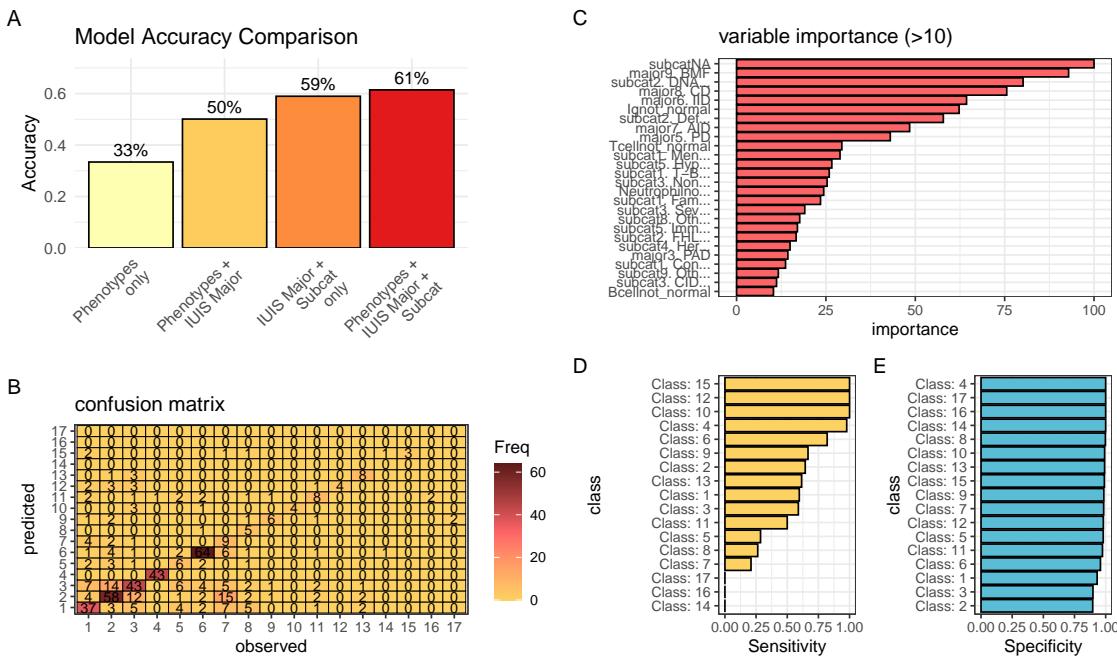


Figure S15: Performance comparison of PID classifiers. Classification predicting PPI cluster membership from IUIS major category, subcategory, and immunological features. (A) Overall accuracy for four rpart models used to predict PPI clustering. The combined model achieves 61.4 % accuracy, exceeding 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.

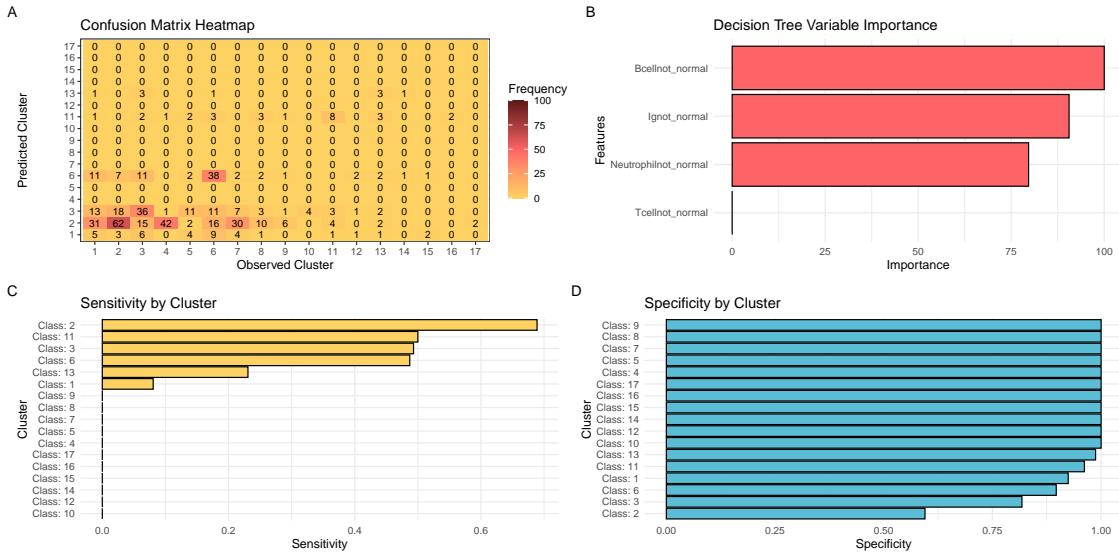


Figure S16: Model performance for fine-tuned PID classification. (A) Confusion matrix heatmap comparing observed and predicted PPI clusters. (B) Variable importance plot ranking immunophenotypic features contributing to the classifier. (C) Per-class sensitivity and (D) per-class specificity bar plots. These panels collectively demonstrate the performance of the decision tree classifier in stratifying PID genes based on immunophenotypic and PPI features.

₁₁₇₂ 6.6 Probability of observing AlphaMissense pathogenicity

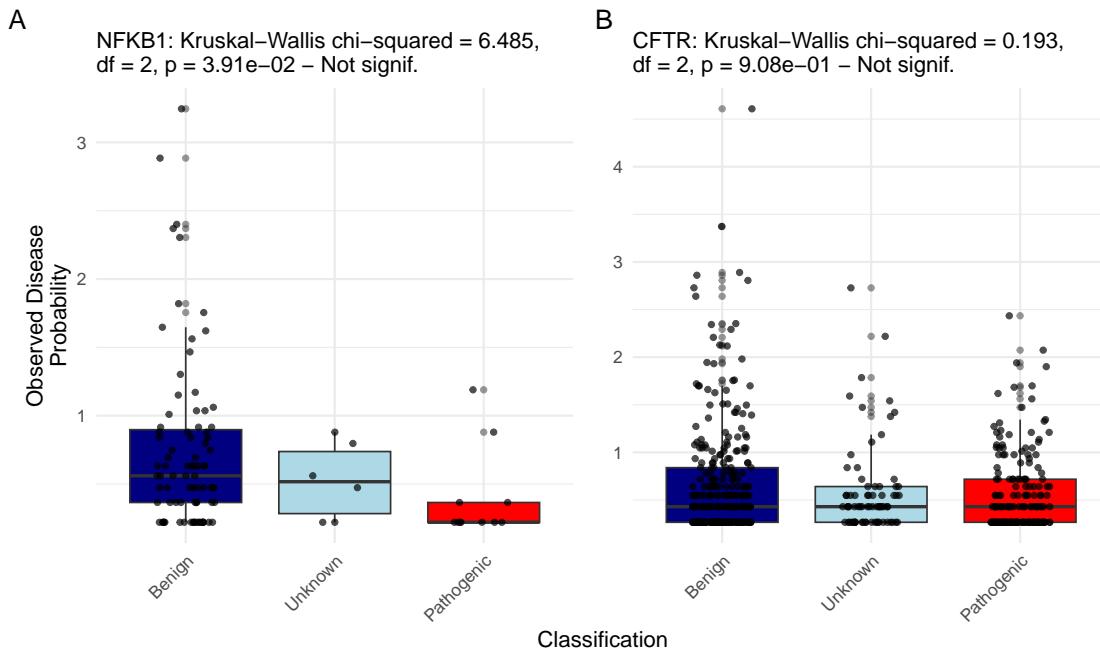


Figure S17: **Observed Disease Probability by Clinical Classification with AlphaMissense.** The figure displays the Kruskal–Wallis test results for NFKB1 and CFTR, showing no significant differences.