

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

Quant Group¹, Simon Boutry², Ali Saadat², Maarja Soomann³, Johannes Trück³, D. Sean Froese⁴, Jacques Fellay², Sinisa Savic⁵, Luregn J. Schlapbach⁶, and Dylan Lawless *⁶

¹The quantitative omic epidemiology group.

²Global Health Institute, School of Life Sciences, École Polytechnique Fédérale de Lausanne, Switzerland.

³Division of Immunology and the Children's Research Center, University Children's Hospital Zurich, University of Zurich, Zurich, Switzerland.

⁴Division of Metabolism and Children's Research Center, University Children's Hospital Zürich, University of Zurich, Zurich, Switzerland.

⁵Leeds Institute of Rheumatic and Musculoskeletal Medicine, University of Leeds, Leeds, UK.

⁶Department of Intensive Care and Neonatology, University Children's Hospital Zurich, University of Zurich, Zurich, Switzerland.

July 23, 2025

Abstract

Background: Accurate interpretation of genetic variants requires quantifying the probability that a variant is disease-causing, including the possibility of alternative or unobserved causal alleles.

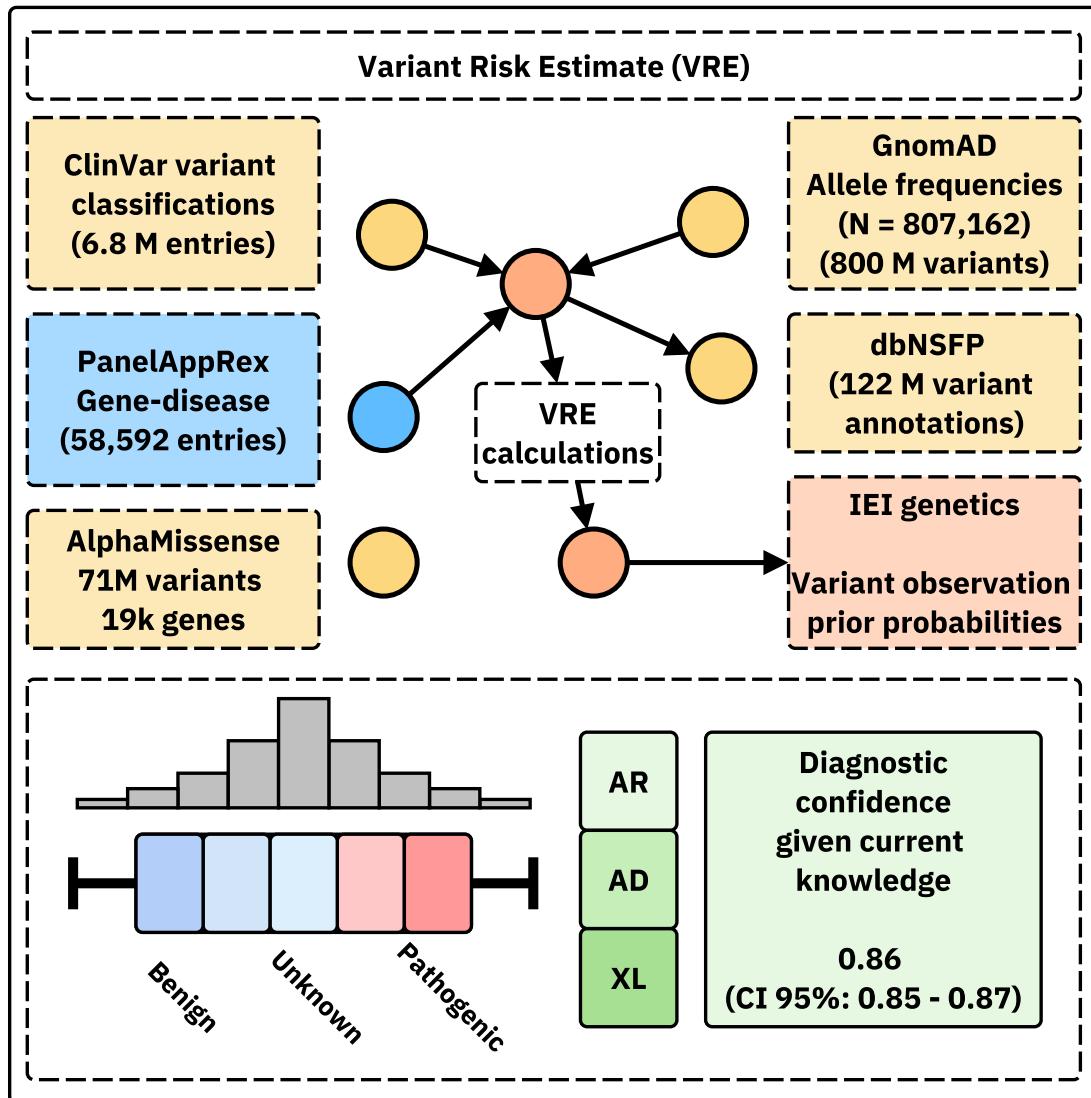
Methods: We developed a statistical framework that computes genome-wide prior probabilities for variant classification, integrating population allele frequencies, disease classifications, and Hardy-Weinberg expectations across inheritance modes. Bayesian modelling then combines these priors with a subject's data to produce credible intervals reflecting diagnostic confidence.

Results: We demonstrated the framework in three known diagnostic scenarios, showing how it quantifies the probability that a candidate variant is causal, including residual uncertainty from unobserved alleles. We then applied it to 557 genes implicated in inborn errors of immunity (IEI), generating variant-level probabilities now publicly accessible. Integration with protein-protein interaction networks and immunophenotypic data revealed patterns of genetic constraint and refined disease classification. Validation in national cohorts showed close agreement between predicted and observed case numbers.

Conclusions: Our method addresses a long-standing gap in clinical genomics by quantifying both observed and unobserved genetic evidence in disease diagnosis. Although demonstrated in IEI, it is broadly applicable and provides a quantitative basis for variant interpretation, clinical decision-making, and future genomic analyses. ¹

¹ * Addresses for correspondence: Dylan.Lawless@kispi.uzh.ch.

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.



25

26 Graphical abstract.

27 Acronyms

28	ACMG American College of Medical Genetics and Genomics	51
29		
30	ACAT Aggregated Cauchy Association Test	51
31		
32	AD Autosomal Dominant.....	7
33		
34	AF Allele Frequency	9
35		
36	ANOVA Analysis of Variance.....	25
37		
38	AR Autosomal Recessive	7
39		
40	BMF Bone Marrow Failure.....	38
41		
42	CD Complement Deficiencies.....	39
43		
44	CI Confidence Interval	36
45		
46	CrI Credible Interval	18
47		
48	CF Cystic Fibrosis.....	22
49		
50	CFTR Cystic Fibrosis Transmembrane Conductance Regulator	10
51		
52	CVID Common Variable Immunodeficiency.....	20
53		
54	DCLRE1C DNA Cross-Link Repair 1C	10
55		
56	dbNSFP database for Non-Synonymous Functional Predictions	10
57		
58	GE Genomics England	8
59		
60	gnomAD Genome Aggregation Database.....	10
61		
62	gVCF genomic variant call format.....	18
63		
64	HGVS Human Genome Variation Society	10
65		
66	HPC High-Performance Computing.....	14
67		

68	HSD Honestly Significant Difference	25
69		
70	HWE Hardy-Weinberg Equilibrium.....	7
71		
72	IEI Inborn Errors of Immunity	8
73		
74	Ig Immunoglobulin	42
75		
76	IL2RG Interleukin 2 Receptor Subunit Gamma.....	10
77		
78	InDel Insertion/Deletion	10
79		
80	IUIS International Union of Immunological Societies	8
81		
82	LD Linkage Disequilibrium.....	41
83		
84	LOEUF Loss-Of-function Observed/Expected Upper bound Fraction.....	25
85		
86	LOF Loss-of-Function.....	25
87		
88	MOI Mode of Inheritance	7
89		
90	NFKB1 Nuclear Factor Kappa B Subunit 1.....	10
91		
92	OMIM Online Mendelian Inheritance in Man	47
93		
94	PID Primary Immunodeficiency	8
95		
96	PPI Protein-Protein Interaction.....	10
97		
98	pLI Probability of being Loss-of-function Intolerant.....	25
99		
100	QC Quality Control	18
101		
102	RAG1 Recombination activating gene 1.....	10
103		
104	SCID Severe Combined Immunodeficiency	10
105		
106	SNV Single Nucleotide Variant	7
107		
108	SKAT Sequence Kernel Association Test	51
109		
110	STRINGdb Search Tool for the Retrieval of Interacting Genes/Proteins....	10

111	TP true positive.....	6
113		
114	FP false positive	6
115		
116	TN true negative.....	6
117		
118	FN false negative	6
119		
120	TNFAIP3 Tumor necrosis factor, alpha-induced protein 3.....	11
121		
122	UMAP Uniform Manifold Approximation and Projection	26
123		
124	UniProt Universal Protein Resource	8
125		
126	VCF variant call format.....	18
127		
128	VEP Variant Effect Predictor	10
129		
130	VRE variant risk estimate.....	11
131		
132	XL X-Linked.....	7
133		

1 Introduction

134 Accurately determining the probability that a patient harbours a disease-causing
135 genetic variant remains a foundational challenge in clinical and statistical ge-
136 netics. For over a century, the primary focus has been on identifying true
137 positive (TP)s, pathogenic causal variants observed in affected individuals.
138 Peer review and classification frameworks also work to suppress false posi-
139 tive (FP)s. However, two critical components of the genetic landscape have
140 received far less attention: false negative (FN)s, where pathogenic variants
141 are missed due to technical or interpretive limitations, and true negative (TN)s,
142 which represent the vast majority of benign or non-causal variants. TNs are
143 more commonly used in contexts such as cancer screening, where a nega-
144 tive result can provide reassurance that a panel of known actionable variants
145 has been checked. Yet outside these specific uses, their broader statistical
146 and clinical value is rarely leveraged. From a statistical perspective, FNs and
147

148 TNs are an untapped goldmine. They hold essential information about what
149 is not observed, what should be expected under baseline assumptions, and
150 how confident one can be in the absence of a pathogenic finding. Yet these di-
151 mensions are rarely quantified, leaving a bias in current variant interpretation
152 frameworks towards known TPs and lacking principled priors for genome-wide
153 disease probability estimation.

154 Quantifying the risk that a patient inherits a disease-causing variant is a
155 fundamental challenge in genomics. Classical statistical approaches grounded
156 in Hardy-Weinberg Equilibrium (HWE) (1; 2) have long been used to calculate
157 genetic probabilities for Single Nucleotide Variant (SNV)s. However, applying
158 these methods becomes more complex when accounting for different Mode
159 of Inheritance (MOI), such as Autosomal Recessive (AR) versus Autosomal
160 Dominant (AD) or X-Linked (XL) disorders. In AR conditions, for example, the
161 occurrence probability must incorporate both the homozygous state and com-
162 pound heterozygosity, whereas for AD and XL disorders, a single pathogenic
163 allele is sufficient to cause disease. Advances in genetic research have re-
164 vealed that MOI can be even more complex (3). Mechanisms such as domi-
165 nant negative effects, haploinsufficiency, mosaicism, and digenic or epistatic
166 interactions can further modulate disease risk and clinical presentation, un-
167 derscoring the need for nuanced approaches in risk estimation. Karczewski
168 et al. (4) made significant advances; however, the remaining challenge lies
169 in applying the necessary statistical genomics data across all MOI for any
170 gene-disease combination. Preliminary approaches have been reported for
171 diseases such Wilson disease, mucopolysaccharidoses, primary ciliary dys-
172 inesia, and treatable metabolic disease, (5; 6), as reviewed by Hannah et al.
173 (7).

174 To our knowledge, all approaches to date have been limited to single MOI,
175 specific to the given disease, or restricted to a small number of genes. We
176 argue that an integrated approach is both necessary and highly powerful be-
177 cause the resulting probabilities can serve as informative priors in a Bayesian
178 framework for variant and disease probability estimation; a perspective that
179 is often overlooked in clinical and statistical genetics. Such a framework not
180 only refines classical HWE-based risk estimates but also has the potential to
181 enrich clinicians' understanding of what to expect in a patient and to enhance
182 the analytical models employed by bioinformaticians.

183 The resulting dataset from these necessary calculations also holds value
184 for AI and reinforcement learning applications, providing an enriched version
185 of the data underpinning frameworks such as AlphaFold (8) and AlphaMis-
186 sense (9).

187 This gap is not only due to conceptual limitations, but to the historical ab-
188 sence of large, harmonised reference datasets. Only recently have resources
189 become available to support rigorous genome-wide probability estimation. These
190 include high-resolution population allele frequencies (e.g. gnomAD v4 (4)),
191 curated variant classifications (e.g. ClinVar (10)), functional annotations (e.g.
192 UniProt (11)), and pathogenicity prediction models (e.g. AlphaMissense (9)).
193 We previously introduced PanelAppRex to aggregate gene panel data from
194 multiple sources, including Genomics England (GE) PanelApp, ClinVar, and
195 Universal Protein Resource (UniProt), thereby enabling advanced natural searches
196 for clinical and research applications (10–13). This earlier work relied on expert-
197 curated panels, such as those from the NHS National Genomic Test Directory
198 and the 100,000 Genomes Project, converted into machine-readable formats
199 for rapid variant discovery and interpretation. Together, these resources now
200 make it possible to model the expected distribution of variant types, frequen-
201 cies, and classifications across the genome.

202 By reframing variant interpretation as a problem of calibrated expectation
203 rather than solely reactive confirmation, our framework empowers clinicians
204 and researchers to anticipate both observed and unobserved pathogenic bur-
205 dens. This scalable, genome-wide approach promises to streamline diagnos-
206 tic workflows, reduce uncertainty in inconclusive cases, inform statistical mod-
207 els and genetic epidemiology studies, and accelerate the integration of ge-
208 netic insights into patient care.

209 In this study, we focused on reporting the probability of disease observa-
210 tion through genome-wide assessments of gene-disease combinations. We
211 then focused on known Inborn Errors of Immunity (IEI) genes, sometimes called
212 the “Primary Immunodeficiency (PID) or Monogenic Inflammatory Bowel Dis-
213 ease genes” (12–14), to validate our approach and demonstrate its clinical
214 relevance. This application to a well-established genotype-phenotype set,
215 comprising over 500 gene-disease associations, underscores its utility. The
216 most recent update on the classification of IEI from the International Union
217 of Immunological Societies (IUIS) expert committee was reported by Poli et al.

²¹⁸ (14), with an accompanying diagnostic guide (15). Our central hypothesis was
²¹⁹ that by using highly curated annotation data including population Allele Fre-
²²⁰ quency (AF)s, disease phenotypes, MOI patterns, and variant classifications
²²¹ and by applying rigorous calculations based on HWE, we could accurately
²²² estimate the expected probabilities of observing disease-associated variants.
²²³ Among other benefits, this knowledge can be used to derive genetic diagnosis
²²⁴ confidence by incorporating these new priors.

225 **2 Methods**

226 **2.1 Dataset**

227 Data from Genome Aggregation Database (gnomAD) v4 comprised 807,162
228 individuals, including 730,947 exomes and 76,215 genomes (4). This dataset
229 provided 786,500,648 SNVs and 122,583,462 Insertion/Deletion (InDel)s, with
230 variant type counts of 9,643,254 synonymous, 16,412,219 missense, 726,924
231 nonsense, 1,186,588 frameshift and 542,514 canonical splice site variants.
232 ClinVar data were obtained from the variant summary dataset (as of: 16 March
233 2025) available from the NCBI FTP site, and included 6,845,091 entries, which
234 were processed into 91,319 gene classification groups and a total of 38,983
235 gene classifications; for example, the gene *A1BG* contained four variants clas-
236 sified as likely benign and 102 total entries (10). For our analysis phase we also
237 used database for Non-Synonymous Functional Predictions (dbNSFP) which
238 consisted of a number of annotations for 121,832,908 SNVs (16). The Pan-
239 elAppRex core model contained 58,592 entries consisting of 52 sets of an-
240 notations, including the gene name, disease-gene panel ID, diseases-related
241 features, confidence measurements. (12) Protein-Protein Interaction (PPI)
242 network data was provided by Search Tool for the Retrieval of Interacting
243 Genes/Proteins (STRINGdb), consisting of 19,566 proteins and 505,968 in-
244 teractions (17). The Human Genome Variation Society (HGVS) nomenclature
245 is used with Variant Effect Predictor (VEP)-based codes for variant IDs. Al-
246 phaMissense includes pathogenicity prediction classifications for 71 million
247 variants in 19 thousand human genes (9; 18). We used these scores to com-
248 pared against the probability of observing the same given variants. **Box 2.1** list
249 the definitions from the IUIS IEI for the major disease categories used through-
250 out this study (14).

251 The following genes were used for disease cohort validations and exam-
252 ples. We used the two most commonly reported genes from the IEI panel Nu-
253 clear Factor Kappa B Subunit 1 (*NFKB1*) (19–22) and Cystic Fibrosis Transmem-
254 brane Conductance Regulator (*CFTR*) (23–25) to demonstrate applications in
255 AD and AR disease genes, respectively. We used Severe Combined Immuno-
256 deficiency (SCID)-specific genes AR DNA Cross-Link Repair 1C (*DCLRE1C*),
257 AR Recombination activating gene 1 (*RAG1*), XL Interleukin 2 Receptor Subunit

258 Gamma (*IL2RG*) to demonstrate a IEI subset disease phenotype of SCID. We
259 also used AD Tumor necrosis factor, alpha-induced protein 3 (*TNFAIP3*) for
260 other examples comparable to *NFKB1* since it is also causes AD pro-inflammatory
261 disease but has more known ClinVar classifications 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

262

2.2 Variant classification occurrence probability

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

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

267 where p is the allele frequency, $q = 1 - p$, p^2 represents the homozygous dominant,
268 $2pq$ the heterozygous, and q^2 the homozygous recessive genotype frequencies.
269 For disease phenotypes, particularly under AR MOI, the risk is traditionally linked
270 to the homozygous state (p^2); however, to account for compound heterozygosity
271 across multiple variants, we allocated the overall gene-level risk proportionally among
272 variants.

273 Our computational pipeline estimated the probability of observing a disease-

275 associated genotype for each variant and aggregated these probabilities by
276 gene and ClinVar classification. This approach included all variant classifi-
277 cations, not limited solely to those deemed “pathogenic”, and explicitly con-
278 ditioned the classification on the given phenotype, recognising that a variant
279 could only be considered pathogenic relative to a defined clinical context. The
280 core calculations proceeded as follows:

281 **1. Allele frequency and total variant frequency.** For each variant i in a
282 gene, the allele frequency was denoted as p_i . For each gene (any genomic
283 region or set), we defined the total variant frequency (summing across all re-
284 ported variants in that gene) as:

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

285 Note that, because each calculation is confined to one gene, no additional
286 scaling was required for our primary analyses (P_{tot}). However, if this same
287 unscaled summation is applied across regions or variant sets of differing size
288 or dosage sensitivity, it can bias burden estimates. In such cases, normali-
289 sation by region length or incorporation of gene- or region-specific dosage
290 constraints is recommended.

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

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

293 where $\max(AN)$ was the maximum allele number observed for that gene.
294 This adjustment ensured that a nonzero risk was incorporated even in the ab-
295 sence of observed variants in the reference database.

296 **2. Occurrence probability based on MOI.** The probability that an individ-
297 ual is affected by a variant depends on the MOI. For **AD** and **XL** variants, a
298 single pathogenic allele suffices:

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

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

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

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

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

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

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

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

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

306 This ensures

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

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

309 More simply, for AD or XL conditions a single pathogenic allele suffices,
310 so the classification risk (e.g. benign, pathogenic) equals its population fre-
311 quency. For AR conditions two pathogenic alleles are required – either two
312 copies of the same variant or one copy each of two different variants, so we
313 divide the overall recessive risk among variants according to each variant's
314 share of the total classification frequency in that gene.

315 **3. Expected case numbers and case detection probability.** Given a
316 population with N births (e.g. as seen in our validation studies, $N = 69\,433\,632$),

317 the expected number of cases attributable to variant i was calculated as:

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

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

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

320 **4. Aggregation by gene and ClinVar classification.** For each gene and
321 for each ClinVar classification (e.g. "Pathogenic", "Likely pathogenic", "Uncer-
322 tain significance", etc.), we aggregated the results across all variants. The
323 classification grouping can be substituted by any alternative score system.
324 The total expected cases for a given group was:

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

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

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

327 **5. Data processing and implementation.** We implemented the calcula-
328 tions within a High-Performance Computing (HPC) pipeline and provided an
329 example for a single dominant disease gene, *TNFAIP3*, in the source code to
330 enhance reproducibility. Variant data were imported in chunks from the anno-
331 tation database for all chromosomes (1-22, X, Y, M).

332 For each data chunk, the relevant fields were gene name, position, allele
333 number, allele frequency, ClinVar classification, and HGVS annotations. Miss-
334 ing classifications (denoted by ".") were replaced with zeros and allele fre-
335 quencies were converted to numeric values. Subsequently, the variant data
336 were merged with gene panel data from PanelAppRex to obtain the disease-
337 related MOI mode for each gene. For each gene, if no variant was observed

338 for a given ClinVar classification (i.e. $p_i = 0$), a minimal risk was assigned as
339 described above. Finally, we computed the occurrence probability, expected
340 cases, and the probability of observing at least one case of disease using the
341 equations presented.

342 The final results were aggregated by gene and ClinVar classification and
343 used to generate summary statistics that reviewed the predicted disease ob-
344 servation probabilities. We define the *VRE* as the prior probability of observing
345 a variant classified as the cause of disease

346 **6. Score-positive-total.** For use as a simple summary statistic on the re-
347 sulting user-interface, we defined the *score-positive-total* as the total num-
348 ber of positively scored variant classifications within a given region (gene,
349 locus, or variant set). Using the ClinVar classification assigned to a scale
350 from -5 (benign) to +5 (pathogenic), we included only scores > 0 , correspond-
351 ing to some evidence of pathogenicity. The score-positive-total yields a non-
352 normalised estimate of the prior probability that a phenotype is explained by
353 known pathogenic variants.

354 **7. Classification scoring system.** Each ClinVar classification was assigned
355 an integer score: pathogenic = +5, likely pathogenic = +4, pathogenic (low
356 penetrance) = +3, likely pathogenic (low penetrance) = +2, conflicting pathogenic-
357 ity = +2, likely risk allele/risk factor/association = +1, drug response/uncertain
358 significance/no classification/affects/other/not provided/uncertain risk allele
359 = 0, protective = -3, likely benign = -4, benign = -5. No further normalisa-
360 tion was applied. The resulting distribution (**Figure S1 A-B**) is naturally com-
361 parable to a zero-centred average rank (**C-D**). This straightforward, modular
362 approach can be readily replaced by any comparable evidence-based classi-
363 fication system. Variants with scores ≤ 0 were omitted, since benign classifi-
364 cations do not inform disease likelihood in the score-positive-total summary.

365 **2.3 Bayesian framework for posterior probability of genetic 366 diagnosis**

367 We developed a Bayesian framework to estimate the probability that at least
368 one variant in a given genome or target set is both present and clinically rel-

369 event in a proband. The method combines variant-specific priors derived
 370 from any classification system or scoring scheme with a probabilistic model of
 371 genotype presence, incorporating both observed and unsequenced genomic
 372 positions.

373 **Prior specification.** Each variant i was assigned a prior probability p_i re-
 374 reflecting its likelihood of being relevant to the user's target interpretation. This
 375 could include, for example, being classified as pathogenic, benign, or uncer-
 376 tain significance, or being prioritised by a quantitative scoring scheme. The
 377 priors may be derived from allele frequency, expert curation, functional predic-
 378 tion, or any other source of variant-level evidence. The framework supports
 379 modular priors without constraint. In our implementation, priors are derived
 380 from ClinVar classification, gnomAD population frequency, and MOI specified
 381 with PanelAppRex gene-disease curation. We modelled the latent causality of
 382 each variant as a Beta distribution parameterised by

$$\alpha_i = \text{round}(p_i \cdot A_N) + w, \quad \beta_i = A_N - \text{round}(p_i \cdot A_N) + 1,$$

383 where A_N is the total number of callable alleles at the site (e.g. $A_N = 2n$
 384 for n diploid individuals), and w is a small stabilising constant (typically $w = 1$).
 385 This construction reflects the expected frequency of the variant in the tested
 386 cohort under the assumption of causal relevance.

387 **Posterior simulation.** We drew $N = 10,000$ samples from the Beta prior to
 388 approximate a posterior distribution for each variant:

$$\theta_i^{(m)} \sim \text{Beta}(\alpha_i, \beta_i), \quad m = 1, \dots, N.$$

389 In each simulation round, these values were normalised across all variants
 390 to reflect the relative share of the total causal signal:

$$\tilde{\theta}_i^{(m)} = \frac{\theta_i^{(m)}}{\sum_j \theta_j^{(m)}}, \quad \hat{P}_i = \frac{1}{N} \sum_{m=1}^N \tilde{\theta}_i^{(m)}.$$

391 This yields a posterior expectation of variant-level contribution, which can
 392 be thresholded or ranked depending on the application.

393 **Genotype presence model.** Independently from the classification or inter-
394 pretive context, each variant was assigned a probability $q_i \in [0, 1]$ of being
395 present in the proband. For observed alternate genotypes, $q_i = 1$; for confi-
396 dently called reference genotypes, $q_i = 0$; and for unsequenced or uncertain
397 sites, q_i was set to the prior probability p_i , or an alternative estimate reflecting
398 presence likelihood. For each simulation round, we drew a genotype indicator:

$$G_i^{(m)} \sim \text{Bernoulli}(q_i),$$

399 and computed the total interpreted variant probability:

$$T^{(m)} = \sum_i G_i^{(m)} \cdot \tilde{\theta}_i^{(m)}.$$

400 The empirical distribution of $\{T^{(m)}\}$ represents the posterior belief that the
401 proband harbours at least one variant matching the specified interpretive crite-
402 ria in the given gene or set. Summary statistics, such as the posterior median
403 and 95% credible interval, are used to report this probability with interval-
404 based accountability.

405 **Flexibility and implementation.** This framework generalises across differ-
406 ent scoring systems, prior models, and sequencing contexts. It can incorpo-
407 rate uncertain classifications, variant weights, unsequenced positions, or al-
408 ternative genotype priors without altering the core inference structure. The
409 method is implemented in R with modular inputs and reproducible simulation
410 logic.

411 **2.4 Application to diagnostic scenarios with observed and 412 unobserved variant data**

413 In this section, we detail our approach to integrating sequencing data with
414 prior classification evidence (e.g. pathogenic on ClinVar) to calculate the pos-
415 terior probability of a complete successful genetic diagnosis. Our method is
416 designed to account for possible outcomes of TP, TN, and FN, by first en-
417 suring that all nucleotides corresponding to known variant classifications (be-
418 nign, pathogenic, etc.) have been accurately sequenced. This implies the

419 use of genomic variant call format (gVCF)-style data which refer to variant
420 call format (VCF)s that contain a record for every position in the genome (or
421 interval of interest) regardless of whether a variant was detected at that site
422 or not. Only after confirming that these positions match the reference alleles
423 (or novel unaccounted variants are classified) do we calculate the probability
424 that additional, alternative pathogenic variants (those not observed in the se-
425 quencing data) could be present. Our Credible Interval (Crl) for pathogenicity
426 thus incorporates uncertainty from the entire process, including the tally of
427 TP, TN, and FN outcomes. We ignore the contribution of FPs as a separate
428 task to be tackled in the future.

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

432 **(i) Data pre-processing.** We synthesized an example patient in a disease
433 cohort of 200 cases. We made several scenarios where a causal genetic diag-
434 nosis based on the available data is either simple, difficult, or impossible. Our
435 example focused on a proband two representative genes for AD IEI: *NFKB1*
436 and *TNFAIP3*. All coding and canonical splice-region variants for *NFKB1* were
437 extracted from the gVCF. We assumed a typical Quality Control (QC) sce-
438 nario, where sites corresponding to previously reported pathogenic alleles
439 were checked for read depth ≥ 10 and genotype quality ≥ 20 . Positions that
440 failed this check were labelled *missing*, thus separating true reference calls
441 from non-sequenced or uninformative sequence.

442 **(ii) Evidence mapping and occurrence probability.** PanelAppRex vari-
443 ants were annotated with ClinVar clinical significance. Each label was con-
444 verted to an ordinal evidence score $S_i \in [-5, 5]$ and rescaled to a pathogenic
445 weight $W_i = \text{rescale}(S_i; -5, 5 \rightarrow 0, 1)$. This scoring system can be replaced
446 with any comparable alternative. The HWE-based pipeline of Section 2.2 sup-
447 plied a per-variant occurrence probability p_i . The adjusted prior was

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

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

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

451 where $\tilde{w}_i = \max(1, S_i + 1)$ contributes an additional pseudo-count whenever
452 $S_i > 0$.

453 **(iv) Posterior simulation and aggregation.** For each variant i we drew
454 $M = 10\,000$ realisations $\pi_i^{(m)}$ and normalised within each iteration,

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

455 Variants with $S_i > 4$ were deemed to have evidence as *causal* (pathogenic
456 or likely pathogenic). We note that an alternative evidence score or condi-
457 tional threshold can be substituted for this step. Their mean posterior share
458 $\bar{\pi}_i = M^{-1} \sum_m \tilde{\pi}_i^{(m)}$ and 95% CrI were retained. The probability that a damaging
459 causal allele is physically present was obtained by a second layer:

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

460 with $g_i = 1$ for present variants, $g_i = 0$ for reference calls, and $g_i = p_i$ for
461 missing variants. The gene-level estimate is the median of $\{P^{(m)}\}_{m=1}^M$ and its
462 2.5th/97.5th percentiles.

463 **(v) Scenario analysis.** The three scenarios were explored for a causal ge-
464 netic diagnosis that is either simple, difficult, or impossible given the exist-
465 ing data. The proband spiked data had either: (1) known classified variants
466 only, including only one known TP pathogenic variant, *NFKB1* p.Ser237Ter,
467 (2) inclusion of an additional plausible yet non-sequenced splice-donor al-
468 lele *NFKB1* c.159+1G>A (likely pathogenic) as a FN, and (3) where no known
469 causal variants were present for a patient, one representative variant from

470 each distinct ClinVar classification was selected and marked as unsequenced
471 to emulate a range of putative FNs. The selected variants were: *TNFAIP3*
472 p.Cys243Arg (pathogenic), p.Tyr246Ter (likely pathogenic), p.His646Pro (con-
473 flicting interpretations of pathogenicity), p.Thr635Ile (uncertain significance),
474 p.Arg162Trp (not provided), p.Arg280Trp (likely benign), p.Ile207Leu (benign/-
475 likely benign), and p.Lys304Glu (benign). All subsequent steps were identical.

476 **2.5 Validation of autosomal dominant estimates using *NFKB1***

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

482 **1. Reference dataset.** We used a reference dataset reported by Tuijnen-
483 burg et al. (19) to build a validation model in an AD disease gene. This study
484 performed whole-genome sequencing of 846 predominantly sporadic, unre-
485 lated PID cases from the NIHR BioResource-Rare Diseases cohort. There were
486 390 CVID cases in the cohort. The study identified *NFKB1* as one of the genes
487 most strongly associated with PID. Sixteen novel heterozygous variants in-
488 cluding truncating, missense, and gene deletion variants, were found in *NFKB1*
489 among the CVID cases.

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

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

492 with a 95% confidence interval (using Wilson's method) of approximately
493 (0.0254, 0.0656).

494 **3. National estimate based on literature.** Based on literature (19; 20; 22),
495 the prevalence of CVID in the general population was estimated as

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

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

$$498 E_{\text{CVID}} \approx \frac{69\,433\,632}{25\,000} \approx 2777.$$

499 Assuming that the proportion of CVID cases attributable to *NFKB1* is equiv-
 500 alent to the cohort estimate, the literature extrapolated estimate is Estimated *NFKB1* cases \approx
 501 $2777 \times 0.041 \approx 114$, with a median value of approximately 118 and a 95% confi-
 502 dence interval of 70 to 181 cases (derived from posterior sampling).

503 **4. Bayesian adjustment.** Recognising that the sequenced cohort cases
 504 likely captures the majority of *NFKB1*-related patients (apart from close rel-
 505 atives), but may still miss rare or geographically dispersed variants, we com-
 506 bined the cohort-based and literature-based estimates using two complemen-
 507 tary Bayesian approaches:

508 1. **Weighted adjustment (emphasising the cohort, $w = 0.9$):** We as-
 509 signed 90% weight to the directly observed cohort count (16) and 10%
 510 to the extrapolated population estimate (114), thereby accounting, illus-
 511 tratively, for a small fraction of unobserved cases while retaining confi-
 512 dence in our well-characterised cohort:

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

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

514 2. **Mixture adjustment (equal weighting, $w = 0.5$):** To reflect greater
 515 uncertainty about how representative the cohort is, we combined cohort
 516 and population prevalences equally. We sampled from the posterior dis-
 517 tribution of the cohort prevalence,

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

518 and mixed this with the literature-based rate at 50% each (19; 20; 22).
 519 This yields a median estimate of 67 cases and a wider 95% CrI of approxi-

520 mately 43 to 99 cases, capturing uncertainty in both under-ascertainment
521 and over-generalisation.

522 **5. Predicted total genotype counts.** The predicted total synthetic geno-
523 type count (before adjustment) was 456, whereas the predicted total geno-
524 types adjusted for synth_flag was 0. This higher synthetic count was set
525 based on a minimal risk threshold, ensuring that at least one genotype is as-
526 sumed to exist (e.g. accounting for a potential unknown de novo variant) even
527 when no variant is observed in gnomAD (as per [section 2.2](#)).

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

532 **2.6 Validation study for autosomal recessive CF using *CFTR***

533 To validate our framework for AR diseases, we focused on Cystic Fibrosis
534 (CF). For comparability sizes between the validation studies, we analysed
535 the most common SNV in the *CFTR* gene, typically reported as p.Arg117His
536 (GRCh38 Chr 7:117530975 G/A, MANE Select HGVSp ENST00000003084.11:
537 p.Arg117His). Our goal was to validate our genome-wide probability estimates
538 by comparing the expected number of CF cases attributable to the p.Arg117His
539 variant in *CFTR* with the nationally reported case count in a well-characterised
540 disease cohort ([23–25](#)).

541 **1. Expected genotype counts.** Let p denote the allele frequency of the
542 p.Arg117His variant and q denote the combined frequency of all other pathogenic
543 *CFTR* variants, such that

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

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

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

546 and

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

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

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

549 **2. Mortality adjustment.** Since CF patients experience increased mortality,
550 we adjusted the expected genotype counts using an exponential survival
551 model (23–25). With an annual mortality rate $\lambda \approx 0.004$ and a median age of
552 22 years, the survival factor was computed as:

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

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

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

554 **3. Bayesian uncertainty simulation.** To incorporate uncertainty in the al-
555 lele frequency p , we modelled p as a beta-distributed random variable:

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

556 using a large effective allele count (AN_{eff}) for illustration. By generating
557 10,000 posterior samples of p , we obtained a distribution of the literature-
558 based adjusted expected counts, E_{adj} .

559 **4. Bayesian Mixture Adjustment.** Since the national registry may not cap-
560 ture all nuances (e.g., reduced penetrance) of *CFTR*-related disease, we fur-
561 ther combined the literature-based estimate with the observed national count
562 (714 cases from the UK Cystic Fibrosis Registry 2023 Annual Data Report)
563 using a 50:50 weighting:

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

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

568 **2.7 Validation of SCID-specific estimates using PID–SCID
569 genes**

570 To validate our genome-wide probability estimates for diagnosing a genetic
571 variant in a patient with an PID phenotype, we focused on a subset of genes
572 implicated in SCID. Given that the overall panel corresponds to PID, but SCID
573 represents a rarer subset, the probabilities were converted to values per mil-
574 lion PID cases.

575 **1. Incidence conversion.** Based on literature, PID occurs in approximately
576 1 in 1,000 births, whereas SCID occurs in approximately 1 in 100,000 births.
577 Consequently, in a population of 1,000,000 births there are about 1,000 PID
578 cases and 10 SCID cases. To express SCID-related variant counts on a per-
579 million PID scale, the observed SCID counts were multiplied by 100. For ex-
580 ample, if a gene is expected to cause SCID in 10 cases within the total PID
581 population, then on a per-million PID basis the count is $10 \times 100 = 1,000$ cases
582 (across all relevant genes).

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

589 **3. Integration with prior probability estimates.** The predicted genotype
590 occurrence probabilities were derived from our framework across the PID gene

591 panel. These probabilities were then converted to expected case counts per
592 million PID cases by multiplying by 1,000,000. For instance, if the probability
593 of observing a pathogenic variant in *IL2RG* is p , the expected SCID-related
594 count becomes $p \times 10^6$. Similar conversions are applied for all relevant SCID
595 genes.

596 **4. Bayesian Uncertainty and Comparison with Observed Data.** To ad-
597 dress uncertainty in the SCID-specific estimates, a Bayesian uncertainty sim-
598 ulation was performed for each gene to generate a distribution of predicted
599 case counts on a per-million PID scale. The resulting median estimates and
600 95% Crls were then compared against known national SCID counts compiled
601 from independent registries. This comparison permuted a direct evaluation
602 of our framework's accuracy in predicting the occurrence of SCID-associated
603 variants within a PID cohort.

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

607 **2.8 Protein network and genetic constraint interpretation**

608 A PPI network was constructed using protein interaction data from STRINGdb
609 (17). We previously prepared and reported on this dataset consisting of 19,566
610 proteins and 505,968 interactions (<https://github.com/DylanLawless/ProteoMCLustR>).
611 Node attributes were derived from log-transformed score-positive-total val-
612 ues, which informed both node size and colour. Top-scoring nodes (top 15
613 based on score) were labelled to highlight prominent interactions. To eval-
614 uate group differences in score-positive-total across major disease categories,
615 one-way Analysis of Variance (ANOVA) was performed followed by Tukey Hon-
616 estly Significant Difference (HSD) post hoc tests (and non-parametric Dunn's
617 test for confirmation). GnomAD v4.1 constraint metrics data was used for
618 the PPI analysis and was sourced from Karczewski et al. (4). This provided
619 transcript-level metrics, such as observed/expected ratios, Loss-Of-function
620 Observed/Expected Upper bound Fraction (LOEUF), Probability of being Loss-
621 of-function Intolerant (pLI), and Z-scores, quantifying Loss-of-Function (LOF)

622 and missense intolerance, along with confidence intervals and related anno-
623 tations for 211,523 observations.

624 **2.9 Gene set enrichment test**

625 To test for overrepresentation of biological functions, the prioritised genes
626 were compared against gene sets from MsigDB (including hallmark, positional,
627 curated, motif, computational, GO, oncogenic, and immunologic signatures)
628 and WikiPathways using hypergeometric tests with FUMA (26; 27). The back-
629 ground set consisted of 24,304 genes. Multiple testing correction was applied
630 per data source using the Benjamini-Hochberg method, and gene sets with an
631 adjusted P-value ≤ 0.05 and more than one overlapping gene are reported.

632 **2.10 Deriving novel PID classifications by genetic PPI and
633 clinical features**

634 We recategorised 315 immunophenotypic features from the original IUIS IEI
635 annotations, reducing the original multi-level descriptors (e.g. "decreased
636 CD8, normal or decreased CD4") first to minimal labels (e.g."low") and second
637 to binary outcomes (normal vs. not-normal) for T cells, B cells, neutrophils,
638 and immunoglobulins Each gene was mapped to its PPI cluster derived from
639 STRINGdb and Uniform Manifold Approximation and Projection (UMAP) em-
640 beddings from previous steps. We first tested for non-random associations
641 between these four binary immunophenotypes and PPI clusters using χ^2 tests.
642 To generate a data-driven PID classification, we trained a decision tree (rpart)
643 to predict PPI cluster membership from the four immunophenotypic features
644 plus the traditional IUIS Major and Subcategory labels. Hyperparameters (com-
645 plexity parameter = 0.001, minimum split = 10, minimum bucket = 5, maximum
646 depth = 30) were optimised via five-fold cross validation using the caret frame-
647 work. Terminal node assignments were then relabelled according to each
648 group's predominant abnormal feature profile.

649 **2.11 Probability of observing AlphaMissense pathogenicity**
650

651 We obtained the subset pathogenicity predictions from AlphaMissense via the
652 AlphaFold database and whole genome data from the studies data repository(9;
653 18). The AlphaMissense data (genome-aligned and amino acid substitutions)
654 were merged with the panel variants based on genomic coordinate and HGVS
655 annotation. Occurrence probabilities were log-transformed and adjusted (y-
656 axis displaying $\log_{10}(\text{occurrence prob} + 1e-5) + 5$), to visualise the distri-
657 bution of pathogenicity scores across the residue sequence. A Kruskal-Wallis
658 test was used to compare the observed disease probability across clinical clas-
659 sification groups.

660 **2.12 Probability model definitions**

661 Estimating disease risk requires accounting for both variant penetrance, $P(D |$
662 $G)$, where D denotes the disease state and G the genotype, and the fraction
663 of cases attributable to a given variant, $P(G | D)$. In a fully penetrant single-
664 variant model ($P(D | G) = 1$), the lifetime risk $P(D)$ equals the genotype fre-
665 quency $P(G)$. For an allele with population frequency p , this gives $P(D) = p^2$
666 for a recessive mode of inheritance and $P(D) = 2p(1 - p) \approx 2p$ for a dominant
667 mode. With incomplete penetrance, $P(D) = P(G) P(D | G)$, and when multiple
668 variants contribute to disease,

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

669 Because both $P(D | G)$ and $P(G | D)$ are often uncertain, we integrate Clin-
670 Var clinical classifications, population allele frequencies and curated gene-
671 disease associations, assuming James-Stein shrinkage to derive robust ag-
672 gregate priors. By focusing on a filtered set of variants \mathcal{V} where each $P(G_i |$
673 $D)$ is the probability that disease D is attributable to variant i and assuming
674 $\sum_{i \in \mathcal{V}} P(G_i | D) \approx 1$, we obtain calibrated estimates of genotype frequency
675 $P(G)$ despite uncertainty in individual parameters.

676 3 Results

677 3.1 Occurrence probability across disease genes

678 Our study integrated large-scale annotation databases with gene panels from
679 PanelAppRex to systematically assess disease genes by MOI (12). By com-
680 bining population allele frequencies with ClinVar clinical classifications, we
681 computed an expected occurrence probability for each SNV, representing the
682 likelihood of encountering a variant of a specific pathogenicity for a given phe-
683 notype. We report these probabilities for 54,814 ClinVar variant classifications
684 across 557 genes (linked dataset (28)).

685 We focused on panels related to Primary Immunodeficiency or Monogenic
686 Inflammatory Bowel Disease, using PanelAppRex panel ID 398. **Figure 1** dis-
687 plays all reported ClinVar variant classifications for this panel. The resulting
688 natural scaling system (-5 to +5) accounts for the frequently encountered
689 combinations of classification labels (e.g. benign to pathogenic). The result-
690 ing dataset (28) is briefly shown in **Table 1** to illustrate that our method yielded
691 estimates of the probability of observing a variant with a particular ClinVar
692 classification.

Table 1: **Example of the first several rows from our main results for 557 genes of PanelAppRex's panel: (ID 398) Primary immunodeficiency or monogenic inflammatory bowel disease.** "ClinVar Significance" indicates the pathogenicity classification assigned by ClinVar, while "Occurrence Prob" represents our 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. (28)

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
...

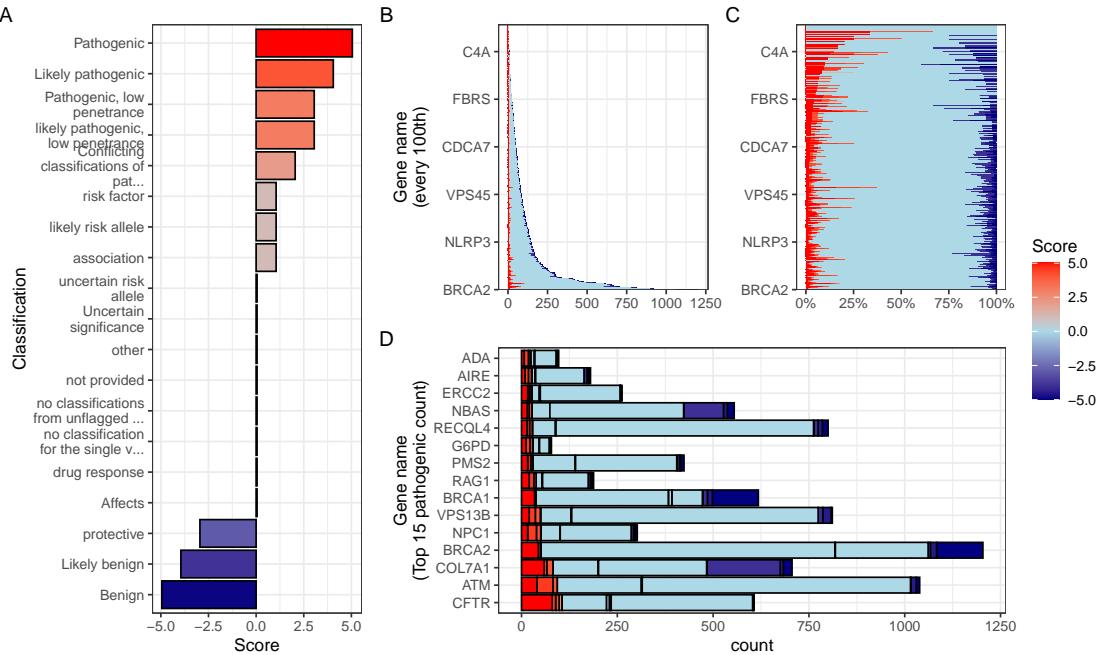


Figure 1: Summary of ClinVar clinical significance classifications in the PID gene panel. (A) Shows the numeric score coding for each classification (i.e. -5 benign to +5 pathogenic) as defined in methods Section 2.2. (B) Displays the stacked absolute count of classifications per gene. The same counts are shown as percentages per gene in (C). (D) For demonstration, the top 15 genes ranked by absolute count of pathogenic (score 5) variant classifications, indicating those most frequently occurring in the population as disease-causing.

693 3.2 Integrating observed true positives and unobserved 694 false negatives into a single, actionable conclusion

695 Having established a probabilistic framework for estimating the prior prob-
696 ability of observing disease-associated variants under different inheritance
697 modes, we then applied this model to an example patient to demonstrate
698 its potential for clinical genetics. The algorithm first verified that all known
699 pathogenic positions have been sequenced and observed as reference (true
700 negatives), and identified any positions that were either observed as vari-
701 ant (true positives) or not assessable due to missing sequence data or failed
702 QC. These missing sites represent potential false negatives. By jointly mod-
703 elling the observed and unobserved space, the method yielded a calibrated,
704 evidence-weighted probability that at least one damaging causal variant could

705 be present within a gene.

706 707 **3.3 Scenario one - complete coverage and simple diagno-** **sis**

708 We present the results from three scenarios for an example single-case pa-
709 tient being investigated for the genetic diagnosis of IEI. **Figure S2** shows
710 the results of the first simple scenario, in which only one known pathogenic
711 variant, *NFKB1* p.Ser237Ter, was observed and all other previously reported
712 pathogenic positions were successfully sequenced and confirmed as refer-
713 ence. In this setting, the model assigned the full posterior probability to the ob-
714 served allele, yielding 100 % confidence that all present evidence supported
715 a single, true positive causal explanation. The most strongly supported ob-
716 served variant was p.Ser237Ter (posterior: 0.594). The strongest (probability
717 of observing) non-sequenced variant was a benign variant p.Thr567Ile (poste-
718 rior: 0). The total probability of a causal diagnosis given the available evidence
719 was 1 (95% CI: 1–1) (**Table S1**).

720 721 **3.4 Scenario two - incomplete coverage and complex di- agnosis**

722 **Figure 2** shows the second more complex scenario, where the same pathogenic
723 variant *NFKB1* p.Ser237Ter was present, but coverage was incomplete at three
724 additional sites of known classified variants. Among these was the likely-
725 pathogenic splice-site variant *NFKB1* c.159+1G>A, which was not captured in
726 the sequencing data. The panels of **Figure 2 (A–F)** illustrate the stepwise
727 integration of observed and missing evidence, culminating in a posterior prob-
728 ability that reflects both confirmed findings and residual uncertainty. **Table**
729 **2** demonstrates our main goal and lists the final conclusion for reporting the
730 clinical genetics results. **Table S2** lists the main metrics used to reach the
731 conclusion (as illustrated in **Figure 2**).

732 Bayesian integration of every annotated allele yielded the quantitative CrIs
733 for pathogenic attribution that (i) preserve Hardy-Weinberg expectations, (ii)
734 accommodate AD, AR, XL inheritance, and (iii) carry explicit uncertainty for

735 non-sequenced (or failed QC) genomic positions. The incremental calculation
736 steps for the variant in question are shown in **Figure 2**.

737 First, **Figure 2 (A)** depicts the prior landscape where occurrence probabilities
738 are partitioned by observed or missing status and by causal or non-causal
739 evidence, with colour reflecting the underlying ClinVar score. **Figure 2 (B)**
740 shows posterior normalisation which concentrates probability density on two
741 high-confidence (high evidence score) alleles since the benign variants are, by
742 definition, non-causal. **Figure 2 (C)** shows the resulting per-variant probability
743 of being simultaneously damaging and causal; only the confirmed present
744 (true positive) nonsense variant p.Ser237Ter and the (false negative) hypothetical
745 splice-donor c.159+1G>A retain substantial support. Restricting the view to
746 causal candidates in **Figure 2 (D)** confirms that posterior mass is distributed
747 across these two variants. **Figure 2 (E)** decomposes the total damaging prob-
748 ability into observed (approximately 40 %) and missing (approximately 34 %)
749 sources, whereas **Figure 2 (F)** summarises the gene-level posterior: inclusion
750 of the splice-site allele (scenario 2) produces a median probability of 0.542
751 with a 95 % CrI of 0.264–0.8.

752 Numerically, the present variant p.Ser237Ter accounts for 0.399 of the
753 posterior share, and the potentially causal but missing splice-donor allele
754 c.159+1G>A contributes 0.339. The remaining alleles together explain a negli-
755 gible share (**Table S2**). Thus, we can report that in this patient's scenario the
756 probability of correct genetic diagnosis due to *NFKB1* p.Ser237Ter is 0.542
757 (95 % CrI of 0.264–0.8) given that a likely alternative remains to be confirmed
758 for this patient. Upon confirmation that the second variant is not present, the
759 confidence will rise to 1 (95 % CrI of 1–1) as 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).	

760 3.5 Scenario three - currently impossible diagnosis

761 **Figure S3** shows the third scenario, in which no observed variants were de-
 762 tected in the proband for *NFKB1*. Instead, a broad range of plausible FN were
 763 detected as missing for the gene *TNFAIP3*. The strongest (probability of ob-
 764 serving and pathogenic) of these non-sequenced variants was p.Cys243Arg
 765 (posterior: 0.347). However, the total probability of a causal diagnosis for the
 766 patient *given the available evidence* was 0 (95% CI: 0–0) since these missing
 767 variants must be accounted for (**Table S3**). Upon confirmation, these proba-
 768 bilities can update, as per scenario two.

769 3.6 Posterior probabilities are calculated across all quali- 770 fying variants

771 While only the top-ranked gene/variant is shown in each of the three scenar-
 772 ios, we emphasise that the same posterior probability and CrI calculations are
 773 performed across all plausible candidates. In real-world diagnostics, we com-
 774 monly find multiple variants to carry non-negligible probabilities. Our frame-

775 work explicitly quantifies these competing hypotheses, enabling a ranked in-
776 terpretation that reflects the totality of evidence. Overlapping CrI do not in-
777 dicate ambiguity in the method, but rather a principled measure of remaining
778 uncertainty. This output can directly inform follow-up actions, such as func-
779 tional validation or treatment trials, and supports the use of CrI thresholds as
780 a transparent decision-making aid when data are incomplete or equivocal.

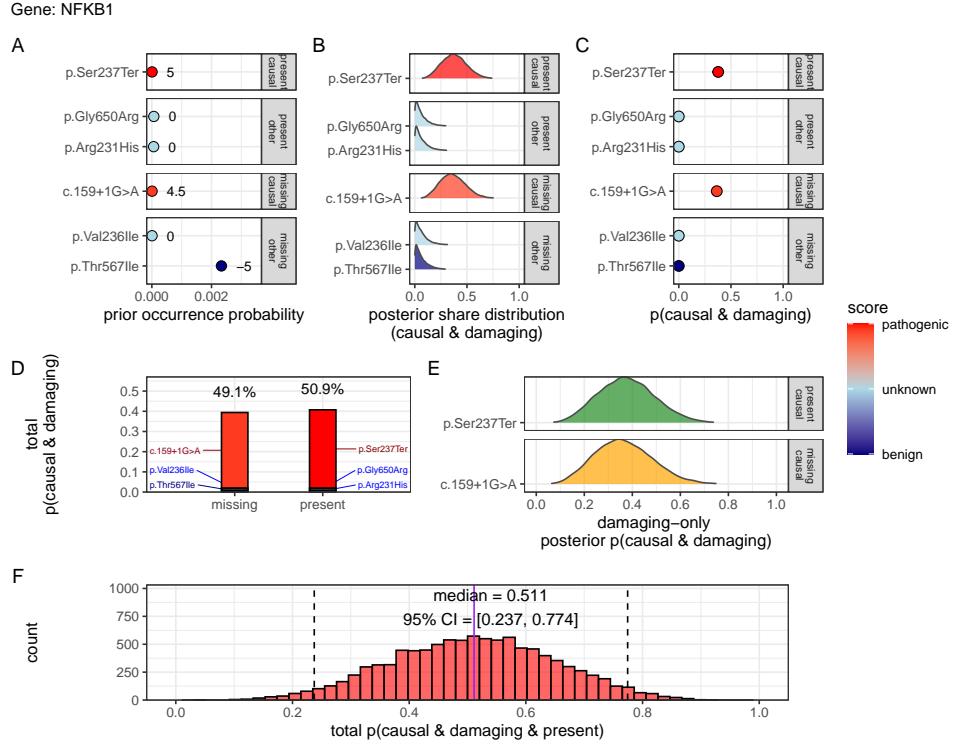


Figure 2: Quantification of present (TP) and missing (FN) causal genetic variants for disease in *NFKB1* (scenario 2). The example proband carried three known heterozygous variants, including pathogenic p.Ser237Ter, and had incomplete coverage at three additional loci, including likely-pathogenic splice-site 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 S2** and **S3**, respectively.

781 **3.7 Validation studies**

782 **3.7.1 Validation of dominant disease occurrence with *NFKB1***

783 To validate our genome-wide probability estimates for AD disorders, we fo-
784 cused on *NFKB1*. We used a reference dataset from Tuijnenburg et al. (19), in
785 which whole-genome sequencing of 846 PID patients identified *NFKB1* as one
786 of the genes most strongly associated with the disease, with 16 *NFKB1*-related
787 CVID cases attributed to AD heterozygous variants. Our goal was to compare
788 the predicted number of *NFKB1*-related CVID cases with the reported count
789 in this well-characterised national-scale cohort.

790 Our model calculated 0 known pathogenic variant *NFKB1*-related CVID cases
791 in the UK with a minimal risk of 456 unknown de novo variants. In the refer-
792 ence cohort, 16 *NFKB1* CVID cases were reported. We additionally wanted to
793 account for potential under-reporting in the reference study. We used an ex-
794 trapolated national CVID prevalence which yielded a median estimate of 118
795 cases (95% CI: 70–181), while a Bayesian-adjusted mixture estimate produced
796 a median of 67 cases (95% CI: 43–99). **Figure S5 (A)** illustrates that our pre-
797 dicted values reflect these ranges and are closer to the observed count. This
798 case supports the validity of our integrated probability estimation framework
799 for AD disorders, and represents a challenging example where pathogenic
800 SNV are not reported in the reference population of gnomAD. Our min-max
801 values successfully contained the true reported values.

802 **3.7.2 Validation of recessive disease occurrence with *CFTR***

803 Our analysis predicted the number of CF cases attributable to carriage of
804 the p.Arg117His variant (either as homozygous or as compound heterozygous
805 with another pathogenic allele) in the UK. Based on HWE calculations and mor-
806 tality adjustments, we predicted approximately 648 cases arising from biallelic
807 variants and 160 cases from homozygous variants, resulting in a total of 808
808 expected cases.

809 In contrast, the nationally reported number of CF cases was 714, as recorded
810 in the UK Cystic Fibrosis Registry 2023 Annual Data Report (23). To account
811 for factors such as reduced penetrance and the mortality-adjusted expected

genotype, we derived a Bayesian-adjusted estimate via posterior simulation. Our Bayesian approach yielded a median estimate of 740 cases (95% Confidence Interval (CI): 696, 786) and a mixture-based estimate of 727 cases (95% CI: 705, 750). **Figure S5 (B)** illustrates the close concordance between the predicted values, the Bayesian-adjusted estimates, and the national report supports the validity of our approach for estimating disease.

Figure S4 shows the final values for these genes of interest in a given population size and phenotype. It reveals that an allele frequency threshold of approximately 0.000007 is required to observe a single heterozygous carrier of a disease-causing variant in the UK population for both genes. However, owing to the AR MOI pattern of *CFTR*, this threshold translates into more than 100,000 heterozygous carriers, compared to only 456 carriers for the AD gene *NFKB1*. Note that this allele frequency threshold, being derived from the current reference population, represents a lower bound that can become more precise as public datasets continue to grow. This marked difference underscores the significant impact of MOI patterns on population carrier frequencies and the observed disease prevalence.

3.7.3 Interpretation of ClinVar variant occurrences

Figure S6 shows the two validation study PID genes, representing AR and dominant MOI. **Figure S6 (A)** illustrates the overall probability of an affected birth by ClinVar variant classification, whereas **Figure S6 (B)** depicts the total expected number of cases per classification for an example population, here the UK, of approximately 69.4 million.

3.7.4 Validation of SCID-specific disease occurrence

Given that SCID is a subset of PID, our probability estimates reflect the likelihood of observing a genetic variant as a diagnosis when the phenotype is PID. However, we additionally tested our results against SCID cohorts in **Figure S8**. The summarised raw cohort data for SCID-specific gene counts are summarised and compared across countries in **Figure S7**. True counts for *IL2RG* and *DCLRE1C* from ten distinct locations yielded 95% CI surrounding our predicted values. For *IL2RG*, the prediction was low (approximately 1 case per 1,000,000 PID), as expected since loss-of-function variants in this XL gene

844 are highly deleterious and rarely observed in gnomAD. In contrast, the pre-
845 dicted value for *RAG1* was substantially higher (553 cases per 1,000,000 PID)
846 than the observed counts (ranging from 0 to 200). We attributed this dis-
847 crepancy to the lower penetrance and higher background frequency of *RAG1*
848 variants in recessive inheritance, whereby reference studies may underreport
849 the true national incidence. Overall, we report that agreement within an order
850 of magnitude was tolerable given the inherent uncertainties from reference
851 studies arising from variable penetrance and allele frequencies.

852 **3.8 Application to inborn errors of immunity**

853 **3.8.1 Genetic constraint in high-impact protein networks**

854 We next applied our framework to IEI, a disease area in which we have exper-
855 tise and which offers a well-curated gene set to validate genome-wide esti-
856 mates and demonstrate potential applications (14).

857 Given that pathogenicity in IEI may reflect shared molecular pathways, we
858 integrated ClinVar-derived variant probability estimates with PPI data to quan-
859 tify pathogenic burden per gene and examine whether genetic constraint ag-
860 gregates within specific networks and corresponds to established IEI classifi-
861 cations and immunophenotypes (4; 17).

862 **3.8.2 Score-positive-total within IEI PPI network**

863 The ClinVar classifications reported in **Figure 1** were scaled -5 to +5 based on
864 their pathogenicity. We were interested in positive (potentially damaging) but
865 not negative (benign) scoring variants, which are statistically incidental in this
866 analysis. We tallied gene-level positive scores to give the score-positive-total
867 metric. **Figure S9 (A)** shows the PPI network of disease-associated genes,
868 where node size and colour encode the score-positive-total (log-transformed).
869 The top 15 genes with the highest total prior probabilities of being observed
870 with disease are labelled (as per **Figure 1**).

871 **3.8.3 Association analysis of score-positive-total across IEI categories**

872 We checked for any statistical enrichment in score-positive-totals, which rep-
873 resents the expected observation of pathogenicity, between the IEI categories.
874 One-way ANOVA revealed an effect of major disease category on score-positive-
875 total ($F(8, 500) = 2.82, p = 0.0046$), indicating that group means were not iden-
876 tical, which we observed in **Figure S9 (B)**. However, despite some apparent
877 differences in median scores across categories (i.e. 9. Bone Marrow Fail-
878 ure (BMF)), the Tukey HSD post hoc comparisons **Figure S9 (C)** showed that
879 all pairwise differences had 95% CIs overlapping zero, suggesting that indi-
880 vidual group differences were not significant.

881 **3.8.4 UMAP embedding of the PPI network**

882 To address the density of the PPI network for the IEI gene panel, we applied
883 UMAP (**Figure 3**). Node sizes reflect interaction degree, a measure of evidence-
884 supported connectivity (17). We tested for a correlation between interaction
885 degree and score-positive-total. In **Figure 3**, gene names with degrees above
886 the 95th percentile are labelled in blue, while the top 15 genes by score-
887 positive-total are labelled in yellow (as per **Figure 1**). Notably, genes with
888 high pathogenic variant loads segregated from highly connected nodes, sug-
889 gesting that LOF in hub genes is selectively constrained, whereas damaging
890 variants in lower-degree genes yield more specific effects. This observation
891 was subsequently tested empirically.

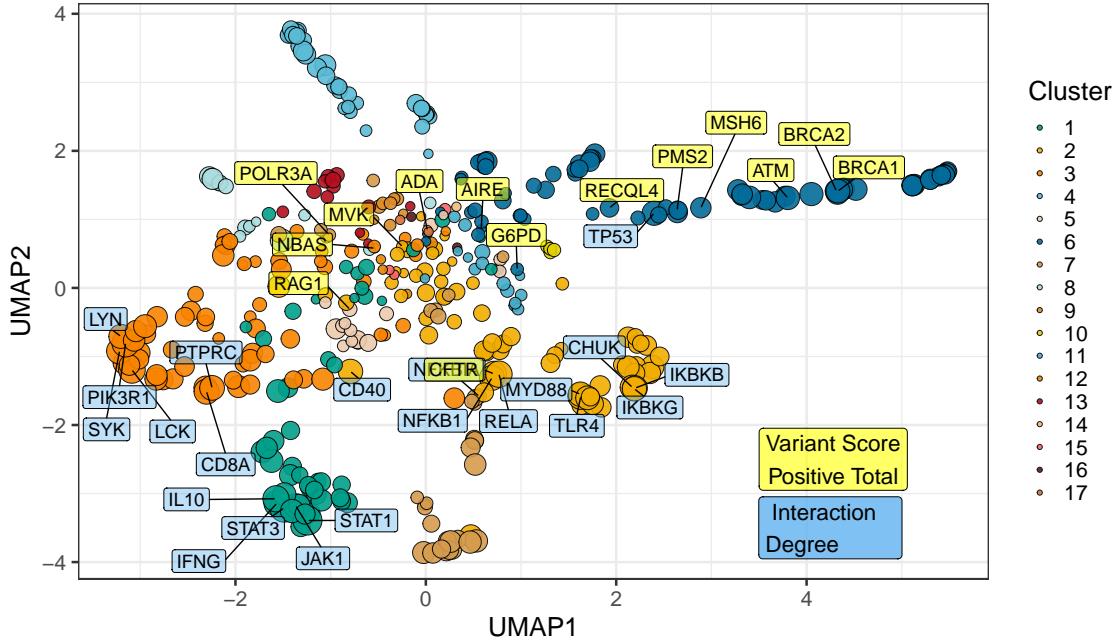


Figure 3: UMAP embedding of the PPI network. The plot projects the high-dimensional protein-protein interaction network into two dimensions, with nodes coloured by cluster and sized by interaction degree. Blue labels indicate hub genes (degree above the 95th percentile) and yellow labels mark the top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

892 3.8.5 Hierarchical clustering of enrichment scores for major disease 893 categories

894 **Figure S10** presents a heatmap of standardised residuals for major disease
895 categories across network clusters, as per **Figure 3**. A dendrogram clusters
896 similar disease categories, while the accompanying bar plot displays the max-
897 imum absolute standardised residual for each category. Notably, (8) Comple-
898 ment Deficiencies (CD) shows the highest maximum enrichment, followed by
899 (9) BMF. While all maximum values exceed 2, the threshold for significance,
900 this likely reflects the presence of protein clusters with strong damaging vari-
901 ant scores rather than uniform significance across all categories (i.e. genes
902 from cluster 4 in 8 CD).

903 **3.8.6 PPI connectivity, LOEUF constraint and enriched network clus-**
904 **ter analysis**

905 Based on the preliminary insight from **Figure S10**, we evaluated the relation-
906 ship between network connectivity (PPI degree) and LOEUF constraint (LOEUF
907 upper rank) Karczewski et al. (4) using Spearman's rank correlation. Overall,
908 there was a weak but significant negative correlation ($\rho = -0.181, p = 0.00024$)
909 at the global scale, indicating that highly connected genes tend to be more
910 constrained. A supplementary analysis (**Figure S11**) did not reveal distinct vi-
911 sual associations between network clusters and constraint metrics, likely due
912 to the high network density. However once stratified by gene clusters, the nat-
913 ural biological scenario based on quantitative PPI evidence (17), some groups
914 showed strong correlations; for instance, cluster 2 ($\rho = -0.375, p = 0.000994$)
915 and cluster 4 ($\rho = -0.800, p < 0.000001$), while others did not. This indicated
916 that shared mechanisms within pathway clusters may underpin genetic con-
917 straints, particularly for LOF intolerance. We observe that the score-positive-
918 total metric effectively summarises the aggregate pathogenic burden across
919 IEI genes, serving as a robust indicator of genetic constraint and highlighting
920 those with elevated disease relevance.

921 **Figure S12 (C, D)** shows the re-plotted PPI networks for clusters with sig-
922 nificant correlations between PPI degree and LOEUF upper rank. In these net-
923 works, node size is scaled by a normalised variant score, while node colour
924 reflects the variant score according to a predefined palette.

925 **3.8.7 New insight from functional enrichment**

926 To interpret the functional relevance of our prioritised IEI gene sets with the
927 highest load of damaging variants (i.e. clusters 2 and 4 in **Figure S12**), we
928 performed functional enrichment analysis for known disease associations us-
929 ing MsigDB with FUMA (i.e. GWAScatalog and Immunologic Signatures) (26).
930 Composite enrichment profiles (**Figure S13**) reveal that our enriched PPI clus-
931 ters were associated with distinct disease-related phenotypes, providing func-
932 tional insights beyond traditional IUIS IEI groupings (14). The gene expres-
933 sion profiles shown in **Figure S14** (GTEx v8 54 tissue types) offer the tissue-
934 specific context for these associations. Together, these results enable the
935 annotation of IEI gene sets with established disease phenotypes, supporting

936 a data-driven classification of IEI.

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

947 **3.8.8 Genome-wide gene distribution and linkage disequilibrium**

948 **Figure 4 (A)** shows a genome-wide karyoplot of all IEI panel genes across
949 GRCh38, with colour-coding based on MOI. Figures **(B)** and **(C)** display zoomed-
950 in locus plots for *NFKB1* and *CFTR*, respectively. In **Figure 4 (B)**, the prob-
951 ability of observing variants with known classifications is high only for vari-
952 ants such as p.Ala475Gly, which are considered benign in the AD *NFKB1* gene
953 that is intolerant to LOF. In **Figure 4 (C)**, high probabilities of observing pa-
954 tients with pathogenic variants in *CFTR* are evident, reproducing this well-
955 established phenomenon. Furthermore, the analysis of Linkage Disequilib-
956 rium (LD) using R^2 shows that high LD regions can be modelled effectively,
957 allowing independent variant signals to be distinguished.

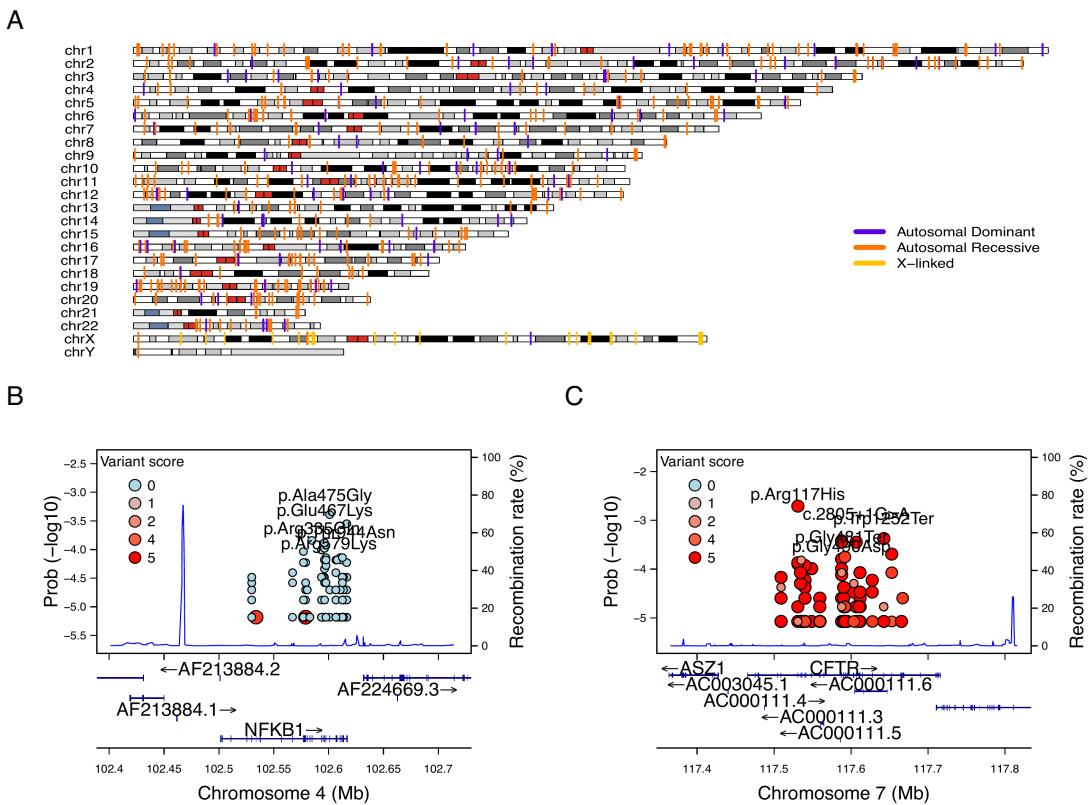


Figure 4: Genome-wide IEI, variant occurrence probability and LD by R^2 .
 (A) Genome-wide karyoplot of all IEI panel genes mapped to GRCh38, with colours indicating MOI. (B) Zoomed-in locus plot example for *NFKB1* showing variant occurrence probabilities; only benign variants such exhibit high probabilities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

958 **3.8.9 Novel PID classifications derived from genetic PPI and clinical
 959 features**

960 We recategorised 315 immunophenotypic features from the original IUIS IEI
 961 annotations, reducing detailed descriptions (e.g. “decreased CD8, normal or
 962 decreased CD4”) to minimal labels (e.g. “low”) and then binarising them (nor-
 963 mal vs. not-normal) for T cells, B cells, Immunoglobulin (Ig) and neutrophils
 964 (**Figure S15**). These simplified profiles were mapped onto STRINGdb PPI clus-
 965 ters, revealing non-random distributions ($\chi^2 < 1e-13$; **Figure S16**), indicating
 966 that network context captures key immunophenotypic variation.

We next compared four classifiers under 5-fold cross-validation to determine which features predicted PPI clustering. As shown in **Figure S17**, the fully combined model achieved the highest accuracy among the four: (i) phenotypes only (33 %) (i.e. T cell, B cell, Ig, Neutrophil); (ii) phenotypes + IUIS major category (50 %) (e.g. CID. See **Box 2.1** for more); (iii) IUIS major + subcategory only (59 %) (e.g. CID, T-B+ SCID); and (iv) phenotypes + IUIS major + subcategory (61 %). This demonstrated that incorporating both traditional IUIS IEI classifications and core immunophenotypic markers into the PPI-based framework yielded the most robust discrimination of PID gene clusters. Variable importance analysis highlighted abnormality status for Ig and T cells were among the top ten features in addition to the other IUIS major and sub categories. Per-class specificity remained uniform across the classes while sensitivity dropped.

The PPI and immunophenotype model yielded 17 data-driven PID groups, whereas incorporating the full complement of IUIS categories expanded this to 33 groups. For clarity, we only demonstrate the decision tree from the smaller 17-group model in **Figure 5**. Each terminal node is annotated by its predominant immunophenotypic signature (for example, "group 65 with abnormal T cell and B cell features"), and the full resulting gene counts per 33 class are plotted in **Figure 5**. Although, less user-friendly, this data-driven taxonomy both aligns with and refines traditional IUIS IEI classifications to provide a scaffold for large-scale computational analyses. Because this framework is fully reproducible, alternative PPI embeddings that incorporate additional molecular annotations can readily swapped to continue building on these IEI classification schemes.

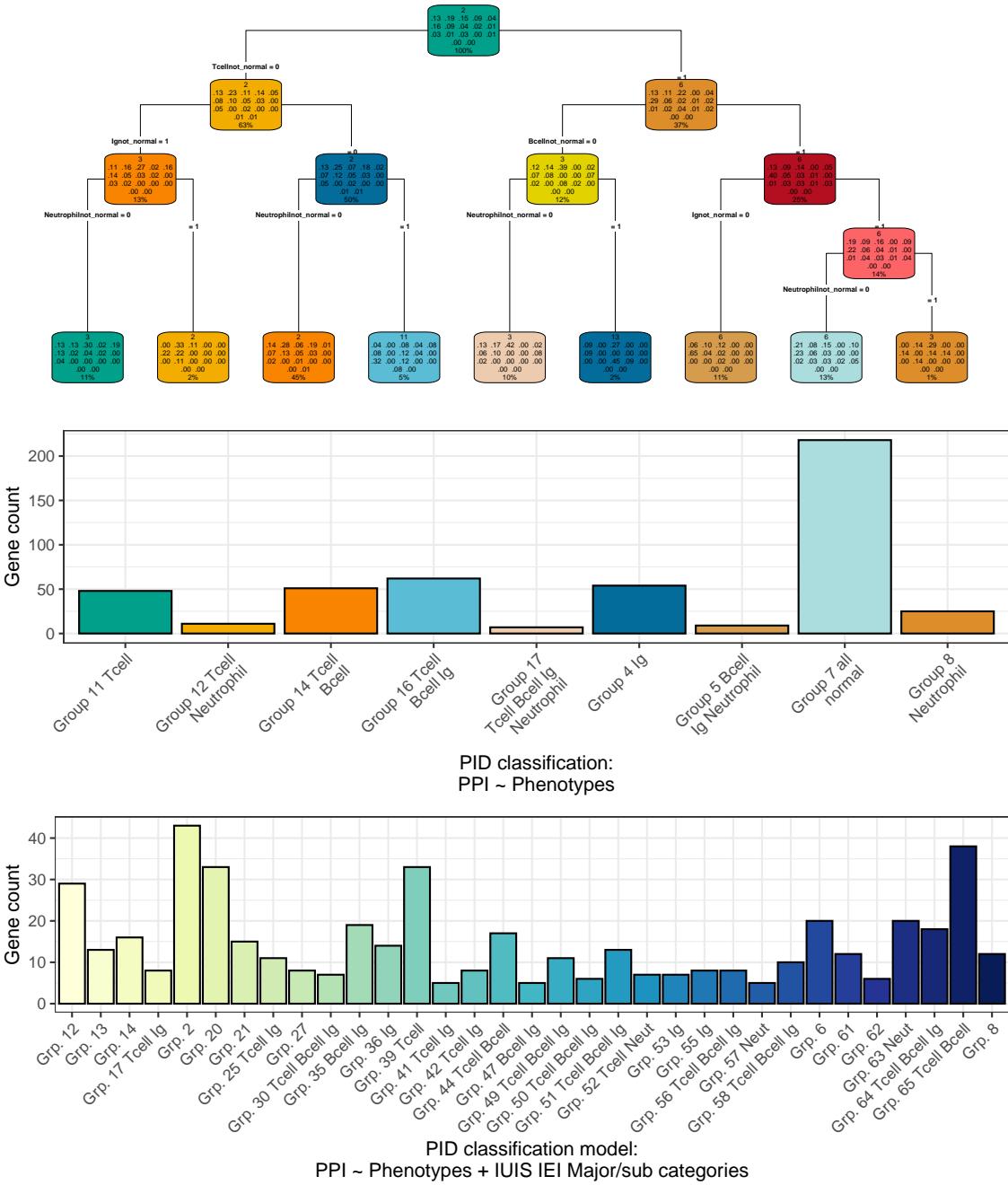


Figure 5: Fine-tuned model for PID classification. (Top) In each terminal node, the top block indicates the number of genes in the node; the middle block shows the fitted class probabilities (which sum to 1); and the bottom block displays the percentage of the total sample in that node. These metrics 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.

992 **3.9 Probability of observing AlphaMissense pathogenic-**
993 **ity**

994 AlphaMissense provides pathogenicity scores for all possible amino acid sub-
995 stitutions; however, our results in **Figure 6** show that the most probable obser-
996 vations in patients occur predominantly for benign or unknown variants. This
997 finding places the likelihood of disease-associated substitutions into perspec-
998 tive and offers a data-driven foundation for future improvements in variant pre-
999 diction. The values in **Figure 6 (A)** can be directly compared to **Figure 1 (D)**
1000 to view the distribution of classifications. A Kruskal-Wallis test was used to
1001 compare the observed disease probability across clinical classification groups
1002 and no significant differences were detected. In general, most variants in pa-
1003 tients are classified as benign or unknown, indicating limited discriminative
1004 power in the current classification, such that pathogenicity prediction does
1005 not infer occurrence prediction (**Figure S18**). Inverse correlation likely de-
1006 pends on factors like MOI and intolerance to LOF.

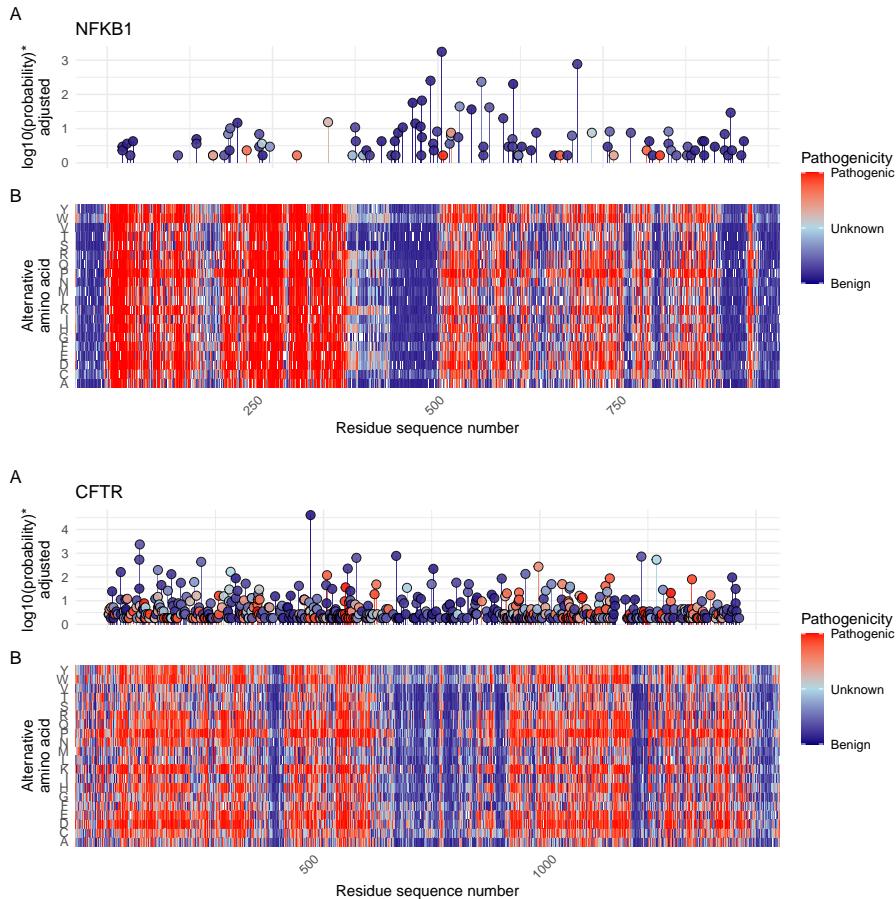


Figure 6: **(A) Probabilities of observing a patient with (B) AlphaMissense-derived pathogenicity scores.** Although AlphaMissense provides scores for every possible amino acid substitution, the most frequently observed variants in patients tend to be classified as benign or of unknown significance. 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.

1007 3.10 Integration of variant probabilities into IEI genetics 1008 data

1009 We integrated the computed prior probabilities for observing variants in all
1010 known genes associated with a given phenotype (14), across AD, AR, and
1011 XL MOI, into our IEI genetics framework. These calculations, derived from
1012 gene panels in PanelAppRex, have yielded novel insights for the IEI disease
1013 panel. The final result comprised of machine- and human-readable datasets,
1014 including the table of variant classifications and priors available via a the linked
1015 repository (28), and a user-friendly web interface that incorporates these new
1016 metrics.

1017 **Figure 7** shows the interface summarising integrated variant data. We in-
1018 clude pre-calculated summary statistics and clinical significance as numer-
1019 ical metrics. Key quantiles (min, Q1, median, Q3, max) for each gene are
1020 rendered as sparkline box plots, and dynamic URLs link table entries to ex-
1021 ternal databases (e.g. ClinVar, Online Mendelian Inheritance in Man (OMIM),
1022 AlphaFold) as per **Section 3.1**. The prepared data are available for bioinfor-
1023 matic application (28) as per **Section 3.2**.

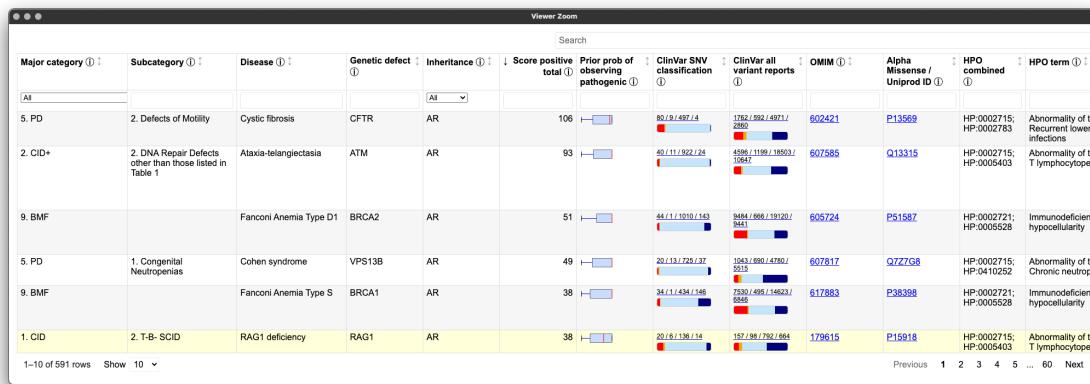


Figure 7: **Integration of variant probabilities into the IEI genetics framework.** The interface summarises the condensed variant data, with pre-calculated summary statistics and dynamic links to external databases. This integration enables immediate access to detailed variant classifications and prior probabilities for each gene.

1024 4 Discussion

1025 Our study presents, to our knowledge, the first comprehensive framework for
1026 calculating prior probabilities of observing disease-associated variants and
1027 the first to demonstrate the method for an evidence-aware genetic diagnosis
1028 with CrI (5; 7). By integrating large-scale genomic annotations, including pop-
1029 ulation allele frequencies from gnomAD (4), variant classifications from Clin-
1030 Var (10), and functional annotations from resources such as dbNSFP, with clas-
1031 sical HWE-based calculations, we derived robust estimates for 54,814 Clin-
1032 Var variant classifications across 557 IEI genes implicated in PID and mono-
1033 genic inflammatory bowel disease (12; 14). Although our results focus on
1034 IEI, the genome-wide framework also supports all inheritance patterns: AD
1035 and XL require a single pathogenic allele, whereas AR demands homozygous
1036 or compound heterozygous states. Classical HWE-based estimates thus fur-
1037 nish baseline occurrence probabilities and serve as robust priors for Bayesian
1038 risk models, a practice underutilised until the advent of large-scale databases
1039 (4; 9; 10; 12).

1040 A major deficit in current clinical genetics is the prevailing focus on con-
1041 firming only the presence of TP variants. Our approach yielded three key
1042 results to overcome this hurdle. We generated per-variant priors across all
1043 MOI. The patient's results of observed and unobserved variants were inte-
1044 grated into a single posterior probability of carrying a damaging causal allele.
1045 As demonstrated in **Table S2** and **Figure 2**, this key result delivers a clinically
1046 applicable, interpretable probability that combines both detected and poten-
1047 tially unobserved variants. When whole-genome sequencing analyses are not
1048 yet available, the score-positive-total metric can serve as an optional decision
1049 aid, enabling manual, evidence-based ranking of candidate genes to prioritise
1050 diagnoses in patients with overlapping phenotypes.

1051 We acknowledge that our framework is currently focused (but not restricted)
1052 on SNVs and does not incorporate numerous other complexities of genetic dis-
1053 ease, such as structural variants, de novo variants, hypomorphic alleles, over-
1054 dominance, variable penetrance, tissue-specific expression, the Wahlund ef-
1055 fect, pleiotropy, and others (3). In certain applications, more refined estimates
1056 would benefit from including factors such as embryonic lethality, condition-
1057 specific penetrance, and age of onset (7). Our analysis also relies on simpli-

1058 fying assumptions of random mating, an effectively infinite population, and
1059 the absence of migration, novel mutations, or natural selection. We demon-
1060 strated the genome-wide gene distribution and MOI for the IEI panel relative
1061 to LD showing that it is an important consideration and is feasible. However,
1062 LD is a challenging feature that requires accurate implementation which de-
1063 pends on the whole genome population-based pairwise genotype matrices for
1064 the given population. We used the reference global population AFs, which is
1065 more generalisable but less accurate than population-specific AF values.

1066 In the example single-case diagnosis scenarios, our approach enabled high-
1067 confidence attribution to a known pathogenic variant while also capturing the
1068 potential impact of a likely-pathogenic splice-site allele that was missed by
1069 sequencing. Scenario two showed a common diagnostic challenge where a
1070 strong candidate exists alongside an unconfirmed but plausible alternative.
1071 Our method distributes confidence across both possibilities. Conventional ap-
1072 proaches focus only on detecting TP and cannot provide this insight. By quan-
1073 tifying residual uncertainty, we can generate structured reports that clearly
1074 distinguish supported, excluded, and plausible-but-unseen variants. We call
1075 this “evidence-aware” interpretation. When combined with genome-wide pri-
1076 ors from the full range of disease-gene panels, this approach applies to any
1077 phenotype from PanelAppRex. By combining variant classification, allele fre-
1078 quency, MOI, and sequencing quality metrics, our method creates a scalable
1079 foundation for evidence-aware diagnostics in clinical genomics.

1080 Estimating disease risk in genetic studies is complicated by uncertainties
1081 in key parameters such as variant penetrance and the fraction of cases at-
1082 tributable to specific variants (3). In the simplest model, where a single, fully
1083 penetrant variant causes disease, the lifetime risk $P(D)$ is equivalent to the
1084 genotype frequency $P(G)$. For an allele with frequency p (ignoring LD for AR),
1085 this translates to:

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

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

1086 When penetrance is incomplete, defined as $P(D | G)$, the risk becomes:
1087 $P(D) = P(G)P(D | G)$. In more realistic scenarios where multiple variants
1088 contribute to disease, $P(G | D)$ denotes the fraction of cases attributable to a

1089 given variant. This leads to:

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

1090 Because both penetrance and $P(G | D)$ are often uncertain, solving this equa-
1091 tion systematically poses a major challenge, which we incidentally tackled in
1092 the validation studies (29; 30).

1093 Our framework addresses this challenge by combining variant classifica-
1094 tions, population allele frequencies, and curated gene-disease associations.
1095 While imperfect on an individual level, these sources exhibit predictable ag-
1096 gregate behaviour, supported by James-Stein estimation principles (31). Cu-
1097 rated gene-disease associations help identify genes that are explainable for
1098 most disease cases, allowing us to approximate $P(G | D)$ close to one. In this
1099 way, we obtain robust estimates of $P(G)$ (the frequency of disease-associated
1100 genotypes), even when exact values of penetrance and case attribution re-
1101 main uncertain.

1102 This approach allows us to pre-calculate priors and summarise the overall
1103 pathogenic burden. By focusing on a subset \mathcal{V} of variants that pass stringent
1104 filtering, where each $P(G_i | D)$ is the probability that a case of disease D is
1105 attributable to variant(s) i , we assume that, in aggregate,

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

1106 Even if the cumulative contribution is slightly less than one, the resultant
1107 risk estimates remain robust within the broad Crls typical of epidemiological
1108 studies. By incorporating these pre-calculated priors into a Bayesian frame-
1109 work, our method refines risk estimates and enhances clinical decision-making
1110 despite inherent uncertainties.

1111 For the IEI-specific investigation, we showed that immunophenotypic and
1112 network-derived features can be used to train and test models that predict
1113 PPIs. From this, we derived a new, simplified classification of immune fea-
1114 tures for IEI genes. We have listed the new immunophenotypic categories
1115 (e.g. T cell low) in the user database, however we have not included the de-
1116 tailed cluster assignments (e.g. PPI groups) because they are too complex for
1117 direct interpretation manually. Instead, our demonstration provides worked

1118 examples that bioinformaticians can use to perform more refined clustering in
1119 larger studies.

1120 Moreover, because variant sets can be collapsed instead of relying on the
1121 gene-level, our method complements existing statistical approaches for ag-
1122 gregating variant effects with methods like Sequence Kernel Association Test
1123 (SKAT) and Aggregated Cauchy Association Test (ACAT) (32–35) and multi-
1124 omics integration techniques (36; 37). It also remains consistent with estab-
1125 lished variant interpretation guidelines from the American College of Med-
1126 ical Genetics and Genomics (ACMG) (38) and complementary frameworks
1127 (39; 40), as well as QC protocols (41; 42). Standardised reporting for qual-
1128 ifying variant sets, such as ACMG Secondary Findings v3.2 (43), further con-
1129 textualises the integration of these probabilities into clinical decision-making.

1130 We compared our occurrence probabilities with AlphaMissense pathogenic-
1131 ity scores and observed that common variants are predominantly scored as
1132 benign or of uncertain significance. While this aligns with their allele frequen-
1133 cies, any pathogenic variant seen in a patient warrants evaluation against its
1134 prior observation probability to assess causality. Predictive tools such as Al-
1135 phaMissense could ostensibly enhance their embedding of variant features
1136 by incorporating gene-disease associations and MOI data, which may not be
1137 fully represented by raw population allele frequencies.

1138 Future work should incorporate the additional variant types and models to
1139 further refine these probability estimates. By continuously updating classical
1140 estimates with emerging data and prior knowledge, we aim to enhance the
1141 precision of genetic diagnostics and ultimately improve patient care.

¹¹⁴² 5 Conclusion

¹¹⁴³ We present a statistical framework that quantifies the probability that any can-
¹¹⁴⁴ didate variant is causally relevant to a genetic disorder, integrating both ob-
¹¹⁴⁵ served genotypes and the possibility of unobserved, disease-causing alleles.
¹¹⁴⁶ The model replaces binary classification schemes with continuous variant-
¹¹⁴⁷ level posterior probabilities and credible intervals that account for evidence
¹¹⁴⁸ strength, genotype availability, and residual uncertainty. By combining classi-
¹¹⁴⁹ cal population genetics, inheritance mode, variant classification, and Bayesian
¹¹⁵⁰ inference, it supports inference across incomplete or ambiguous sequencing
¹¹⁵¹ data by modelling the full space of plausible causal variants. This enables
¹¹⁵² principled interpretation in scenarios involving compound heterozygosity, low
¹¹⁵³ coverage, or uncertain variant evidence, and offers a reproducible metric for
¹¹⁵⁴ estimating the probability of a genetic diagnosis at the level of a gene, panel,
¹¹⁵⁵ or genome. Although demonstrated here in the context of inborn errors of
¹¹⁵⁶ immunity, the framework is generalisable to other rare disease cohorts and
¹¹⁵⁷ supports future work on quantitative variant interpretation at scale.

¹¹⁵⁸ Acknowledgements

¹¹⁵⁹ We would like to thank all the patients and families who have been providing
¹¹⁶⁰ advice on SwissPedHealth and its projects, as well as the clinical and research
¹¹⁶¹ teams at the participating institutions. We acknowledge Genomics England for
¹¹⁶² providing public access to the PanelApp data. The use of data from Genomics
¹¹⁶³ England panelapp was licensed under the Apache License 2.0. The use of
¹¹⁶⁴ data from UniProt was licensed under Creative Commons Attribution 4.0 In-
¹¹⁶⁵ ternational (CC BY 4.0). ClinVar asks its users who distribute or copy data to
¹¹⁶⁶ provide attribution to them as a data source in publications and websites ([10](#)).
¹¹⁶⁷ dbNSFP version 4.4a is licensed under the Creative Commons Attribution-
¹¹⁶⁸ NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0); while we
¹¹⁶⁹ cite this dataset as used our research publication, it is not used for the final
¹¹⁷⁰ version which instead used ClinVar and gnomAD directly. GnomAD is licensed
¹¹⁷¹ under Creative Commons Zero Public Domain Dedication (CC0 1.0 Universal).
¹¹⁷² GnomAD request that usages cites the gnomAD flagship paper ([4](#)) and any
¹¹⁷³ online resources that include the data set provide a link to the browser, and

¹¹⁷⁴ note that tool includes data from the gnomAD v4.1 release. AlphaMissense
¹¹⁷⁵ asks to cite Cheng et al. (9) for usage in research, with data available from
¹¹⁷⁶ Cheng et al. (18).

¹¹⁷⁷ Contributions

¹¹⁷⁸ DL performed main analyses and wrote the manuscript. SB, AS, MS, and JT de-
¹¹⁷⁹ signed analysis and wrote the manuscript. DSF, and SS wrote the manuscript.
¹¹⁸⁰ JF, LJS supervised the work, and applied for funding. The Quant Group is
¹¹⁸¹ a collaboration across multiple institutions where authors contribute equally;
¹¹⁸² the members on this project were DL, SB, AS, and MS.

¹¹⁸³ Competing interest

¹¹⁸⁴ The authors declare no competing interest.

¹¹⁸⁵ Ethics statement

¹¹⁸⁶ This study only used data which was previously published and publicly avail-
¹¹⁸⁷ able, as cited in the manuscript. This SwissPedHealth study, under which this
¹¹⁸⁸ work was carried out, was approved based on the advice of the ethical commit-
¹¹⁸⁹ tee Northwest and Central Switzerland (EKNZ, AO_2022-00018). The study
¹¹⁹⁰ was conducted in accordance with the Declaration of Helsinki.

¹¹⁹¹ Funding

¹¹⁹² This project was supported through the grant Swiss National Science Foun-
¹¹⁹³ dation (SNF) 320030_201060, and NDS-2021-911 (SwissPedHealth) from the
¹¹⁹⁴ Swiss Personalized Health Network and the Strategic Focal Area 'Personalized
¹¹⁹⁵ Health and Related Technologies' of the ETH Domain (Swiss Federal Institutes
¹¹⁹⁶ of Technology).

1197 References

- 1198 [1] Oliver Mayo. A Century of Hardy–Weinberg Equilibrium. *Twin Research and Human Genetics*, 11(3):249–256, June 2008. ISSN
1199 1832-4274, 1839-2628. doi: 10.1375/twin.11.3.249. URL https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal_article.
- 1200
1201
1202
- 1203 [2] Nikita Abramovs, Andrew Brass, and May Tassabehji. Hardy–Weinberg
1204 Equilibrium in the Large Scale Genomic Sequencing Era. *Frontiers in Ge-*
1205 *netics*, 11:210, March 2020. ISSN 1664-8021. doi: 10.3389/fgene.2020.
1206 00210. URL <https://www.frontiersin.org/article/10.3389/fgene.2020.00210/full>.
- 1207
1208 [3] Johannes Zschocke, Peter H. Byers, and Andrew O. M. Wilkie. Mendelian
1209 inheritance revisited: dominance and recessiveness in medical genet-
1210 ics. *Nature Reviews Genetics*, 24(7):442–463, July 2023. ISSN 1471-
1211 0056, 1471-0064. doi: 10.1038/s41576-023-00574-0. URL <https://www.nature.com/articles/s41576-023-00574-0>.
- 1212
1213 [4] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings,
1214 Jessica Alföldi, Qingbo Wang, Ryan L Collins, Kristen M Laricchia, Andrea
1215 Ganna, Daniel P Birnbaum, et al. The mutational constraint spectrum
1216 quantified from variation in 141,456 humans. *Nature*, 581(7809):434–
1217 443, 2020.
- 1218 [5] Sarah L. Bick, Aparna Nathan, Hannah Park, Robert C. Green, Monica H.
1219 Wojcik, and Nina B. Gold. Estimating the sensitivity of genomic new-
1220 born screening for treatable inherited metabolic disorders. *Genetics in*
1221 *Medicine*, 27(1):101284, January 2025. ISSN 10983600. doi: 10.1016/
1222 j.gim.2024.101284. URL <https://linkinghub.elsevier.com/retrieve/pii/S1098360024002181>.
- 1223
1224 [6] Benjamin D. Evans, Piotr Słowiński, Andrew T. Hattersley, Samuel E.
1225 Jones, Seth Sharp, Robert A. Kimmitt, Michael N. Weedon, Richard A.
1226 Oram, Krasimira Tsaneva-Atanasova, and Nicholas J. Thomas. Estimat-
1227 ing disease prevalence in large datasets using genetic risk scores. *Na-*
1228 *ture Communications*, 12(1):6441, November 2021. ISSN 2041-1723. doi:

- 1229 10.1038/s41467-021-26501-7. URL [https://www.nature.com/articles/
1230 s41467-021-26501-7](https://www.nature.com/articles/s41467-021-26501-7).
- 1231 [7] William B. Hannah, Mitchell L. Drumm, Keith Nykamp, Tiziano Prampano,
1232 Robert D. Steiner, and Steven J. Schrodi. Using genomic databases
1233 to determine the frequency and population-based heterogeneity of au-
1234 tosomal recessive conditions. *Genetics in Medicine Open*, 2:101881,
1235 2024. ISSN 29497744. doi: 10.1016/j.gimo.2024.101881. URL <https://linkinghub.elsevier.com/retrieve/pii/S2949774424010276>.
- 1236
- 1237 [8] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Fig-
1238 urnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Au-
1239 gustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Si-
1240 mon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-
1241 Paredes, Stanislav Nikolov, Rishabh Jain, Jonas Adler, Trevor Back, Stig
1242 Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steineg-
1243 ger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein,
1244 David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Push-
1245 meet Kohli, and Demis Hassabis. Highly accurate protein structure
1246 prediction with AlphaFold. *Nature*, 596(7873):583–589, August 2021.
1247 ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL
1248 <https://www.nature.com/articles/s41586-021-03819-2>.
- 1249
- 1250 [9] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė,
1251 Taylor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski,
1252 Tobias Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper,
1253 Demis Hassabis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-
1254 wide missense variant effect prediction with AlphaMissense. *Science*,
1255 381(6664):eadg7492, September 2023. ISSN 0036-8075, 1095-9203.
1256 doi: 10.1126/science.adg7492. URL <https://www.science.org/doi/10.1126/science.adg7492>.
- 1257
- 1258 [10] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen
1259 Chao, Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoff-
1260 man, Wonhee Jang, Karen Karapetyan, Kenneth Katz, Chunlei Liu, Zenith
1261 Maddipatla, Adriana Malheiro, Kurt McDaniel, Michael Ovetsky, George
Riley, George Zhou, J Bradley Holmes, Brandi L Kattman, and Donna R

- 1262 Maglott. ClinVar: improving access to variant interpretations and sup-
1263 porting evidence. *Nucleic Acids Research*, 46(D1):D1062–D1067, Jan-
1264 uary 2018. ISSN 0305-1048, 1362-4962. doi: 10.1093/nar/gkx1153. URL
1265 <http://academic.oup.com/nar/article/46/D1/D1062/4641904>.
- 1266 [11] The UniProt Consortium, Alex Bateman, Maria-Jesus Martin, Sandra Or-
1267 chard, Michele Magrane, Aduragbemi Adesina, Shadab Ahmad, Bowler-
1268 Barnett, and Others. UniProt: the Universal Protein Knowledgebase
1269 in 2025. *Nucleic Acids Research*, 53(D1):D609–D617, January 2025.
1270 ISSN 0305-1048, 1362-4962. doi: 10.1093/nar/gkae1010. URL <https://academic.oup.com/nar/article/53/D1/D609/7902999>.
- 1271 [12] Dylan Lawless. PanelAppRex aggregates disease gene panels and fa-
1272 cilitates sophisticated search. March 2025. doi: 10.1101/2025.03.20.
1273 25324319. URL <http://medrxiv.org/lookup/doi/10.1101/2025.03.20.25324319>.
- 1274 [13] Antonio Rueda Martin, Eleanor Williams, Rebecca E. Foulger, Sarah
1275 Leigh, Louise C. Daugherty, Olivia Niblock, Ivone U. S. Leong, Katherine
1276 R. Smith, Oleg Gerasimenko, Eik Haraldsdottir, Ellen Thomas,
1277 Richard H. Scott, Emma Baple, Arianna Tucci, Helen Brittain, Anna
1278 De Burca, Kristina Ibañez, Dalia Kasperaviciute, Damian Smedley, Mark
1279 Caulfield, Augusto Rendon, and Ellen M. McDonagh. PanelApp crowd-
1280 sources expert knowledge to establish consensus diagnostic gene pan-
1281 els. *Nature Genetics*, 51(11):1560–1565, November 2019. ISSN 1061-
1282 4036, 1546-1718. doi: 10.1038/s41588-019-0528-2. URL <https://www.nature.com/articles/s41588-019-0528-2>.
- 1283 [14] M. Cecilia Poli, Ivona Aksentijevich, Ahmed Aziz Bousfiha, Charlotte
1284 Cunningham-Rundles, Sophie Hambleton, Christoph Klein, Tomohiro
1285 Morio, Capucine Picard, Anne Puel, Nima Rezaei, Mikko R.J. Seppänen,
1286 Raz Somech, Helen C. Su, Kathleen E. Sullivan, Troy R. Torgerson, Is-
1287 abelle Meyts, and Stuart G. Tangye. Human inborn errors of immu-
1288 nity: 2024 update on the classification from the International Union of
1289 Immunological Societies Expert Committee. *Journal of Human Immu-
1290 nity*, 1(1):e20250003, May 2025. ISSN 3065-8993. doi: 10.70962/
1291 jhi.20250003. URL <https://rupress.org/jhi/article/1/1/e20250003/277390/Human-inborn-errors-of-immunity-2024-update-on-the>.
- 1292
1293
1294
1295

- 1296 [15] Ahmed Aziz Bousfiha, Leïla Jeddane, Abderrahmane Moundir, M. Ce-
1297 cilia Poli, Ivona Aksentijevich, Charlotte Cunningham-Rundles, Sophie
1298 Hambleton, Christoph Klein, Tomohiro Morio, Capucine Picard, Anne
1299 Puel, Nima Rezaei, Mikko R.J. Seppänen, Raz Somech, Helen C. Su,
1300 Kathleen E. Sullivan, Troy R. Torgerson, Stuart G. Tangye, and Is-
1301 abelle Meyts. The 2024 update of IUIS phenotypic classification
1302 of human inborn errors of immunity. *Journal of Human Immunity*,
1303 1(1):e20250002, May 2025. ISSN 3065-8993. doi: 10.70962/
1304 jhi.20250002. URL <https://rupress.org/jhi/article/1/1/e20250002/277374/The-2024-update-of-IUIS-phenotypic-classification>.
1305
- 1306 [16] Xiaoming Liu, Chang Li, Chengcheng Mou, Yibo Dong, and Yicheng Tu.
1307 dbNSFP v4: a comprehensive database of transcript-specific functional
1308 predictions and annotations for human nonsynonymous and splice-site
1309 SNVs. *Genome Medicine*, 12(1):103, December 2020. ISSN 1756-994X.
1310 doi: 10.1186/s13073-020-00803-9. URL <https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-020-00803-9>.
1311
- 1312 [17] Damian Szkłarczyk, Katerina Nastou, Mikaela Koutrouli, Rebecca Kirsch,
1313 Farrokh Mehryary, Radja Hachilif, Dewei Hu, Matteo E Peluso, Qingyao
1314 Huang, Tao Fang, et al. The string database in 2025: protein networks
1315 with directionality of regulation. *Nucleic Acids Research*, 53(D1):D730–
1316 D737, 2025.
- 1317 [18] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė,
1318 Taylor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski,
1319 Tobias Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper,
1320 Demis Hassabis, Pushmeet Kohli, and Žiga Avsec. Predictions for al-
1321 phamissense, September 2023. URL <https://doi.org/10.5281/zenodo.8208688>.
1322
- 1323 [19] Paul Tuijnenburg, Hana Lango Allen, Siobhan O. Burns, Daniel Greene,
1324 Machiel H. Jansen, and Others. Loss-of-function nuclear factor κB sub-
1325 unit 1 (NFKB1) variants are the most common monogenic cause of com-
1326 mon variable immunodeficiency in Europeans. *Journal of Allergy and*
1327 *Clinical Immunology*, 142(4):1285–1296, October 2018. ISSN 00916749.
1328 doi: 10.1016/j.jaci.2018.01.039. URL <https://linkinghub.elsevier.com/retrieve/pii/S0091674918302860>.
1329

- 1330 [20] WHO Scientific Group et al. Primary immunodeficiency diseases: report
1331 of a who scientific group. *Clin. Exp. Immunol.*, 109(1):1–28, 1997.
- 1332 [21] Charlotte Cunningham-Rundles and Carol Bodian. Common variable im-
1333 munodeficiency: clinical and immunological features of 248 patients.
1334 *Clinical immunology*, 92(1):34–48, 1999.
- 1335 [22] Eric Oksenhendler, Laurence Gérard, Claire Fieschi, Marion Malphettes,
1336 Gael Mouillot, Roland Jaussaud, Jean-François Viallard, Martine
1337 Gardembas, Lionel Galicier, Nicolas Schleinitz, et al. Infections in 252
1338 patients with common variable immunodeficiency. *Clinical Infectious Dis-*
1339 *eases*, 46(10):1547–1554, 2008.
- 1340 [23] Y Naito, F Adams, S Charman, J Duckers, G Davies, and S Clarke. Uk
1341 cystic fibrosis registry 2023 annual data report. *London: Cystic Fibrosis*
1342 *Trust*, 2023.
- 1343 [24] Carlo Castellani, CFTR2 team, et al. Cftr2: how will it help care? *Paedi-*
1344 *atric respiratory reviews*, 14:2–5, 2013.
- 1345 [25] Hartmut Grasemann and Felix Ratjen. Cystic fibrosis. *New Eng-*
1346 *land Journal of Medicine*, 389(18):1693–1707, 2023. doi: 10.
1347 1056/NEJMra2216474. URL <https://www.nejm.org/doi/full/10.1056/NEJMra2216474>.
- 1349 [26] Kyoko Watanabe, Erdogan Taskesen, Arjen Van Bochoven, and Danielle
1350 Posthuma. Functional mapping and annotation of genetic associa-
1351 tions with FUMA. *Nature Communications*, 8(1):1826, November 2017.
1352 ISSN 2041-1723. doi: 10.1038/s41467-017-01261-5. URL <https://www.nature.com/articles/s41467-017-01261-5>.
- 1354 [27] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thor-
1355 valdsdóttir, Pablo Tamayo, and Jill P. Mesirov. Molecular signa-
1356 tures database (MSigDB) 3.0. *Bioinformatics*, 27(12):1739–1740, June
1357 2011. ISSN 1367-4811, 1367-4803. doi: 10.1093/bioinformatics/btr260. URL <https://academic.oup.com/bioinformatics/article/27/12/1739/257711>.

- 1360 [28] Dylan Lawless. Variant risk estimate probabilities for iei genes. March
1361 2025. doi: 10.5281/zenodo.15111584. URL <https://doi.org/10.5281/zenodo.15111584>.
- 1363 [29] Eric Vallabh Minikel, Sonia M. Vallabh, Monkol Lek, Karol Estrada,
1364 Kaitlin E. Samocha, J. Fah Sathirapongsasuti, Cory Y. McLean, Joyce Y.
1365 Tung, Linda P. C. Yu, Pierluigi Gambetti, Janis Blevins, Shulin Zhang,
1366 Yvonne Cohen, Wei Chen, Masahito Yamada, Tsuyoshi Hamaguchi,
1367 Nobuo Sanjo, Hidehiro Mizusawa, Yosikazu Nakamura, Tetsuyuki Kitamoto,
1368 Steven J. Collins, Alison Boyd, Robert G. Will, Richard Knight, Claudia Ponto,
1369 Inga Zerr, Theo F. J. Kraus, Sabina Eigenbrod, Armin Giese,
1370 Miguel Calero, Jesús De Pedro-Cuesta, Stéphane Haïk, Jean-Louis Laplanche,
1371 Elodie Bouaziz-Amar, Jean-Philippe Brandel, Sabina Capellari,
1372 Piero Parchi, Anna Poleggi, Anna Ladogana, Anne H. O'Donnell-Luria,
1373 Konrad J. Karczewski, Jamie L. Marshall, Michael Boehnke, Markku
1374 Laakso, Karen L. Mohlke, Anna Kähler, Kimberly Chambert, Steven Mc-
1375 Carroll, Patrick F. Sullivan, Christina M. Hultman, Shaun M. Purcell,
1376 Pamela Sklar, Sven J. Van Der Lee, Annemieke Rozemuller, Casper
1377 Jansen, Albert Hofman, Robert Kraaij, Jeroen G. J. Van Rooij, M. Ar-
1378 fan Ikram, André G. Uitterlinden, Cornelia M. Van Duijn, Exome Ag-
1379 gregation Consortium (ExAC), Mark J. Daly, and Daniel G. MacArthur.
1380 Quantifying prion disease penetrance using large population control co-
1381 cohorts. *Science Translational Medicine*, 8(322), January 2016. ISSN
1382 1946-6234, 1946-6242. doi: 10.1126/scitranslmed.aad5169. URL <https://www.science.org/doi/10.1126/scitranslmed.aad5169>.
- 1384 [30] Nicola Whiffin, Eric Minikel, Roddy Walsh, Anne H O'Donnell-Luria, Kon-
1385 rad Karczewski, Alexander Y Ing, Paul J R Barton, Birgit Funke, Stu-
1386 art A Cook, Daniel MacArthur, and James S Ware. Using high-resolution
1387 variant frequencies to empower clinical genome interpretation. *Genet-
ics in Medicine*, 19(10):1151–1158, October 2017. ISSN 10983600. doi:
1388 10.1038/gim.2017.26. URL <https://linkinghub.elsevier.com/retrieve/pii/S1098360021013678>.
- 1391 [31] Bradley Efron and Carl Morris. Stein's Estimation Rule and Its
1392 Competitors—An Empirical Bayes Approach. *Journal of the American Sta-
1393 tistical Association*, 68(341):117, March 1973. ISSN 01621459. doi: 10.

- 1394 2307/2284155. URL <https://www.jstor.org/stable/2284155?origin=crossref>.
- 1395
- 1396 [32] Yaowu Liu, Sixing Chen, Zilin Li, Alanna C Morrison, Eric Boerwinkle, and
1397 Xihong Lin. Acat: a fast and powerful p value combination method for
1398 rare-variant analysis in sequencing studies. *The American Journal of*
1399 *Human Genetics*, 104(3):410–421, 2019.
- 1400 [33] Xihao Li, Zilin Li, Hufeng Zhou, Sheila M Gaynor, Yaowu Liu, Han Chen,
1401 Ryan Sun, Rounak Dey, Donna K Arnett, Stella Aslibekyan, et al. Dynamic
1402 incorporation of multiple in silico functional annotations empowers rare
1403 variant association analysis of large whole-genome sequencing studies
1404 at scale. *Nature genetics*, 52(9):969–983, 2020.
- 1405 [34] Michael C Wu, Seunggeun Lee, Tianxi Cai, Yun Li, Michael Boehnke, and
1406 Xihong Lin. Rare-variant association testing for sequencing data with
1407 the sequence kernel association test. *The American Journal of Human*
1408 *Genetics*, 89(1):82–93, 2011.
- 1409 [35] Seunggeun Lee, Mary J Emond, Michael J Bamshad, Kathleen C Barnes,
1410 Mark J Rieder, Deborah A Nickerson, David C Christiani, Mark M Wur-
1411 fel, and Xihong Lin. Optimal unified approach for rare-variant associa-
1412 tion testing with application to small-sample case-control whole-exome
1413 sequencing studies. *The American Journal of Human Genetics*, 91(2):
1414 224–237, 2012.
- 1415 [36] Augustine Kong, Gudmar Thorleifsson, Michael L Frigge, Bjarni J Vilh-
1416 jalmsson, Alexander I Young, Thorgeir E Thorgeirsson, Stefania Benonis-
1417 dottir, Asmundur Oddsson, Bjarni V Halldorsson, Gisli Masson, et al. The
1418 nature of nurture: Effects of parental genotypes. *Science*, 359(6374):
1419 424–428, 2018.
- 1420 [37] Laurence J Howe, Michel G Nivard, Tim T Morris, Ailin F Hansen, Hu-
1421 maira Rasheed, Yoonsu Cho, Geetha Chittoor, Penelope A Lind, Teemu
1422 Palviainen, Matthijs D van der Zee, et al. Within-sibship gwas improve
1423 estimates of direct genetic effects. *BioRxiv*, pages 2021–03, 2021.
- 1424 [38] Sue Richards, Nazneen Aziz, Sherri Bale, David Bick, Soma Das, Julie
1425 Gastier-Foster, Wayne W Grody, Madhuri Hegde, Elaine Lyon, Elaine

- 1426 Spector, et al. Standards and guidelines for the interpretation of se-
1427 quence variants: a joint consensus recommendation of the american col-
1428 lege of medical genetics and genomics and the association for molecular
1429 pathology. *Genetics in medicine*, 17(5):405–423, 2015.
- 1430 [39] Sean V Tavtigian, Steven M Harrison, Kenneth M Boucher, and Leslie G
1431 Biesecker. Fitting a naturally scaled point system to the acmg/amp vari-
1432 ant classification guidelines. *Human mutation*, 41(10):1734–1737, 2020.
- 1433 [40] Quan Li and Kai Wang. Intervar: clinical interpretation of genetic variants
1434 by the 2015 acmg-amp guidelines. *The American Journal of Human Ge-
1435 netics*, 100(2):267–280, 2017.
- 1436 [41] Brent S Pedersen, Joe M Brown, Harriet Dashnow, Amelia D Wal-
1437 lace, Matt Velinder, Martin Tristani-Firouzi, Joshua D Schiffman, Tatiana
1438 Tvrzik, Rong Mao, D Hunter Best, et al. Effective variant filtering and
1439 expected candidate variant yield in studies of rare human disease. *NPJ
1440 Genomic Medicine*, 6(1):1–8, 2021.
- 1441 [42] Carl A Anderson, Fredrik H Pettersson, Geraldine M Clarke, Lon R Car-
1442 don, Andrew P Morris, and Krina T Zondervan. Data quality control in
1443 genetic case-control association studies. *Nature protocols*, 5(9):1564–
1444 1573, 2010. URL <https://doi.org/10.1038/nprot.2010.116>.
- 1445 [43] David T Miller, Kristy Lee, Noura S Abul-Husn, Laura M Amendola, Kyle
1446 Brothers, Wendy K Chung, Michael H Gollob, Adam S Gordon, Steven M
1447 Harrison, Ray E Hershberger, et al. Acmg sf v3. 2 list for reporting of
1448 secondary findings in clinical exome and genome sequencing: a policy
1449 statement of the american college of medical genetics and genomics
1450 (acmg). *Genetics in Medicine*, 25(8):100866, 2023.

6 Supplemental

1452 Supplemental data are presented under the same headings that correspond
 1453 to their relevant main text sections.

1454 **6.1 Variant class occurrence probability**

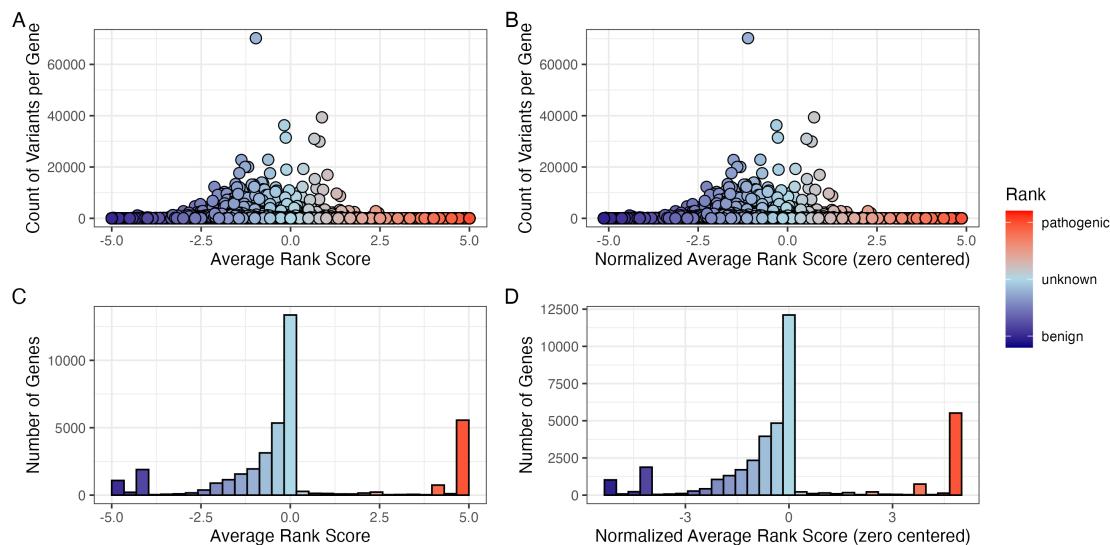


Figure S1: **Global distribution of ClinVar clinical-significance classification scoring.** (A) Number of variants per gene containing the assigned score for each ClinVar classification term (-5 to +5). (B) The same data after normalisation by zero centring the average rank score. (C) The tally of genes for their average rank and (D) after normalisation. No normalisation was required for the scoring system as shown by comparison of A-C and B-D.

1455 **6.2 Integrating observed true positives and unobserved**
 1456 **false 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: 0.381). The strongest unsequenced variant was c.159+1G>A (posterior: 0.353). The total probability of a causal diagnosis given the available evidence was 0.52 (95% CI: 0.248–0.787).

Variant	Flag	Class	Evidence Score	Occurrence Prob	Adj Occ Prob	Alpha	Beta	Lower	Median	Upper	Posterior Share	Prob Causal
p.Ser237Ter	present	causal	5.0	0.000	0	6.0	371	0.003	0.096	0.557	0.381	0.381
c.159+1G>A	missing	causal	4.5	0.000	0	5.5	367	NA	NA	NA	0.353	0.353
p.Thr567Ile	missing	other	-5.0	0.002	0	1.0	365	NA	NA	NA	0.000	0.000
p.Arg231His	present	other	0.0	0.000	0	1.0	359	0.003	0.096	0.557	0.000	0.000
p.Gly650Arg	present	other	0.0	0.000	0	1.0	349	0.003	0.096	0.557	0.000	0.000
p.Val236Ile	missing	other	0.0	0.000	0	1.0	363	NA	NA	NA	0.000	0.000
Total	NA	NA	NA	NA	NA	NA	NA	0.248	0.520	0.787	NA	0.520

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

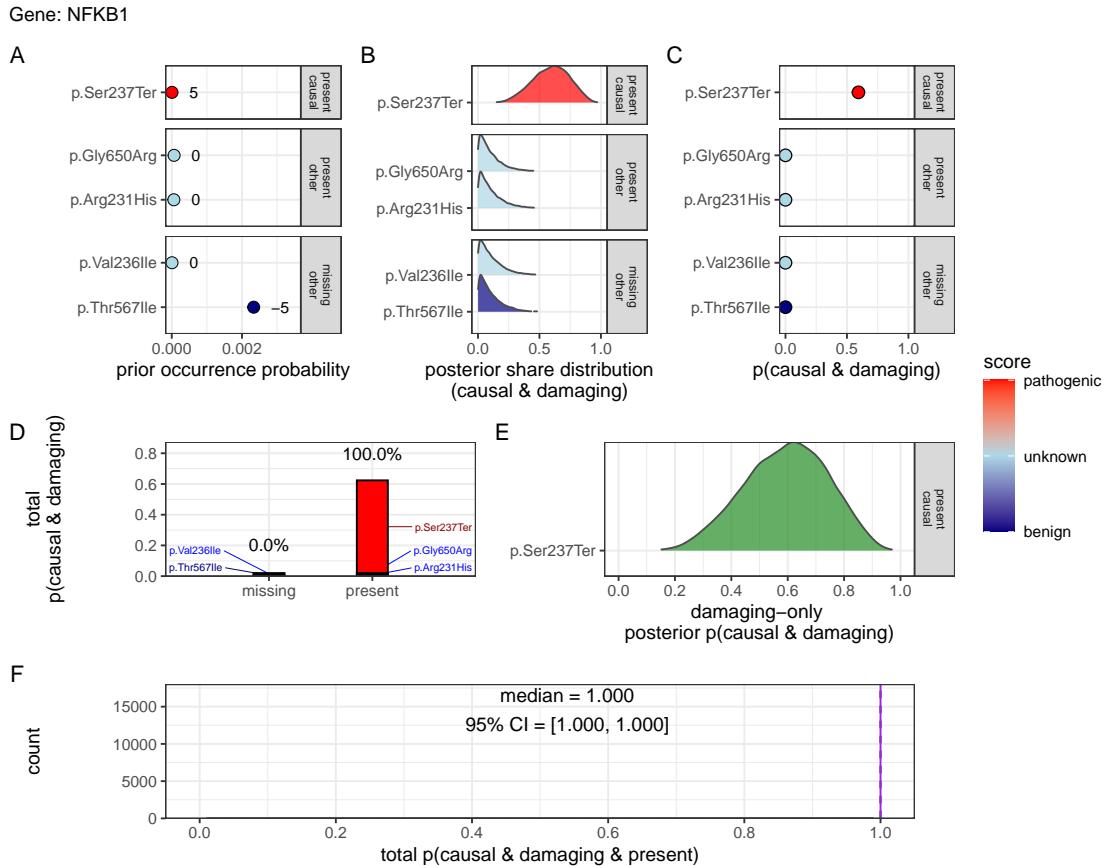


Figure S2: Quantification of present (TP) and no missing (FN) causal genetic variants for disease in *NFKB1* (scenario 1). Only one known pathogenic variant, p.Ser237Ter, was observed and all previously reported pathogenic positions were successfully sequenced and confirmed as reference (true negatives). 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.

Gene: TNFAIP3

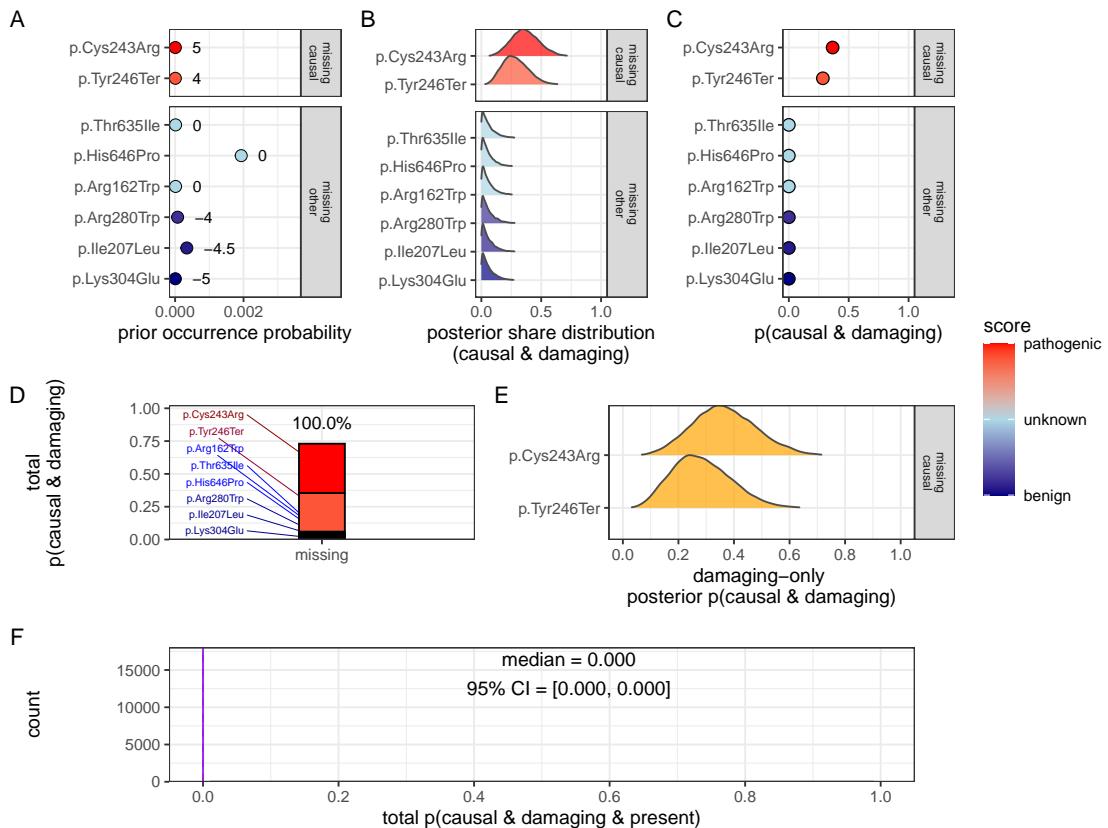
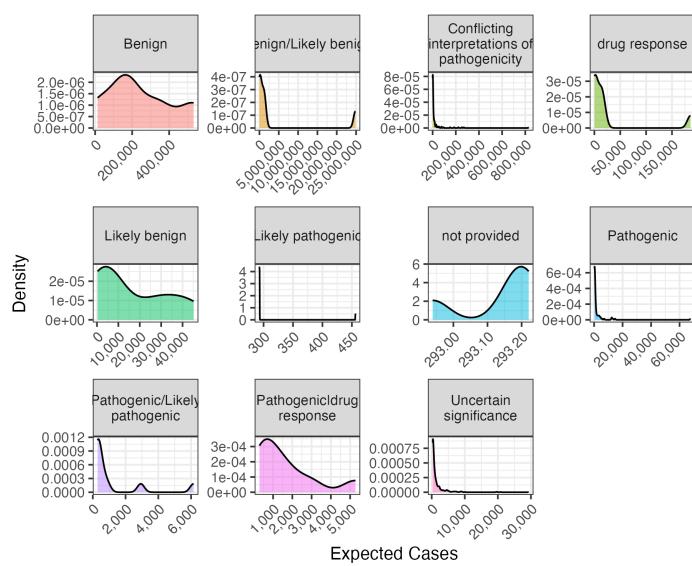


Figure S3: Quantification of no present (TP) in *NFKB1* and only missing (FN) causal genetic variants for disease in *TNFAIP3* (scenario 3). No known causal variants were observed in *NFKB1*, but one representative unsequenced allele was selected from each distinct ClinVar classification and treated 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).

6.3 Validation studies

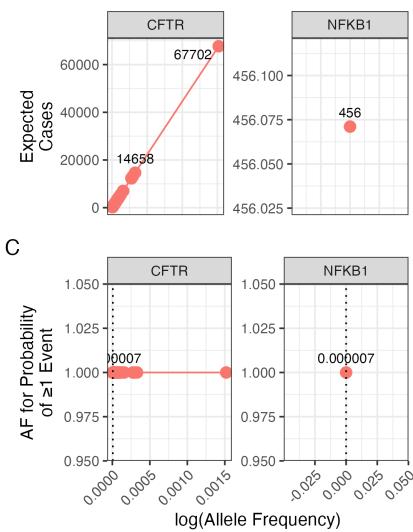
Condition: population size 69433632, phenotype PID-related, genes CFTR and NFKB1.

A



B

clinvar_clnsig • Pathogenic



C

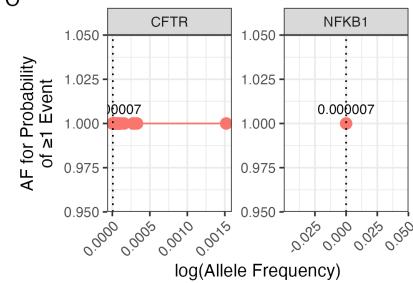


Figure S4: Interpretation of probability of observing a variant classification. The result from the chosen validation genes *CFTR* and *NFKB1* are shown. Case counts are dependant on the population size and phenotype. (A) The density plots of expected observations by ClinVar clinical significance. We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

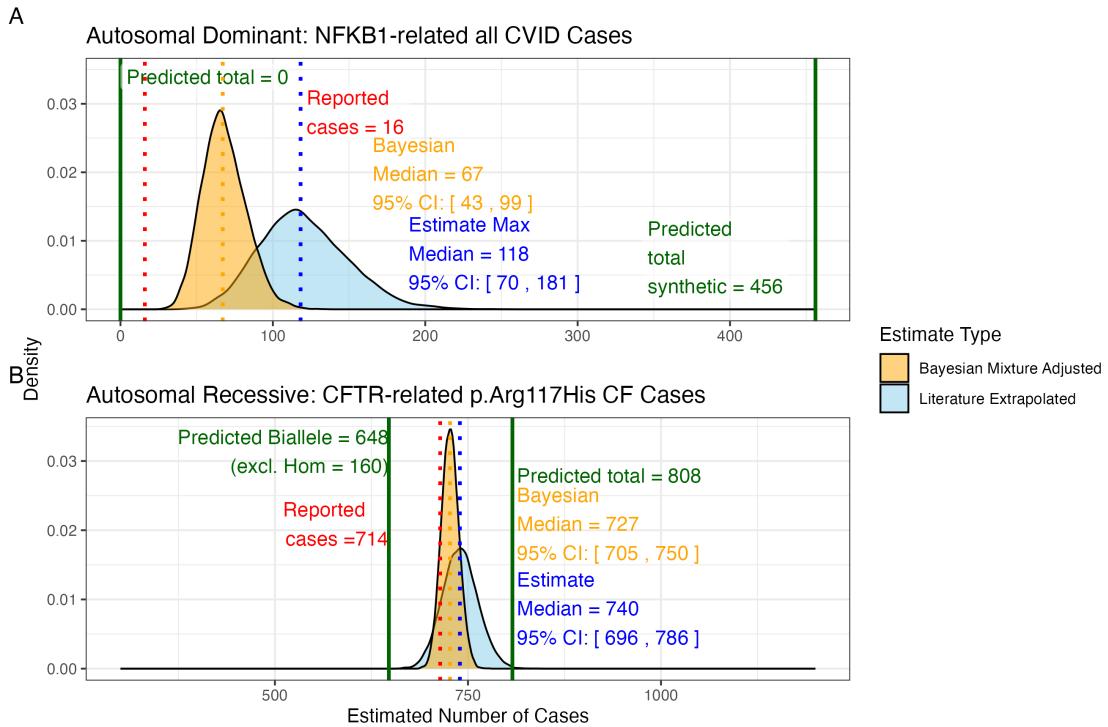


Figure S5: Prior probabilities compared to validation disease cohort metrics. (A) Density distributions for the number of *NFKB1*-related CVID cases in the UK. Our model (green) predicted 456 cases, which falls between the observed cohort count (red) of 390 and the upper extrapolated values. The blue 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.

6.3.1 Interpretation of ClinVar variant observations

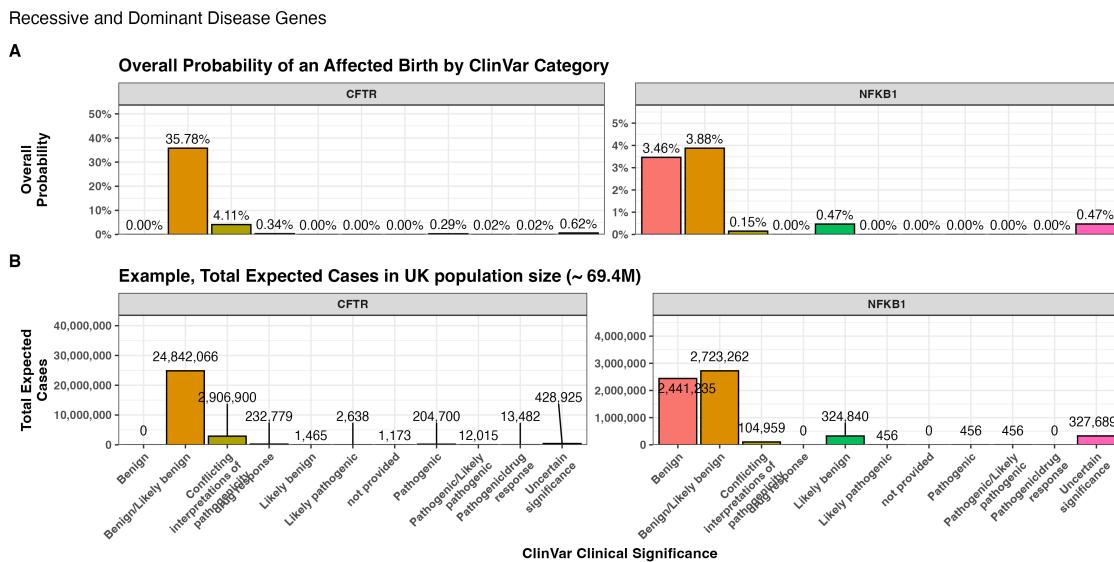


Figure S6: Combined bar charts summarising the genome-wide analysis of ClinVar clinical significance for the PID gene panel. Panel (A) shows the overall probability of an affected birth by variant classification, and (B) displays the total expected number of cases per classification, both stratified by gene. These integrated results illustrate the variability in variant observations across genes and underpin our validation of the probability estimation framework.

6.3.2 Validation of SCID-specific disease occurrence

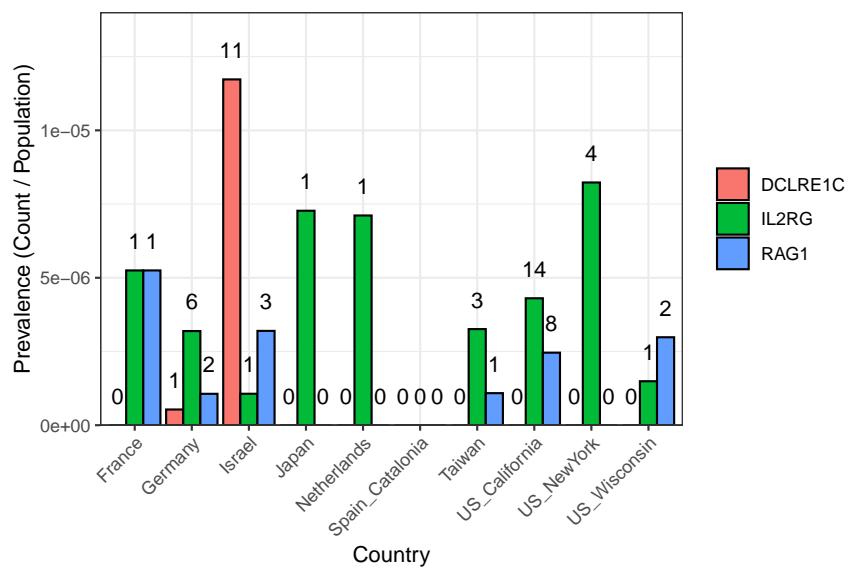


Figure S7: SCID-specific gene comparison across regions. The bar plot shows the prevalence of SCID-related cases (count divided by population) for each gene and country (or region), with numbers printed above the bars representing the actual counts in the original cohort (ranging from 0 to 11 per region and gene).

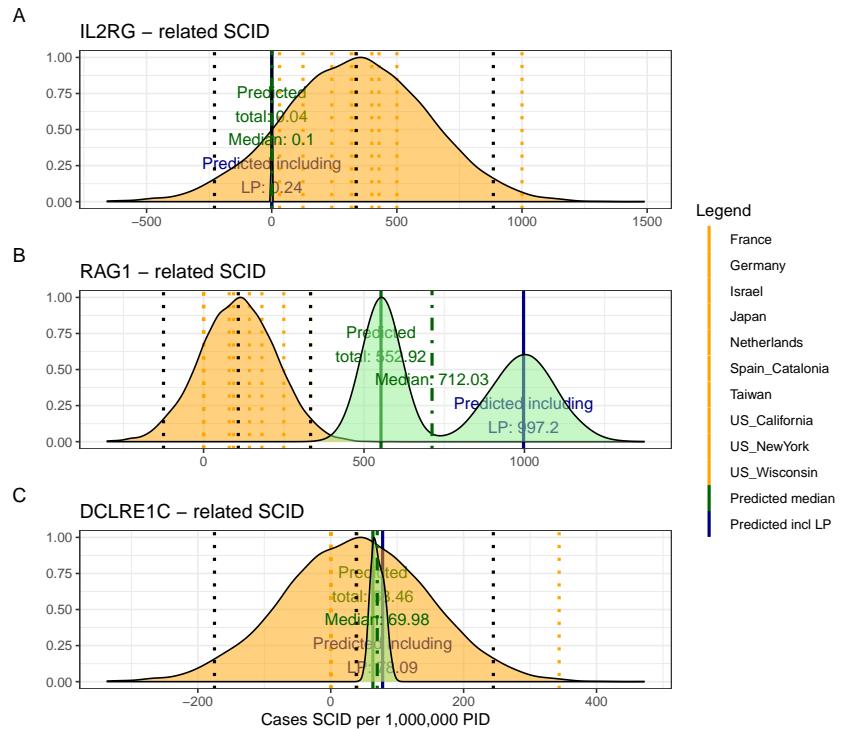
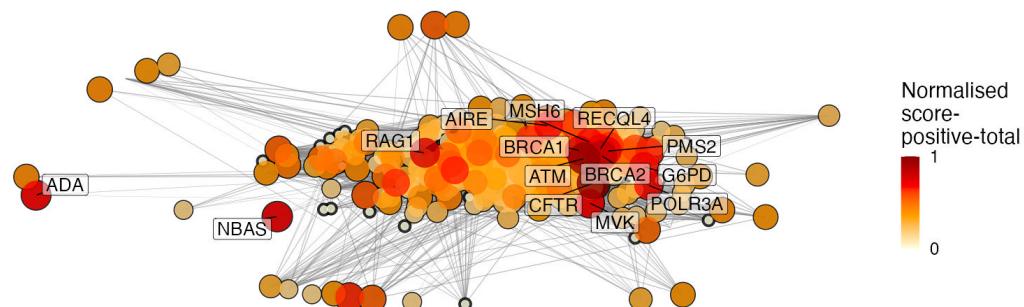


Figure S8: Combined SCID-specific Predictions and Observed Rates per 1,000,000 PID. The figure presents density distributions for the predicted SCID case counts (per 1,000,000 PID) for three genes: *IL2RG*, *RAG1*, and *DCLRE1C*. Country-specific rates (displayed as dotted vertical lines) are overlaid 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.

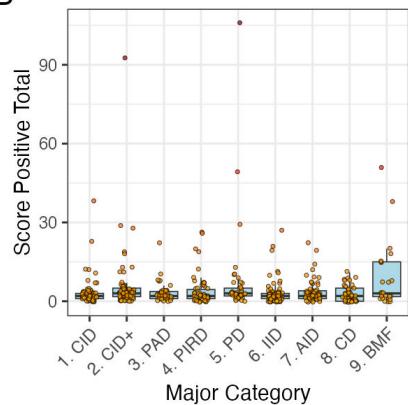
6.4 Genetic constraint in high-impact protein networks

6.4.1 Score-positive-total within IEI PPI network

A



B



C

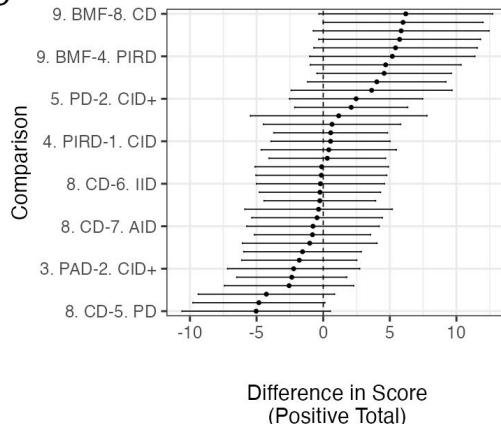


Figure S9: PPI network and score-positive-total ClinVar significance variants. (A) PPI network of disease-associated genes. Node size and colour represent the log-transformed score-positive-total, the top 15 genes/proteins with the highest probability of being observed in disease are labelled. (B) Distribution 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.

6.4.2 Hierarchical Clustering of Enrichment Scores for Major Disease Categories

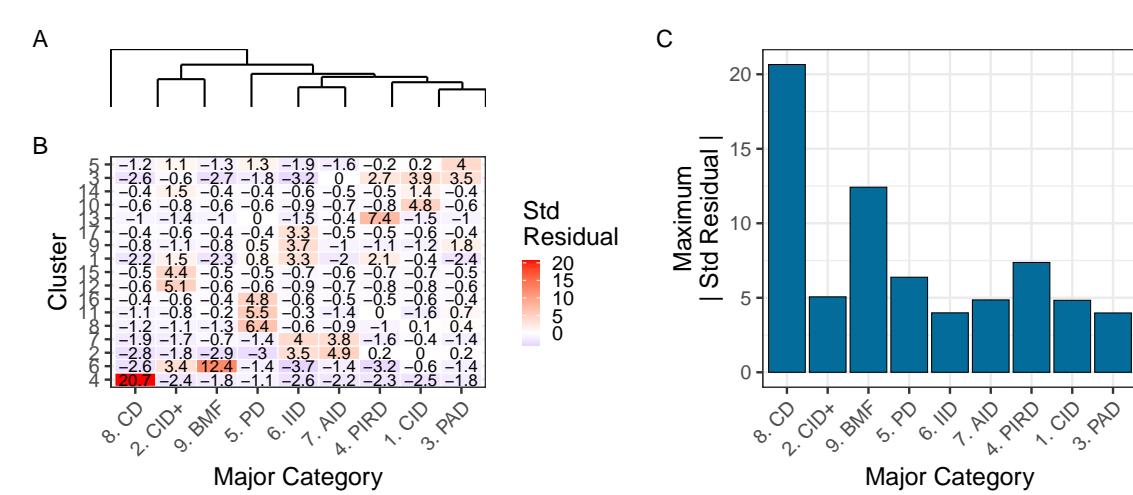


Figure S10: Hierarchical clustering of enrichment scores. The heatmap displays standardised residuals for major disease categories (x-axis) across network clusters (y-axis). A dendrogram groups similar disease categories, and the bar plot shows the maximum absolute residual per category. (8) CD and (9)BMF show the highest values, indicating significant enrichment or depletion (residuals $> |2|$). Definitions in **Box 2.1**.

1464

6.4.3 PPI connectivity, LOEUF constraint and enriched network cluster analysis

1465

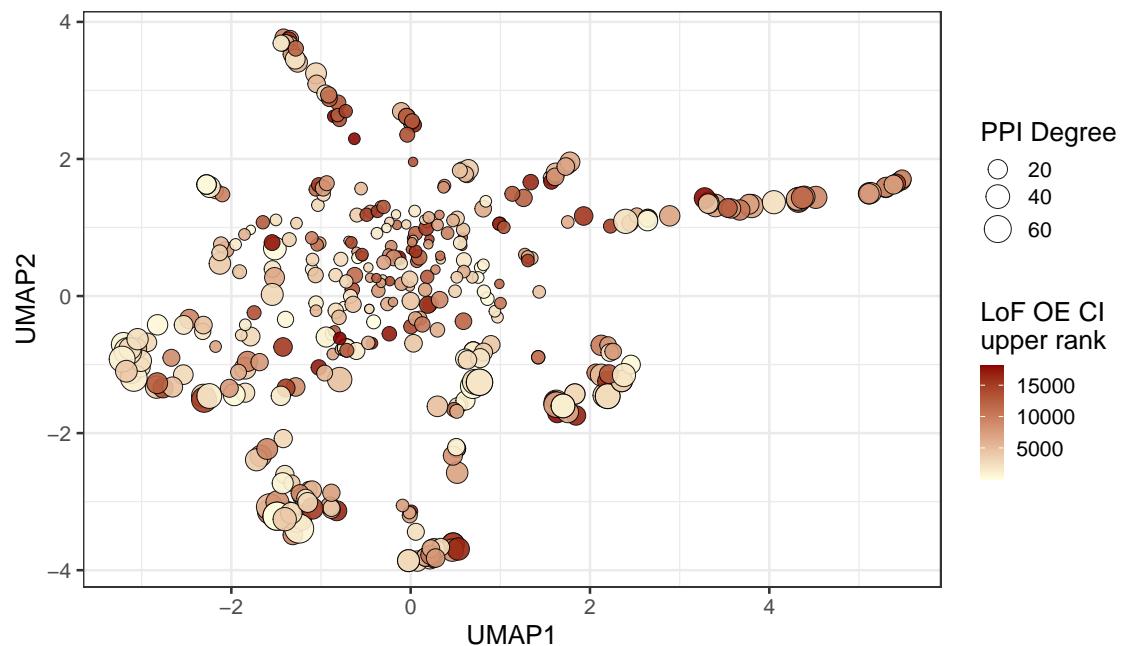


Figure S11: **Analysis of PPI degree versus LOEUF upper rank with UMAP embedding of the PPI network.** The relationship between PPI degree (size) and LOEUF upper rank (color) across gene clusters. No clear patterns are evident.

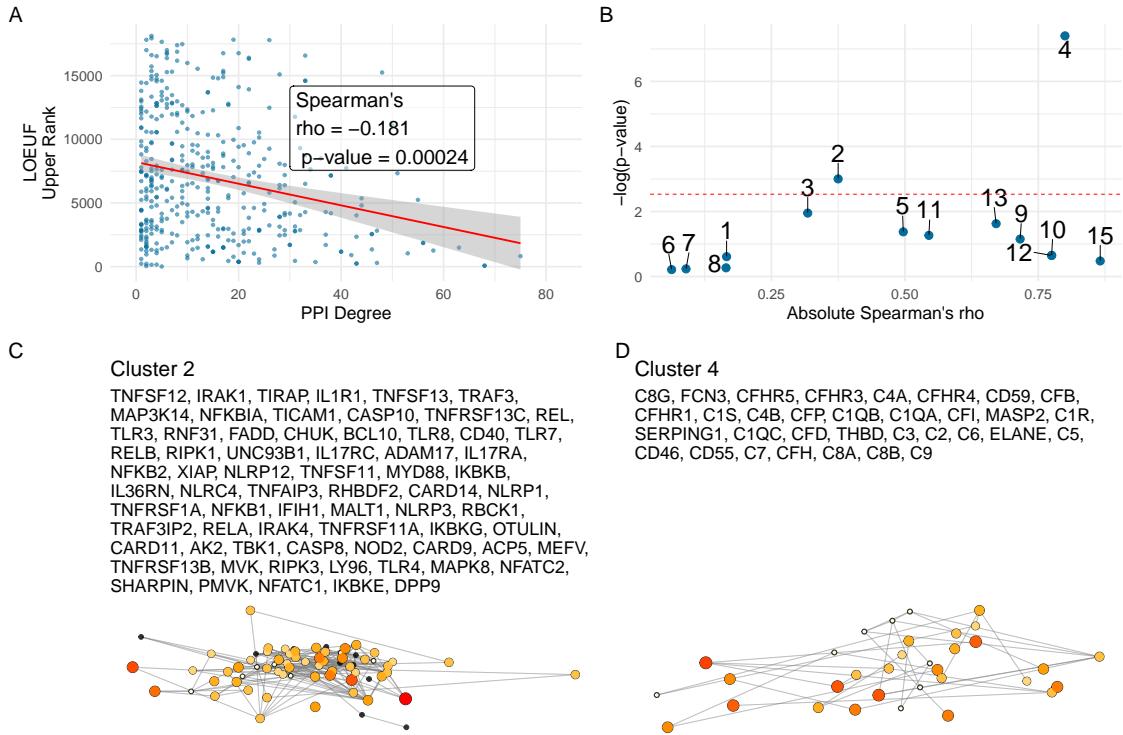


Figure S12: Correlation between PPI degree and LOEUF upper rank. **(A)** Ananlysis across all genes revealed a weak, significant negative correlation between PPI degree and LOEUF upper rank. **(B)** The cluster-wise analysis showed that clusters 2 and 4 exhibited moderate to strong correlations, while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

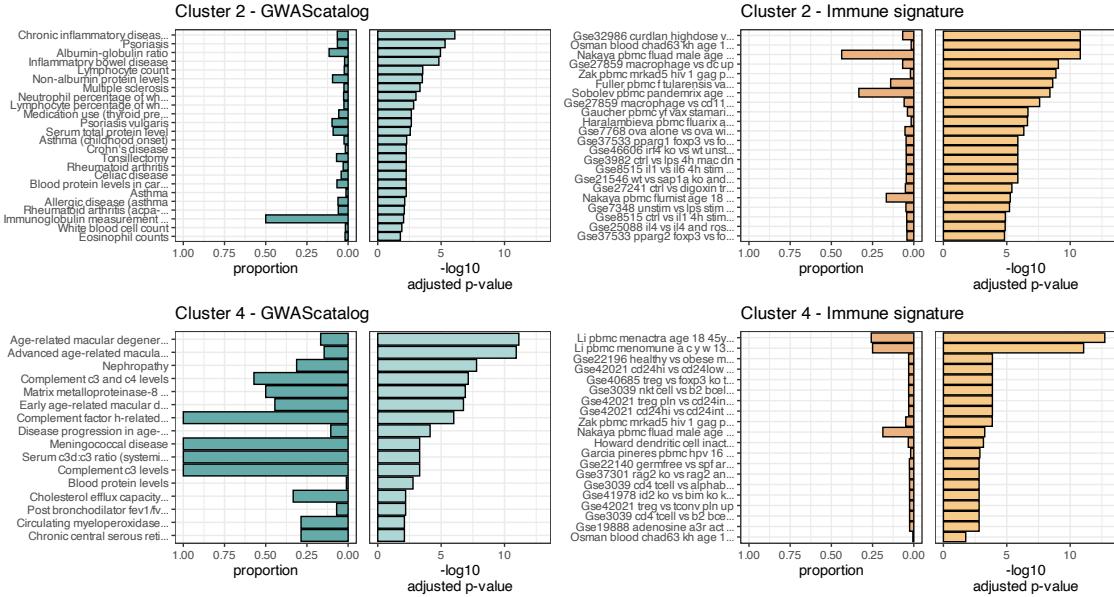


Figure S13: Composite Enrichment Profiles for IEI Gene Sets. We selected the top two enriched clusters (as per **Figure S12**) and performed functional enrichment analysis derived from known disease associations. For each gene set, the left panel displays the proportion of input genes overlapping with 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.

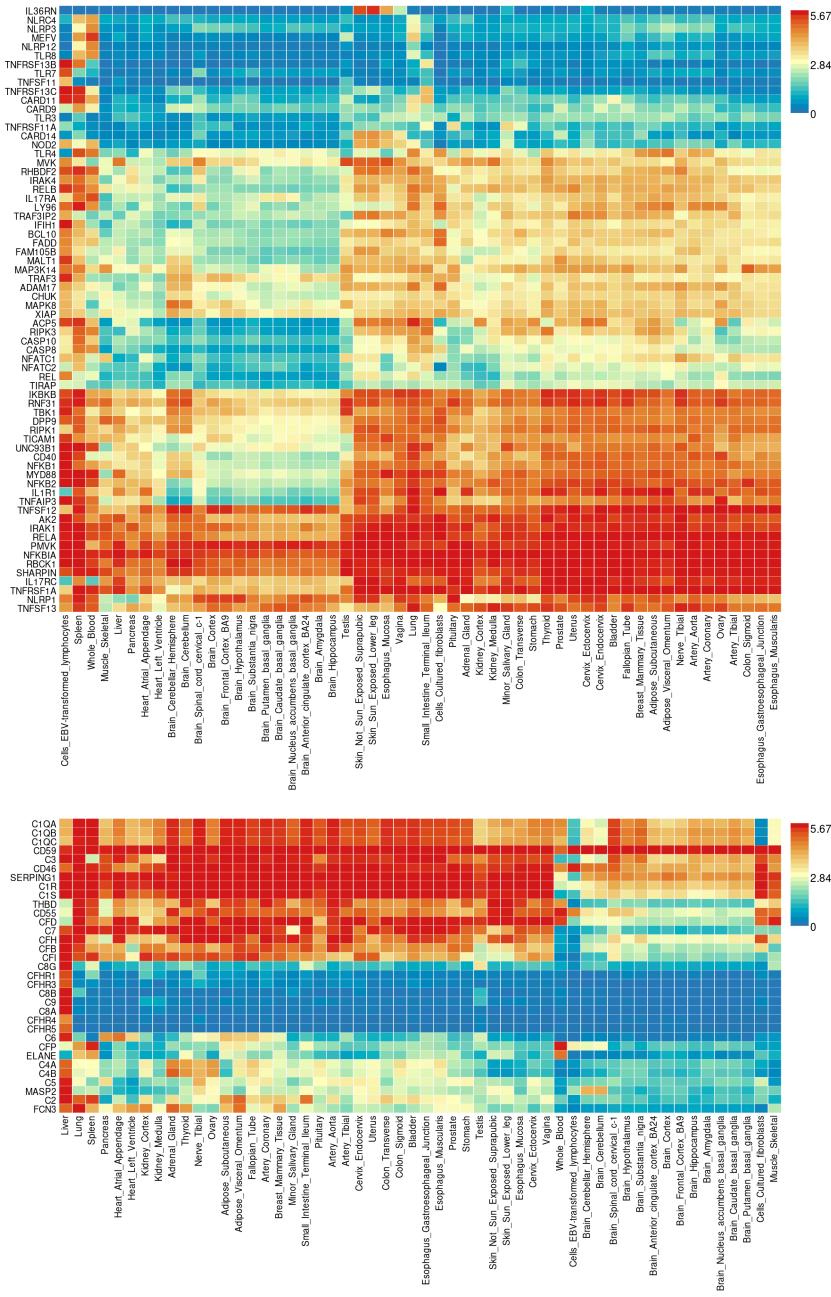


Figure S14: Gene Expression Heatmaps for IEI Genes. GTEx v8 data from 54 tissue types display the average expression per tissue label (\log_2 transformed) for the IEI gene panels. Top: Cluster 2; Bottom: Cluster 4.

6.5 Novel PID classifications derived from genetic PPI and clinical features

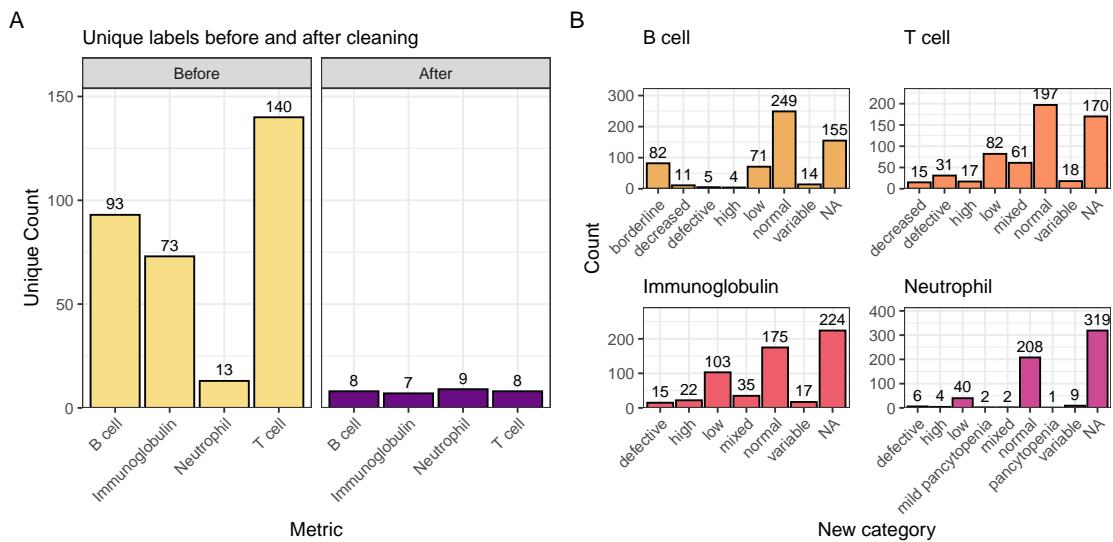


Figure S15: Distribution of immunophenotypic features before and after recategorisation. The original IUIS IEI descriptions contain information such as T cell-related “decreased CD8, normal or decreased CD4 cells” which we recategorise as “low”. The bar plot shows the count of unique labels for each status (normal, not_normal) across the T cell, B cell, Ig, and Neutrophil features.

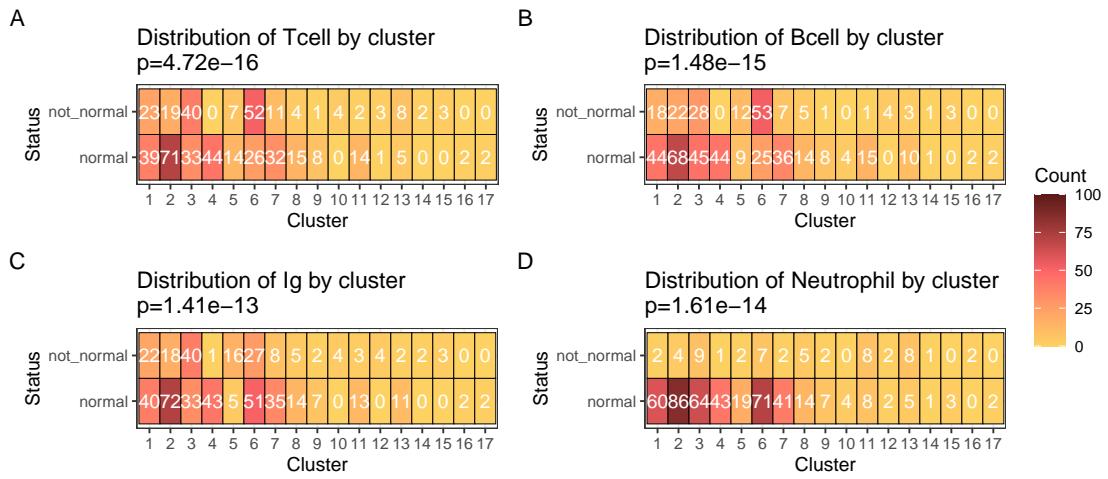


Figure S16: Heatmaps of clinical feature distributions by PPI cluster.
The heatmaps display the count of observations for abnormality of each clinical feature (A) T cell, (B) B cell, (C) Immunoglobulin, (D) Neutrophil, in relation to the PPI clusters, with p-values from chi-square tests annotated in the titles.

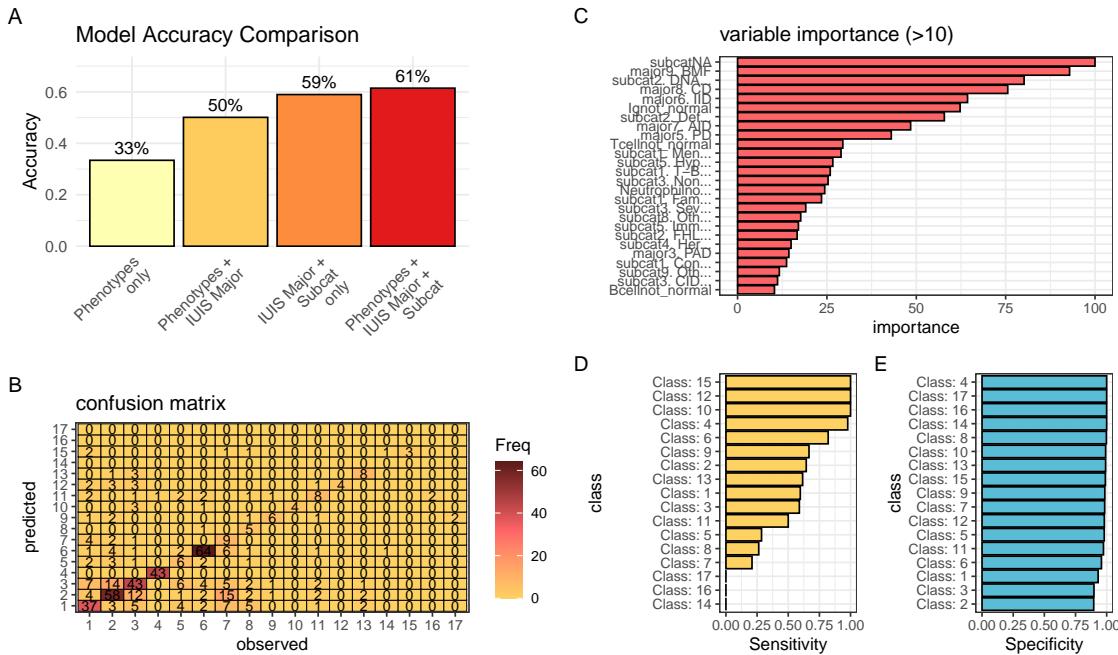


Figure S17: Performance comparison of PID classifiers. Classification predicting PPI cluster membership from IUIS major category, subcategory, and immunological features. (A) Overall accuracy for four rpart models used to predict PPI clustering. The combined model achieves 61.4 % accuracy, exceeding 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.

6.6 Probability of observing AlphaMissense pathogenicity

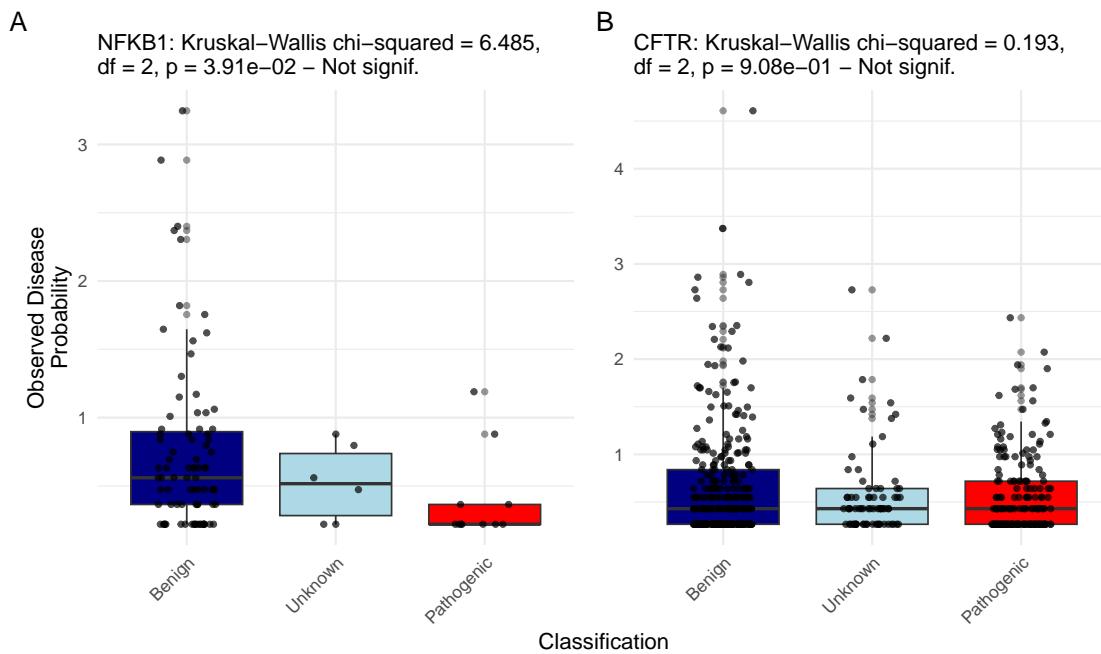


Figure S18: **Observed Disease Probability by Clinical Classification with AlphaMissense.** The figure displays the Kruskal-Wallis test results for NFKB1 and CFTR, showing no significant differences.