

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

Quant Group<sup>1</sup>, Simon Boutry<sup>2</sup>, Ali Saadat<sup>2</sup>, Maarja Soomann<sup>3</sup>, Johannes Trück<sup>3</sup>, D. Sean Froese<sup>4</sup>, Jacques Fellay<sup>2</sup>, Sinisa Savic<sup>5</sup>, Luregn J. Schlapbach<sup>6</sup>, and Dylan Lawless \*<sup>6</sup>

<sup>1</sup>The quantitative omic epidemiology group.

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

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

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

<sup>5</sup>Leeds Institute of Rheumatic and Musculoskeletal Medicine, University of Leeds, Leeds, UK.

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

July 1, 2025

## Abstract

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

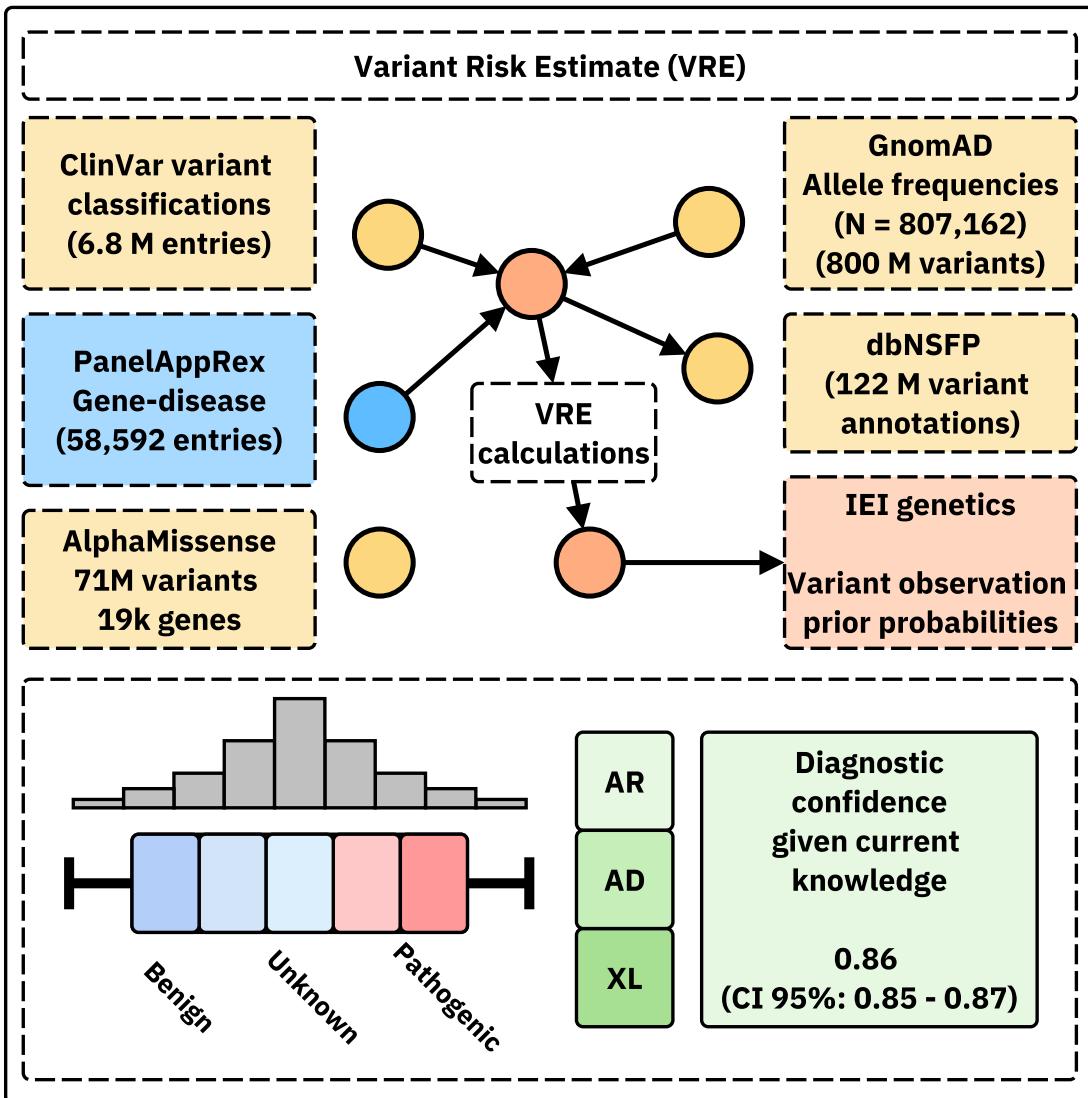
**Methods:** We developed a statistical framework that computes genome-wide prior probabilities for variant pathogenicity, integrating population allele frequencies, variant classifications, and Hardy-Weinberg expectations across inheritance modes. Bayesian modelling then combines these priors with patient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

**Results:** We applied the framework to 557 genes implicated in inborn errors of immunity (IEI), generating variant-level probabilities now publicly accessible. Integrating these data with protein-protein interaction networks and immunophenotypic features revealed patterns of genetic constraint and refined disease classification. Validation in national cohorts demonstrated 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. <sup>1</sup>

<sup>1</sup> \* 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](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.



21

22 Graphical abstract.

<sup>23</sup> **Acronyms**

<sup>24</sup> <b>ACMG</b>	American College of Medical Genetics and Genomics.....	<sup>40</sup>
<sup>25</sup> <b>ACAT</b>	Aggregated Cauchy Association Test .....	<sup>40</sup>
<sup>27</sup> <b>AD</b>	Autosomal Dominant.....	<sup>6</sup>
<sup>29</sup> <b>AF</b>	Allele Frequency.....	<sup>7</sup>
<sup>31</sup> <b>ANOVA</b>	Analysis of Variance .....	<sup>19</sup>
<sup>33</sup> <b>AR</b>	Autosomal Recessive .....	<sup>6</sup>
<sup>35</sup> <b>BMF</b>	Bone Marrow Failure.....	<sup>29</sup>
<sup>37</sup> <b>CD</b>	Complement Deficiencies .....	<sup>30</sup>
<sup>39</sup> <b>CI</b>	Confidence Interval.....	<sup>27</sup>
<sup>41</sup> <b>CrI</b>	Credible Interval .....	<sup>13</sup>
<sup>43</sup> <b>CF</b>	Cystic Fibrosis .....	<sup>16</sup>
<sup>45</sup> <b>CFTR</b>	Cystic Fibrosis Transmembrane Conductance Regulator.....	<sup>8</sup>
<sup>47</sup> <b>CVID</b>	Common Variable Immunodeficiency.....	<sup>15</sup>
<sup>49</sup> <b>DCLRE1C</b>	DNA Cross-Link Repair 1C .....	<sup>8</sup>
<sup>51</sup> <b>dbNSFP</b>	database for Non-Synonymous Functional Predictions .....	<sup>8</sup>
<sup>53</sup> <b>GE</b>	Genomics England .....	<sup>7</sup>
<sup>55</sup> <b>gnomAD</b>	Genome Aggregation Database .....	<sup>8</sup>
<sup>57</sup> <b>gVCF</b>	genomic variant call format .....	<sup>13</sup>
<sup>59</sup> <b>HGVS</b>	Human Genome Variation Society.....	<sup>8</sup>
<sup>61</sup> <b>HPC</b>	High-Performance Computing.....	<sup>12</sup>
<sup>63</sup> <b>HSD</b>	Honestly Significant Difference .....	<sup>19</sup>
<sup>65</sup> <b>HWE</b>	Hardy-Weinberg Equilibrium .....	<sup>6</sup>
<sup>67</sup> <b>IEI</b>	Inborn Errors of Immunity .....	<sup>7</sup>
<sup>69</sup> <b>Ig</b>	Immunoglobulin .....	<sup>33</sup>
<sup>71</sup> <b>IL2RG</b>	Interleukin 2 Receptor Subunit Gamma.....	<sup>8</sup>
<sup>73</sup> <b>InDel</b>	Insertion/Deletion .....	<sup>8</sup>
<sup>75</sup> <b>IUIS</b>	International Union of Immunological Societies .....	<sup>7</sup>

<sup>77</sup>

78	<b>LD</b> Linkage Disequilibrium .....	32
79	<b>LOEUF</b> Loss-Of-function Observed/Expected Upper bound Fraction .....	19
81	<b>LOF</b> Loss-of-Function .....	19
83	<b>MOI</b> Mode of Inheritance .....	6
85	<b>NFKB1</b> Nuclear Factor Kappa B Subunit 1 .....	8
87	<b>OMIM</b> Online Mendelian Inheritance in Man .....	37
89	<b>PID</b> Primary Immunodeficiency .....	7
91	<b>PPI</b> Protein-Protein Interaction .....	8
93	<b>pLI</b> Probability of being Loss-of-function Intolerant .....	19
95	<b>QC</b> Quality Control .....	13
97	<b>RAG1</b> Recombination activating gene 1 .....	8
99	<b>SCID</b> Severe Combined Immunodeficiency .....	8
101	<b>SNV</b> Single Nucleotide Variant .....	6
103	<b>SKAT</b> Sequence Kernel Association Test.....	40
105	<b>STRINGdb</b> Search Tool for the Retrieval of Interacting Genes/Proteins.....	8
107	<b>TP</b> true positive.....	6
109	<b>FP</b> false positive.....	6
111	<b>TN</b> true negative .....	6
113	<b>FN</b> false negative.....	6
115	<b>TNFAIP3</b> Tumor necrosis factor, alpha-induced protein 3 .....	8
117	<b>UMAP</b> Uniform Manifold Approximation and Projection .....	20
119	<b>UniProt</b> Universal Protein Resource.....	7
121	<b>VCF</b> variant call format .....	13
123	<b>VEP</b> Variant Effect Predictor.....	8
125	<b>VRE</b> variant risk estimate .....	9
127	<b>XL</b> X-Linked .....	6
129		

# 130 1 Introduction

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

148 Quantifying the risk that a patient inherits a disease-causing variant is a fun-  
149 damental challenge in genomics. Classical statistical approaches grounded in Hardy-  
150 Weinberg Equilibrium (HWE) (1; 2) have long been used to calculate genetic probabili-  
151 ties for Single Nucleotide Variant (SNV)s. However, applying these methods becomes  
152 more complex when accounting for different Mode of Inheritance (MOI), such as Auto-  
153 somal Recessive (AR) versus Autosomal Dominant (AD) or X-Linked (XL) disorders.  
154 In AR conditions, for example, the occurrence probability must incorporate both the  
155 homozygous state and compound heterozygosity, whereas for AD and XL disorders,  
156 a single pathogenic allele is sufficient to cause disease. Advances in genetic research  
157 have revealed that MOI can be even more complex (3). Mechanisms such as dominant  
158 negative effects, haploinsufficiency, mosaicism, and digenic or epistatic interactions  
159 can further modulate disease risk and clinical presentation, underscoring the need  
160 for nuanced approaches in risk estimation. Karczewski et al. (4) made significant  
161 advances; however, the remaining challenge lies in applying the necessary statistical  
162 genomics data across all MOI for any gene-disease combination. Preliminary ap-  
163 proaches have been reported for diseases such Wilson disease, mucopolysaccharidoses,  
164 primary ciliary dyskinesia, and treatable metabolic disease, (5; 6), as reviewed by  
165 Hannah et al. (7).

166 To our knowledge, all approaches to date have been limited to single MOI, spe-  
167 cific to the given disease, or restricted to a small number of genes. We argue that  
168 an integrated approach is both necessary and highly powerful because the resulting  
169 probabilities can serve as informative priors in a Bayesian framework for variant and  
170 disease probability estimation; a perspective that is often overlooked in clinical and  
171 statistical genetics. Such a framework not only refines classical HWE-based risk esti-  
172 mates but also has the potential to enrich clinicians' understanding of what to expect

173 in a patient and to enhance the analytical models employed by bioinformaticians.

174 The resulting dataset from these necessary calculations also holds value for AI  
175 and reinforcement learning applications, providing an enriched version of the data  
176 underpinning frameworks such as AlphaFold (8) and AlphaMissense (9).

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

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

197 In this study, we focused on reporting the probability of disease observation  
198 through genome-wide assessments of gene-disease combinations. We then focused  
199 on known Inborn Errors of Immunity (IEI) genes, sometimes called the “Primary  
200 Immunodeficiency (PID) or Monogenic Inflammatory Bowel Disease genes” (12–14),  
201 to validate our approach and demonstrate its clinical relevance. This application to  
202 a well-established genotype-phenotype set, comprising over 500 gene-disease associa-  
203 tions, underscores its utility. The most recent update on the classification of IEI  
204 from the International Union of Immunological Societies (IUIS) expert committee  
205 was reported by Poli et al. (14), with an accompanying diagnostic guide (15). Our  
206 central hypothesis was that by using highly curated annotation data including popu-  
207 lation Allele Frequency (AF)s, disease phenotypes, MOI patterns, and variant classi-  
208 fications and by applying rigorous calculations based on HWE, we could accurately  
209 estimate the expected probabilities of observing disease-associated variants. Among  
210 other benefits, this knowledge can be used to derive genetic diagnosis confidence by  
211 incorporating these new priors.

212 **2 Methods**

213 **2.1 Dataset**

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

236 The following genes were used for disease cohort validations and examples. We  
237 used the two most commonly reported genes from the IEI panel Nuclear Factor  
238 Kappa B Subunit 1 (*NFKB1*) (19–22) and Cystic Fibrosis Transmembrane Conduc-  
239 tance Regulator (*CFTR*) (23–25) to demonstrate applications in AD and AR disease  
240 genes, respectively. We used Severe Combined Immunodeficiency (SCID)-specific  
241 genes AR DNA Cross-Link Repair 1C (*DCLRE1C*), AR Recombination activating  
242 gene 1 (*RAG1*), XL Interleukin 2 Receptor Subunit Gamma (*IL2RG*) to demonstrate  
243 a IEI subset disease phenotype of SCID. We also used AD Tumor necrosis factor,  
244 alpha-induced protein 3 (*TNFAIP3*) for other examples comparable to *NFKB1* since  
245 it is also causes AD pro-inflammatory disease but has more known ClinVar classifica-  
246 tions at higher AF then *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	

247

## 2.2 Variant classification occurrence probability

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

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

249 where  $p$  is the allele frequency,  $q = 1 - p$ ,  $p^2$  represents the homozygous dominant,  
 250  $2pq$  the heterozygous, and  $q^2$  the homozygous recessive genotype frequencies. For  
 251 disease phenotypes, particularly under AR MOI, the risk is traditionally linked to  
 252 the homozygous state ( $p^2$ ); however, to account for compound heterozygosity across  
 253 multiple variants, we allocated the overall gene-level risk proportionally among vari-  
 254 ants.

255 Our computational pipeline estimated the probability of observing a disease-associated  
 256 genotype for each variant and aggregated these probabilities by gene and ClinVar  
 257 classification. This approach included all variant classifications, not limited solely to  
 258 those deemed “pathogenic”, and explicitly conditioned the classification on the given  
 259 phenotype, recognising that a variant could only be considered pathogenic relative to  
 260 a defined clinical context. The core calculations proceeded as follows:

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

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

265 Note that, because each calculation is confined to one gene, no additional scaling  
 266 was required for our primary analyses ( $P_{\text{tot}}$ ). However, if this same unscaled  
 267 summation is applied across regions or variant sets of differing size or dosage sensitivity,  
 268 it can bias burden estimates. In such cases, normalisation by region length or  
 269 incorporation of gene- or region-specific dosage constraints is recommended.

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

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

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

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

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

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

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

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

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

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

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

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

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

285 This ensures

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

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

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

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

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

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

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

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

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

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

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

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

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

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

324 **6. Score-positive-total.** For use as a simple summary statistic on the resulting  
325 user-interface, we defined the *score-positive-total* as the total number of positively  
326 scored variant classifications within a given region (gene, locus, or variant set). Using  
327 the ClinVar classification assigned to a scale from -5 (benign) to +5 (pathogenic),  
328 we included only scores  $> 0$ , corresponding to some evidence of pathogenicity. The  
329 score-positive-total yields a non-normalised estimate of the prior probability that a  
330 phenotype is explained by known pathogenic variants.

331 **7. Classification scoring system.** Each ClinVar classification was assigned an  
332 integer score: pathogenic = +5, likely pathogenic = +4, pathogenic (low penetrance)  
333 = +3, likely pathogenic (low penetrance) = +2, conflicting pathogenicity = +2, likely  
334 risk allele/risk factor/association = +1, drug response/uncertain significance/no clas-  
335 sification/affects/other/not provided/uncertain risk allele = 0, protective = -3, likely  
336 benign = -4, benign = -5. No further normalisation was applied. The resulting dis-  
337 tribution (**Figure S1 A-B**) is naturally comparable to a zero-centred average rank  
338 (**C-D**). This straightforward, modular approach can be readily replaced by any com-  
339 parable evidence-based classification system. Variants with scores  $\leq 0$  were omitted,

340 since benign classifications do not inform disease likelihood in the score-positive-total  
341 summary.

342 **2.3 Integrating observed true positives and unobserved false**  
343 **negatives into a single, actionable conclusion**

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

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

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

370 **(ii) Evidence mapping and occurrence probability.** PanelAppRex variants  
371 were annotated with ClinVar clinical significance. Each label was converted to an ordi-  
372 nal evidence score  $S_i \in [-5, 5]$  and rescaled to a pathogenic weight  $W_i = \text{rescale}(S_i; -5, 5 \rightarrow$   
373  $0, 1)$ . This scoring system can be replaced with any comparable alternative. The  
374 HWE-based pipeline of Section 2.2 supplied a per-variant occurrence probability  $p_i$ .  
375 The adjusted prior was

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

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

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

379 where  $\tilde{w}_i = \max(1, S_i + 1)$  contributes an additional pseudo-count whenever  $S_i >$   
 380 0.

381 **(iv) Posterior simulation and aggregation.** For each variant  $i$  we drew  $M =$   
 382 10 000 realisations  $\pi_i^{(m)}$  and normalised within each iteration,

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

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

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

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

391 **(v) Scenario analysis.** The three scenarios were explored for a causal genetic di-  
 392 agnosis that is either simple, difficult, or impossible given the existing data. The  
 393 proband spiked data had either: (1) known classified variants only, including only  
 394 one known TP pathogenic variant, *NFKB1* p.Ser237Ter, (2) inclusion of an ad-  
 395 dditional plausible yet non-sequenced splice-donor allele *NFKB1* c.159+1G>A (likely  
 396 pathogenic) as a FN, and (3) where no known causal variants were present for a pa-  
 397 tient, one representative variant from each distinct ClinVar classification was selected  
 398 and marked as unsequenced to emulate a range of putative FNs. The selected vari-  
 399 ants were: *TNFAIP3* p.Cys243Arg (pathogenic), p.Tyr246Ter (likely pathogenic),  
 400 p.His646Pro (conflicting interpretations of pathogenicity), p.Thr635Ile (uncertain  
 401 significance), p.Arg162Trp (not provided), p.Arg280Trp (likely benign), p.Ile207Leu  
 402 (benign/likely benign), and p.Lys304Glu (benign). All subsequent steps were identi-  
 403 cal.

## 404 2.4 Validation of autosomal dominant estimates using *NFKB1*

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

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

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

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

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

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

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

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

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

424 Assuming that the proportion of CVID cases attributable to *NFKB1* is equivalent  
425 to the cohort estimate, the literature extrapolated estimate is Estimated *NFKB1* cases  $\approx$   
426  $2777 \times 0.041 \approx 114$ , with a median value of approximately 118 and a 95% confidence  
427 interval of 70 to 181 cases (derived from posterior sampling).

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

- 432     1. **Weighted adjustment (emphasising the cohort,  $w = 0.9$ ):** We assigned  
 433       90% weight to the directly observed cohort count (16) and 10% to the extrapolated  
 434       population estimate (114), thereby accounting, illustratively, for a small fraction of unobserved cases while retaining confidence in our well-characterised cohort:

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

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

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

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

442       and mixed this with the literature-based rate at 50% each (19; 20; 22). This  
 443       yields a median estimate of 67 cases and a wider 95% CrI of approximately  
 444       43 to 99 cases, capturing uncertainty in both under-ascertainment and over-generalisation.

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

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

## 456     2.5 Validation study for autosomal recessive CF using *CFTR*

457       To validate our framework for AR diseases, we focused on Cystic Fibrosis (CF).  
 458       For comparability sizes between the validation studies, we analysed the most common SNV in the *CFTR* gene, typically reported as p.Arg117His (GRCh38 Chr  
 459       7:117530975 G/A, MANE Select HGVS ENST00000003084.11: p.Arg117His). Our  
 460       goal was to validate our genome-wide probability estimates by comparing the ex-  
 461       pected number of CF cases attributable to the p.Arg117His variant in *CFTR* with  
 462       the nationally reported case count in a well-characterised disease cohort (23–25).

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

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

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

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

and

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

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

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

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

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

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

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

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

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

471 using a large effective allele count ( $\text{AN}_{\text{eff}}$ ) for illustration. By generating 10,000  
472 posterior samples of  $p$ , we obtained a distribution of the literature-based adjusted  
473 expected counts,  $E_{\text{adj}}$ .

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

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

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

## 482 **2.6 Validation of SCID-specific estimates using PID–SCID 483 genes**

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

488 **1. Incidence conversion.** Based on literature, PID occurs in approximately 1 in  
489 1,000 births, whereas SCID occurs in approximately 1 in 100,000 births. Consequently,  
490 in a population of 1,000,000 births there are about 1,000 PID cases and 10 SCID cases.  
491 To express SCID-related variant counts on a per-million PID scale, the observed SCID  
492 counts were multiplied by 100. For example, if a gene is expected to cause SCID in  
493 10 cases within the total PID population, then on a per-million PID basis the count  
494 is  $10 \times 100 = 1,000$  cases (across all relevant genes).

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

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

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

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

## 517 **2.7 Protein network and genetic constraint interpretation**

518 A PPI network was constructed using protein interaction data from STRINGdb (17).  
519 We previously prepared and reported on this dataset consisting of 19,566 proteins and  
520 505,968 interactions (<https://github.com/DylanLawless/ProteoMCLustR>). Node  
521 attributes were derived from log-transformed score-positive-total values, which in-  
522 formed both node size and colour. Top-scoring nodes (top 15 based on score) were  
523 labelled to highlight prominent interactions. To evaluate group differences in score-  
524 positive-total across major disease categories, one-way Analysis of Variance (ANOVA)  
525 was performed followed by Tukey Honestly Significant Difference (HSD) post hoc  
526 tests (and non-parametric Dunn's test for confirmation). GnomAD v4.1 constraint  
527 metrics data was used for the PPI analysis and was sourced from Karczewski et al.  
528 (4). This provided transcript-level metrics, such as observed/expected ratios, Loss-Of-  
529 function Observed/Expected Upper bound Fraction (LOEUF), Probability of being  
530 Loss-of-function Intolerant (pLI), and Z-scores, quantifying Loss-of-Function (LOF)  
531 and missense intolerance, along with confidence intervals and related annotations for  
532 211,523 observations.

## 533 **2.8 Gene set enrichment test**

534 To test for overrepresentation of biological functions, the prioritised genes were com-  
535 pared against gene sets from MsigDB (including hallmark, positional, curated, motif,  
536 computational, GO, oncogenic, and immunologic signatures) and WikiPathways using  
537 hypergeometric tests with FUMA (26; 27). The background set consisted of 24,304  
538 genes. Multiple testing correction was applied per data source using the Benjamini-  
539 Hochberg method, and gene sets with an adjusted P-value  $\leq 0.05$  and more than one  
540 overlapping gene are reported.

541 **2.9 Deriving novel PID classifications by genetic PPI and**  
542 **clinical features**

543 We recategorised 315 immunophenotypic features from the original IUIS IEI annotations,  
544 reducing the original multi-level descriptors (e.g. “decreased CD8, normal or  
545 decreased CD4”) first to minimal labels (e.g.“low”) and second to binary outcomes  
546 (normal vs. not-normal) for T cells, B cells, neutrophils, and immunoglobulins Each  
547 gene was mapped to its PPI cluster derived from STRINGdb and Uniform Manifold  
548 Approximation and Projection (UMAP) embeddings from previous steps. We first  
549 tested for non-random associations between these four binary immunophenotypes and  
550 PPI clusters using  $\chi^2$  tests. To generate a data-driven PID classification, we trained  
551 a decision tree (rpart) to predict PPI cluster membership from the four immunophe-  
552 notypic features plus the traditional IUIS Major and Subcategory labels. Hyperpa-  
553 rameters (complexity parameter = 0.001, minimum split = 10, minimum bucket = 5,  
554 maximum depth = 30) were optimised via five-fold cross validation using the caret  
555 framework. Terminal node assignments were then relabelled according to each group’s  
556 predominant abnormal feature profile.

557 **2.10 Probability of observing AlphaMissense pathogenicity**

558 We obtained the subset pathogenicity predictions from AlphaMissense via the Al-  
559 phaFold database and whole genome data from the studies data repository(9; 18). The  
560 AlphaMissense data (genome-aligned and amino acid substitutions) were merged with  
561 the panel variants based on genomic coordinate and HGVS annotation. Occurrence  
562 probabilities were log-transformed and adjusted (y-axis displaying  $\log_{10}(\text{occurrence}$   
563 prob + 1e-5) + 5), to visualise the distribution of pathogenicity scores across the  
564 residue sequence. A Kruskal-Wallis test was used to compare the observed disease  
565 probability across clinical classification groups.

566 **2.11 Probability model definitions**

567 Estimating disease risk requires accounting for both variant penetrance,  $P(D | G)$ ,  
568 where  $D$  denotes the disease state and  $G$  the genotype, and the fraction of cases  
569 attributable to a given variant,  $P(G | D)$ . In a fully penetrant single-variant model  
570 ( $P(D | G) = 1$ ), the lifetime risk  $P(D)$  equals the genotype frequency  $P(G)$ . For an  
571 allele with population frequency  $p$ , this gives  $P(D) = p^2$  for a recessive mode of inheri-  
572 tance and  $P(D) = 2p(1 - p) \approx 2p$  for a dominant mode. With incomplete penetrance,  
573  $P(D) = P(G) P(D | G)$ , and when multiple variants contribute to disease,

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

574 Because both  $P(D | G)$  and  $P(G | D)$  are often uncertain, we integrate ClinVar  
575 clinical classifications, population allele frequencies and curated gene–disease associa-  
576 tions, assuming James–Stein shrinkage to derive robust aggregate priors. By focusing  
577 on a filtered set of variants  $\mathcal{V}$  where each  $P(G_i | D)$  is the probability that disease  $D$   
578 is attributable to variant  $i$  and assuming  $\sum_{i \in \mathcal{V}} P(G_i | D) \approx 1$ , we obtain calibrated  
579 estimates of genotype frequency  $P(G)$  despite uncertainty in individual parameters.

## 580 3 Results

### 581 3.1 Occurrence probability across disease genes

582 Our study integrated large-scale annotation databases with gene panels from PanelAp-  
583 pRex to systematically assess disease genes by MOI (12). By combining population  
584 allele frequencies with ClinVar clinical classifications, we computed an expected occur-  
585 rence probability for each SNV, representing the likelihood of encountering a variant  
586 of a specific pathogenicity for a given phenotype. We report these probabilities for  
587 54,814 ClinVar variant classifications across 557 genes (linked dataset (28)).

588 We focused on panels related to Primary Immunodeficiency or Monogenic Inflam-  
589 matory Bowel Disease, using PanelAppRex panel ID 398. **Figure 1** displays all  
590 reported ClinVar variant classifications for this panel. The resulting natural scaling  
591 system (-5 to +5) accounts for the frequently encountered combinations of classifica-  
592 tion labels (e.g. benign to pathogenic). The resulting dataset (28) is briefly shown in  
593 **Table 1** to illustrate that our method yielded estimates of the probability of observing  
594 a variant with a particular ClinVar classification.

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

### 595 3.2 Integrating observed true positives and unobserved false 596 negatives into a single, actionable conclusion

597 Having established a probabilistic framework for estimating the prior probability of  
598 observing disease-associated variants under different inheritance modes, we then ap-  
599 plied this model to an example patient to demonstrate its potential for clinical ge-  
600 netics. The algorithm first verified that all known pathogenic positions have been  
601 sequenced and observed as reference (true negatives), and identified any positions  
602 that were either observed as variant (true positives) or not assessable due to missing

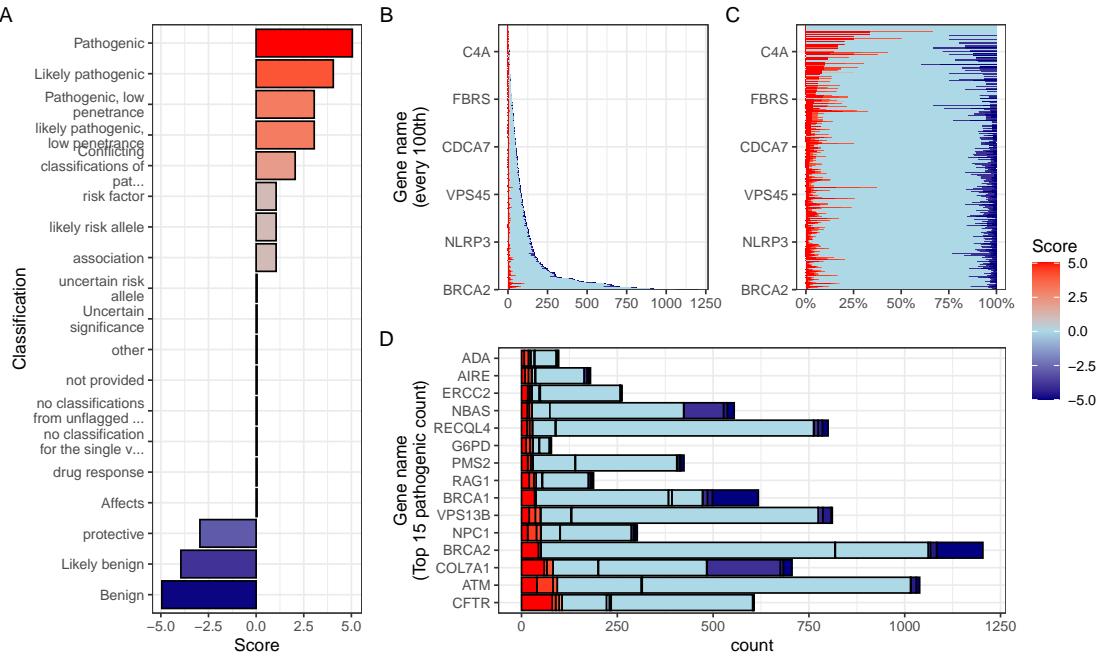


Figure 1: **Summary of ClinVar clinical significance classifications in the PID gene panel.** (A) Shows the numeric score coding for each classification (i.e. -5 benign to +5 pathogenic) as defined in methods Section 2.2. (B) Displays the stacked absolute count of classifications per gene. The same counts 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.

603 sequence data of failed QC. These missing sites represent potential false negatives.  
 604 By jointly modelling the observed and unobserved space, the method yielded a cali-  
 605 brated, evidence-weighted probability that at least one damaging causal variant could  
 606 be present within a gene.

### 607 3.3 Scenario one - complete coverage and simple diagnosis

608 We present the results from three scenarios for an example single-case patient being  
 609 investigated for the genetic diagnosis of IEI. **Figure S2** shows the results of the first  
 610 simple scenario, in which only one known pathogenic variant, *NFKB1* p.Ser237Ter,  
 611 was observed and all other previously reported pathogenic positions were successfully  
 612 sequenced and confirmed as reference. In this setting, the model assigned the full  
 613 posterior probability to the observed allele, yielding 100 % confidence that all present  
 614 evidence supported a single, true positive causal explanation. The most strongly sup-  
 615 ported observed variant was p.Ser237Ter (posterior: 0.594). The strongest (probabil-  
 616 ity of observing) non-sequenced variant was a benign variant p.Thr567Ile (posterior:  
 617 0). The total probability of a causal diagnosis given the available evidence was 1 (95%)

618 CI: 1–1) (**Table S1**).

619 **3.4 Scenario two - incomplete coverage and complex diagno-**  
620 **sis**

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

630 Bayesian integration of every annotated allele yielded the quantitative CrIs for  
631 pathogenic attribution that (i) preserve Hardy-Weinberg expectations, (ii) accommo-  
632 date AD, AR, XL inheritance, and (iii) carry explicit uncertainty for non-sequenced  
633 (or failed QC) genomic positions. The incremental calculation steps for the variant  
634 in question are shown in **Figure 2**.

635 First, **Figure 2 (A)** depicts the prior landscape where occurrence probabilities  
636 are partitioned by observed or missing status and by causal or non-causal evidence,  
637 with colour reflecting the underlying ClinVar score. **Figure 2 (B)** shows posterior  
638 normalisation which concentrates probability density on two high-confidence (high ev-  
639 idence score) alleles since the benign variants are, by definition, non-causal. **Figure 2**  
640 (**C**) shows the resulting per-variant probability of being simultaneously damaging and  
641 causal; only the confirmed present (true positive) nonsense variant p.Ser237Ter and  
642 the (false negative) hypothetical splice-donor c.159+1G>A retain substantial support.  
643 Restricting the view to causal candidates in **Figure 2 (D)** confirms that posterior  
644 mass is distributed across these two variants. **Figure 2 (E)** decomposes the total  
645 damaging probability into observed (approximately 40 %) and missing (approximately  
646 34 %) sources, whereas **Figure 2 (F)** summarises the gene-level posterior: inclusion  
647 of the splice-site allele (scenario 2) produces a median probability of 0.542 with a  
648 95 % CrI of 0.264–0.8.

649 Numerically, the present variant p.Ser237Ter accounts for 0.399 of the posterior  
650 share, and the potentially causal but missing splice-donor allele c.159+1G>A con-  
651 tributes 0.339. The remaining alleles together explain a negligible share (**Table S2**).  
652 Thus, we can report that in this patient’s scenario the probability of correct genetic  
653 diagnosis due to *NFKB1* p.Ser237Ter is 0.542 (95 % CrI of 0.264–0.8) given that a  
654 likely alternative remains to be confirmed for this patient. Upon confirmation that  
655 the second variant is not present, the confidence will rise to 1 (95 % CrI of 1–1) as  
656 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
<b>Summary</b>	Overall probability of correct causal diagnosis due to SNV/INDEL given the currently available evidence: 0.511 (95% CI 0.237–0.774).	

### 657 3.5 Scenario three - currently impossible diagnosis

658 **Figure S3** shows the third scenario, in which no observed variants were detected in  
 659 the proband for *NFKB1*. Instead, a broad range of plausible FN were detected as missing  
 660 for the gene *TNFAIP3*. The strongest (probability of observing and pathogenic)  
 661 of these non-sequenced variants was p.Cys243Arg (posterior: 0.347). However, the  
 662 total probability of a causal diagnosis for the patient *given the available evidence* was  
 663 0 (95% CI: 0–0) since these missing variants must be accounted for (**Table S3**). Upon  
 664 confirmation, these probabilities can update, as per scenario two.

### 665 3.6 Posterior probabilities are calculated across all qualifying 666 variants

667 While only the top-ranked gene/variant is shown in each of the three scenarios, we  
 668 emphasise that the same posterior probability and CrI calculations are performed  
 669 across all plausible candidates. In real-world diagnostics, we commonly find multiple  
 670 variants to carry non-negligible probabilities. Our framework explicitly quantifies  
 671 these competing hypotheses, enabling a ranked interpretation that reflects the totality  
 672 of evidence. Overlapping CrI do not indicate ambiguity in the method, but rather a  
 673 principled measure of remaining uncertainty. This output can directly inform follow-  
 674 up actions, such as functional validation or treatment trials, and supports the use  
 675 of CrI thresholds as a transparent decision-making aid when data are incomplete or

676 equivocal.

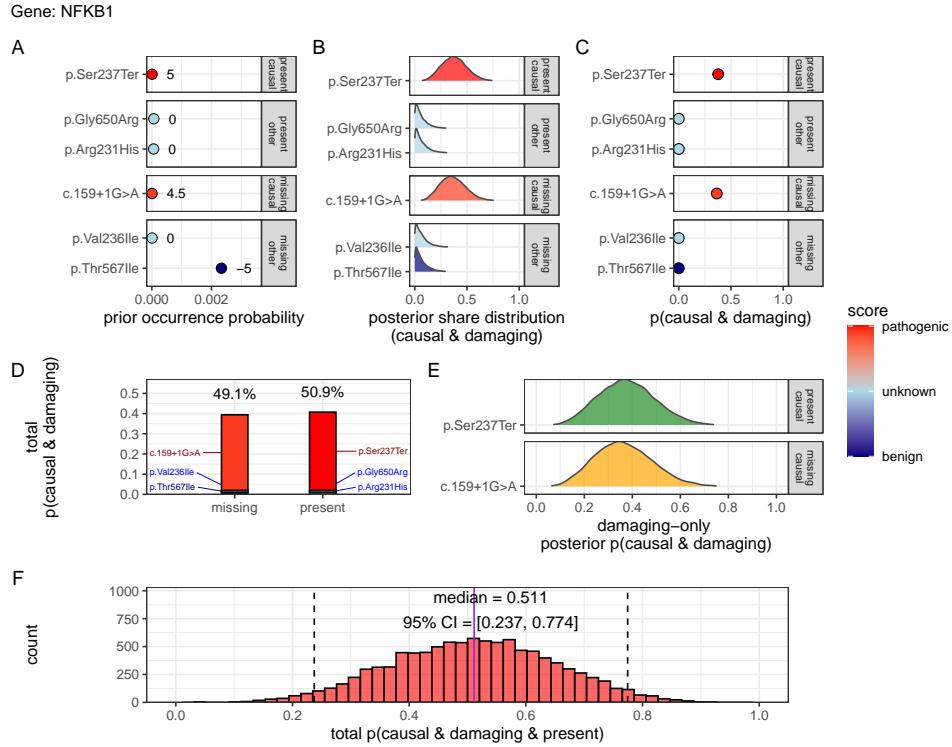


Figure 2: Quantification of present (TP) and missing (FN) causal genetic variants for disease in *NFKB1* (scenario 2). The example proband carried three known heterozygous variants, including pathogenic p.Ser237Ter, and had incomplete coverage at three additional loci, including likely-pathogenic splice-site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

### 3.7 Validation studies

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

To validate our genome-wide probability estimates for AD disorders, we focused on *NFKB1*. We used a reference dataset from Tuijnjenburg et al. (19), in which whole-

681 genome sequencing of 846 PID patients identified *NFKB1* as one of the genes most  
682 strongly associated with the disease, with 16 *NFKB1*-related CVID cases attributed to  
683 AD heterozygous variants. Our goal was to compare the predicted number of *NFKB1*-  
684 related CVID cases with the reported count in this well-characterised national-scale  
685 cohort.

686 Our model calculated 0 known pathogenic variant *NFKB1*-related CVID cases  
687 in the UK with a minimal risk of 456 unknown de novo variants. In the reference  
688 cohort, 16 *NFKB1* CVID cases were reported. We additionally wanted to account for  
689 potential under-reporting in the reference study. We used an extrapolated national  
690 CVID prevalence which yielded a median estimate of 118 cases (95% CI: 70–181),  
691 while a Bayesian-adjusted mixture estimate produced a median of 67 cases (95% CI:  
692 43–99). **Figure S5 (A)** illustrates that our predicted values reflect these ranges and  
693 are closer to the observed count. This case supports the validity of our integrated  
694 probability estimation framework for AD disorders, and represents a challenging ex-  
695 ample where pathogenic SNV are not reported in the reference population of gnomAD.  
696 Our min-max values successfully contained the true reported values.

### 697 3.7.2 Validation of recessive disease occurrence with *CFTR*

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

703 In contrast, the nationally reported number of CF cases was 714, as recorded in the  
704 UK Cystic Fibrosis Registry 2023 Annual Data Report (23). To account for factors  
705 such as reduced penetrance and the mortality-adjusted expected genotype, we derived  
706 a Bayesian-adjusted estimate via posterior simulation. Our Bayesian approach yielded  
707 a median estimate of 740 cases (95% Confidence Interval (CI): 696, 786) and a  
708 mixture-based estimate of 727 cases (95% CI: 705, 750). **Figure S5 (B)** illustrates  
709 the close concordance between the predicted values, the Bayesian-adjusted estimates,  
710 and the national report supports the validity of our approach for estimating disease.

711 **Figure S4** shows the final values for these genes of interest in a given population  
712 size and phenotype. It reveals that an allele frequency threshold of approximately  
713 0.000007 is required to observe a single heterozygous carrier of a disease-causing  
714 variant in the UK population for both genes. However, owing to the AR MOI pattern  
715 of *CFTR*, this threshold translates into more than 100,000 heterozygous carriers,  
716 compared to only 456 carriers for the AD gene *NFKB1*. Note that this allele frequency  
717 threshold, being derived from the current reference population, represents a lower  
718 bound that can become more precise as public datasets continue to grow. This marked  
719 difference underscores the significant impact of MOI patterns on population carrier  
720 frequencies and the observed disease prevalence.

721 **3.7.3 Interpretation of ClinVar variant occurrences**

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

727 **3.7.4 Validation of SCID-specific disease occurrence**

728 Given that SCID is a subset of PID, our probability estimates reflect the likelihood of  
729 observing a genetic variant as a diagnosis when the phenotype is PID. However, we  
730 additionally tested our results against SCID cohorts in **Figure S8**. The summarised  
731 raw cohort data for SCID-specific gene counts are summarised and compared across  
732 countries in **Figure S7**. True counts for *IL2RG* and *DCLRE1C* from ten distinct lo-  
733 cations yielded 95% CI surrounding our predicted values. For *IL2RG*, the prediction  
734 was low (approximately 1 case per 1,000,000 PID), as expected since loss-of-function  
735 variants in this XL gene are highly deleterious and rarely observed in gnomAD. In con-  
736 trast, the predicted value for *RAG1* was substantially higher (553 cases per 1,000,000  
737 PID) than the observed counts (ranging from 0 to 200). We attributed this discrep-  
738 ancies to the lower penetrance and higher background frequency of *RAG1* variants in  
739 recessive inheritance, whereby reference studies may underreport the true national  
740 incidence. Overall, we report that agreement within an order of magnitude was tol-  
741 erable given the inherent uncertainties from reference studies arising from variable  
742 penetrance and allele frequencies.

743 **3.8 Application to inborn errors of immunity**

744 **3.8.1 Genetic constraint in high-impact protein networks**

745 We next applied our framework to IEI, a disease area in which we have expertise and  
746 which offers a well-curated gene set to validate genome-wide estimates and demon-  
747 strate potential applications (14).

748 Given that pathogenicity in IEI may reflect shared molecular pathways, we in-  
749 tegrated ClinVar-derived variant probability estimates with PPI data to quantify  
750 pathogenic burden per gene and examine whether genetic constraint aggregates within  
751 specific networks and corresponds to established IEI classifications and immunophe-  
752 notypes (4; 17).

753 **3.8.2 Score-positive-total within IEI PPI network**

754 The ClinVar classifications reported in **Figure 1** were scaled -5 to +5 based on  
755 their pathogenicity. We were interested in positive (potentially damaging) but not

756 negative (benign) scoring variants, which are statistically incidental in this analysis.  
757 We tallied gene-level positive scores to give the score-positive-total metric. **Figure S9**  
758 (**A**) shows the PPI network of disease-associated genes, where node size and colour  
759 encode the score-positive-total (log-transformed). The top 15 genes with the highest  
760 total prior probabilities of being observed with disease are labelled (as per **Figure**  
761 **1**).

762 **3.8.3 Association analysis of score-positive-total across IEI categories**

763 We checked for any statistical enrichment in score-positive-totals, which represents the  
764 expected observation of pathogenicity, between the IEI categories. One-way ANOVA  
765 revealed an effect of major disease category on score-positive-total ( $F(8, 500) =$   
766 2.82,  $p = 0.0046$ ), indicating that group means were not identical, which we ob-  
767 served in **Figure S9 (B)**. However, despite some apparent differences in median  
768 scores across categories (i.e. 9. Bone Marrow Failure (BMF)), the Tukey HSD post  
769 hoc comparisons **Figure S9 (C)** showed that all pairwise differences had 95% CIs  
770 overlapping zero, suggesting that individual group differences were not significant.

771 **3.8.4 UMAP embedding of the PPI network**

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

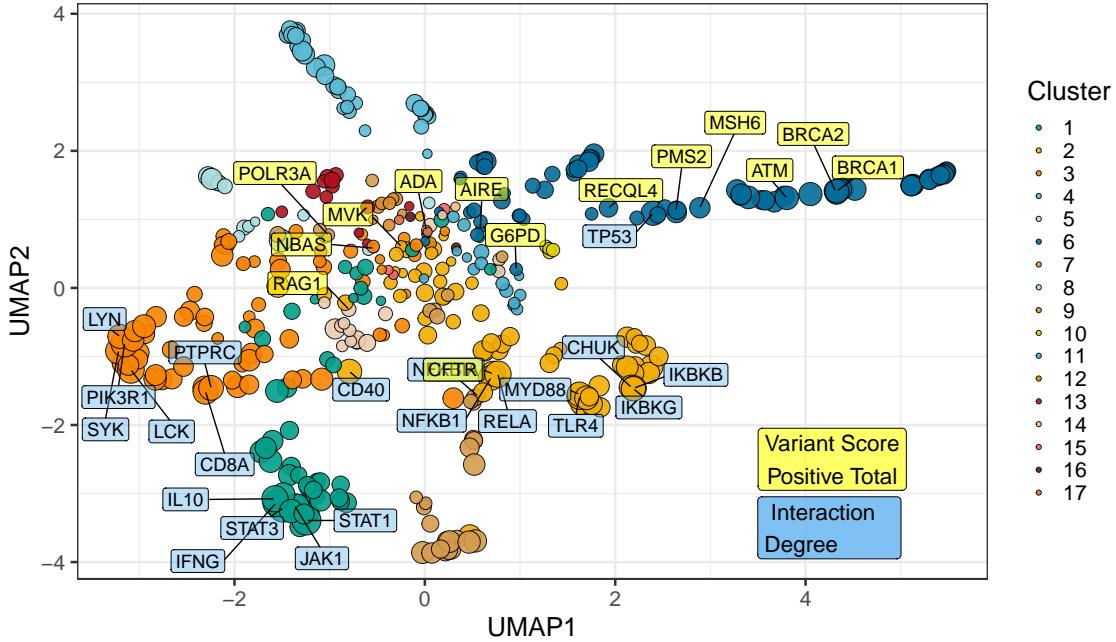


Figure 3: **UMAP embedding of the PPI network.** The plot projects the high-dimensional protein-protein interaction network into two dimensions, with nodes coloured by cluster and sized by interaction degree. Blue labels indicate hub genes (degree above the 95th percentile) and yellow labels mark the top 15 genes by score-positive-total (damaging ClinVar classifications). The spatial segregation suggests that genes with high pathogenic variant loads are distinct from highly connected nodes.

781    **3.8.5 Hierarchical clustering of enrichment scores for major disease cate-**  
 782    **gories**

783    **Figure S10** presents a heatmap of standardised residuals for major disease categories  
 784 across network clusters, as per **Figure 3**. A dendrogram clusters similar disease cate-  
 785 gories, while the accompanying bar plot displays the maximum absolute standardised  
 786 residual for each category. Notably, (8) Complement Deficiencies (CD) shows the  
 787 highest maximum enrichment, followed by (9) BMF. While all maximum values  
 788 exceed 2, the threshold for significance, this likely reflects the presence of protein  
 789 clusters with strong damaging variant scores rather than uniform significance across  
 790 all categories (i.e. genes from cluster 4 in 8 CD).

791    **3.8.6 PPI connectivity, LOEUF constraint and enriched network cluster**  
 792    **analysis**

793    Based on the preliminary insight from **Figure S10**, we evaluated the relationship  
 794 between network connectivity (PPI degree) and LOEUF constraint (LOEUF upper rank)  
 795 Karczewski et al. (4) using Spearman's rank correlation. Overall, there was a weak

796 but significant negative correlation ( $\rho = -0.181$ ,  $p = 0.00024$ ) at the global scale,  
797 indicating that highly connected genes tend to be more constrained. A supplementary  
798 analysis (**Figure S11**) did not reveal distinct visual associations between network  
799 clusters and constraint metrics, likely due to the high network density. However  
800 once stratified by gene clusters, the natural biological scenario based on quantitative  
801 PPI evidence (17), some groups showed strong correlations; for instance, cluster 2  
802 ( $\rho = -0.375$ ,  $p = 0.000994$ ) and cluster 4 ( $\rho = -0.800$ ,  $p < 0.000001$ ), while others did  
803 not. This indicated that shared mechanisms within pathway clusters may underpin  
804 genetic constraints, particularly for LOF intolerance. We observe that the score-  
805 positive-total metric effectively summarises the aggregate pathogenic burden across  
806 IEI genes, serving as a robust indicator of genetic constraint and highlighting those  
807 with elevated disease relevance.

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

### 812 3.8.7 New insight from functional enrichment

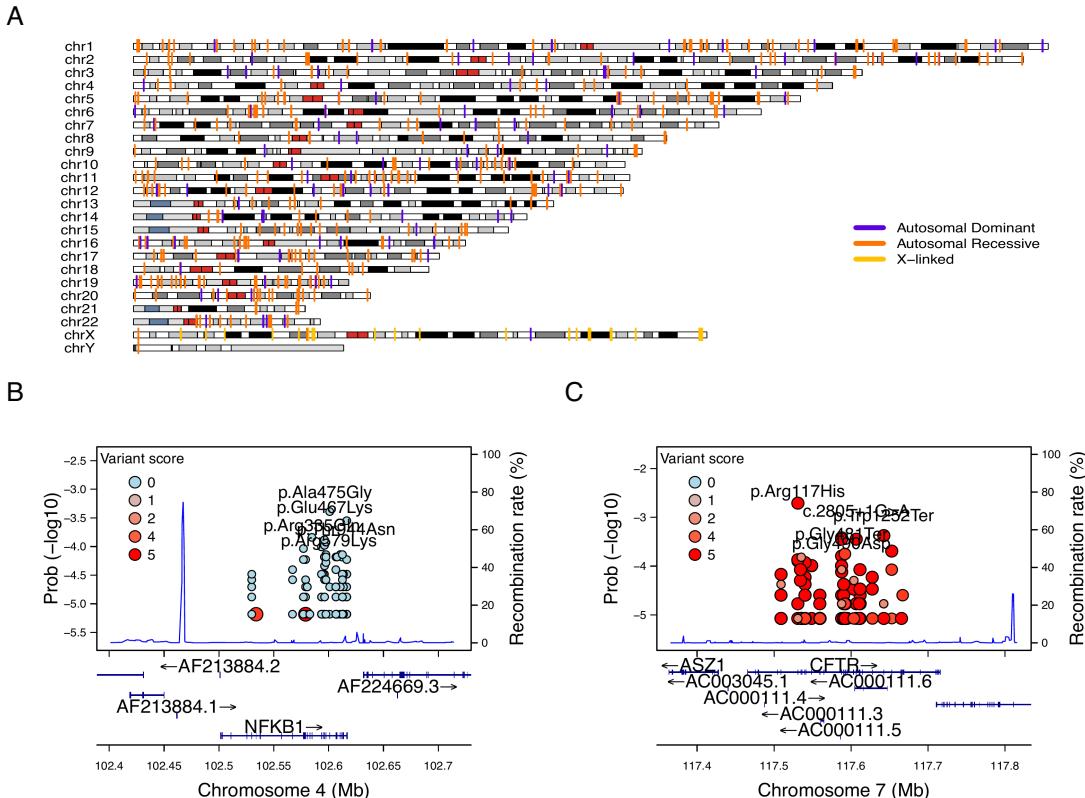
813 To interpret the functional relevance of our prioritised IEI gene sets with the highest  
814 load of damaging variants (i.e. clusters 2 and 4 in **Figure S12**), we performed  
815 functional enrichment analysis for known disease associations using MsigDB with  
816 FUMA (i.e. GWAScatalog and Immunologic Signatures) (26). Composite enrichment  
817 profiles (**Figure S13**) reveal that our enriched PPI clusters were associated with  
818 distinct disease-related phenotypes, providing functional insights beyond traditional  
819 IUIS IEI groupings (14). The gene expression profiles shown in **Figure S14** (GTEx v8  
820 54 tissue types) offer the tissue-specific context for these associations. Together, these  
821 results enable the annotation of IEI gene sets with established disease phenotypes,  
822 supporting a data-driven classification of IEI.

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

### 832 3.8.8 Genome-wide gene distribution and linkage disequilibrium

833 **Figure 4 (A)** shows a genome-wide karyoplot of all IEI panel genes across GRCh38,  
834 with colour-coding based on MOI. Figures (B) and (C) display zoomed-in locus plots

for *NFKB1* and *CFTTR*, respectively. In **Figure 4 (B)**, the probability of observing variants with known classifications is high only for variants such as p.Ala475Gly, which are considered benign in the AD *NFKB1* gene that is intolerant to LOF. In **Figure 4 (C)**, high probabilities of observing patients with pathogenic variants in *CFTTR* are evident, reproducing this well-established phenomenon. Furthermore, the analysis of Linkage Disequilibrium (LD) using  $R^2$  shows that high LD regions can be modelled effectively, allowing independent variant signals to be distinguished.



**Figure 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 *CFTTR* 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.

### 842 3.8.9 Novel PID classifications derived from genetic PPI and clinical fea- 843 tures

844 We recategorised 315 immunophenotypic features from the original IUIS IEI annotations,  
 845 reducing detailed descriptions (e.g. “decreased CD8, normal or decreased CD4”)

846 to minimal labels (e.g. “low”) and then binarising them (normal vs. not-normal) for  
847 T cells, B cells, Immunoglobulin (Ig) and neutrophils (**Figure S15**). These sim-  
848 plified profiles were mapped onto STRINGdb PPI clusters, revealing non-random  
849 distributions ( $\chi^2 < 1e-13$ ; **Figure S16**), indicating that network context captures  
850 key immunophenotypic variation.

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

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

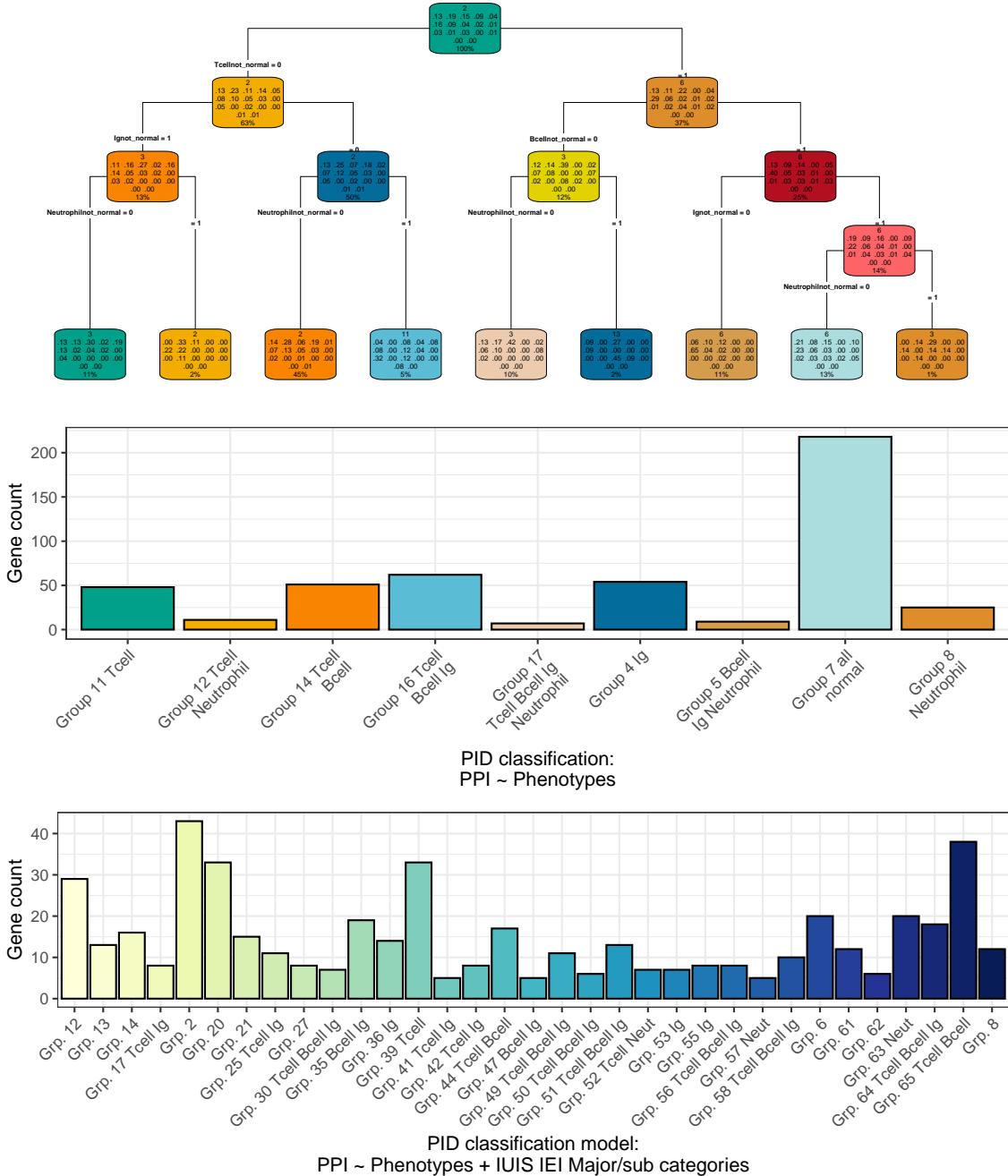


Figure 5: **Fine-tuned model for PID classification.** (Top) In each terminal node, the top block indicates the number of genes in the node; the middle block shows the fitted class probabilities (which sum to 1); and the bottom block displays the percentage of the total sample in that node. These metrics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

874 **3.9 Probability of observing AlphaMissense pathogenicity**

875 AlphaMissense provides pathogenicity scores for all possible amino acid substitutions;  
876 however, our results in **Figure 6** show that the most probable observations in pa-  
877 tients occur predominantly for benign or unknown variants. This finding places the  
878 likelihood of disease-associated substitutions into perspective and offers a data-driven  
879 foundation for future improvements in variant prediction. The values in **Figure 6**  
880 (**A**) can be directly compared to **Figure 1 (D)** to view the distribution of classifi-  
881 cations. A Kruskal-Wallis test was used to compare the observed disease probability  
882 across clinical classification groups and no significant differences were detected. In  
883 general, most variants in patients are classified as benign or unknown, indicating  
884 limited discriminative power in the current classification, such that pathogenicity pre-  
885 diction does not infer occurrence prediction (**Figure S18**). Inverse correlation likely  
886 depends on factors like MOI and intolerance to LOF.

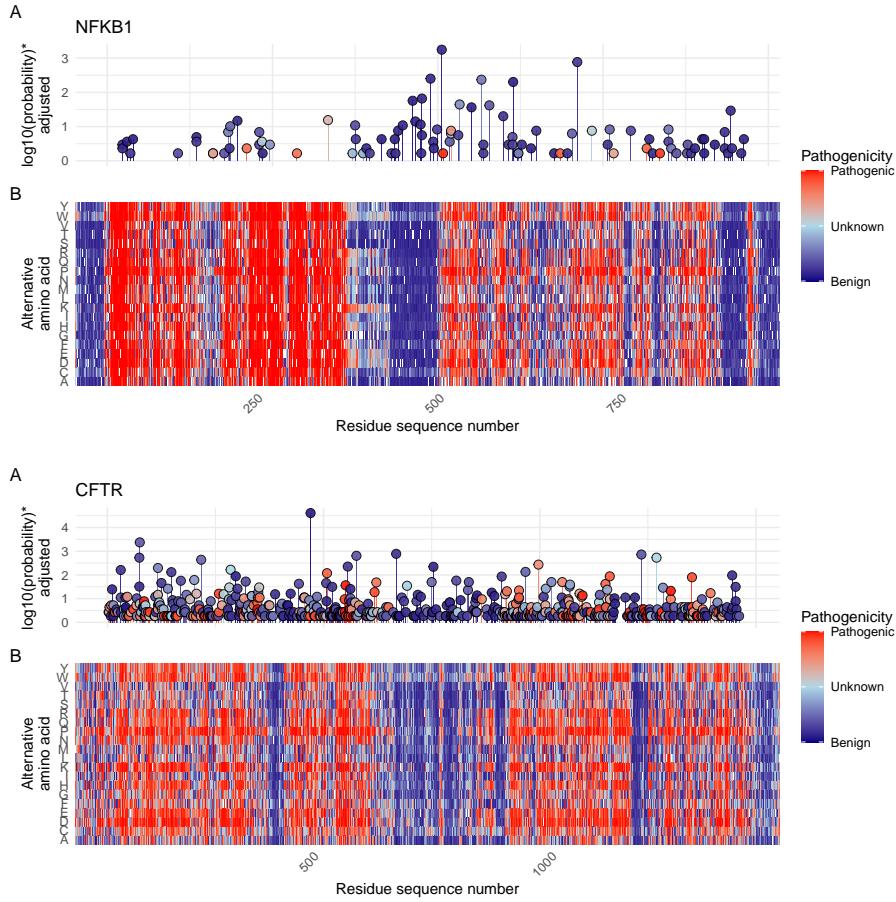


Figure 6: **(A) Probabilities of observing a patient with (B) AlphaMissense-derived pathogenicity scores.** Although AlphaMissense provides scores for every possible amino acid substitution, the most frequently observed variants in patients tend to be classified as benign or of unknown significance.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

### 3.10 Integration of variant probabilities into IEI genetics data

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

**Figure 7** shows the interface summarising integrated variant data. We include pre-calculated summary statistics and clinical significance as numerical metrics. Key quantiles (min, Q1, median, Q3, max) for each gene are rendered as sparkline box plots, and dynamic URLs link table entries to external databases (e.g. ClinVar, Online Mendelian Inheritance in Man (OMIM), AlphaFold) as per **Section 3.1**. The prepared data are available for bioinformatic application (28) as per **Section 3.2**.

Major category	Subcategory	Disease	Genetic defect	Inheritance	Score positive total observing pathogenic	Prior prob of observing pathogenic	ClinVar SNV classification	ClinVar all variant reports	OMIM	Alpha Missense / Uniprot ID	HPO combined	HPO term
All				All								
5. PD	2. Defects of Motility	Cystic fibrosis	CFTR	AR	106	80 / 9 / 497 / 4	1782 / 592 / 4971 / 200	602421	P13569	HP:0002715; HP:0002783	Abnormality of II Recurrent lower infections	
2. CID+	2. DNA Repair Defects other than those listed in Table 1	Ataxia-telangiectasia	ATM	AR	93	40 / 11 / 922 / 24	4196 / 1109 / 18503 / 1084	607585	Q13315	HP:0002715; HP:005403	Abnormality of II T lymphocyte	
9. BMF		Fanconi Anemia Type D1	BRCA2	AR	51	44 / 1 / 1010 / 143	9484 / 866 / 19120 / 841	605724	PS1587	HP:0002721; HP:0005528	Immunodeficiency hypcellularity	
5. PD	1. Congenital Neutropenias	Cohen syndrome	VPS13B	AR	49	20 / 13 / 725 / 37	1043 / 890 / 4790 / 5515	607817	Q7Z7G8	HP:0002715; HP:0410252	Abnormality of II Chronic neutrop	
9. BMF		Fanconi Anemia Type S	BRCA1	AR	38	24 / 1 / 1434 / 148	7320 / 405 / 16933 / 6868	617883	P38398	HP:0002721; HP:0005528	Immunodeficiency hypcellularity	
1. CID	2. T-B- SCID	RAG1 deficiency	RAG1	AR	38	20 / 6 / 139 / 14	157 / 98 / 792 / 184	179615	P15918	HP:0002715; HP:005403	Abnormality of II T lymphocyte	

**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.

901 **4 Discussion**

902 Our study presents, to our knowledge, the first comprehensive framework for cal-  
903 culating prior probabilities of observing disease-associated variants and the first to  
904 demonstrate the method for an evidence-aware genetic diagnosis with Cri (5; 7). By  
905 integrating large-scale genomic annotations, including population allele frequencies  
906 from gnomAD (4), variant classifications from ClinVar (10), and functional annota-  
907 tions from resources such as dbNSFP, with classical HWE-based calculations, we de-  
908 rived robust estimates for 54,814 ClinVar variant classifications across 557 IEI genes  
909 implicated in PID and monogenic inflammatory bowel disease (12; 14). Although  
910 our results focus on IEI, the genome-wide framework also supports all inheritance  
911 patterns: AD and XL require a single pathogenic allele, whereas AR demands ho-  
912 mozygous or compound heterozygous states. Classical HWE-based estimates thus  
913 furnish baseline occurrence probabilities and serve as robust priors for Bayesian risk  
914 models, a practice underutilised until the advent of large-scale databases (4; 9; 10; 12).

915 A major deficit in current clinical genetics is the prevailing focus on confirming  
916 only the presence of TP variants. Our approach yielded three key results to overcome  
917 this hurdle. We generated per-variant priors across all MOI. The patient's results of  
918 observed and unobserved variants were integrated into a single posterior probability  
919 of carrying a damaging causal allele. As demonstrated in **Table S2** and **Figure 2**,  
920 this key result delivers a clinically applicable, interpretable probability that combines  
921 both detected and potentially unobserved variants. When whole-genome sequencing  
922 analyses are not yet available, the score-positive-total metric can serve as an optional  
923 decision aid, enabling manual, evidence-based ranking of candidate genes to prioritise  
924 diagnoses in patients with overlapping phenotypes.

925 We acknowledge that our framework is currently focused (but not restricted) on  
926 SNVs and does not incorporate numerous other complexities of genetic disease, such  
927 as structural variants, de novo variants, hypomorphic alleles, overdominance, variable  
928 penetrance, tissue-specific expression, the Wahlund effect, pleiotropy, and others (3).  
929 In certain applications, more refined estimates would benefit from including factors  
930 such as embryonic lethality, condition-specific penetrance, and age of onset (7). Our  
931 analysis also relies on simplifying assumptions of random mating, an effectively infinite  
932 population, and the absence of migration, novel mutations, or natural selection. We  
933 demonstrated the genome-wide gene distribution and MOI for the IEI panel relative  
934 to LD showing that it is an important consideration and is feasible. However, LD  
935 is a challenging feature that requires accurate implementation which depends on the  
936 whole genome population-based pairwise genotype matrices for the given population.  
937 We used the reference global population AFs, which is more generalisable but less  
938 accurate than population-specific AF values.

939 In the example single-case diagnosis scenarios, our approach enabled high-confidence  
940 attribution to a known pathogenic variant while also capturing the potential impact  
941 of a likely-pathogenic splice-site allele that was missed by sequencing. Scenario two  
942 showed a common diagnostic challenge where a strong candidate exists alongside an

943 unconfirmed but plausible alternative. Our method distributes confidence across both  
 944 possibilities. Conventional approaches focus only on detecting TP and cannot provide  
 945 this insight. By quantifying residual uncertainty, we can generate structured reports  
 946 that clearly distinguish supported, excluded, and plausible-but-unseen variants. We  
 947 call this “evidence-aware” interpretation. When combined with genome-wide priors  
 948 from the full range of disease-gene panels, this approach applies to any phenotype  
 949 from PanelAppRex. By combining variant classification, allele frequency, MOI, and  
 950 sequencing quality metrics, our method creates a scalable foundation for evidence-  
 951 aware diagnostics in clinical genomics.

952 Estimating disease risk in genetic studies is complicated by uncertainties in key  
 953 parameters such as variant penetrance and the fraction of cases attributable to specific  
 954 variants (3). In the simplest model, where a single, fully penetrant variant causes  
 955 disease, the lifetime risk  $P(D)$  is equivalent to the genotype frequency  $P(G)$ . For an  
 956 allele with frequency  $p$  (ignoring LD for AR), this translates to:

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

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

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

961 Because both penetrance and  $P(G | D)$  are often uncertain, solving this equation  
 962 systematically poses a major challenge, which we incidentally tackled in the validation  
 963 studies (29; 30).

964 Our framework addresses this challenge by combining variant classifications, pop-  
 965 ulation allele frequencies, and curated gene-disease associations. While imperfect on  
 966 an individual level, these sources exhibit predictable aggregate behaviour, supported  
 967 by James-Stein estimation principles (31). Curated gene-disease associations help  
 968 identify genes that are explainable for most disease cases, allowing us to approximate  
 969  $P(G | D)$  close to one. In this way, we obtain robust estimates of  $P(G)$  (the fre-  
 970 quency of disease-associated genotypes), even when exact values of penetrance and  
 971 case attribution remain uncertain.

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

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

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

For the IEI-specific investigation, we showed that immunophenotypic and network-derived features can be used to train and test models that predict PPIs. From this, we derived a new, simplified classification of immune features for IEI genes. We have listed the new immunophenotypic categories (e.g. T cell low) in the user database, however we have not included the detailed cluster assignments (e.g. PPI groups) because they are too complex for direct interpretation manually. Instead, our demonstration provides worked examples that bioinformaticians can use to perform more refined clustering in larger studies.

Moreover, because variant sets can be collapsed instead of relying on the gene-level, our method complements existing statistical approaches for aggregating variant effects with methods like Sequence Kernel Association Test (SKAT) and Aggregated Cauchy Association Test (ACAT) (32–35) and multi-omics integration techniques (36; 37). It also remains consistent with established variant interpretation guidelines from the American College of Medical Genetics and Genomics (ACMG) (38) and complementary frameworks (39; 40), as well as QC protocols (41; 42). Standardised reporting for qualifying variant sets, such as ACMG Secondary Findings v3.2 (43), further contextualises the integration of these probabilities into clinical decision-making.

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

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

## 1008 5 Conclusion

1009 We present a statistical framework that resolves the long-standing challenge of quantifying  
1010 the probability that any candidate variant is causal for genetic disease. Unlike  
1011 traditional methods that stop at pathogenic or benign labels, our approach integrates  
1012 both observed variants and the possibility of unobserved causal alleles into  
1013 a single probabilistic model. By combining classical population genetics, large-scale  
1014 variant data, and Bayesian inference, we generate genome-wide priors across inheritance  
1015 modes and provide credible intervals that reflect both evidence and residual uncertainty.  
1016 Although demonstrated here in inborn errors of immunity, the framework  
1017 is broadly applicable and establishes a quantitative foundation for variant interpretation,  
1018 clinical decision-making, and future genomic analyses.

## 1019 Acknowledgements

1020 We would like to thank all the patients and families who have been providing advice  
1021 on SwissPedHealth and its projects, as well as the clinical and research teams at the  
1022 participating institutions. We acknowledge Genomics England for providing public  
1023 access to the PanelApp data. The use of data from Genomics England panelapp was  
1024 licensed under the Apache License 2.0. The use of data from UniProt was licensed  
1025 under Creative Commons Attribution 4.0 International (CC BY 4.0). ClinVar asks  
1026 its users who distribute or copy data to provide attribution to them as a data source  
1027 in publications and websites ([10](#)). dbNSFP version 4.4a is licensed under the Creative  
1028 Commons Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-  
1029 ND 4.0); while we cite this dataset as used our research publication, it is not used  
1030 for the final version which instead used ClinVar and gnomAD directly. GnomAD  
1031 is licensed under Creative Commons Zero Public Domain Dedication (CC0 1.0 Uni-  
1032 versal). GnomAD request that usages cites the gnomAD flagship paper ([4](#)) and any  
1033 online resources that include the data set provide a link to the browser, and note that  
1034 tool includes data from the gnomAD v4.1 release. AlphaMissense asks to cite Cheng  
1035 et al. ([9](#)) for usage in research, with data available from Cheng et al. ([18](#)).

## 1036 Contributions

1037 DL performed main analyses and wrote the manuscript. SB, AS, MS, and JT de-  
1038 signed analysis and wrote the manuscript. JF, LJS supervised the work, and applied  
1039 for funding. The Quant Group is a collaboration across multiple institutions where  
1040 authors contribute equally; the members on this project were DL, SB, AS, and MS.

<sup>1041</sup> **Competing interest**

<sup>1042</sup> The authors declare no competing interest.

<sup>1043</sup> **Ethics statement**

<sup>1044</sup> This study only used data which was previously published and publicly available,  
<sup>1045</sup> as cited in the manuscript. This SwissPedHealth study, under which this work was  
<sup>1046</sup> carried out, was approved based on the advice of the ethical committee Northwest  
<sup>1047</sup> and Central Switzerland (EKNZ, AO\_2022-00018). The study was conducted in  
<sup>1048</sup> accordance with the Declaration of Helsinki.

<sup>1049</sup> **Funding**

<sup>1050</sup> This project was supported through the grant Swiss National Science Foundation  
<sup>1051</sup> (SNF) 320030\_201060, and NDS-2021-911 (SwissPedHealth) from the Swiss Personalized  
<sup>1052</sup> Health Network and the Strategic Focal Area ‘Personalized Health and Related Technologies’  
<sup>1053</sup> of the ETH Domain (Swiss Federal Institutes of Technology).

<sup>1054</sup> **References**

- <sup>1055</sup> [1] Oliver Mayo. A Century of Hardy–Weinberg Equilibrium. *Twin Research and Human Genetics*, 11(3):249–256, June 2008. ISSN 1832-4274, 1839-2628. doi: 10.1375/twin.11.3.249. URL [https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal\\_article](https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal_article).
- <sup>1059</sup> [2] Nikita Abramovs, Andrew Brass, and May Tassabehji. Hardy–Weinberg Equilibrium in the Large Scale Genomic Sequencing Era. *Frontiers in Genetics*, 11:210, March 2020. ISSN 1664-8021. doi: 10.3389/fgene.2020.00210. URL <https://www.frontiersin.org/article/10.3389/fgene.2020.00210/full>.
- <sup>1063</sup> [3] Johannes Zschocke, Peter H. Byers, and Andrew O. M. Wilkie. Mendelian inheritance revisited: dominance and recessiveness in medical genetics. *Nature Reviews Genetics*, 24(7):442–463, July 2023. ISSN 1471-0056, 1471-0064. doi: 10.1038/s41576-023-00574-0. URL <https://www.nature.com/articles/s41576-023-00574-0>.
- <sup>1068</sup> [4] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings, Jessica Alföldi, Qingbo Wang, Ryan L Collins, Kristen M Laricchia, Andrea Ganna, Daniel P Birnbaum, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*, 581(7809):434–443, 2020.

- 1072 [5] Sarah L. Bick, Aparna Nathan, Hannah Park, Robert C. Green, Monica H. Wo-  
1073 jcik, and Nina B. Gold. Estimating the sensitivity of genomic newborn screen-  
1074 ing for treatable inherited metabolic disorders. *Genetics in Medicine*, 27(1):  
1075 101284, January 2025. ISSN 10983600. doi: 10.1016/j.gim.2024.101284. URL  
1076 <https://linkinghub.elsevier.com/retrieve/pii/S1098360024002181>.
- 1077 [6] Benjamin D. Evans, Piotr Słowiński, Andrew T. Hattersley, Samuel E. Jones,  
1078 Seth Sharp, Robert A. Kimmitt, Michael N. Weedon, Richard A. Oram,  
1079 Krasimira Tsaneva-Atanasova, and Nicholas J. Thomas. Estimating disease  
1080 prevalence in large datasets using genetic risk scores. *Nature Communications*,  
1081 12(1):6441, November 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-26501-7.  
1082 URL <https://www.nature.com/articles/s41467-021-26501-7>.
- 1083 [7] William B. Hannah, Mitchell L. Drumm, Keith Nykamp, Tiziano Prampano,  
1084 Robert D. Steiner, and Steven J. Schrodi. Using genomic databases to de-  
1085 termine the frequency and population-based heterogeneity of autosomal reces-  
1086 sive conditions. *Genetics in Medicine Open*, 2:101881, 2024. ISSN 29497744.  
1087 doi: 10.1016/j.gimo.2024.101881. URL <https://linkinghub.elsevier.com/retrieve/pii/S2949774424010276>.
- 1088 [8] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov,  
1089 Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek,  
1090 Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J.  
1091 Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh  
1092 Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy,  
1093 Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer,  
1094 Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray  
1095 Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate pro-  
1096 tein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August  
1097 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL  
1098 <https://www.nature.com/articles/s41586-021-03819-2>.
- 1099 [9] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-  
1100 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias  
1101 Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-  
1102 abis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense vari-  
1103 ant effect prediction with AlphaMissense. *Science*, 381(6664):eadg7492, Septem-  
1104 ber 2023. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.adg7492. URL  
1105 <https://www.science.org/doi/10.1126/science.adg7492>.
- 1106 [10] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen Chao,  
1107 Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Wonhee  
1108 Jang, Karen Karapetyan, Kenneth Katz, Chunlei Liu, Zenith Maddipatla, Adri-  
1109 ana Malheiro, Kurt McDaniel, Michael Ovetsky, George Riley, George Zhou,  
1110 J Bradley Holmes, Brandi L Kattman, and Donna R Maglott. ClinVar: im-  
1111 proving access to variant interpretations and supporting evidence. *Nucleic Acids*

- 1113      *Research*, 46(D1):D1062–D1067, January 2018. ISSN 0305-1048, 1362-4962. doi:  
1114      10.1093/nar/gkx1153. URL <http://academic.oup.com/nar/article/46/D1/D1062/4641904>.
- 1115
- 1116 [11] The UniProt Consortium, Alex Bateman, Maria-Jesus Martin, Sandra Orchard,  
1117 Michele Magrane, Aduragbemi Adesina, Shadab Ahmad, Bowler-Barnett, and  
1118 Others. UniProt: the Universal Protein Knowledgebase in 2025. *Nucleic  
1119 Acids Research*, 53(D1):D609–D617, January 2025. ISSN 0305-1048, 1362-4962.  
1120 doi: 10.1093/nar/gkae1010. URL <https://academic.oup.com/nar/article/53/D1/D609/7902999>.
- 1121
- 1122 [12] Dylan Lawless. PanelAppRex aggregates disease gene panels and facilitates  
1123 sophisticated search. March 2025. doi: 10.1101/2025.03.20.25324319. URL  
1124 <http://medrxiv.org/lookup/doi/10.1101/2025.03.20.25324319>.
- 1125
- 1126 [13] Antonio Rueda Martin, Eleanor Williams, Rebecca E. Foulger, Sarah Leigh,  
1127 Louise C. Daugherty, Olivia Niblock, Ivone U. S. Leong, Katherine R. Smith,  
1128 Oleg Gerasimenko, Eik Haraldsdottir, Ellen Thomas, Richard H. Scott, Emma  
1129 Baple, Arianna Tucci, Helen Brittain, Anna De Burca, Kristina Ibañez, Dalia  
1130 Kasperaviciute, Damian Smedley, Mark Caulfield, Augusto Rendon, and Ellen M.  
1131 McDonagh. PanelApp crowdsources expert knowledge to establish consensus  
1132 diagnostic gene panels. *Nature Genetics*, 51(11):1560–1565, November 2019.  
1133 ISSN 1061-4036, 1546-1718. doi: 10.1038/s41588-019-0528-2. URL <https://www.nature.com/articles/s41588-019-0528-2>.
- 1134
- 1135 [14] M. Cecilia Poli, Ivona Aksentijevich, Ahmed Aziz Bousfiha, Charlotte  
1136 Cunningham-Rundles, Sophie Hambleton, Christoph Klein, Tomohiro Morio,  
1137 Capucine Picard, Anne Puel, Nima Rezaei, Mikko R.J. Seppänen, Raz  
1138 Somech, Helen C. Su, Kathleen E. Sullivan, Troy R. Torgerson, Is-  
1139 abelle Meyts, and Stuart G. Tangye. Human inborn errors of immu-  
1140 nity: 2024 update on the classification from the International Union of  
1141 Immunological Societies Expert Committee. *Journal of Human Immu-  
1142 nity*, 1(1):e20250003, May 2025. ISSN 3065-8993. doi: 10.70962/jhi.  
1143 20250003. URL <https://rupress.org/jhi/article/1/1/e20250003/277390/Human-inborn-errors-of-immunity-2024-update-on-the>.
- 1144
- 1145 [15] Ahmed Aziz Bousfiha, Leïla Jeddane, Abderrahmane Moundir, M. Cecilia  
1146 Poli, Ivona Aksentijevich, Charlotte Cunningham-Rundles, Sophie Hambleton,  
1147 Christoph Klein, Tomohiro Morio, Capucine Picard, Anne Puel, Nima Rezaei,  
1148 Mikko R.J. Seppänen, Raz Somech, Helen C. Su, Kathleen E. Sullivan, Troy R.  
1149 Torgerson, Stuart G. Tangye, and Isabelle Meyts. The 2024 update of IUIS  
1150 phenotypic classification of human inborn errors of immunity. *Journal of Hu-  
1151 man Immunity*, 1(1):e20250002, May 2025. ISSN 3065-8993. doi: 10.70962/jhi.  
1152 20250002. URL <https://rupress.org/jhi/article/1/1/e20250002/277374/The-2024-update-of-IUIS-phenotypic-classification>.
- 1153

- 1153 [16] Xiaoming Liu, Chang Li, Chengcheng Mou, Yibo Dong, and Yicheng Tu.  
1154 dbNSFP v4: a comprehensive database of transcript-specific functional pre-  
1155 dictions and annotations for human nonsynonymous and splice-site SNVs.  
1156 *Genome Medicine*, 12(1):103, December 2020. ISSN 1756-994X. doi: 10.  
1157 1186/s13073-020-00803-9. URL <https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-020-00803-9>.  
1158
- 1159 [17] Damian Szklarczyk, Katerina Nastou, Mikaela Koutrouli, Rebecca Kirsch, Far-  
1160 rokh Mehryary, Radja Hachilif, Dewei Hu, Matteo E Peluso, Qingyao Huang,  
1161 Tao Fang, et al. The string database in 2025: protein networks with directional-  
1162 ity of regulation. *Nucleic Acids Research*, 53(D1):D730–D737, 2025.
- 1163 [18] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Tay-  
1164 lor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias  
1165 Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hass-  
1166 abis, Pushmeet Kohli, and Žiga Avsec. Predictions for alphanonsense, September  
1167 2023. URL <https://doi.org/10.5281/zenodo.8208688>.
- 1168 [19] Paul Tijnenburg, Hana Lango Allen, Siobhan O. Burns, Daniel Greene,  
1169 Machiel H. Jansen, and Others. Loss-of-function nuclear factor B subunit  
1170 1 (NFKB1) variants are the most common monogenic cause of common vari-  
1171 able immunodeficiency in Europeans. *Journal of Allergy and Clinical Im-*  
1172 *munology*, 142(4):1285–1296, October 2018. ISSN 00916749. doi: 10.1016/  
1173 j.jaci.2018.01.039. URL <https://linkinghub.elsevier.com/retrieve/pii/S0091674918302860>.  
1174
- 1175 [20] WHO Scientific Group et al. Primary immunodeficiency diseases: report of a  
1176 who scientific group. *Clin. Exp. Immunol.*, 109(1):1–28, 1997.
- 1177 [21] Charlotte Cunningham-Rundles and Carol Bodian. Common variable immunod-  
1178 eficiency: clinical and immunological features of 248 patients. *Clinical immunol-*  
1179 *ogy*, 92(1):34–48, 1999.
- 1180 [22] Eric Oksenhendler, Laurence Gérard, Claire Fieschi, Marion Malphettes, Gael  
1181 Mouillot, Roland Jaussaud, Jean-François Viallard, Martine Gardembas, Lionel  
1182 Galicier, Nicolas Schleinitz, et al. Infections in 252 patients with common variable  
1183 immunodeficiency. *Clinical Infectious Diseases*, 46(10):1547–1554, 2008.
- 1184 [23] Y Naito, F Adams, S Charman, J Duckers, G Davies, and S Clarke. Uk cystic  
1185 fibrosis registry 2023 annual data report. *London: Cystic Fibrosis Trust*, 2023.
- 1186 [24] Carlo Castellani, CFTR2 team, et al. Cftr2: how will it help care? *Paediatric*  
1187 *respiratory reviews*, 14:2–5, 2013.
- 1188 [25] Hartmut Grasemann and Felix Ratjen. Cystic fibrosis. *New England Journal*  
1189 *of Medicine*, 389(18):1693–1707, 2023. doi: 10.1056/NEJMra2216474. URL  
1190 <https://www.nejm.org/doi/full/10.1056/NEJMra2216474>.

- 1191 [26] Kyoko Watanabe, Erdogan Taskesen, Arjen Van Bochoven, and Danielle  
1192 Posthuma. Functional mapping and annotation of genetic associations with  
1193 FUMA. *Nature Communications*, 8(1):1826, November 2017. ISSN 2041-1723.  
1194 doi: 10.1038/s41467-017-01261-5. URL <https://www.nature.com/articles/s41467-017-01261-5>.
- 1195
- 1196 [27] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir,  
1197 Pablo Tamayo, and Jill P. Mesirov. Molecular signatures database (MSigDB)  
1198 3.0. *Bioinformatics*, 27(12):1739–1740, June 2011. ISSN 1367-4811, 1367-  
1199 4803. doi: 10.1093/bioinformatics/btr260. URL <https://academic.oup.com/bioinformatics/article/27/12/1739/257711>.
- 1200
- 1201 [28] Dylan Lawless. Variant risk estimate probabilities for iei genes. March 2025. doi:  
1202 10.5281/zenodo.15111584. URL <https://doi.org/10.5281/zenodo.15111584>.
- 1203
- 1204 [29] Eric Vallabh Minikel, Sonia M. Vallabh, Monkol Lek, Karol Estrada, Kaitlin E.  
1205 Samocha, J. Fah Sathirapongsasuti, Cory Y. McLean, Joyce Y. Tung, Linda  
1206 P. C. Yu, Pierluigi Gambetti, Janis Blevins, Shulin Zhang, Yvonne Cohen,  
1207 Wei Chen, Masahito Yamada, Tsuyoshi Hamaguchi, Nobuo Sanjo, Hidehiro  
1208 Mizusawa, Yosikazu Nakamura, Tetsuyuki Kitamoto, Steven J. Collins, Alison  
1209 Boyd, Robert G. Will, Richard Knight, Claudia Ponto, Inga Zerr, Theo  
1210 F. J. Kraus, Sabina Eigenbrod, Armin Giese, Miguel Calero, Jesús De Pedro-  
1211 Cuesta, Stéphane Haïk, Jean-Louis Laplanche, Elodie Bouaziz-Amar, Jean-  
1212 Philippe Brandel, Sabina Capellari, Piero Parchi, Anna Poleggi, Anna Ladogana,  
1213 Anne H. O’Donnell-Luria, Konrad J. Karczewski, Jamie L. Marshall, Michael  
1214 Boehnke, Markku Laakso, Karen L. Mohlke, Anna Kähler, Kimberly Chambert,  
1215 Steven McCarroll, Patrick F. Sullivan, Christina M. Hultman, Shaun M. Purcell,  
1216 Pamela Sklar, Sven J. Van Der Lee, Annemieke Rozemuller, Casper Jansen, Albert  
1217 Hofman, Robert Kraaij, Jeroen G. J. Van Rooij, M. Arfan Ikram, André G.  
1218 Uitterlinden, Cornelia M. Van Duijn, Exome Aggregation Consortium (ExAC),  
1219 Mark J. Daly, and Daniel G. MacArthur. Quantifying prion disease penetrance  
1220 using large population control cohorts. *Science Translational Medicine*, 8(322),  
1221 January 2016. ISSN 1946-6234, 1946-6242. doi: 10.1126/scitranslmed.aad5169.  
1222 URL <https://www.science.org/doi/10.1126/scitranslmed.aad5169>.
- 1223
- 1224 [30] Nicola Whiffin, Eric Minikel, Roddy Walsh, Anne H O’Donnell-Luria, Konrad  
1225 Karczewski, Alexander Y Ing, Paul J R Barton, Birgit Funke, Stuart A Cook,  
1226 Daniel MacArthur, and James S Ware. Using high-resolution variant frequencies  
1227 to empower clinical genome interpretation. *Genetics in Medicine*, 19(10):1151–  
1158, October 2017. ISSN 10983600. doi: 10.1038/gim.2017.26. URL <https://linkinghub.elsevier.com/retrieve/pii/S1098360021013678>.
- 1228
- 1229 [31] Bradley Efron and Carl Morris. Stein’s Estimation Rule and Its Competitors—  
1230 An Empirical Bayes Approach. *Journal of the American Statistical Association*,  
1231 68(341):117, March 1973. ISSN 01621459. doi: 10.2307/2284155. URL <https://www.jstor.org/stable/2284155?origin=crossref>.
- 1232

- 1232 [32] Yaowu Liu, Sixing Chen, Zilin Li, Alanna C Morrison, Eric Boerwinkle, and  
1233 Xihong Lin. Acat: a fast and powerful p value combination method for rare-  
1234 variant analysis in sequencing studies. *The American Journal of Human Genetics*,  
1235 104(3):410–421, 2019.
- 1236 [33] Xiacao Li, Zilin Li, Hufeng Zhou, Sheila M Gaynor, Yaowu Liu, Han Chen, Ryan  
1237 Sun, Rounak Dey, Donna K Arnett, Stella Aslibekyan, et al. Dynamic incorpora-  
1238 tion of multiple in silico functional annotations empowers rare variant association  
1239 analysis of large whole-genome sequencing studies at scale. *Nature genetics*, 52  
1240 (9):969–983, 2020.
- 1241 [34] Michael C Wu, Seunggeun Lee, Tianxi Cai, Yun Li, Michael Boehnke, and Xi-  
1242 hong Lin. Rare-variant association testing for sequencing data with the sequence  
1243 kernel association test. *The American Journal of Human Genetics*, 89(1):82–93,  
1244 2011.
- 1245 [35] Seunggeun Lee, Mary J Emond, Michael J Bamshad, Kathleen C Barnes, Mark J  
1246 Rieder, Deborah A Nickerson, David C Christiani, Mark M Wurfel, and Xihong  
1247 Lin. Optimal unified approach for rare-variant association testing with applica-  
1248 tion to small-sample case-control whole-exome sequencing studies. *The American  
1249 Journal of Human Genetics*, 91(2):224–237, 2012.
- 1250 [36] Augustine Kong, Gudmar Thorleifsson, Michael L Frigge, Bjarni J Vilhjalmsson,  
1251 Alexander I Young, Thorgeir E Thorgeirsson, Stefania Benonisdottir, Asmundur  
1252 Oddsson, Bjarni V Halldorsson, Gisli Masson, et al. The nature of nurture:  
1253 Effects of parental genotypes. *Science*, 359(6374):424–428, 2018.
- 1254 [37] Laurence J Howe, Michel G Nivard, Tim T Morris, Ailin F Hansen, Humaira  
1255 Rasheed, Yoonsoo Cho, Geetha Chittoor, Penelope A Lind, Teemu Palviainen,  
1256 Matthijs D van der Zee, et al. Within-sibship gwas improve estimates of direct  
1257 genetic effects. *BioRxiv*, pages 2021–03, 2021.
- 1258 [38] Sue Richards, Nazneen Aziz, Sherri Bale, David Bick, Soma Das, Julie Gastier-  
1259 Foster, Wayne W Grody, Madhuri Hegde, Elaine Lyon, Elaine Spector, et al.  
1260 Standards and guidelines for the interpretation of sequence variants: a joint  
1261 consensus recommendation of the american college of medical genetics and ge-  
1262 nomics and the association for molecular pathology. *Genetics in medicine*, 17  
1263 (5):405–423, 2015.
- 1264 [39] Sean V Tavtigian, Steven M Harrison, Kenneth M Boucher, and Leslie G  
1265 Biesecker. Fitting a naturally scaled point system to the acmng/amp variant  
1266 classification guidelines. *Human mutation*, 41(10):1734–1737, 2020.
- 1267 [40] Quan Li and Kai Wang. Intervar: clinical interpretation of genetic variants by  
1268 the 2015 acmng-amp guidelines. *The American Journal of Human Genetics*, 100  
1269 (2):267–280, 2017.

- 1270 [41] Brent S Pedersen, Joe M Brown, Harriet Dashnow, Amelia D Wallace, Matt  
1271 Velinder, Martin Tristani-Firouzi, Joshua D Schiffman, Tatiana Tvardik, Rong  
1272 Mao, D Hunter Best, et al. Effective variant filtering and expected candidate  
1273 variant yield in studies of rare human disease. *NPJ Genomic Medicine*, 6(1):1–8,  
1274 2021.
- 1275 [42] Carl A Anderson, Fredrik H Pettersson, Geraldine M Clarke, Lon R Cardon,  
1276 Andrew P Morris, and Krina T Zondervan. Data quality control in genetic  
1277 case-control association studies. *Nature protocols*, 5(9):1564–1573, 2010. URL  
1278 <https://doi.org/10.1038/nprot.2010.116>.
- 1279 [43] David T Miller, Kristy Lee, Noura S Abul-Husn, Laura M Amendola, Kyle Broth-  
1280 ers, Wendy K Chung, Michael H Gollob, Adam S Gordon, Steven M Harrison,  
1281 Ray E Hershberger, et al. Acmg sf v3. 2 list for reporting of secondary findings  
1282 in clinical exome and genome sequencing: a policy statement of the american  
1283 college of medical genetics and genomics (acmg). *Genetics in Medicine*, 25(8):  
1284 100866, 2023.

1285 **6 Supplemental**

1286 Supplemental data are presented under the same headings that correspond to their  
1287 relevant main text sections.

1288 **6.1 Variant class occurrence probability**

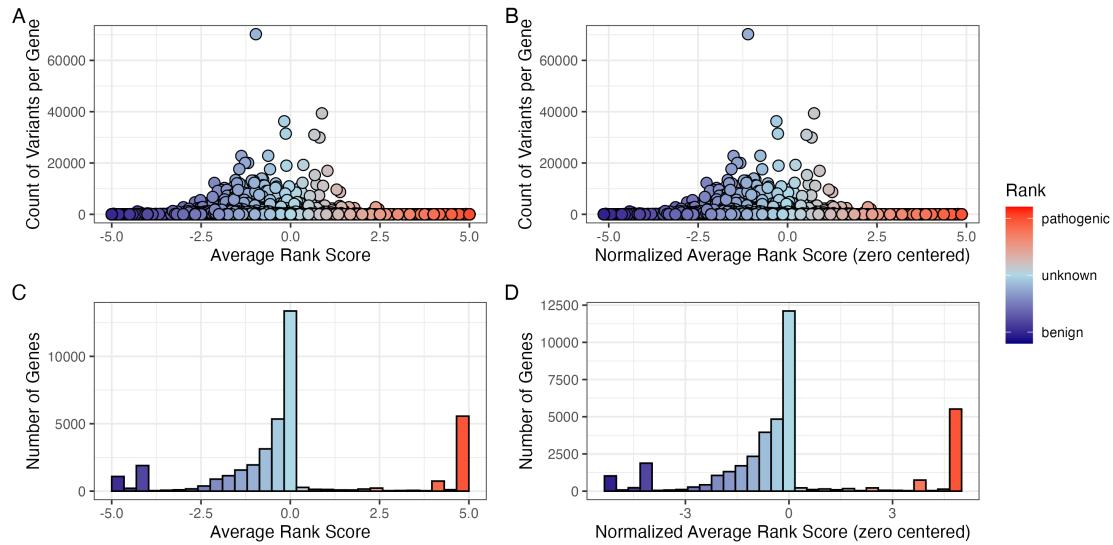


Figure S1: **Global distribution of ClinVar clinical-significance classification scoring.** (A) Number of variants per gene containing the assigned score for each ClinVar classification term ( $-5$  to  $+5$ ). (B) The same data after normalisation by zero centring the average rank score. (C) The tally 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.

1289    **6.2 Integrating observed true positives and unobserved false**  
 1290    **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	1.000	1.000	1.000	NA	1.000

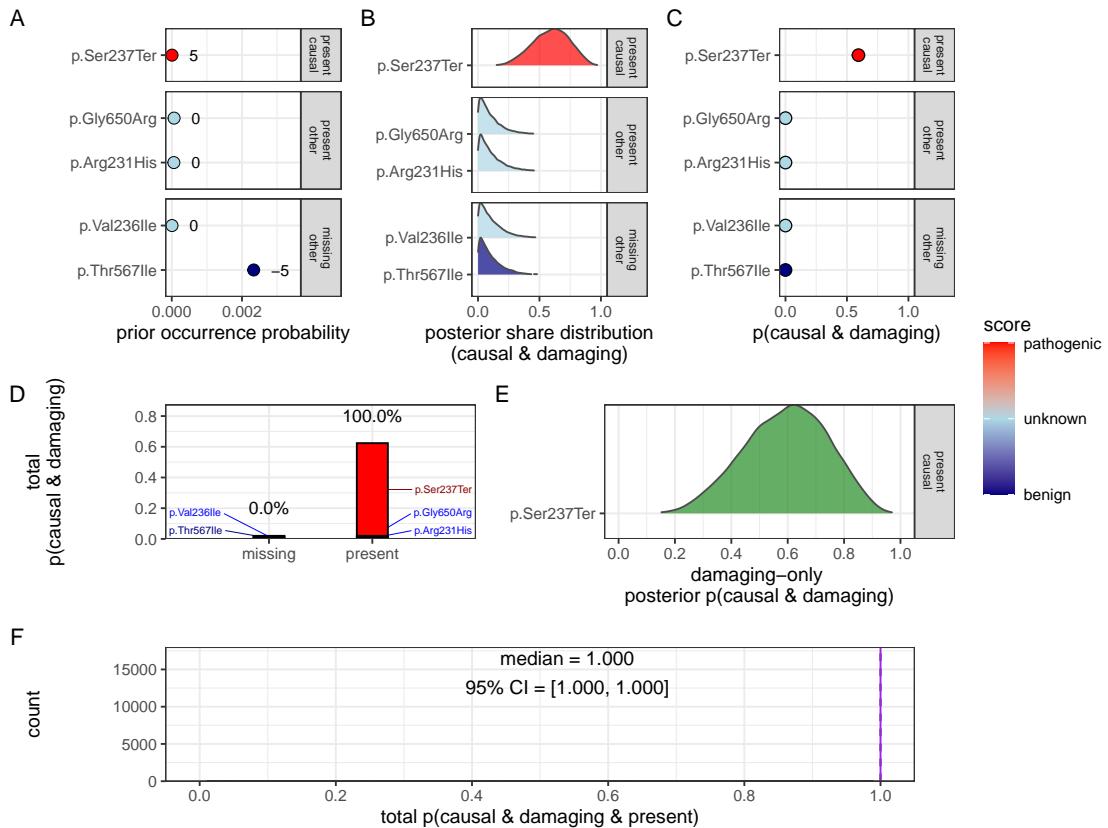
Table S2: Result of clinical genetics diagnosis scenario 2 including metadata. The most strongly supported observed variant was **p.Ser237Ter** (posterior: 0.381). The strongest unsequenced variant was **c.159+1G>A** (posterior: 0.353). The total probability of a causal diagnosis given the available 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	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	0	0	0	NA	0.000

Gene: *NFKB1*



**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

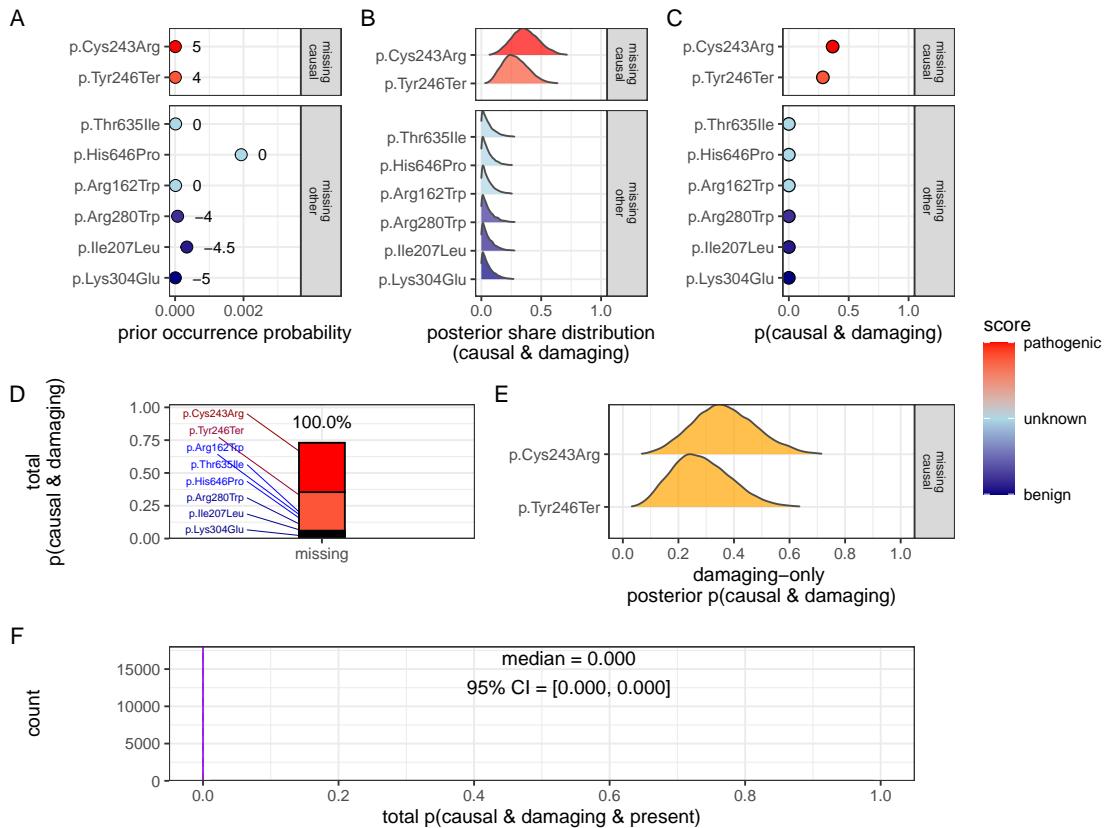
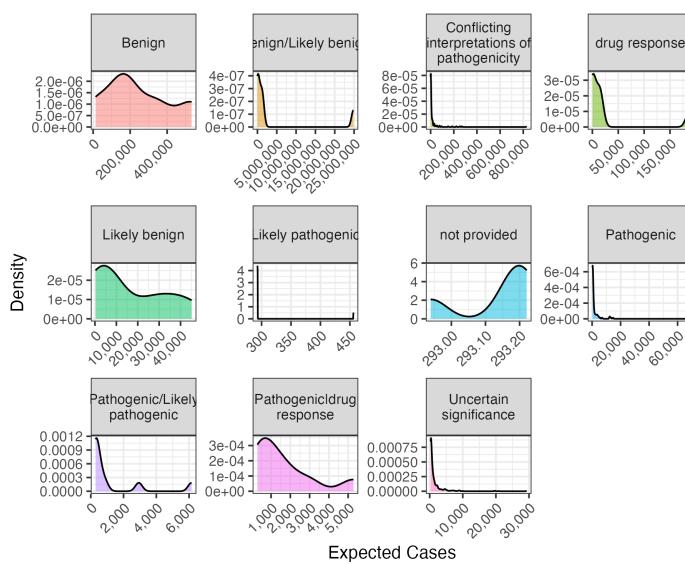


Figure S3: Quantification of no present (TP) in *NFKB1* and only missing (FN) causal genetic variants for disease in *TNFAIP3* (scenario 3). No known causal variants were observed in *NFKB1*, but one representative unsequenced allele was selected from each distinct ClinVar classification and treated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

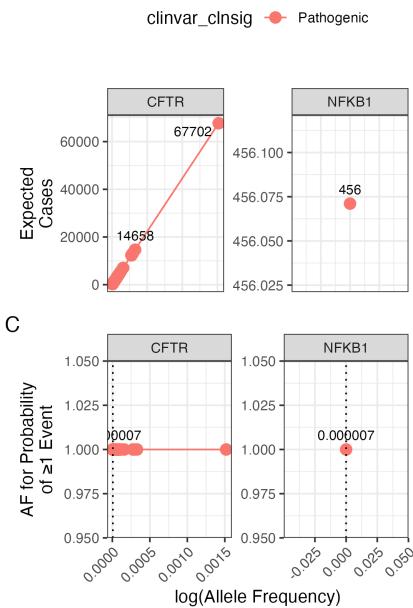
### 6.3 Validation studies

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

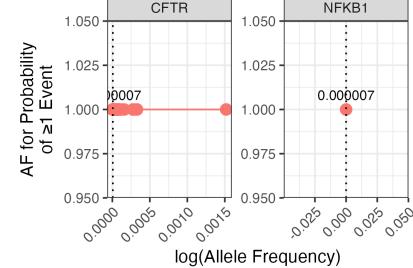
A



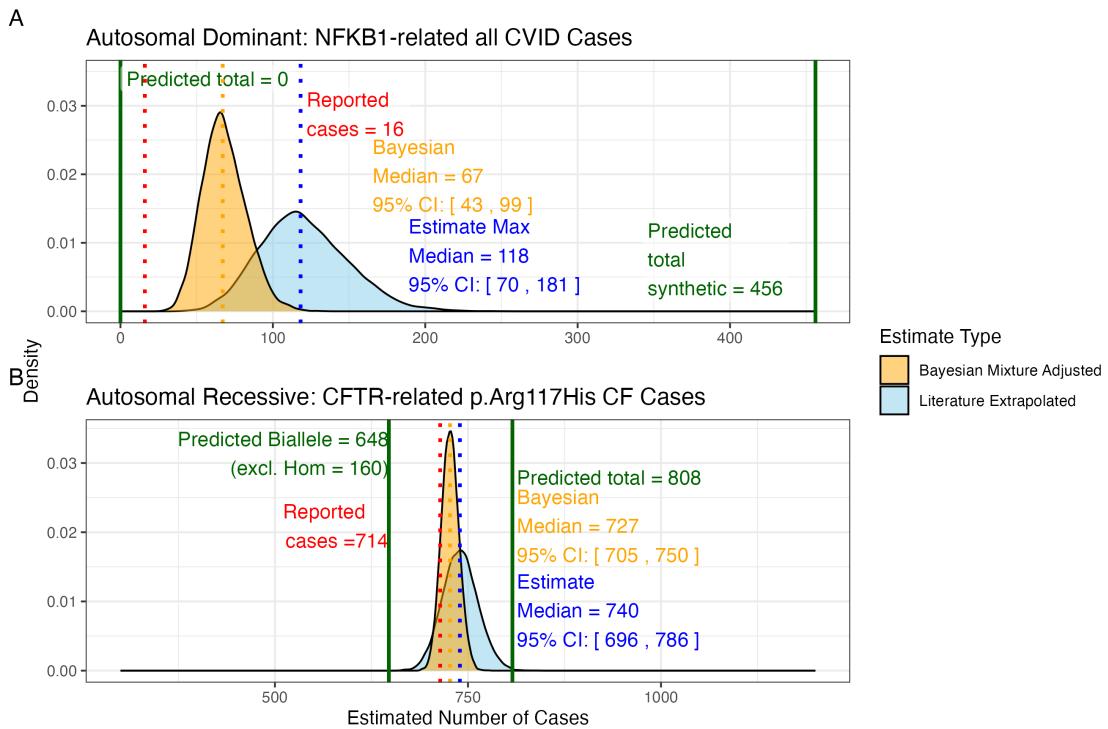
B



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.

1292 6.3.1 Interpretation of ClinVar variant observations

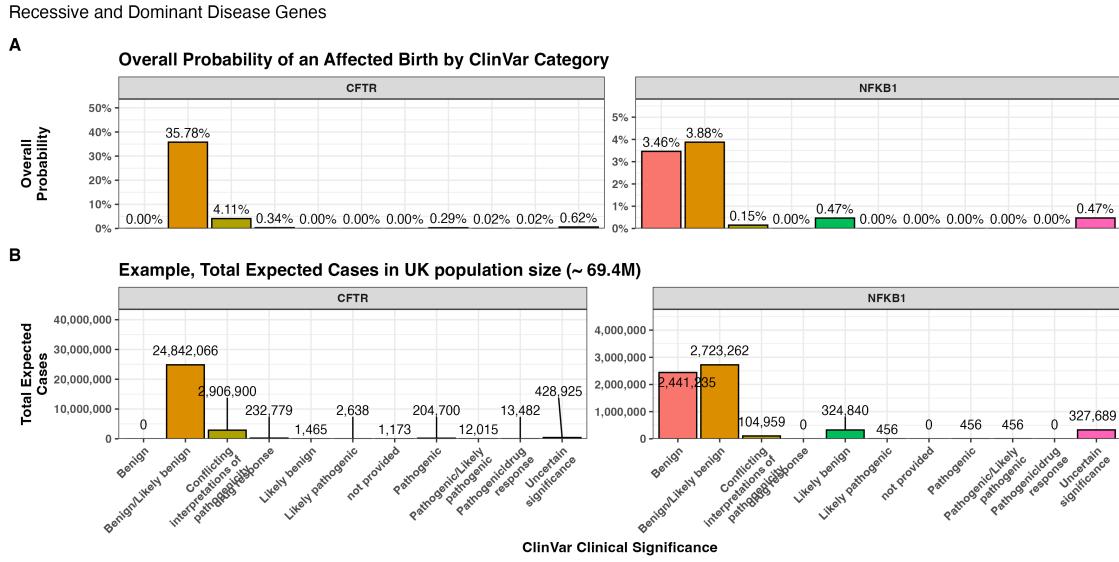
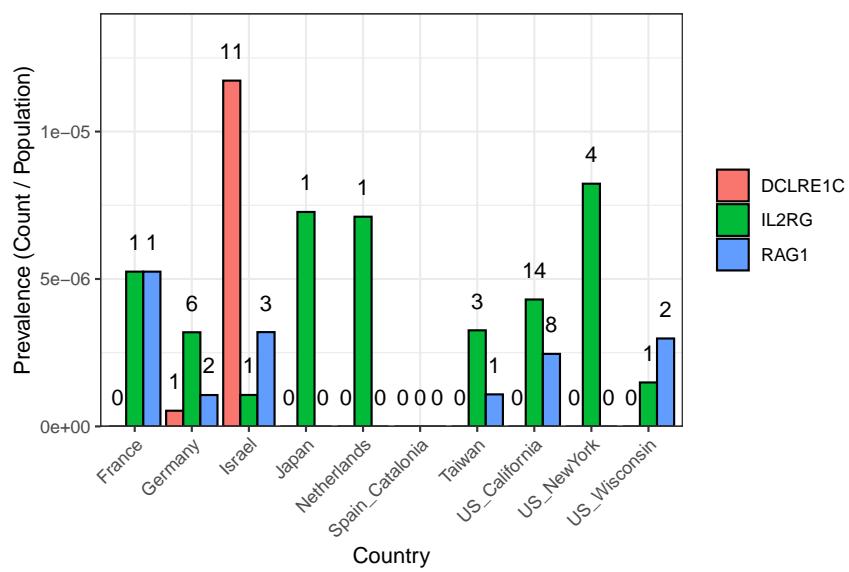
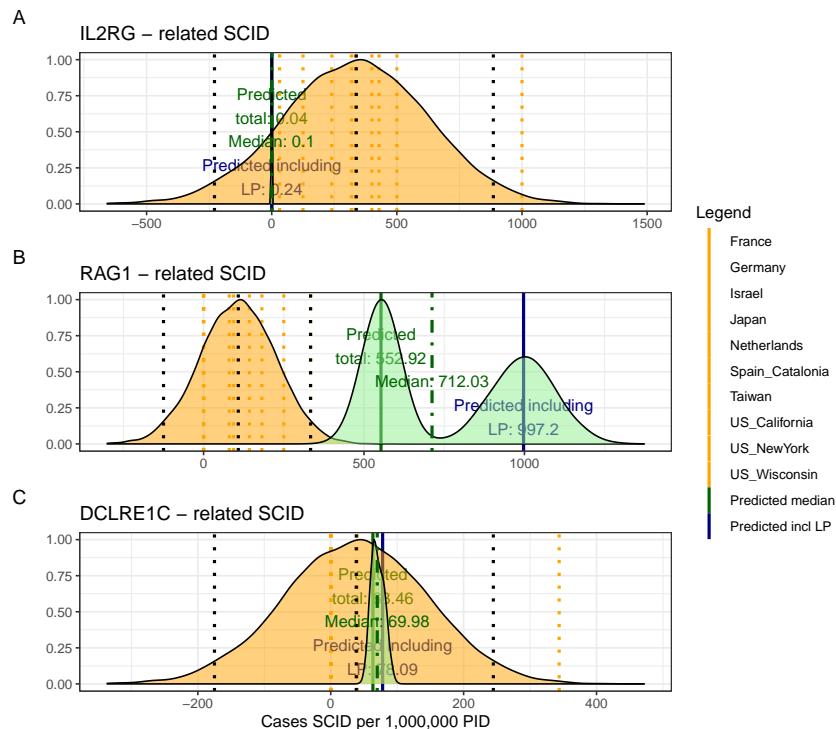


Figure S6: Combined bar charts summarising the genome-wide analysis of ClinVar clinical significance for the PID gene panel. Panel (A) shows the overall probability of an affected birth by variant classification, and (B) displays the total expected number of cases per classification, both stratified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

1293 6.3.2 Validation of SCID-specific disease occurrence



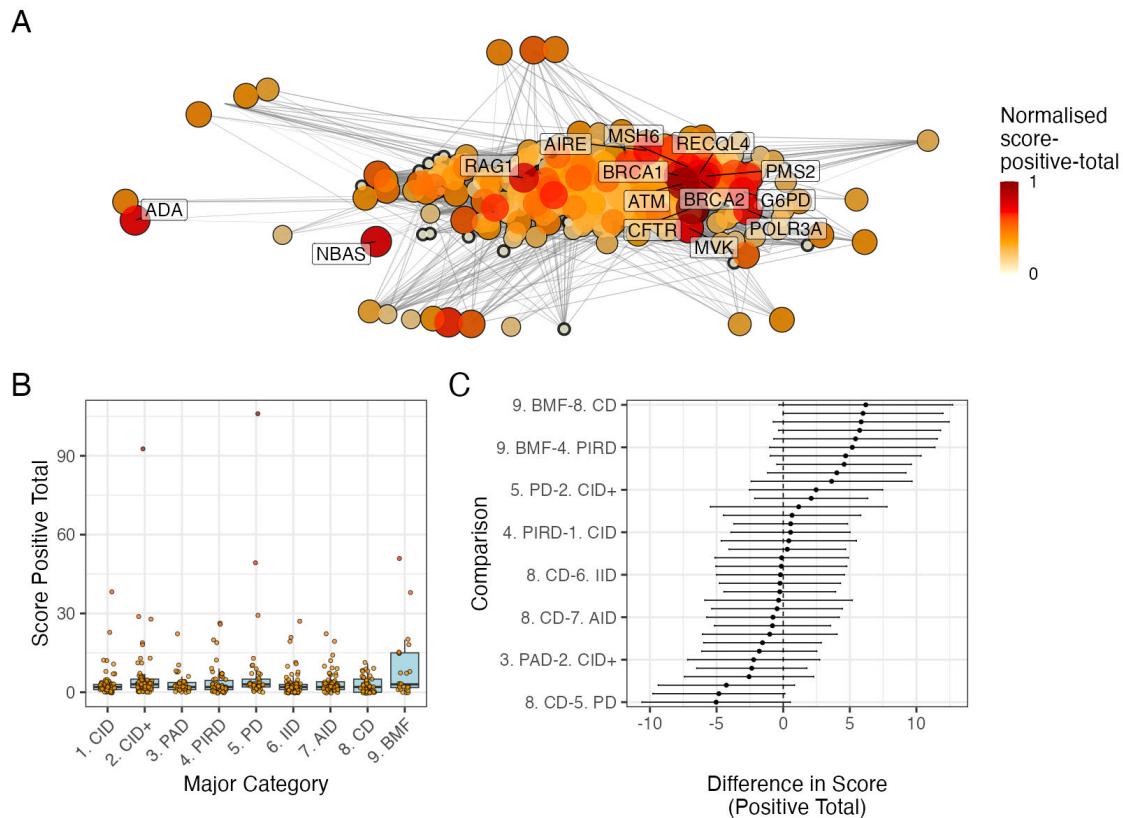
**Figure S7: SCID-specific gene comparison across regions.** The bar plot shows the prevalence of SCID-related cases (count divided by population) for each gene and country (or region), with numbers printed above the bars representing the actual counts in the original cohort (ranging from 0 to 11 per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

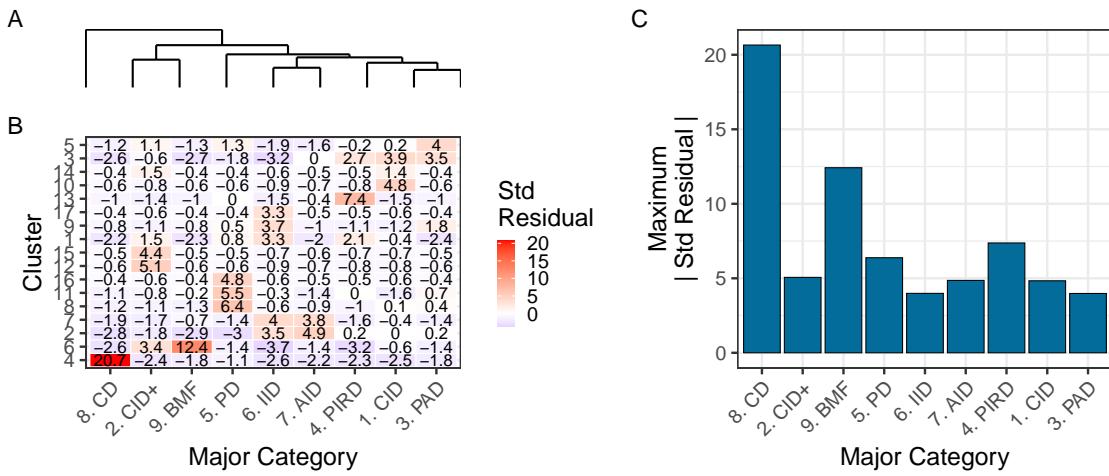
<sup>1294</sup> **6.4 Genetic constraint in high-impact protein networks**

<sup>1295</sup> **6.4.1 Score-positive-total within IEI PPI network**



**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**.

1298    6.4.3 PPI connectivity, LOEUF constraint and enriched network cluster  
1299    analysis

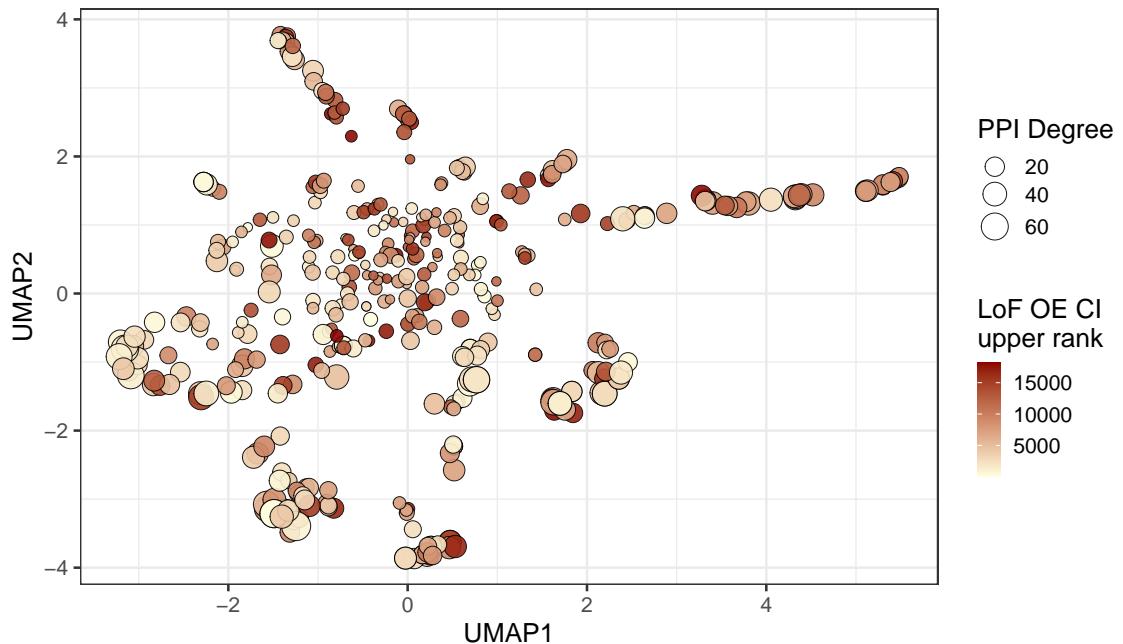
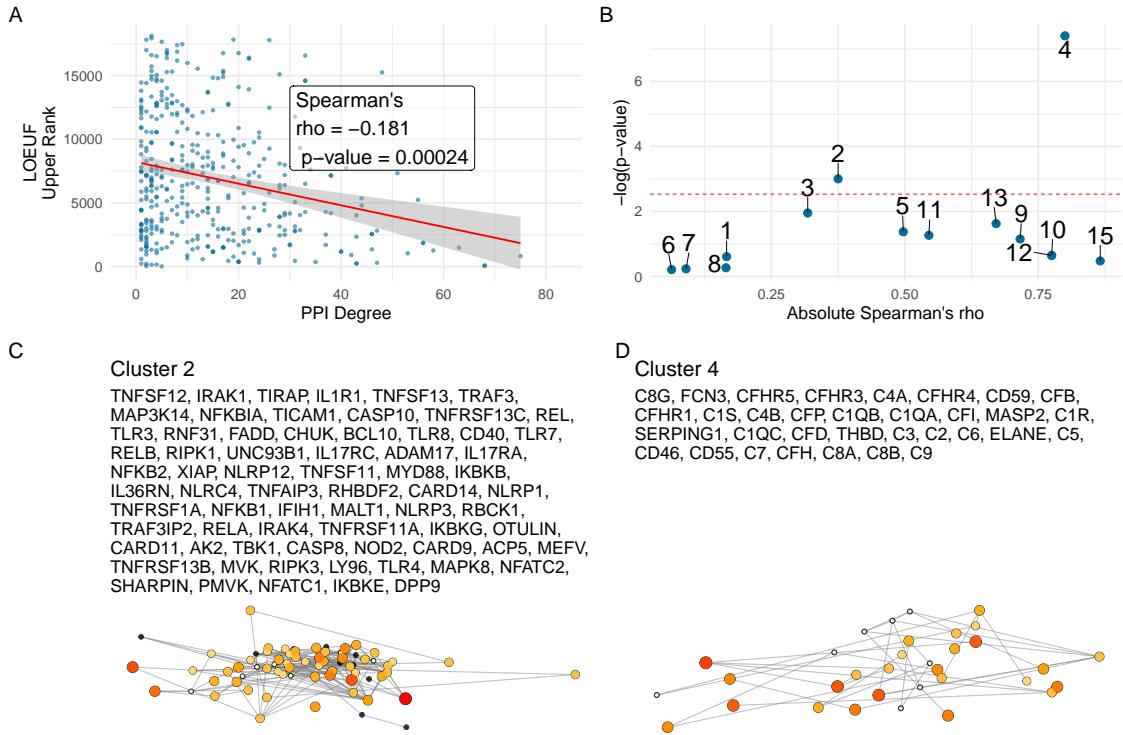
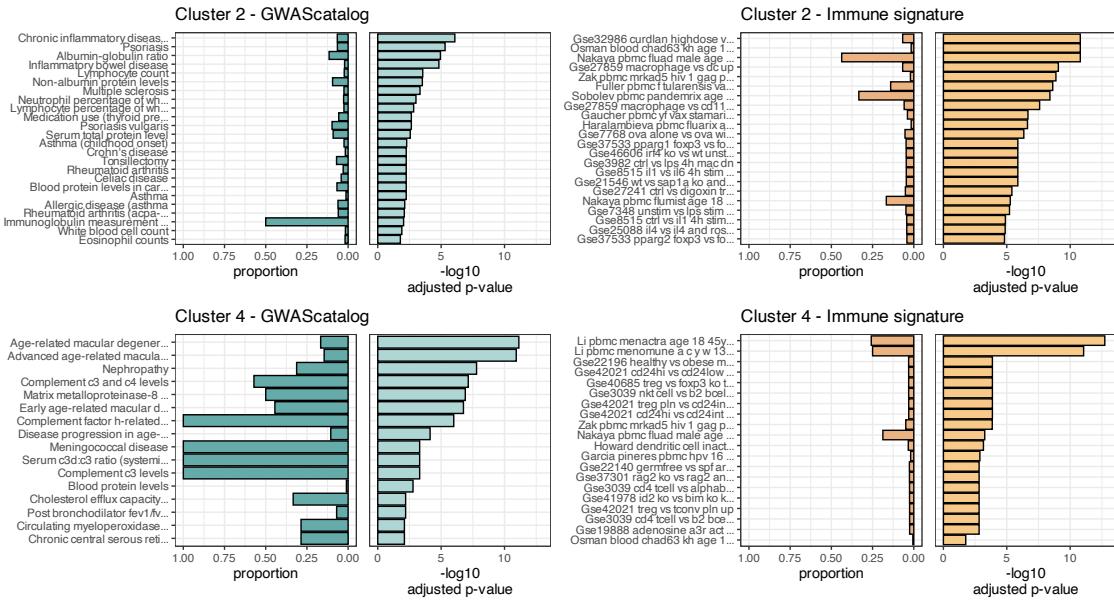


Figure S11: **Analysis of PPI degree versus LOEUF upper rank with UMAP embedding of the PPI network.** The relationship between PPI degree (size) and LOEUF upper rank (color) across gene clusters. No clear patterns are evident.



**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.

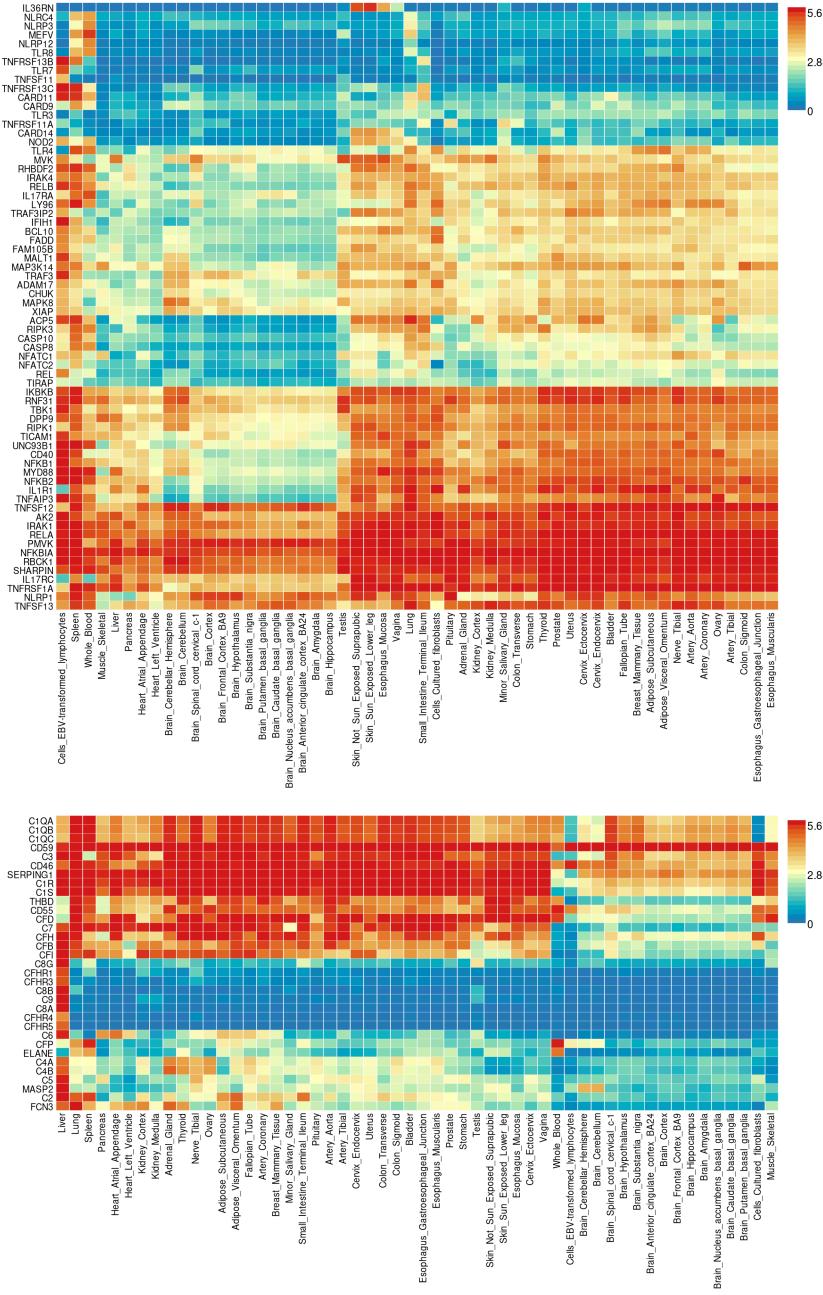
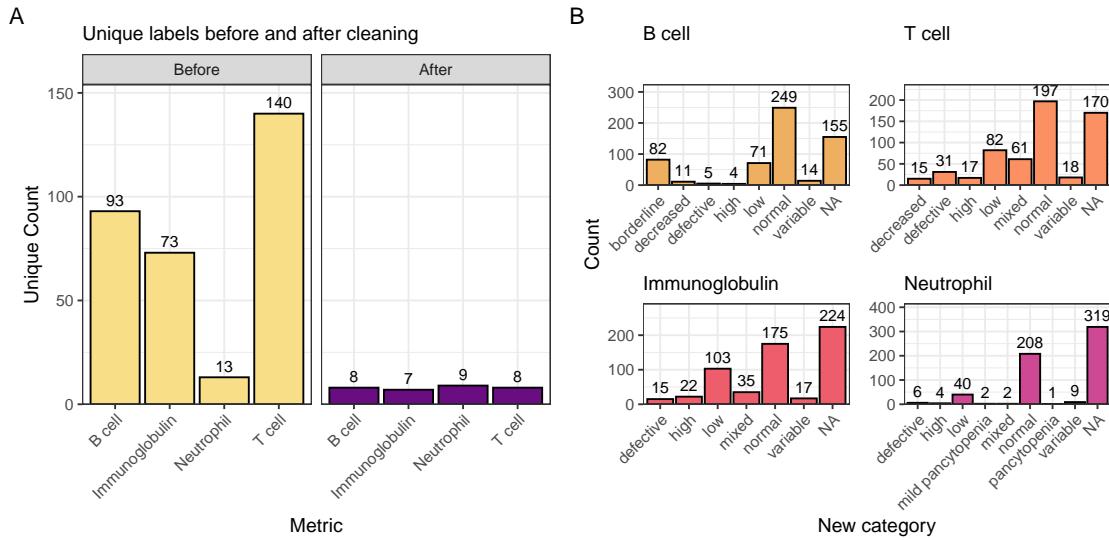
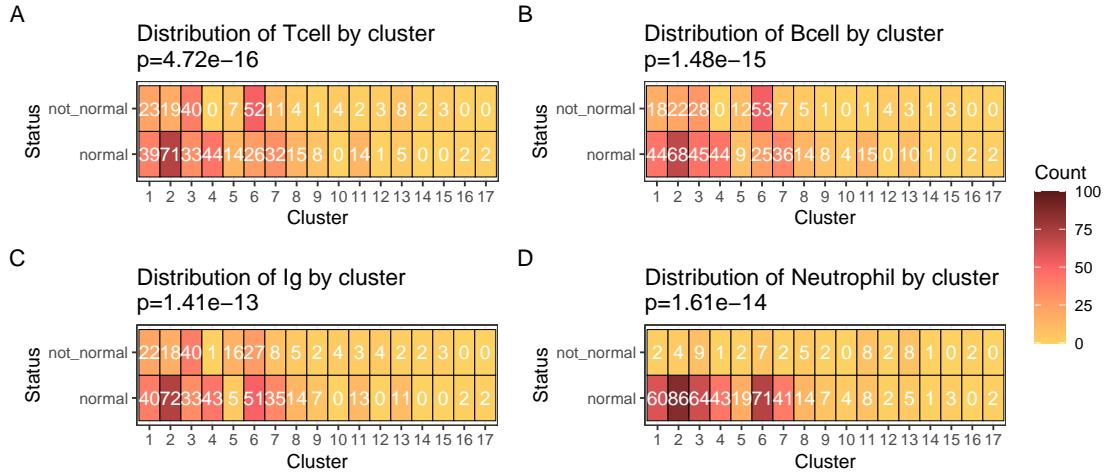


Figure S14: **Gene Expression Heatmaps for IEI Genes.** GTEx v8 data from 54 tissue types display the average expression per tissue label (log<sub>2</sub> transformed) for the IEI gene panels. Top: Cluster 2; Bottom: Cluster 4.

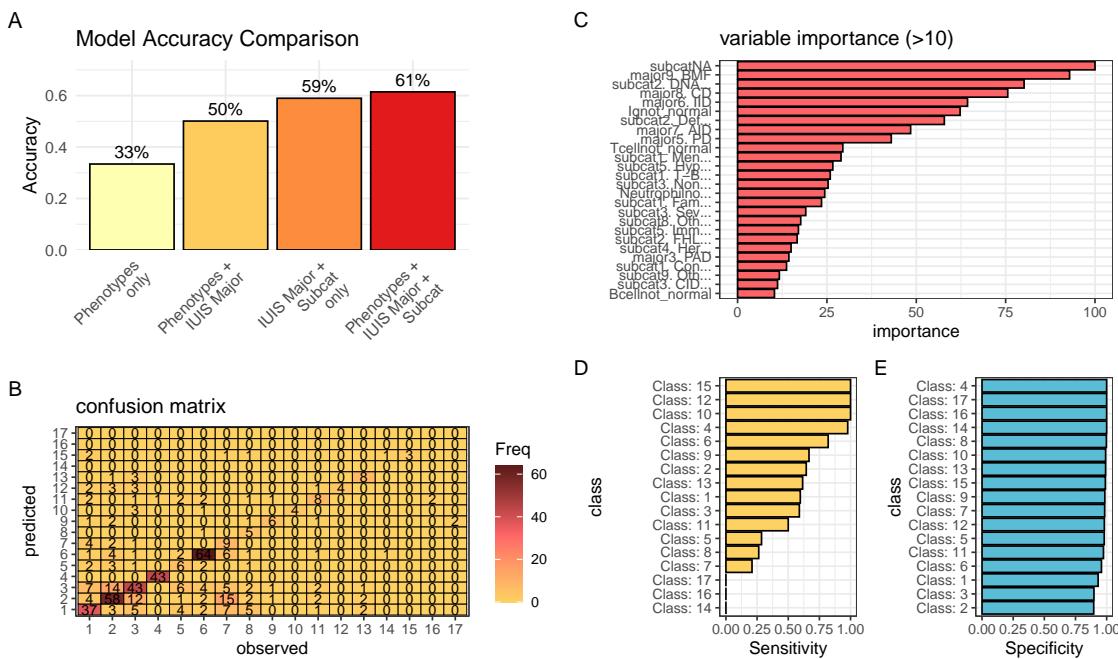
1300    **6.5 Novel PID classifications derived from genetic PPI and  
1301 clinical features**



**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.

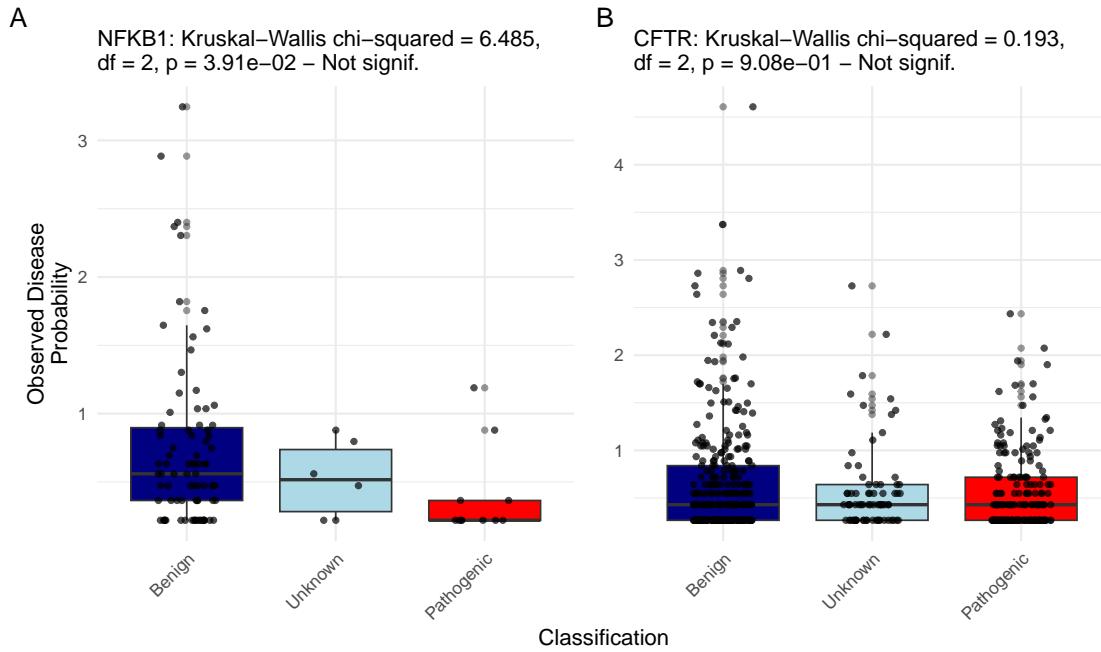


**Figure S16: Heatmaps of clinical feature distributions by PPI cluster.** The heatmaps display the count of observations for abnormality of each clinical feature (A) T cell, (B) B cell, (C) Immunoglobulin, (D) Neutrophil, in relation to the PPI clusters, with p-values from chi-square tests annotated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