

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

Dylan Lawless^{*1}

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

April 3, 2025

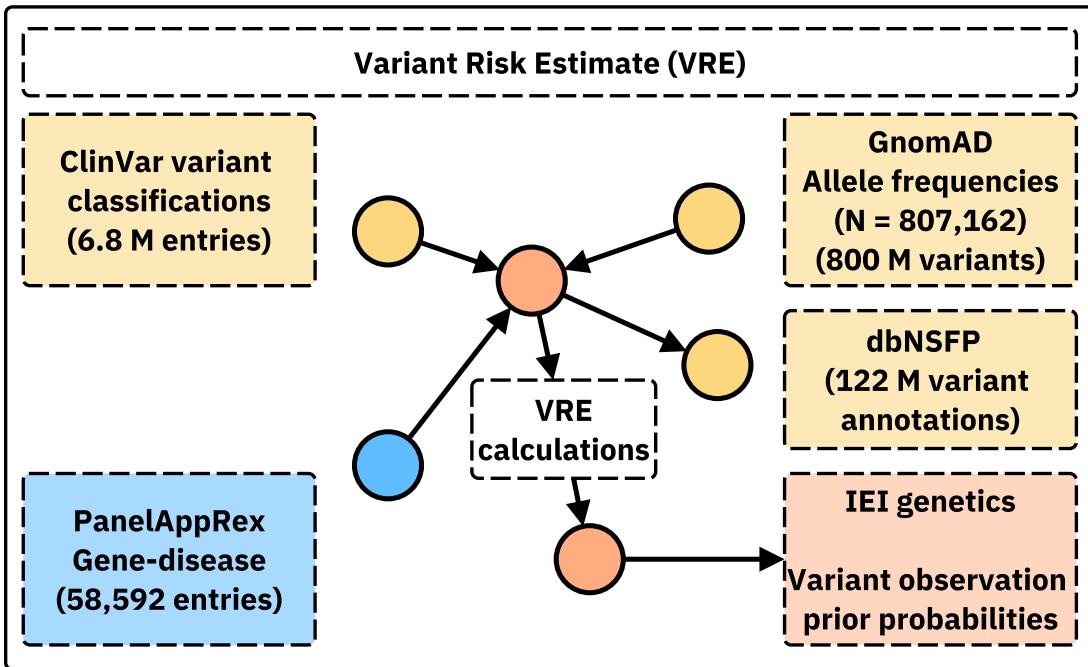
¹

Abstract

We present a novel framework for quantifying the prior probability of observing disease-associated variants in any gene for a given phenotype. By integrating large-scale genomic annotations, including population allele frequencies and ClinVar variant classifications, with Hardy-Weinberg-based calculations, our method estimates per-variant observation probabilities under autosomal dominant (AD), autosomal recessive (AR), and X-linked modes of inheritance. Applied to 557 genes implicated in primary immunodeficiency and inflammatory disease, our approach generated 54,814 variant probabilities. First, these detailed, pre-calculated results provide robust priors for any gene-disease combination. Second, a score positive total metric summarises the aggregate pathogenic burden, serving as an indicator of the likelihood of observing a patient with the disease and reflecting genetic constraint. Validation in *NFKB1* (AD) and *CFTR* (AR) disorders confirmed close concordance between predicted and observed case counts. The resulting datasets, available in both machine-readable and human-friendly formats, support Bayesian variant interpretation and clinical decision-making.¹

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

¹ **Availability:** This data is integrated in public panels at <https://iei-genetics.github.io>. The source code and data are accessible as part of the variant risk estimation project at https://github.com/DylanLawless/var_risk_est. The variant-level data is available from the Zenodo repository: <https://doi.org/10.5281/zenodo.15111583> (VarRiskEst PanelAppRex ID 398 gene variants.tsv). VarRiskEst is available under the MIT licence.



18

¹⁹ **Acronyms**

²⁰ ACMG American College of Medical Genetics and Genomics.....	²⁴
²¹ ACAT Aggregated Cauchy Association Test	²⁴
²² AD Autosomal Dominant.....	⁴
²³ ANOVA Analysis of Variance	¹¹
²⁴ AR Autosomal Recessive	⁴
²⁵ BMF Bone Marrow Failure.....	¹⁸
²⁶ CD Complement Deficiencies	¹⁹
²⁷ CI Confidence Interval.....	¹³
²⁸ CF Cystic Fibrosis	¹⁰
²⁹ CFTR Cystic Fibrosis Transmembrane Conductance Regulator.....	⁵
³⁰ CVID Common Variable Immunodeficiency	⁸
³¹ dbNSFP database for Non-Synonymous Functional Predictions	⁵
³² GE Genomics England	⁵
³³ gnomAD Genome Aggregation Database	⁵
³⁴ HGVS Human Genome Variation Society	⁵
³⁵ HPC High-Performance Computing.....	⁸
³⁶ HWE Hardy-Weinberg Equilibrium	⁴
³⁷ IEI Inborn Errors of Immunity.....	⁴
³⁸ InDel Insertion/Deletion	⁵
³⁹ IUIS International Union of Immunological Societies	⁶
⁴⁰ LD Linkage Disequilibrium	²¹
⁴¹ LOEUF Loss-Of-function Observed/Expected Upper bound Fraction	¹¹
⁴² LOF Loss-of-Function	¹⁸
⁴³ MOI Mode of Inheritance	⁴
⁴⁴ NFKB1 Nuclear Factor Kappa B Subunit 1	⁵
⁴⁵ OMIM Online Mendelian Inheritance in Man	²²
⁴⁶ PID Primary Immunodeficiency	⁴
⁴⁷ PPI Protein-Protein Interaction	⁵
⁴⁸ SNV Single Nucleotide Variant	⁴
⁴⁹ SKAT Sequence Kernel Association Test.....	²⁴
⁵⁰ STRINGdb Search Tool for the Retrieval of Interacting Genes/Proteins.....	⁵
⁵¹ HSD Honestly Significant Difference	¹¹
⁵² UMAP Uniform Manifold Approximation and Projection	¹⁸
⁵³ UniProt Universal Protein Resource	⁵
⁵⁴ VEP Variant Effect Predictor.....	⁵
⁵⁵ XL X-Linked	⁴

92 1 Introduction

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

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

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

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

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

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

To cite: <https://doi.org/10.1016/j.gimo.2024.101881> <https://doi.org/10.1016/j.gim.2024.101284> and some from Eric’s <https://www.cureffi.org/2019/06/05/using-genetic-data-to-estimate-disease-prevalence/>.

2 Methods

2.1 Dataset

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

¹⁷⁴ from the International Union of Immunological Societies (IUIS) IEI for the major
¹⁷⁵ disease categories used throughout this study (1).

Box 2.1 Definitions for IEI Major Disease Categories

Major Category	Description
1. CID	Immunodeficiencies affecting cellular and humoral immunity
2. CID+	Combined immunodeficiencies with associated or syndromic features
3. PAD	- Predominantly Antibody Deficiencies
4. PIRD	- Diseases of Immune Dysregulation
5. PD	- Congenital defects of phagocyte number or function
6. IID	- Defects in intrinsic and innate immunity
7. AID	- Autoinflammatory Disorders
8. CD	- Complement Deficiencies
9. BMF	- Bone marrow failure

¹⁷⁶

¹⁷⁷ 2.2 Variant Class Observation Probability

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

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

¹⁷⁸ where p is the allele frequency, $q = 1 - p$, p^2 represents the homozygous dominant,
¹⁷⁹ $2pq$ the heterozygous, and q^2 the homozygous recessive genotype frequencies. For dis-
¹⁸⁰ ease phenotypes, particularly under AR MOI, the risk is traditionally linked to the
¹⁸¹ homozygous state (p^2); however, to account for compound heterozygosity across mul-
¹⁸² tiple variants, we extend this by incorporating the contribution from other pathogenic
¹⁸³ alleles.

¹⁸⁴ Our computational pipeline estimated the probability of observing a disease-associated
¹⁸⁵ genotype for each variant and aggregated these probabilities by gene and ClinVar
¹⁸⁶ classification. This approach included all variant classifications, not limited solely to
¹⁸⁷ those deemed “pathogenic”, and explicitly conditioned the classification on the given
¹⁸⁸ phenotype, recognising that a variant could only be considered pathogenic relative to
¹⁸⁹ a defined clinical context. The core calculations proceeded as follows:

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

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

If a variant had no observed allele ($p_i = 0$), we assigned a minimal risk:

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

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

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

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

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

- For **AR** variants, disease manifested when two pathogenic alleles were present. In this case, we accounted for both the homozygous state and the possibility of compound heterozygosity:

$$p_{\text{disease},i} = p_i^2 + 2p_i(P_{\text{tot}} - p_i).$$

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

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

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

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

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

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

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

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

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

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

212 The final results were aggregated by gene and ClinVar classification and used to
213 generate summary statistics that reviewed the predicted disease observation proba-
214bilities.

215 **2.3 Validation of Autosomal Dominant Estimates Using *NFKB1***

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

220 **1. Reference Dataset.** We used a reference dataset reported by Tuijnenburg
221 et al. (17) to build a validation model in an AD disease gene. A whole-genome se-
222 quencing study of 846 predominantly sporadic, unrelated PID cases from the NIHR
223 BioResource-Rare Diseases cohort identified *NFKB1* as one of the genes most strongly
224 associated with PID. Sixteen novel heterozygous variants-including truncating, mis-
225 sense, and gene deletion variants-in *NFKB1* were found, accounting for 46% of CVID
226 cases ($n = 390$) in the cohort.

227 Functional analyses, including structural protein evaluation, immunophenotyping,
228 immunoblotting, and ex vivo lymphocyte stimulation, revealed that all carriers exhib-
229 ited deficiencies in B-lymphocyte differentiation, particularly an increased CD21low
230 B-cell population. These findings had established heterozygous loss-of-function vari-
231 ants in *NFKB1* as the most common monogenic cause of CVID, with significant
232 prognostic implications.

233 **2. Cohort Prevalence Calculation.** Therefore, we used this UK-based cohort
234 of 846 unrelated PID patients where 390 cases of CVID were attributed to *NFKB1*,

yielding an observed cohort prevalence of

$$\text{Prevalence}_{\text{cohort}} = \frac{390}{846} \approx 0.461.$$

3. National Estimate Based on Literature. Based on literature, the prevalence of CVID in the general population was estimated at approximately 1/25 000 (17–20). For a UK population of $N_{\text{UK}} \approx 69\,433\,632$, the expected number of CVID cases was calculated as

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

Thus, the maximum expected number of *NFKB1*-related CVID cases in the entire population was estimated as

$$\text{Estimated } NFKB1 \text{ cases} \approx 2777 \times 0.461 \approx 1280,$$

²³³ with an approximate 95% confidence interval (derived from Wilson's method) of 1188
²³⁴ to 1374 cases.

4. Bayesian Adjustment. Given that the clinical cohort was derived from a specialized setting-likely capturing nearly all PID cases-the observed 390 cases may have better represented the true burden. To reconcile these perspectives, we performed a Bayesian adjustment by combining the known cohort data with the national estimate. Specifically, we computed a weighted average to symbolically acknowledge potential uncertainty:

$$\text{Adjusted Estimate} = w \cdot 390 + (1 - w) \cdot 1280,$$

with w set to 0.9 to reflect a strong preference for the observed data. Additionally, we modelled the uncertainty in the observed prevalence using a beta distribution:

$$p \sim \text{Beta}(390 + 1, 846 - 390 + 1),$$

²³⁵ and generated 10 000 posterior samples to obtain a density distribution for the ad-
²³⁶ justed estimate.

²³⁷ **5. Validation test.** Thus, the expected number of *NFKB1*-related CVID cases
²³⁸ derived from our genome-wide probability estimates was compared with the observed
²³⁹ counts from the UK-based PID cohort. This comparison validated our framework for
²⁴⁰ estimating disease incidence in AD disorders.

241 **2.4 Validation Study for Autosomal Recessive CF Using CFTR**

242 To validate our framework for AR diseases, we focused on Cystic Fibrosis (CF).
243 For comparability sizes between the validation studies, we analysed the most com-
244 mon SNV in the *CFTR* gene, typically reported as “p.Arg117His” (GRCh38 Chr
245 7:117530975 G/A, MANE Select HGVS p.ENST00000003084.11: p.Arg117His). Our
246 goal was to validate our genome-wide probability estimates by comparing the ex-
247 pected number of CF cases attributable to the p.Arg117His variant in *CFTR* with
248 the nationally reported case count in a well-characterised disease cohort (21–23).

1. Expected Genotype Counts. Let p denote the allele frequency of the p.Arg117His variant and q denote the combined frequency of all other pathogenic *CFTR* variants, such that

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

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

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

and

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

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

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

2. Mortality Adjustment. Since CF patients experience increased mortality, we adjusted the expected genotype counts using an exponential survival model (21–23). With an annual mortality rate $\lambda \approx 0.004$ and a median age of 22 years, the survival factor was computed as:

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

Thus, the mortality-adjusted expected genotype count became:

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

3. Bayesian Uncertainty Simulation. To incorporate uncertainty in the allele frequency p , we modelled p as a beta-distributed random variable:

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

249 using a large effective allele count (AN_{eff}) for illustration. By generating 10,000 poste-
250 rior samples of p , we obtained a distribution of the literature-based adjusted expected
251 counts, E_{adj} .

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

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

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

256 2.5 Protein Network and Genetic Constraint Interpretation

257 A PPI network was constructed using protein interaction data from STRINGdb (16).
258 We previously prepared and reported on this dataset consisting of 19,566 proteins and
259 505,968 interactions (<https://github.com/DylanLawless/ProteoMCLustR>). Node
260 attributes were derived from log-transformed score-positive-total values, which in-
261 formed both node size and colour. Top-scoring nodes (top 15 based on score) were
262 labelled to highlight prominent interactions. To evaluate group differences in score-
263 positive-total across major disease categories, one-way Analysis of Variance (ANOVA)
264 was performed followed by Tukey Honestly Significant Difference (HSD) post hoc tests
265 (and non-parametric Dunn's test for confirmation). GnomAD v4.1 constraint metrics
266 data was used for the PPI analysis and was sourced from Karczewski et al. (7). This
267 provided transcript-level metrics, such as observed/expected ratios, Loss-Of-function
268 Observed/Expected Upper bound Fraction (LOEUF), pLI, and Z-scores, quantifying
269 loss-of-function and missense intolerance, along with confidence intervals and related
270 annotations for 211,523 observations.

271 2.6 Gene Set Enrichment Test

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

279

3 Results

280

3.1 Observation Probability Across Disease Genes

281 Our study integrated large-scale annotation databases with gene panels from PanelAppRex to systematically assess disease genes by MOI. By combining population
282 allele frequencies with ClinVar clinical classifications, we computed an expected obser-
283 vation probability for each SNV, representing the likelihood of encountering a variant
284 of a specific pathogenicity for a given phenotype. We report these probabilities for
285 54,814 ClinVar variant classifications across 557 genes (linked dataset (26)).

287 In practice, our approach computed a simple observation probability for every
288 SNV across the genome and was applicable to any disease-gene panel. Here, we fo-
289 cused on panels related to Primary Immunodeficiency or Monogenic Inflammatory
290 Bowel Disease, using PanelAppRex panel ID 398 as a case study. **Figure 1** dis-
291 plays all reported ClinVar variant classifications for this panel. The resulting natural
292 scaling system (-5 to +5) accounts for the frequently encountered combinations of
293 classification labels (e.g. benign to pathogenic). The resulting data set (26) is briefly
294 shown in **Table 1** to illustrate that our method yielded estimations of the probability
295 of observing a variant with a particular ClinVar classification.

Table 1: Example of the first several rows from our main results for 557 genes of PanelAppRex’s panel: (ID 398) Primary immunodeficiency or monogenic inflammatory bowel disease. “ClinVar Significance” indicates the pathogenicity classification assigned by ClinVar, while “inVar Significance” indicates the pathogenicity classification assigned by ClinVar, while “Occurrence Prob” represents our calculated probability of observing the corresponding variant class for a given phenotype. Additional columns, such as population allele frequency, are not shown. (26)

Gene	Panel ID	ClinVar Clinical Significance	GRCh38 Pos	HGVSc (VEP)	HGVSp (VEP)	Inheritance	Occurrence Probability
ABI3	398	Uncertain significance	49210771	c.47G>A	p.Arg16Gln	AR	0.000000007
ABI3	398	Uncertain significance	49216678	c.265C>T	p.Arg89Cys	AR	0.000000005
ABI3	398	Uncertain significance	49217742	c.289G>A	p.Val97Met	AR	0.000000002
ABI3	398	Uncertain significance	49217781	c.328G>A	p.Gly110Ser	AR	0.000000002
ABI3	398	Uncertain significance	49217844	c.391C>T	p.Pro131Ser	AR	0.000000015
ABI3	398	Uncertain significance	49220257	c.733C>G	p.Pro245Ala	AR	0.000000022

296

3.2 Validation studies

297

3.2.1 Validation of Dominant Disease Occurrence with *NFKB1*

298 To validate our genome-wide probability estimates for AD disorders, we focused
299 on *NFKB1*. We used a reference dataset from Tuijnernburg et al. (17), in which
300 whole-genome sequencing of 846 PID patients identified *NFKB1* as one of the genes
301 most strongly associated with the disease, with 390 CVID cases attributed to het-
302 erozygous variants. Our goal was to compare the predicted number of *NFKB1*-related
303 CVID cases with the reported count in this well-characterised national-scale cohort.

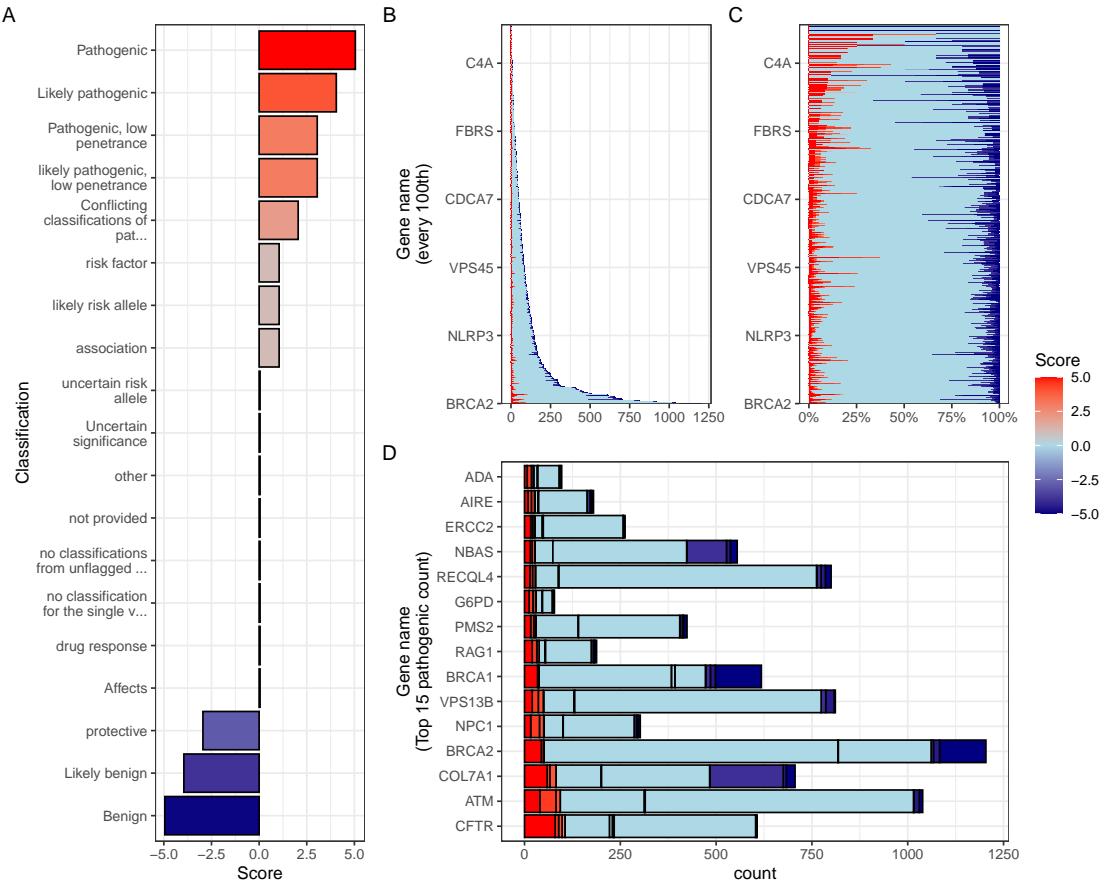


Figure 1: Summary of ClinVar clinical significance classifications in the PID gene panel. (A) Shows the numeric score coding for each classification. Panels (B) and (C) display the tally of classifications per gene as absolute counts and as percentages, respectively. (D) Highlights the top 15 genes with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

304 Our model calculated 456 *NFKB1*-related CVID cases in the UK. In the reference
 305 cohort, 390 *NFKB1* CVID cases were reported. We additionally wanted to account for
 306 potential under-reporting in the reference study. We used an extrapolated national
 307 CVID prevalence to yield an upper bound maximum of 1280 cases (95% Confidence
 308 Interval (CI): 1188–1374), while a Bayesian-adjusted mixture estimate produced a
 309 median of 835 cases (95% CI: 789–882). **Figure 2 (A)** illustrates that our predicted
 310 value of 456 lies within these ranges and is closer to the observed count, thereby
 311 supporting the validity of our integrated probability estimation framework for AD
 312 disorders.

313 **3.2.2 Validation of Recessive Disease Occurrence with *CFTR***

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

319 In contrast, the nationally reported number of CF cases was 714, as recorded in the
320 UK Cystic Fibrosis Registry 2023 Annual Data Report (21). To account for factors
321 such as reduced penetrance and the mortality-adjusted expected genotype, we derived
322 a Bayesian-adjusted estimate via posterior simulation. Our Bayesian approach yielded
323 a median estimate of 740 cases (95% CI: 696, 786) and a mixture-based estimate of 727
324 cases (95% CI: 705, 750). **Figure 2 (B)** illustrates the close concordance between the
325 predicted values, the Bayesian-adjusted estimates, and the national report supports
326 the validity of our approach for estimating disease.

327 **Figure S1** shows the final values for these genes of interest in a given population
328 size and phenotype. It reveals that an allele frequency threshold of approximately
329 0.000007 is required to observe a single heterozygous disease-causing variant carrier in
330 the UK population for both genes. However, owing to the AR MOI pattern of *CFTR*,
331 this threshold translates into more than 100,000 heterozygous carriers, compared to
332 only 456 carriers for the AD gene *NFKB1*. Note that this allele frequency threshold,
333 being derived from the current reference population, represents a lower bound that
334 can become more precise as public datasets continue to grow. This marked difference
335 underscores the significant impact of MOI patterns on population carrier frequencies
336 and the observed disease prevalence.

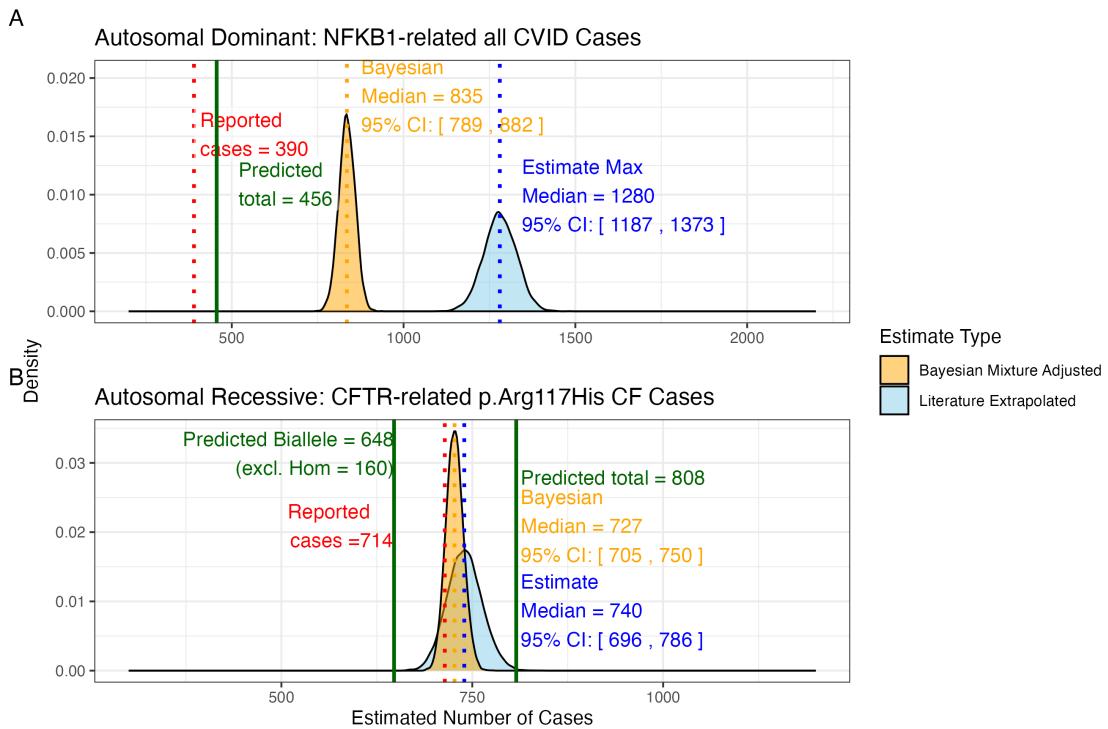


Figure 2: Prior probabilities compared to validation disease cohort metrics.

(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.

337 **3.2.3 Interpretation of ClinVar Variant Observations**

338 **Figure 3** shows the two validation study PID genes, representing AR and dominant
 339 MOI. **Figure 3 (A)** illustrates the overall probability of an affected birth by ClinVar
 340 variant classification, whereas **Figure 3 (B)** depicts the total expected number of
 341 cases per classification for an example population, here the UK, of approximately 69.4
 342 million.

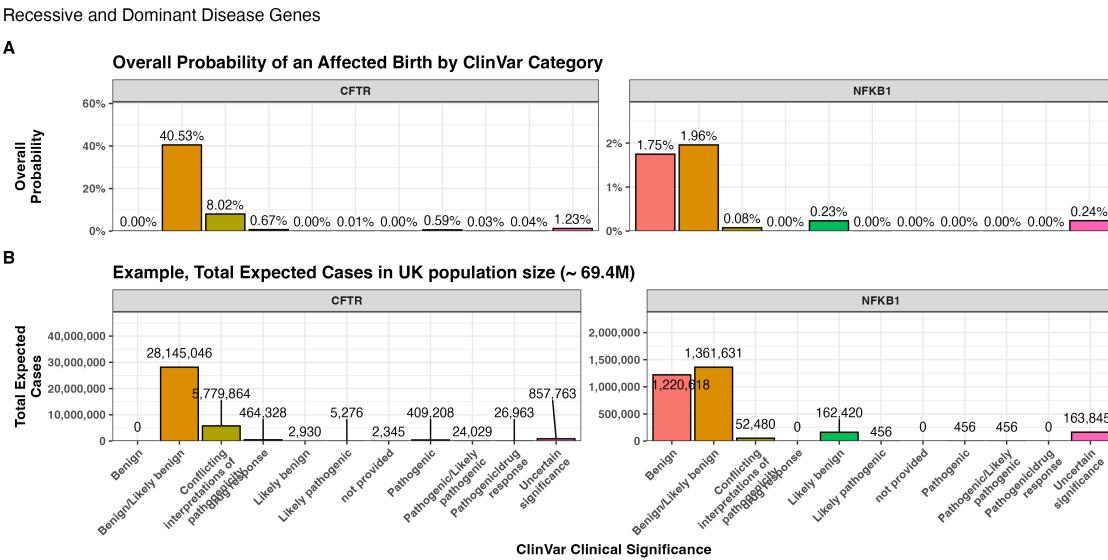


Figure 3: **Combined bar charts summarizing the genome-wide analysis of ClinVar clinical significance for the PID gene panel.** Panel (A) shows the overall probability of an affected birth by variant classification, and (B) displays the total expected number of cases per classification, both stratified by gene. These integrated results illustrate the variability in variant observations across genes and underpin our validation of the probability estimation framework.

343 **3.3 Genetic constraint in high-impact protein networks**

344 We next examined genetic constraint in high-impact protein networks across the whole
 345 IEI gene set of over 500 known disease-gene phenotypes (1). By integrating ClinVar
 346 variant classification scores with PPI data, we quantified the pathogenic burden per
 347 gene and assessed its relationship with network connectivity and genetic constraint
 348 (7; 16).

349 **3.3.1 Score-Positive-Total within IEI PPI network**

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

(benign) scoring variants, which are statistically incidental in this analysis. We tallied gene-level positive scores to give the score positive total metric. **Figure 4 (A)** shows the PPI network of disease-associated genes, where node size and colour encode the score positive total (log-transformed). The top 15 genes/proteins with the highest total prior probabilities of being observed with disease are labelled (as per **Figure 1**).

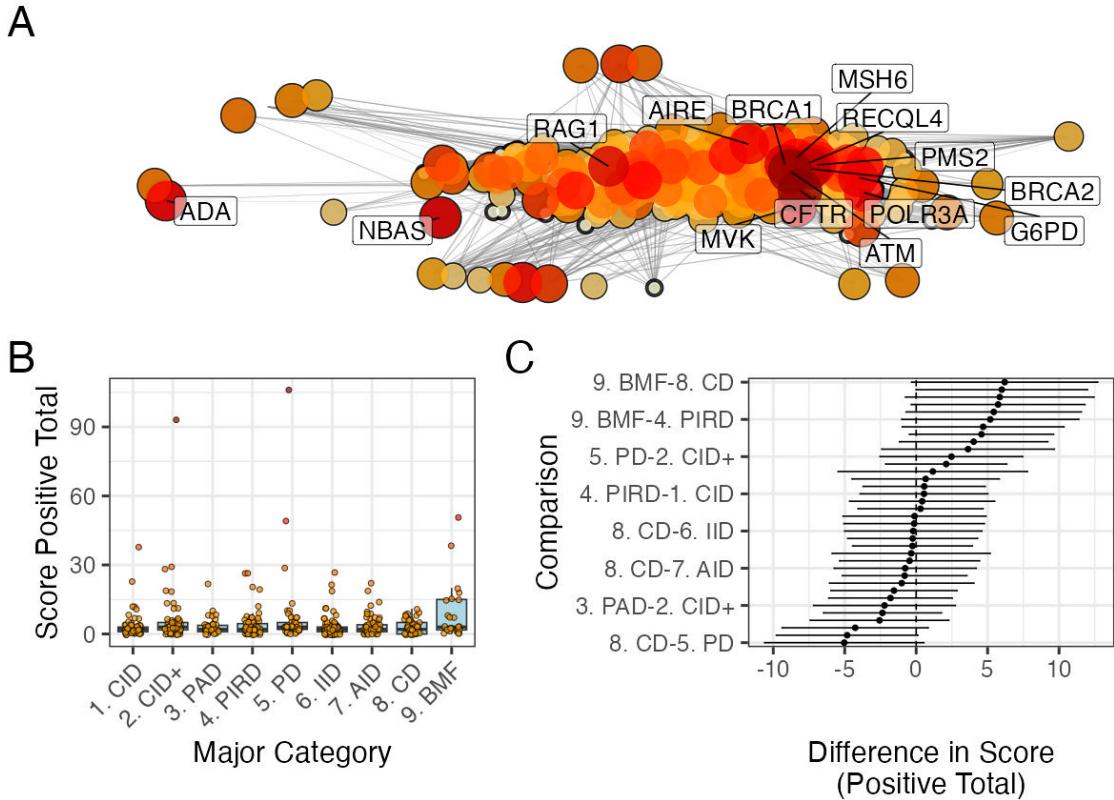


Figure 4: **PPI network and score positive total ClinVar significance variants.** (A) PPI network of disease-associated genes. Node size and colour represent the log-transformed score positive total, the top 15 genes/proteins with the highest probability of being observed in disease are labelled. (B) 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.

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

We checked for any statistical enrichment in score positive totals, which represents the expected observation of pathogenicity, between the IEI categories. The one-way ANOVA revealed an effect of major disease category on score positive total ($F(8, 500) = 2.82, p = 0.0046$), indicating that group means were not identical, which we observed in **Figure 4 (B)**. However, despite some apparent differences in median scores across

363 categories (i.e. 9. Bone Marrow Failure (BMF)), the Tukey HSD post hoc comparisons
 364 **Figure 4 (C)** showed that all pairwise differences had 95% confidence intervals
 365 overlapping zero, suggesting that individual group differences were not significant.

366 3.3.3 UMAP Embedding of the PPI Network

367 To address the density of the PPI network for the IEI gene panel, we applied Uniform
 368 Manifold Approximation and Projection (UMAP) (**Figure 5**). Node sizes reflect
 369 interaction degree, a measure of evidence-supported connectivity (16). We tested
 370 for a correlation between interaction degree and score positive total. In **Figure**
 371 **5**, gene names with degrees above the 95th percentile are labelled in blue, while
 372 the top 15 genes by score positive total are labelled in yellow (as per **Figure 1**).
 373 Notably, genes with high pathogenic variant loads segregated from highly connected
 374 nodes, suggesting that Loss-of-Function (LOF) in hub genes is selectively constrained,
 375 whereas damaging variants in lower-degree genes yield more specific effects. This
 376 observation was subsequently tested empirically.

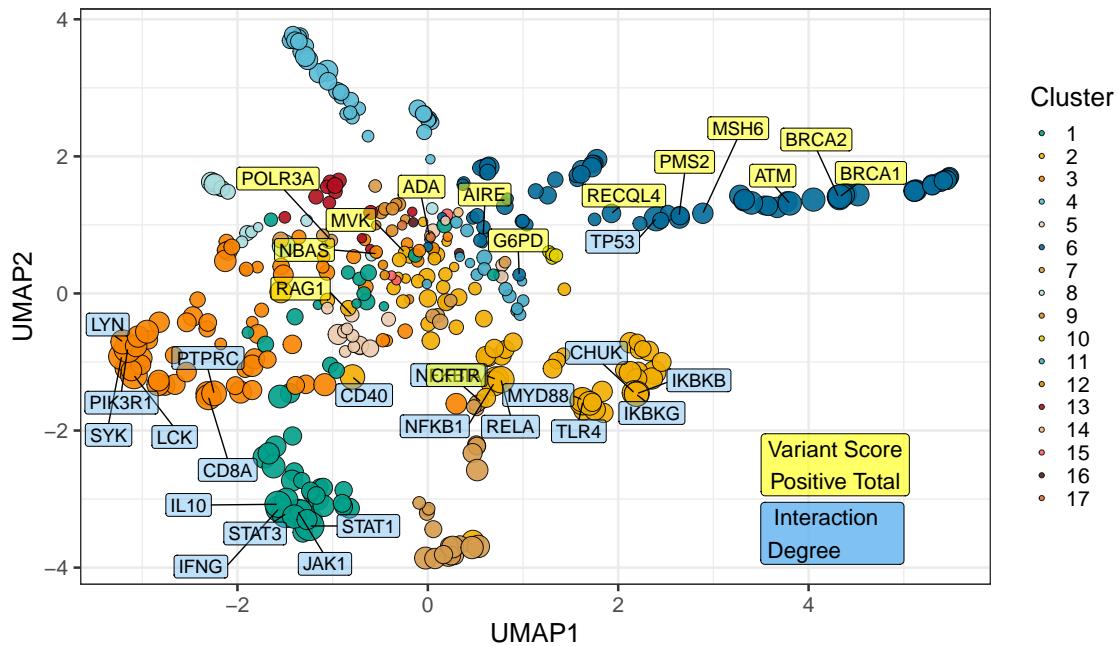


Figure 5: **UMAP embedding of the PPI network (p_umap).** The plot projects the high-dimensional protein-protein interaction network into two dimensions, with nodes coloured by cluster and sized by interaction degree. Blue labels indicate hub genes (degree above the 95th percentile) and yellow labels mark the top 15 genes by score positive total (damaging ClinVar classifications). The spatial segregation suggests that genes with high pathogenic variant loads are distinct from highly connected nodes.

377 **3.3.4 Hierarchical Clustering of Enrichment Scores for Major Disease Cat-**
378 **egories**

379 **Figure S2** presents a heatmap of standardised residuals for major disease categories
380 across network clusters, as per **Figure 5**. A dendrogram clusters similar disease cate-
381 gories, while the accompanying bar plot displays the maximum absolute standardised
382 residual for each category. Notably, (8) Complement Deficiencies (CD) shows the
383 highest maximum enrichment, followed by (9) BMF. While all maximum values
384 exceed 2, the threshold for significance, this likely reflects the presence of protein
385 clusters with strong damaging variant scores rather than uniform significance across
386 all categories (i.e. genes from cluster 4 in 8 CD).

387 **3.3.5 PPI Connectivity, LOEUF Constraint and Enriched Network Clus-**
388 **ter Analysis**

389 Based on the preliminary insight from **Figure S2**, we evaluated the relationship
390 between network connectivity (PPI degree) and LOEUF constraint (LOEUF upper rank)
391 Karczewski et al. (7) using Spearman's rank correlation. Overall, there was a weak
392 but significant negative correlation ($\rho = -0.181, p = 0.00024$) at the global scale,
393 indicating that highly connected genes tend to be more constrained. A supplementary
394 analysis (see **Figure 6**) did not reveal distinct visual associations between network
395 clusters and constraint metrics, likely due to the high network density. However
396 once stratified by gene clusters, the natural biological scenario based on quantitative
397 PPI evidence (16), some groups showed strong correlations; for instance, cluster 2
398 ($\rho = -0.375, p = 0.000994$) and cluster 4 ($\rho = -0.800, p < 0.000001$), while others did
399 not. This indicated that shared mechanisms within pathway clusters may underpin
400 genetic constraints, particularly for LOF intolerance. We observe that the score
401 positive total metric effectively summarises the aggregate pathogenic burden across
402 IEI genes, serving as a robust indicator of genetic constraint and highlighting those
403 with elevated disease relevance.

404 **Figure 6 (C, D)** shows the re-plotted PPI networks for clusters with significant
405 correlations between PPI degree and LOEUF upper rank. In these networks, node
406 size is scaled by a normalised variant score, while node colour reflects the variant
407 score according to a predefined palette.

408 **3.4 New Insight from Functional Enrichment**

409 To interpret the functional relevance of our prioritised IEI gene sets with the highest
410 load of damaging variants (i.e. clusters 2 and 4 in **Figure 6**), we performed func-
411 tional enrichment analysis for known disease associations using MsigDB with FUMA
412 (i.e. GWAScatalog and Immunologic Signatures) (24). Composite enrichment pro-
413 files (**Figure S4**) reveal that our enriched PPI clusters were associated with distinct
414 disease-related phenotypes, providing functional insights beyond traditional IUIS IEI

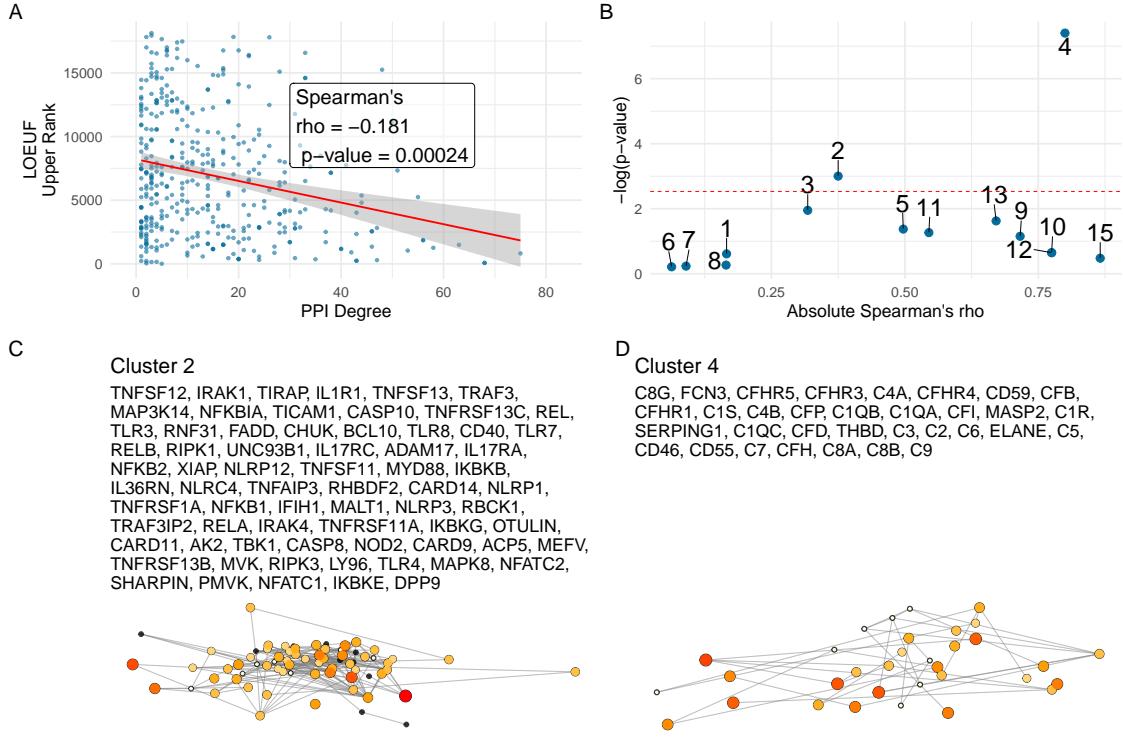


Figure 6: **Correlation between PPI degree and LOEUF upper rank.** (A) Ananlysis across all genes revealed a weak, significant negative correlation between PPI degree and LOEUF upper rank. (B) The cluster-wise analysis showed that clusters 2 and 4 exhibited moderate to strong correlations, while other 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.

groupings (1). The gene expression profiles shown in **Figure S5** (GTEx v8 54 tissue types) offer the tissue-specific context for these associations. Together, these results enable the annotation of IEI gene sets with established disease phenotypes, supporting a data-driven classification of IEI.

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

428 **3.5 Genome-wide Gene Distribution and Locus-specific Vari-**
 429 **ant Occurrence**

430 **Figure 7 (A)** shows a genome-wide karyoplot of all IEI panel genes across GRCh38,
 431 with colour-coding based on MOI. Figures (B) and (C) display zoomed-in locus plots
 432 for *NFKB1* and *CFTTR*, respectively. In **Figure 7 (B)**, the probability of observing
 433 variants with known classifications is high only for variants such as p.Ala475Gly,
 434 which are considered benign in the AD *NFKB1* gene that is intolerant to LOF. In
 435 **Figure 7 (C)**, high probabilities of observing patients with pathogenic variants in
 436 *CFTTR* are evident, reproducing this well-established phenomenon. Furthermore, the
 437 analysis of Linkage Disequilibrium (LD) using R^2 shows that high LD regions can be
 438 modelled effectively, allowing independent variant signals to be distinguished.

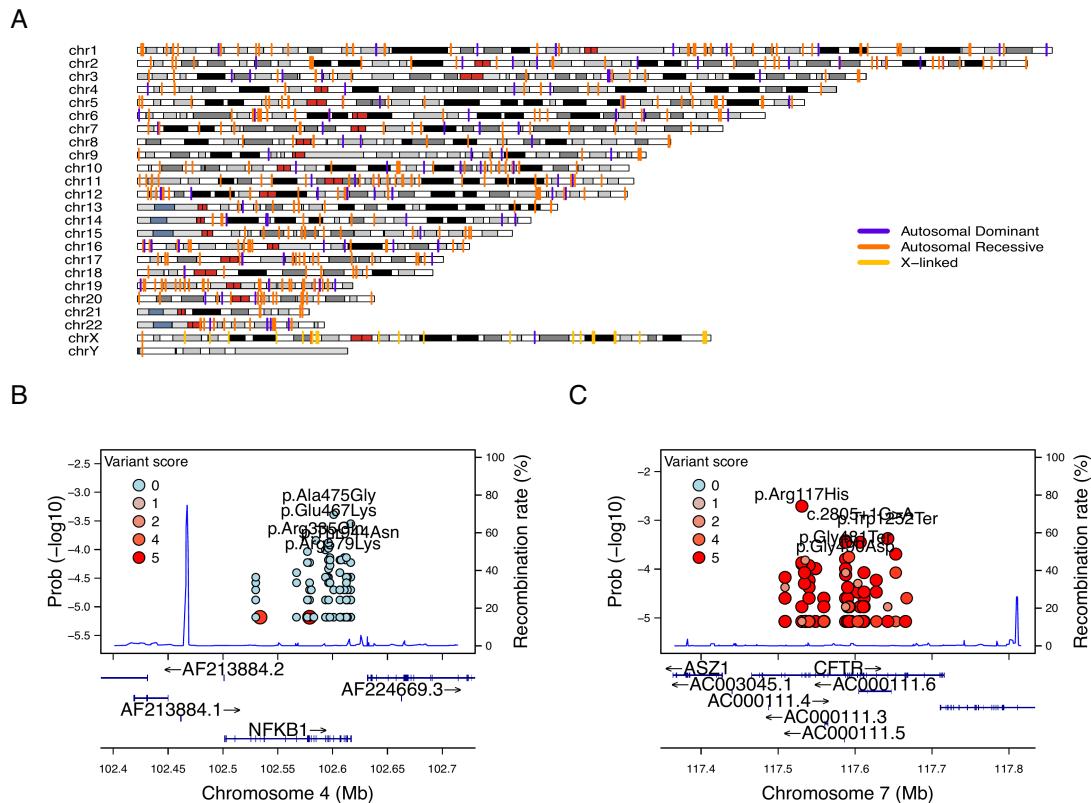


Figure 7: Genome-wide IEI, variant occurrence probability and LD by R^2 . (A) Genome-wide karyoplot of all IEI panel genes mapped to GRCh38, with colours indicating MOI. (B) Zoomed-in locus plot for *NFKB1* showing variant observation probabilities; only benign variants such exhibit high probabilities in this AD gene intolerant to LOF. (C) Locus plot for *CFTTR* displaying high probabilities for pathogenic variants; due to the dense clustering of pathogenic variants, score filter >0 was applied. Top five variant are labelled per gene.

3.6 Deriving new IEI classifications and insight

This section looks at the current 9 major IEI classes, and 45 subclasses. These were designed over last few decades based on phenotype observation and subjective decisions. We can use the new information be potentially define better, new data-based classifications: clustering by PPI, severity/probability, and markers (T, B, NK, Ig), phenotype CVID, etc. Potential but maybe it will be too noisy. From the previous section, cross compare the clusters from PPI with the IUIS categories. This section will be complete by 2025-04-07.

3.6.1 Integration of Variant Probabilities into IEI Genetics Data

We integrated the computed prior probabilities for observing variants in all known genes associated with a given phenotype (1), across AD, AR, and XL MOI, into our IEI genetics framework. These calculations, derived from gene panels in PanelAppRex, have yielded novel insights for the IEI disease panel. The final result comprised of machine- and human-readable datasets, including the table of variant classifications and priors available via a the linked repository (26), and a user-friendly web interface that incorporates these new metrics.

Figure 8 shows the interface summarising integrated variant data. Server-side pre-calculation of summary statistics minimises browser load, while clinical significance is converted to 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).

The screenshot displays a web-based application for managing genetic variant data. The interface includes a header with 'Viewer Zoom' and 'Search' functions. Below is a detailed table of variants:

Major category	Subcategory	Disease	Genetic defect	Inheritance	Gene score ClinVar pathogenicity	Prior prob of observing pathogenic	ClinVar SNV classification	ClinVar all variant reports	OMIM	Alpha Missense / Uniprot ID	HPO combined	HPO term
All				All								
1. CID	1. T+B+ SCID	CD3z deficiency	CD247	AR	2	0/0 - 1	1/0 / 33/1	15/0 / 133 / 218	186780	P20963	HP-0002715; Abnormalit	
1. CID	1. T+B+ SCID	CD3d deficiency	CD3D	AR	3	1 - 1	2/1 / 34 / 0	20/19 / 162 / 234	186790	P04234	HP-0002715; Abnormalit	
1. CID	1. T+B+ SCID	CD3e deficiency	CD3E	AR	3	1 - 1	1/2 / 29 / 2	29/14 / 173 / 346	186830	P07766	HP-0002715; Abnormalit	
1. CID	1. T+B+ SCID	Coronin-1A deficiency	CORO1A	AR	2	1	1/1 / 43 / 2	19/14 / 99 / 376	605000	P31146	HP-0002715; Abnormalit	
1. CID	1. T+B+ SCID	go-deficiency (common gamma chain SCID, CD132 deficiency)	IL2RG	XL	3	1 - 1	1/1 / 16 / 28	594/106 / 364 / 414	308380	P31785	HP-0002715; Abnormalit	
1. CID	1. T+B+ SCID	IL7Ra deficiency	IL7R	AR	12	1 - 1	6/2 / 81 / 14	81/26 / 458 / 586	146661	P16871	HP-0002715; Abnormalit	
1. CID	1. T+B+ SCID	ITPKB deficiency	ITPKB	AR	3	1 - 1	1/2 / 15 / 9	0/12 / 130 / 40	query	P27987	HP-0002715; Abnormalit	
1. CID	1. T+B+ SCID	JAK3 deficiency	JAK3	AR	12	1 - 1	9/15 / 131 / 13	152/19 / 427 / 1528	600173	P52333	HP-0002715; Abnormalit	
1. CID	1. T+B+ SCID	LAT deficiency	LAT	AR	1	1	1/0 / 39 / 3	14/2 / 138 / 242	602354	Q43561	HP-0002715; Abnormalit	

At the bottom, there are navigation links: '1-10 of 591 rows', 'Show 10', 'Previous', '1 2 3 4 5 ... 60 Next'.

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

461 **4 Discussion**

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

469 Our approach yielded two key results. First, our detailed, per-variant pre-calculated
470 results provide prior probabilities of observing disease-associated variants across all
471 MOI for any gene-disease combination. Second, the score positive total metric effec-
472 tively summarises the aggregate pathogenic burden across genes, serving as a robust
473 indicator of genetic constraint and highlighting those with elevated disease relevance.

Estimating disease risk in genetic studies is complicated by uncertainties in key parameters such as variant penetrance and the fraction of cases attributable to specific variants (6). In the simplest model, where a single, fully penetrant variant causes disease, the lifetime risk $P(D)$ is equivalent to the genotype frequency $P(G)$. For an allele with frequency p , this translates to:

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

When penetrance is incomplete, defined as $P(D | G)$, the risk becomes:

$$P(D) = P(G) P(D | G).$$

In more realistic scenarios where multiple variants contribute to disease, $P(G | D)$ denotes the fraction of cases attributable to a given variant. This leads to:

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

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

476 Our framework addresses this challenge by combining variant classifications, pop-
477 ulation allele frequencies, and curated gene-disease associations. While imperfect on
478 an individual level, these sources exhibit predictable aggregate behaviour, supported
479 by James-Stein estimation principles (27). Curated gene-disease associations help
480 identify genes that explainable for most disease cases, allowing us to approximate
481 $P(G | D)$ close to one. In this way, we obtain robust estimates of $P(G)$ (the fre-
482 quency of disease-associated genotypes), even when exact values of penetrance and
483 case attribution remain uncertain.

This approach allows us to pre-calculate priors and summarise the overall pathogenic burden using our *score positive total* metric. By focusing on a subset \mathcal{V} of variants

that pass stringent filtering, where each $P(G_i | D)$ is the probability that a case of disease D is attributable to variant i , we assume that, in aggregate,

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

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

Our results focused on IEI, but the genome-wide approach accommodates the distinct MOI patterns of AD, AR, and XL disorders. Whereas AD and XL conditions require only a single pathogenic allele, AR disorders necessitate the consideration of both homozygous and compound heterozygous states. These classical HWE-based estimates provide an informative baseline for predicting variant occurrence and serve as robust priors for Bayesian models of variant and disease risk estimation. This is an approach that has been underutilised in clinical and statistical genetics. As such, our framework refines risk calculations by incorporating MOI complexities and enhances clinicians' understanding of expected variant occurrences, thereby improving diagnostic precision.

Moreover, our method complements existing statistical approaches for aggregating variant effects with methods like Sequence Kernel Association Test (SKAT) and Aggregated Cauchy Association Test (ACAT) (28–31)) and multi-omics integration techniques (32; 33), while remaining consistent with established variant interpretation guidelines from the American College of Medical Genetics and Genomics (ACMG) (34) and complementary frameworks (35; 36), as well as quality control protocols (37; 38). Standardised reporting for qualifying variant sets, such as ACMG Secondary Findings v3.2 (39), further contextualises the integration of these probabilities into clinical decision-making.

We acknowledge that our current framework is restricted to SNVs and does not incorporate numerous other complexities of genetic disease, such as structural variants, de novo variants, hypomorphic alleles, overdominance, variable penetrance, tissue-specific expression, the Wahlund effect, pleiotropy, and others (6). In certain applications, more refined estimates would benefit from including factors such as embryonic lethality, condition-specific penetrance, and age of onset (10). Our analysis also relies on simplifying assumptions of random mating, an effectively infinite population, and the absence of migration, novel mutations, or natural selection.

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

520 **5 Conclusion**

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

524 **Acknowledgements**

525 We acknowledge Genomics England for providing public access to the PanelApp data.
526 The use of data from Genomics England panelapp was licensed under the Apache
527 License 2.0. The use of data from UniProt was licensed under Creative Commons
528 Attribution 4.0 International (CC BY 4.0). ClinVar asks its users who distribute or
529 copy data to provide attribution to them as a data source in publications and websites
530 (13). dbNSFP version 4.4a is licensed under the Creative Commons Attribution-
531 NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0); while we cite
532 this dataset as used our research publication, it is not used for the final version which
533 instead used ClinVar and gnomAD directly. GnomAD is licensed under Creative
534 Commons Zero Public Domain Dedication (CC0 1.0 Universal). GnomAD request
535 that usages cites the gnomAD flagship paper (7) and any online resources that include
536 the data set provide a link to the browser, and note that tool includes data from the
537 gnomAD v4.1 release.

538 **Competing interest**

539 We declare no competing interest.

540 **References**

- 541 [1] Stuart G. Tangye, Waleed Al-Herz, Aziz Bousfiha, Charlotte Cunningham-
542 Rundles, Jose Luis Franco, Steven M. Holland, Christoph Klein, Tomohiro Morio,
543 Eric Oksenhendler, Capucine Picard, Anne Puel, Jennifer Puck, Mikko R. J.
544 Seppänen, Raz Somech, Helen C. Su, Kathleen E. Sullivan, Troy R. Torger-
545 son, and Isabelle Meyts. Human Inborn Errors of Immunity: 2022 Update
546 on the Classification from the International Union of Immunological Societies
547 Expert Committee. *Journal of Clinical Immunology*, 42(7):1473–1507, October
548 2022. ISSN 0271-9142, 1573-2592. doi: 10.1007/s10875-022-01289-3. URL
549 <https://link.springer.com/10.1007/s10875-022-01289-3>.
- 550 [2] Dylan Lawless. PanelAppRex aggregates disease gene panels and facilitates
551 sophisticated search. March 2025. doi: 10.1101/2025.03.20.25324319. URL
552 <http://medrxiv.org/lookup/doi/10.1101/2025.03.20.25324319>.

- 553 [3] Antonio Rueda Martin, Eleanor Williams, Rebecca E. Foulger, Sarah Leigh,
554 Louise C. Daugherty, Olivia Niblock, Ivone U. S. Leong, Katherine R. Smith,
555 Oleg Gerasimenko, Eik Haraldsdottir, Ellen Thomas, Richard H. Scott, Emma
556 Baple, Arianna Tucci, Helen Brittain, Anna De Burca, Kristina Ibañez, Dalia
557 Kasperaviciute, Damian Smedley, Mark Caulfield, Augusto Rendon, and Ellen M.
558 McDonagh. PanelApp crowdsources expert knowledge to establish consensus
559 diagnostic gene panels. *Nature Genetics*, 51(11):1560–1565, November 2019.
560 ISSN 1061-4036, 1546-1718. doi: 10.1038/s41588-019-0528-2. URL <https://www.nature.com/articles/s41588-019-0528-2>.
- 561
- 562 [4] Oliver Mayo. A Century of Hardy–Weinberg Equilibrium. *Twin Research
563 and Human Genetics*, 11(3):249–256, June 2008. ISSN 1832-4274, 1839-
564 2628. doi: 10.1375/twin.11.3.249. URL https://www.cambridge.org/core/product/identifier/S1832427400009051/type/journal_article.
- 565
- 566 [5] Nikita Abramovs, Andrew Brass, and May Tassabehji. Hardy-Weinberg Equi-
567 librium in the Large Scale Genomic Sequencing Era. *Frontiers in Genetics*,
568 11:210, March 2020. ISSN 1664-8021. doi: 10.3389/fgene.2020.00210. URL
569 <https://www.frontiersin.org/article/10.3389/fgene.2020.00210/full>.
- 570
- 571 [6] Johannes Zschocke, Peter H. Byers, and Andrew O. M. Wilkie. Mendelian
572 inheritance revisited: dominance and recessiveness in medical genetics. *Nature
573 Reviews Genetics*, 24(7):442–463, July 2023. ISSN 1471-0056, 1471-0064.
574 doi: 10.1038/s41576-023-00574-0. URL <https://www.nature.com/articles/s41576-023-00574-0>.
- 575
- 576 [7] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings,
577 Jessica Alfoldi, Qingbo Wang, Ryan L Collins, Kristen M Laricchia, Andrea
578 Ganna, Daniel P Birnbaum, et al. The mutational constraint spectrum quantified
from variation in 141,456 humans. *Nature*, 581(7809):434–443, 2020.
- 579
- 580 [8] Sarah L. Bick, Aparna Nathan, Hannah Park, Robert C. Green, Monica H. Wo-
581 jcik, and Nina B. Gold. Estimating the sensitivity of genomic newborn screen-
582 ing for treatable inherited metabolic disorders. *Genetics in Medicine*, 27(1):
583 101284, January 2025. ISSN 10983600. doi: 10.1016/j.gim.2024.101284. URL
584 <https://linkinghub.elsevier.com/retrieve/pii/S1098360024002181>.
- 585
- 586 [9] Benjamin D. Evans, Piotr Słowiński, Andrew T. Hattersley, Samuel E. Jones,
587 Seth Sharp, Robert A. Kimmitt, Michael N. Weedon, Richard A. Oram,
588 Krasimira Tsaneva-Atanasova, and Nicholas J. Thomas. Estimating disease
589 prevalence in large datasets using genetic risk scores. *Nature Communications*,
12(1):6441, November 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-26501-7.
URL <https://www.nature.com/articles/s41467-021-26501-7>.
- 590
- 591 [10] William B. Hannah, Mitchell L. Drumm, Keith Nykamp, Tiziano Prampano,
592 Robert D. Steiner, and Steven J. Schrodi. Using genomic databases to de-
593 termine the frequency and population-based heterogeneity of autosomal reces-
594 sive conditions. *Genetics in Medicine Open*, 2:101881, 2024. ISSN 29497744.

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

- 595 [11] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov,
596 Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek,
597 Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J.
598 Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh
599 Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy,
600 Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer,
601 Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray
602 Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein
603 structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August
604 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL
605 <https://www.nature.com/articles/s41586-021-03819-2>.
- 606 [12] Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Taylor
607 Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hassabis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Science*, 381(6664):eadg7492, September
608 2023. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.adg7492. URL
609 <https://www.science.org/doi/10.1126/science.adg7492>.
- 610 [13] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen Chao,
611 Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Wonhee Jang, Karen Karapetyan, Kenneth Katz, Chunlei Liu, Zenith Maddipatla, Adriana Malheiro, Kurt McDaniel, Michael Ovetsky, George Riley, George Zhou, J Bradley Holmes, Brandi L Kattman, and Donna R Maglott. ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Research*, 46(D1):D1062–D1067, January 2018. ISSN 0305-1048, 1362-4962. doi:
612 10.1093/nar/gkx1153. URL <http://academic.oup.com/nar/article/46/D1/D1062/4641904>.
- 613 [14] The UniProt Consortium, Alex Bateman, Maria-Jesus Martin, Sandra Orchard,
614 Michele Magrane, Aduragbemi Adesina, Shadab Ahmad, Bowler-Barnett, and Others. UniProt: the Universal Protein Knowledgebase in 2025. *Nucleic Acids Research*, 53(D1):D609–D617, January 2025. ISSN 0305-1048, 1362-4962. doi:
615 10.1093/nar/gkae1010. URL <https://academic.oup.com/nar/article/53/D1/D609/7902999>.
- 616 [15] Xiaoming Liu, Chang Li, Chengcheng Mou, Yibo Dong, and Yicheng Tu. dbNSFP v4: a comprehensive database of transcript-specific functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Genome Medicine*, 12(1):103, December 2020. ISSN 1756-994X. doi: 10.1186/s13073-020-00803-9. URL <https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-020-00803-9>.

- 635 [16] Damian Szklarczyk, Katerina Nastou, Mikaela Koutrouli, Rebecca Kirsch, Far-
636 rokh Mehryary, Radja Hachilif, Dewei Hu, Matteo E Peluso, Qingyao Huang,
637 Tao Fang, et al. The string database in 2025: protein networks with directional-
638 ity of regulation. *Nucleic Acids Research*, 53(D1):D730–D737, 2025.
- 639 [17] Paul Tuijnenburg, Hana Lango Allen, Siobhan O. Burns, Daniel Greene,
640 Machiel H. Jansen, and Others. Loss-of-function nuclear factor B subunit
641 1 (NFKB1) variants are the most common monogenic cause of common vari-
642 able immunodeficiency in Europeans. *Journal of Allergy and Clinical Im-*
643 *munology*, 142(4):1285–1296, October 2018. ISSN 00916749. doi: 10.1016/
644 j.jaci.2018.01.039. URL <https://linkinghub.elsevier.com/retrieve/pii/S0091674918302860>.
- 645 [18] WHO Scientific Group et al. Primary immunodeficiency diseases: report of a
646 who scientific group. *Clin. Exp. Immunol.*, 109(1):1–28, 1997.
- 647 [19] Charlotte Cunningham-Rundles and Carol Bodian. Common variable immunod-
648 eficiency: clinical and immunological features of 248 patients. *Clinical immunol-*
649 *ogy*, 92(1):34–48, 1999.
- 650 [20] Eric Oksenhendler, Laurence Gérard, Claire Fieschi, Marion Malphettes, Gael
651 Mouillot, Roland Jaussaud, Jean-François Viallard, Martine Gardembas, Lionel
652 Galicier, Nicolas Schleinitz, et al. Infections in 252 patients with common variable
653 immunodeficiency. *Clinical Infectious Diseases*, 46(10):1547–1554, 2008.
- 654 [21] Y Naito, F Adams, S Charman, J Duckers, G Davies, and S Clarke. Uk cystic
655 fibrosis registry 2023 annual data report. *London: Cystic Fibrosis Trust*, 2023.
- 656 [22] Carlo Castellani, CFTR2 team, et al. Cftr2: how will it help care? *Paediatric*
657 *respiratory reviews*, 14:2–5, 2013.
- 658 [23] Hartmut Grasemann and Felix Ratjen. Cystic fibrosis. *New England Journal*
659 *of Medicine*, 389(18):1693–1707, 2023. doi: 10.1056/NEJMra2216474. URL
660 <https://www.nejm.org/doi/full/10.1056/NEJMra2216474>.
- 661 [24] Kyoko Watanabe, Erdogan Taskesen, Arjen Van Bochoven, and Danielle
662 Posthuma. Functional mapping and annotation of genetic associations with
663 FUMA. *Nature Communications*, 8(1):1826, November 2017. ISSN 2041-1723.
664 doi: 10.1038/s41467-017-01261-5. URL <https://www.nature.com/articles/s41467-017-01261-5>.
- 665 [25] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir,
666 Pablo Tamayo, and Jill P. Mesirov. Molecular signatures database (MSigDB)
667 3.0. *Bioinformatics*, 27(12):1739–1740, June 2011. ISSN 1367-4811, 1367-
668 4803. doi: 10.1093/bioinformatics/btr260. URL <https://academic.oup.com/bioinformatics/article/27/12/1739/257711>.
- 669
- 670
- 671

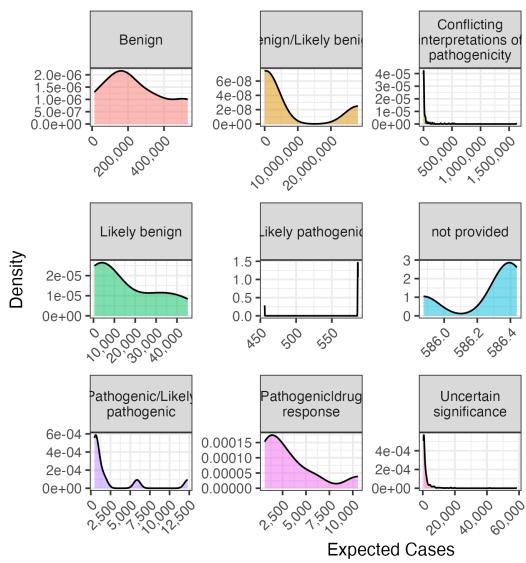
- 672 [26] Dylan Lawless. Variant risk estimate probabilities for iei genes. March 2025. doi:
673 10.5281/zenodo.15111584. URL <https://doi.org/10.5281/zenodo.15111584>.
- 674 [27] Bradley Efron and Carl Morris. Stein’s Estimation Rule and Its Competitors—
675 An Empirical Bayes Approach. *Journal of the American Statistical Association*,
676 68(341):117, March 1973. ISSN 01621459. doi: 10.2307/2284155. URL <https://www.jstor.org/stable/2284155?origin=crossref>.
- 678 [28] Yaowu Liu, Sixing Chen, Zilin Li, Alanna C Morrison, Eric Boerwinkle, and
679 Xihong Lin. Acat: a fast and powerful p value combination method for rare-
680 variant analysis in sequencing studies. *The American Journal of Human Genetics*,
681 104(3):410–421, 2019.
- 682 [29] Xihao Li, Zilin Li, Hufeng Zhou, Sheila M Gaynor, Yaowu Liu, Han Chen, Ryan
683 Sun, Rounak Dey, Donna K Arnett, Stella Aslibekyan, et al. Dynamic incorpora-
684 tion of multiple in silico functional annotations empowers rare variant association
685 analysis of large whole-genome sequencing studies at scale. *Nature genetics*, 52
686 (9):969–983, 2020.
- 687 [30] Michael C Wu, Seunggeun Lee, Tianxi Cai, Yun Li, Michael Boehnke, and Xi-
688 hong Lin. Rare-variant association testing for sequencing data with the sequence
689 kernel association test. *The American Journal of Human Genetics*, 89(1):82–93,
690 2011.
- 691 [31] Seunggeun Lee, Mary J Emond, Michael J Bamshad, Kathleen C Barnes, Mark J
692 Rieder, Deborah A Nickerson, David C Christiani, Mark M Wurfel, and Xihong
693 Lin. Optimal unified approach for rare-variant association testing with applica-
694 tion to small-sample case-control whole-exome sequencing studies. *The American
695 Journal of Human Genetics*, 91(2):224–237, 2012.
- 696 [32] Augustine Kong, Gudmar Thorleifsson, Michael L Frigge, Bjarni J Vilhjalmsson,
697 Alexander I Young, Thorgeir E Thorgeirsson, Stefania Benonisdottir, Asmundur
698 Oddsson, Bjarni V Halldorsson, Gisli Masson, et al. The nature of nurture:
699 Effects of parental genotypes. *Science*, 359(6374):424–428, 2018.
- 700 [33] Laurence J Howe, Michel G Nivard, Tim T Morris, Ailin F Hansen, Humaira
701 Rasheed, Yoonsoo Cho, Geetha Chittoor, Penelope A Lind, Teemu Palviainen,
702 Matthijs D van der Zee, et al. Within-sibship gwas improve estimates of direct
703 genetic effects. *BioRxiv*, pages 2021–03, 2021.
- 704 [34] Sue Richards, Nazneen Aziz, Sherri Bale, David Bick, Soma Das, Julie Gastier-
705 Foster, Wayne W Grody, Madhuri Hegde, Elaine Lyon, Elaine Spector, et al.
706 Standards and guidelines for the interpretation of sequence variants: a joint
707 consensus recommendation of the american college of medical genetics and ge-
708 nomics and the association for molecular pathology. *Genetics in medicine*, 17
709 (5):405–423, 2015.

- 710 [35] Sean V Tavtigian, Steven M Harrison, Kenneth M Boucher, and Leslie G
711 Biesecker. Fitting a naturally scaled point system to the acmg/amp variant
712 classification guidelines. *Human mutation*, 41(10):1734–1737, 2020.
- 713 [36] Quan Li and Kai Wang. Intervar: clinical interpretation of genetic variants by
714 the 2015 acmg-amp guidelines. *The American Journal of Human Genetics*, 100
715 (2):267–280, 2017.
- 716 [37] Brent S Pedersen, Joe M Brown, Harriet Dashnow, Amelia D Wallace, Matt
717 Velinder, Martin Tristani-Firouzi, Joshua D Schiffman, Tatiana Tvardik, Rong
718 Mao, D Hunter Best, et al. Effective variant filtering and expected candidate
719 variant yield in studies of rare human disease. *NPJ Genomic Medicine*, 6(1):1–8,
720 2021.
- 721 [38] Carl A Anderson, Fredrik H Pettersson, Geraldine M Clarke, Lon R Cardon,
722 Andrew P Morris, and Krina T Zondervan. Data quality control in genetic
723 case-control association studies. *Nature protocols*, 5(9):1564–1573, 2010. URL
724 <https://doi.org/10.1038/nprot.2010.116>.
- 725 [39] David T Miller, Kristy Lee, Noura S Abul-Husn, Laura M Amendola, Kyle Broth-
726 ers, Wendy K Chung, Michael H Gollob, Adam S Gordon, Steven M Harrison,
727 Ray E Hershberger, et al. Acmg sf v3. 2 list for reporting of secondary findings
728 in clinical exome and genome sequencing: a policy statement of the american
729 college of medical genetics and genomics (acmg). *Genetics in Medicine*, 25(8):
730 100866, 2023.

6 Supplemental

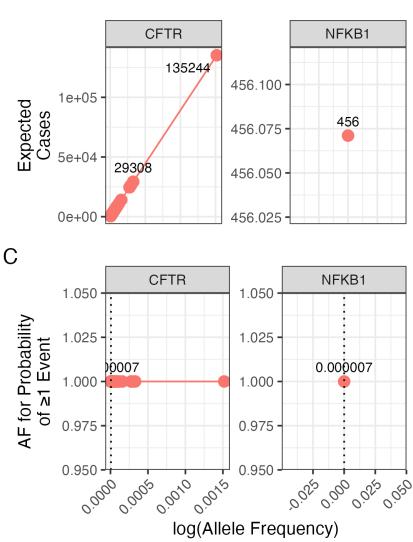
Condition: population size 69433632, phenotype PID-related, genes CFTR and NFKB1.

A



B

clinvar_clnsig • Pathogenic



C

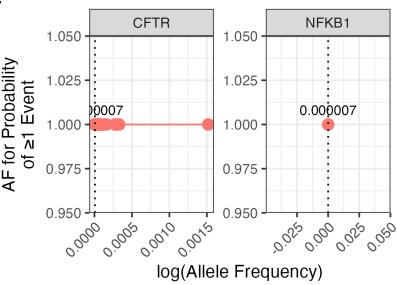


Figure S1: Interpretation of probability of observing a variant classification.
The result from the chosen validation genes *CFTR* and *NFKB1* are shown. Case counts are dependant on the population size and phenotype. (A) The density plots of expected observations by ClinVar clinical significance. We 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.

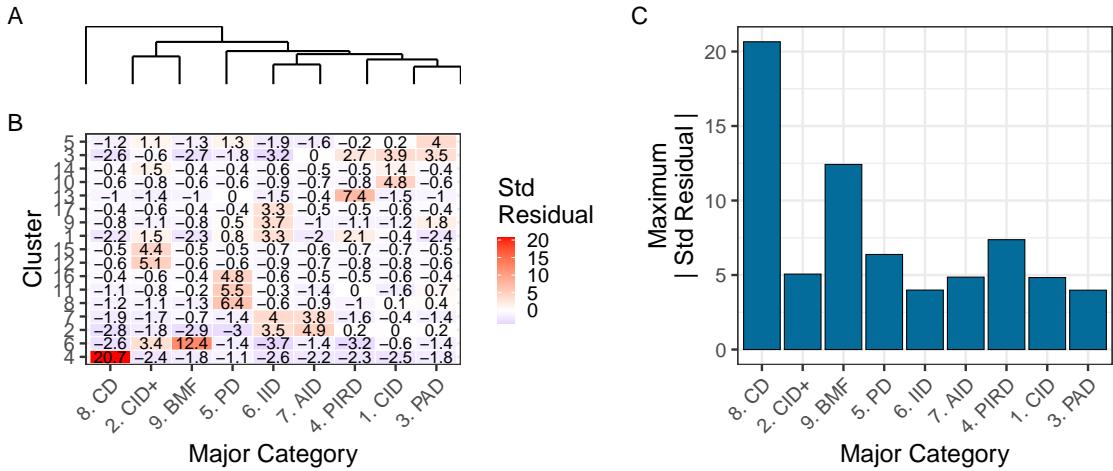


Figure S2: Hierarchical clustering of enrichment scores. The heatmap displays standardised residuals for major disease categories (x-axis) across network clusters (y-axis). A dendrogram groups similar disease categories, and the bar plot shows the maximum absolute residual per category. (8) CD and (9)BMF show the highest values, indicating significant enrichment or depletion (residuals $> |2|$). Definitions in **Box 2.1**.

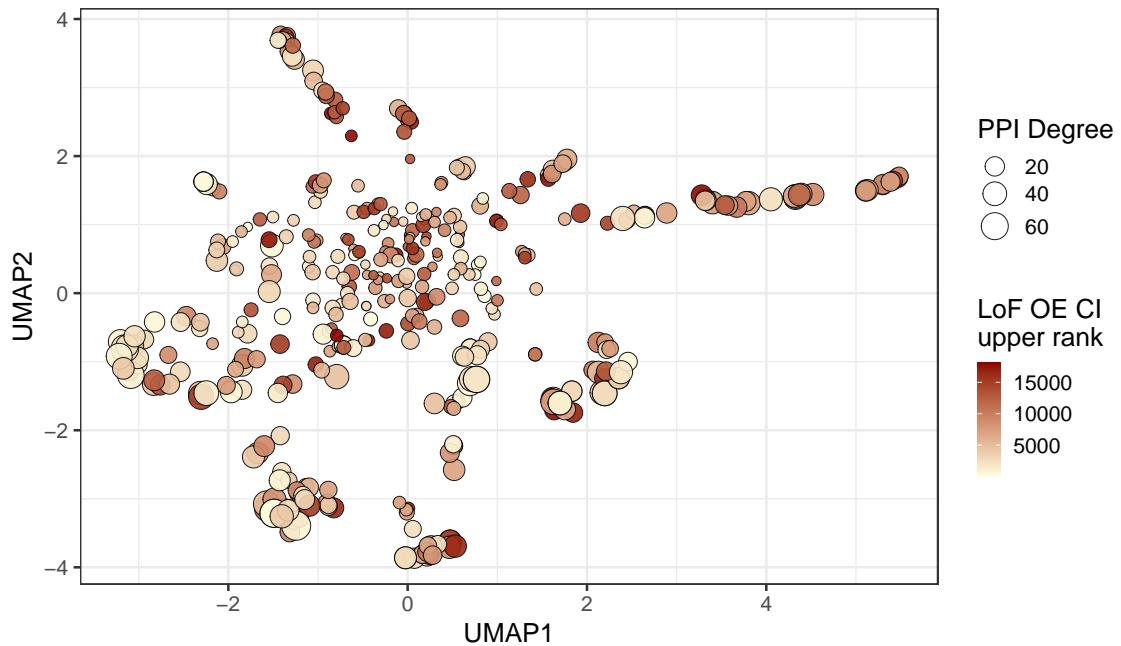


Figure S3: Supplementary analysis of PPI degree versus LOEUF upper rank with UMAP embedding of the PPI network. The relationship between PPI degree (size) and LOEUF upper rank (color) across gene clusters. No clear patterns are evident.

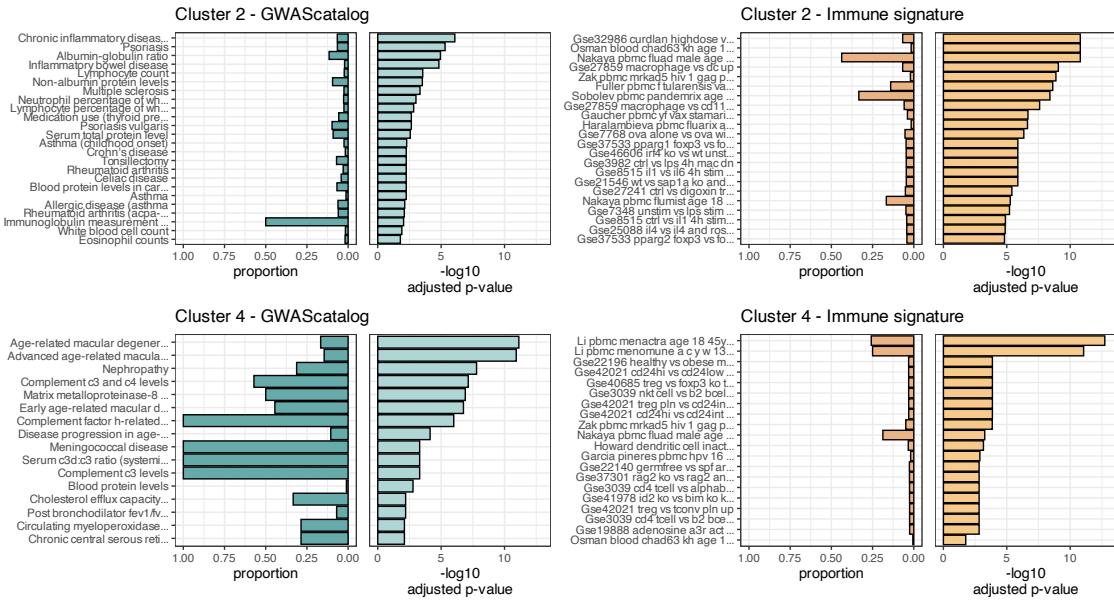


Figure S4: Composite Enrichment Profiles for IEI Gene Sets. We selected the top two enriched clusters (as per **Figure 6**) and performed functional enrichment analysis derived from known disease associations. For each gene set, the left panel displays the proportion of input genes overlapping with a 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.

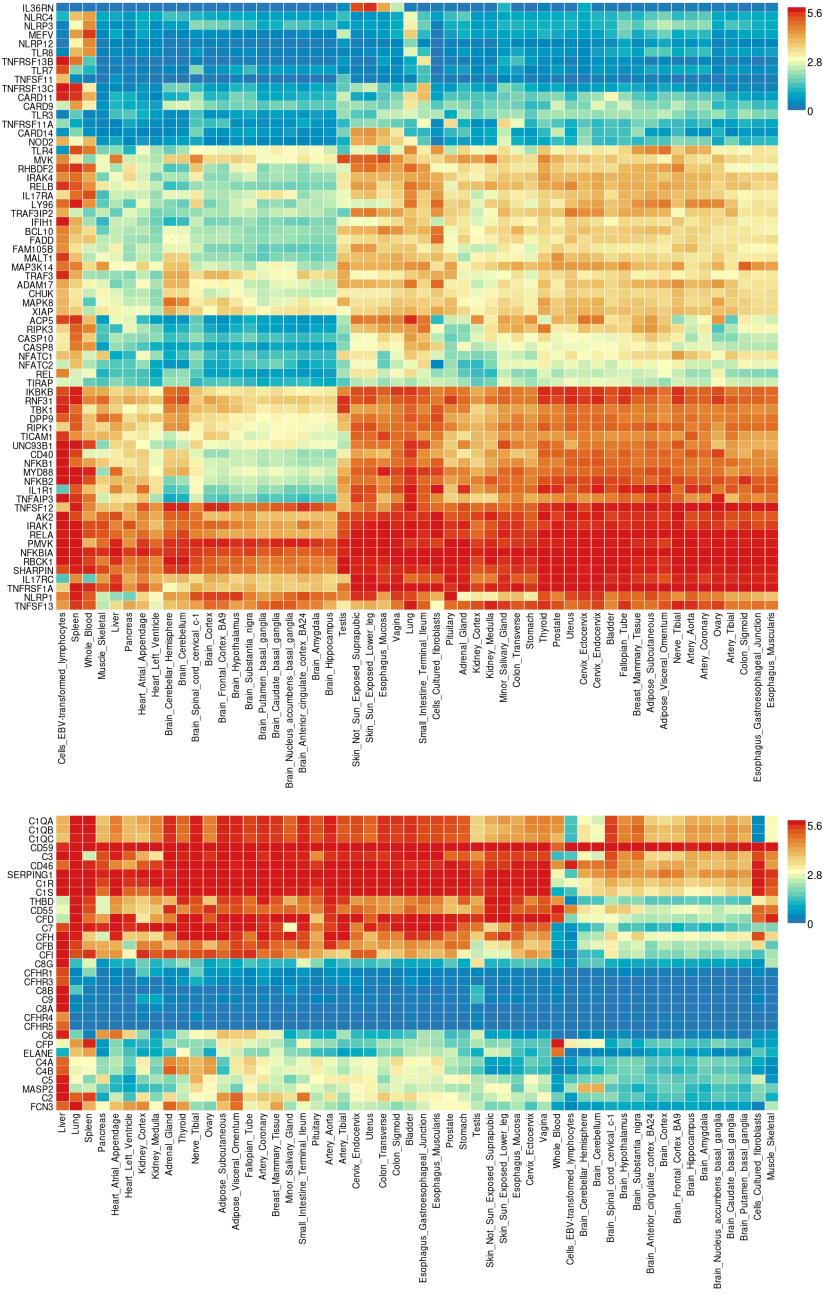


Figure S5: **Gene Expression Heatmaps for IEI Genes.** GTEx v8 data from 54 tissue types display the average expression per tissue label (log₂ transformed) for the IEI gene panels. Top: Cluster 2; Bottom: Cluster 4.