

1
2

1. Extended Data

Figure or Table # Please group Extended Data items by type, in sequential order. Total number of items (Figs. + Tables) must not exceed 10.	Figure/Table title One sentence only	Filename Whole original file name including extension. i.e.: Smith_ED_Fig1.jpg	Figure/Table Legend If you are citing a reference for the first time in these legends, please include all new references in the main text Methods References section, and carry on the numbering from the main References section of the paper. If your paper does not have a Methods section, include all new references at the end of the main Reference list.
Extended Data Fig. 1	Reduction in the number of genotypes stored per sample.	ext-data-fig1.eps	For 100 randomly chosen 100KGP participants belonging to each ancestry group (taken from amongst those with an inferred probability >0.9 of belonging): a , boxplots showing the distribution of the number of non-homozygous reference PASSing genotypes for variants on chromosomes 1–22 and X which meet the default Rareservoir MAF filtering criteria (i.e. a PMAF score >0 using gnomAD v3.0 and internal MAF <0.002); b , boxplots showing the distribution of the proportion of all PASSing non-homozygous reference genotypes that meet the default Rareservoir MAF filtering criteria. In both plots, the lower, centre and upper lines respectively indicate the lower quartile, median and upper quartile. Whiskers are drawn up to the most extreme points that are less than 1.5× the interquartile range away from the nearest quartile.
Extended Data Fig. 2	General schematic of the database build procedure and contents.	ext-data-fig2.eps	Variants are extracted from VCF files, filtered on internal cohort allele frequency, encoded as 64-bit RSVR IDs and loaded into a table containing the corresponding genotypes. The variants are annotated with scores reflecting their predicted deleteriousness (in this case, CADD scores) and probabilistic minor allele frequency scores (PMAF) from gnomAD. The consequences of each variant with respect to a reference set of transcripts are generated and loaded into a table. Sample information including pedigree membership

			and membership of an MSUP is loaded into a table. The case groupings for case/control association analyses are stored in a table.
Extended Data Fig. 3	Detailed schematic of the database build procedure.	ext-data-fig3.eps	Variants may be imported to a Rareservoir from either single gVCF files or a merged VCF file, following the procedures indicated by red and blue arrows respectively.
Extended Data Fig. 4	Schematic showing the variant data in the 100KGP Main Programme Rareservoir.	ext-data-fig4.eps	The number of variant/transcript pairs, the distribution of CADD scores and a breakdown of gnomAD frequency classes is shown for each annotated SO term in the context of the structure of the ontology.
Extended Data Fig. 5	The 269 case sets, Disease Groups A–I.	ext-data-fig5.eps	The names and sizes of the case sets used for the genetic association analyses, grouped by Disease Group and coloured by type (Disease Sub Group or Specific Disease). Disease Sub Groups with only one Specific Disease were excluded to avoid repeating identical analyses. Case sets smaller than 5 are labelled '<5' and shown as having size 4 to comply with 100KGP policy on limiting participant identifiability. For legibility, only Disease Groups starting with the letters A–I are shown here.
Extended Data Fig. 6	The 269 case sets, Disease Groups M–Z.	ext-data-fig6.eps	An extension of Extended Data Fig. 5 showing the case sets in Disease Groups starting with the letters M–Z.
Extended Data Fig. 7	Breakdown of cases attributable to associations with 'Posterior segment abnormalities' by Specific Disease.	ext-data-fig7.eps	For each gene associated with the Disease Sub Group 'Posterior segment abnormalities', a bar plot showing the number of cases having each of the different Specific Diseases who have an inferred pathogenic configuration of alleles in the gene. This example illustrates that sets of cases with the same aetiological gene may be assigned different Specific Diseases. Consequently, pooling cases within Disease Sub Group can boost power.
Extended Data Fig. 8	Microscopy images of	ext-data-fig8.eps	Exemplar immunofluorescence microscopy images of HEK293 cells overexpressing wild

	HEK293 cells overexpressing ERG.		type ERG (from 20 replicates) and each of the p.S182Afs*22, p.T224Rfs*15 and p.A447Cfs*19 variants of ERG (each from 17 replicates). Cells were stained for ERG (green) and nuclear marker DAPI (blue). Scale bar, 20µm.
Extended Data Fig. 9	Illustrative audiograms for GPR156 cases.	ext-data-fig9.eps	Air and bone conduction audiograms for the two affected daughters of the family with compound heterozygous <i>GPR156</i> truncating alleles.

3

2. Supplementary Information:

4

A. PDF Files

5

6

Item	Present?	Filename Whole original file name including extension. i.e.: Smith_SI.pdf. The extension must be .pdf	A brief, numerical description of file contents. i.e.: <i>Supplementary Figures 1-4, Supplementary Discussion, and Supplementary Tables 1-4.</i>
Supplementary Information	No		
Reporting Summary	Yes.	nr-reporting-summary-comments-addressed.pdf	
Peer Review Information	No	OFFICE USE ONLY	

7

3. Source Data

8

9

Parent Figure or Table	Filename Whole original file name including extension. i.e.: <i>Smith_SourceData_Fig1.xls</i> , or <i>Smith_Unmodified_Gels_Fig1.pdf</i>	Data description i.e.: Unprocessed western Blots and/or gels, Statistical Source Data, etc.
Source Data Fig. 1	source-data-fig1b.xlsx	Sheet 1 , Table of associations shown in Fig. 1b annotated with BeviMed PPAs (PPA), the level of the case set in the disease label hierarchy (Level), the inferred variant class and MOI for the association, the matched

		PanelApp panel for the association, the method that was used to find the match (Match method, either 'Automatic' or 'Manual'), the associated evidence level for the match, and the notes on the consistency between the MOI listed by PanelApp for the association and the inferred MOI (MOI match comment). Sheet 2 , Table of variants having a probability of pathogenicity >0.8 conditional on the modal model and forming a pathogenic configuration of alleles in at least one case. While these variants contributed to the reported statistical associations, they have not been individually scrutinised according to ACMG guidelines.
Source Data Fig. 2	source-data-fig2e.jpg	Uncropped western blot images corresponding to Fig. 2e .
Source Data Fig. 3	source-data-fig4e.jpg	Uncropped western blot images corresponding to Fig. 4e .

10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

Genetic association analysis of 77,539 genomes reveals rare disease etiologies

Daniel Greene^{1,2}, Genomics England Research Consortium, Daniela Pirri³, Karen Frudd^{3,4}, Ege Sackey⁵, Mohammed Al-Owain⁶, Arnaud P.J. Giese⁷, Khushnooda Ramzan⁸, Sehar Riaz^{7,9}, Itaru Yamanaka¹⁰, Nele Boeckx¹¹, Chantal Thys¹², Bruce D. Gelb^{2,13,14}, Paul Brennan¹⁵, Verity Hartill^{16,17}, Julie Harvengt¹⁸, Tomoki Koshi^{19,20}, Sahar Mansour^{5,21}, Mitsuo Masuno²², Takako Ohata²³, Helen Stewart²⁴, Khalid Taibah²⁵, Claire L.S. Turner²⁶, Faiqa Imtiaz⁸, Saima Riazuddin^{7,9}, Takayuki Morisaki^{10,27}, Pia Ostergaard⁵, Bart L. Loeys^{11,28}, Hiroko Morisaki^{10,29}, Zubair M. Ahmed^{7,9}, Graeme M. Birdsey³, Kathleen Freson¹², Andrew Mumford^{30,31} & Ernest Turro^{2,14,32,33*}

*Corresponding author. Contact: ernest.turro@mssm.edu

¹Department of Medicine, University of Cambridge, Cambridge, UK. ²Mindich Child Health and Development Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA. ³National Heart and Lung Institute, Imperial College London, London, UK. ⁴UCL Institute of Ophthalmology, University College London, London, UK. ⁵Molecular and Clinical Sciences Institute, St George's University of London, London, UK. ⁶Department of Medical Genomics, Centre for Genomic Medicine, King Faisal Specialist Hospital & Research Centre, Riyadh, Saudi Arabia. ⁷Department of Otorhinolaryngology Head and Neck Surgery, School of Medicine, University of Maryland Baltimore, Baltimore, MD, USA. ⁸Department of Clinical Genomics, Centre for Genomic Medicine, King Faisal Specialist Hospital & Research Centre, Riyadh, Saudi Arabia. ⁹Department of Biochemistry and Molecular Biology, School of Medicine, University of Maryland, Baltimore, MD, USA. ¹⁰Department of Bioscience and Genetics, National Cerebral and Cardiovascular Center, Osaka, Japan. ¹¹Center for Medical Genetics, Antwerp University Hospital/University of Antwerp, Antwerp, Belgium. ¹²Department of Cardiovascular Sciences, Center for Molecular and Vascular Biology, KU Leuven, Leuven, Belgium. ¹³Department of Pediatrics, Icahn School of Medicine at Mount Sinai, New York, NY,

38 USA. ¹⁴Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
39 ¹⁵Northern Genetics Service, Newcastle upon Tyne Hospitals NHS Trust, International Centre for Life, Newcastle upon Tyne,
40 UK. ¹⁶Department of Clinical Genetics, Chapel Allerton Hospital, Leeds Teaching Hospitals NHS Trust, Leeds, UK. ¹⁷Leeds
41 Institute of Medical Research, University of Leeds, Leeds, UK. ¹⁸Center for Medical Genetics, CHU Liege, Belgium.
42 ¹⁹Department of Medical Genetics, Shinshu University School of Medicine, Matsumoto, Japan. ²⁰Center for Medical Genetics,
43 Shinshu University Hospital. ²¹SW Thames Regional Genetics Service, St George's University Hospitals NHS Foundation Trust,
44 London, UK. ²²Department of Medical Genetics, Kawasaki Medical School Hospital, Okayama, Japan. ²³Okinawa Chubu
45 Hospital, Okinawa, Japan. ²⁴Oxford University Hospitals NHS Foundation Trust, Oxford, UK. ²⁵ENT Medical Centre, Riyadh,
46 Saudi Arabia. ²⁶Peninsula Clinical Genetics Service, Royal Devon & Exeter Hospital, Exeter, UK. ²⁷Division of Molecular
47 Pathology and Department of Internal Medicine, Institute of Medical Science, The University of Tokyo, Tokyo, Japan.
48 ²⁸Department of Human Genetics, Radboud University Medical Center, Nijmegen, The Netherlands. ²⁹Department of Medical
49 Genetics, Sakakibara Heart Institute, Tokyo, Japan. ³⁰School of Cellular and Molecular Medicine, University of Bristol, Bristol,
50 UK. ³¹South West NHS Genomic Medicine Service Alliance, Bristol, UK. ³²Department of Haematology, University of
51 Cambridge, Cambridge Biomedical Campus, Cambridge, UK. ³³Charles Bronfman Institute for Personalized Medicine, Icahn
52 School of Medicine at Mount Sinai, New York, NY, USA.
53

54 **Abstract**

55 **The genetic aetiologies of more than half of rare diseases remain unknown. Standardised**
56 **genome sequencing (GS) and phenotyping of large patient cohorts provides an opportunity for**
57 **discovering the unknown aetiologies, but this depends on efficient and powerful analytical**
58 **methods. We built a compact database, the 'Rareservoir,' containing the rare variant genotypes**
59 **and phenotypes of 77,539 participants sequenced by the 100,000 Genomes Project (100KGP). We**
60 **then used the Bayesian genetic association method BeviMed to infer associations between genes**
61 **and each of 269 rare disease classes assigned by clinicians to the participants. We identified 241**
62 **known and 19 previously unidentified associations. We validated associations with *ERG*,**
63 ***PMEPA1* and *GPR156* by searching for pedigrees in other cohorts and using bioinformatic and**
64 **experimental approaches. We provide evidence that (1) loss-of-function variants in the ETS-**
65 **family transcription factor encoding gene *ERG* lead to primary lymphoedema, (2) truncating**
66 **variants in the last exon of TGF β regulator *PMEPA1* result in Loews-Dietz syndrome, and (3) loss-**
67 **of-function variants in *GPR156* give rise to recessive congenital hearing impairment. The**
68 **Rareservoir provides a lightweight, flexible and portable system for synthesising the genetic and**
69 **phenotypic data required to study rare disease cohorts with tens of thousands of participants.**
70

71 **Introduction**

72 Collectively, rare diseases affect 1 in 20 people¹, but fewer than half of the approximately 10,000
73 catalogued rare diseases have a resolved genetic aetiology². Standardised GS of large, phenotypically
74 diverse collections of rare disease patients enables aetiological discovery across a wide range of
75 pathologies^{3,4,5} while boosting genetic diagnostic rates for patients. The 100KGP, the largest GS study of
76 rare disease patients to date, sequenced 34,523 United Kingdom National Health Service patients with
77 rare diseases and 43,016 of their unaffected relatives. The linked genetic and phenotypic data of
78 100KGP participants were then made available to researchers through a web portal called the Genomics
79 England Research Environment. The scale and complexity of such large GS datasets and the
80 hierarchical nature of patient phenotype coding⁶ induce numerous bioinformatic and statistical
81 challenges. Most importantly, the full genotype data from GS studies of tens of thousands of individuals
82 are typically stored in unmodifiable files many terabytes in size, leading to high storage and processing
83 costs. Recently developed frameworks such as Hail⁷ and OpenCGA⁸ afford greater flexibility. However,
84 they are designed to capture genotypes for variants across the full minor allele frequency (MAF)

85 spectrum, from rare (MAF<0.1%) to common (MAF>5%) variants. To accommodate large numbers of
86 genotypes, they depend on distributed storage systems and require numerous software packages,
87 hindering deployment. We developed a database schema, the 'Rareservoir,' for working with rare variant
88 genotypes and patient phenotypes flexibly and efficiently. We deployed a Rareservoir only 5.5GB in size
89 of the 100KGP data and applied the Bayesian statistical method BeviMed⁹ to identify genetic
90 associations between coding genes and each of the 269 rare disease classes assigned to patients by
91 clinicians. Out of the previously unknown associations that we identified, we followed up the most
92 plausible subset in confirmatory analytical and experimental work.

93

94 **Results**

95 **The 'Rareservoir'**

96 Relational databases (RDBs) provide a unified, centralised structure for storing, querying and modifying
97 data of multiple underlying types. In principle, an RDB could provide a convenient foundation for the
98 analysis of genotypes, variants, genes, participants and statistical results, but they cannot accommodate
99 tables of the scale required to store exome or genome-wide genotypes in a moderately sized cohort. An
100 RDB can, however, accommodate a sparse representation of genotypes corresponding to rare variants
101 only, which encompass almost all variants having a large effect on rare disease risk. We developed an
102 RDB schema, the 'Rareservoir', and build procedure for the analysis of rare diseases, which, by default,
103 stores genotypes corresponding to variants for which all population-specific MAFs are likely to be <0.1%.
104 This reduces the number of stored genotypes in a large study by about 99% (**Extended Data Fig. 1**).
105 The Rareservoir encodes variants as 64-bit integers ('RSVR IDs', **Extended Data Fig. 2**), which can
106 represent 99.3% of variants encountered in practice without loss of information. RSVR IDs occupy a
107 single column and increase numerically with respect to genomic position, allowing fast location-based
108 queries within a simple database structure. To support the build process of a Rareservoir, we developed
109 a complementary software package called 'rsvr' (**Extended Data Figs. 2–3**). The package includes tools
110 to annotate variants with MAF information from control databases (e.g. gnomAD¹⁰), pathogenicity scores
111 (e.g. combined annotation dependent depletion (CADD) scores¹¹) and predicted Sequence Ontology
112 (SO)¹² consequences with respect to a set of transcripts. We use a 64-bit integer ('CSQ ID') to record the
113 consequences for interacting variant/transcript pairs, where each bit encodes one of the possible
114 consequences, ordered by severity. Encoding the consequences in this way is efficient and enables
115 succinct queries that threshold or sort based on severity of impact. The Rareservoir also contains a table
116 with genetically derived data for each sample (including ancestry, sex and membership of a maximal set
117 of unrelated participants (MSUP)), and a table of 'case sets' storing the rare disease classes assigned to
118 each participant.

119

120 **BeviMed infers 241 known and 19 unknown genetic associations**

121 We built a Rareservoir, 5.5GB in size, containing 11.9 million rare exonic and splicing single nucleotide
122 variants (SNVs) and short insertions or deletions (indels) affecting canonical transcripts of protein-coding
123 genes in Ensembl v104¹³ from a merged variant call format file (VCF) containing genotype calls for
124 77,539 participants, including 29,741 probands, in the Rare Diseases Main Programme of the 100KGP
125 (Data Release Version 13) (**Extended Data Fig. 4**). During enrolment to the 100KGP, expert clinicians
126 used the clinical characteristics of each affected participant to assign them to one or more of 220 *Specific*
127 *Diseases*. The Specific Diseases are hierarchically arranged into 88 *Disease Sub Groups*, each of which
128 belongs to one of 20 *Disease Groups*. Whereas the eligibility criteria for many specific diseases aligned

129 to the same or closely related rare diseases, for others such as 'intellectual disability,' the criteria were
130 broader and encompassed diverse genetic aetiologies. We generated 269 analytical case sets
131 corresponding to all distinct Specific Diseases and Disease Sub Groups, ranging in size from 5,809 to
132 one proband, and stored them in the Rareservoir (**Fig. 1a, Extended Data Figs. 5–6**). We included
133 these two levels of the phenotyping hierarchy to account for heterogeneity in presentation or diagnosis
134 among cases sharing the same genetic aetiology, with the aim of boosting power to identify statistical
135 genetic associations.

136
137 Using the Bayesian statistical method BeviMed⁹, we obtained a posterior probability of association (PPA)
138 between each of the 19,663 protein-coding genes and each of the 269 rare disease classes. BeviMed
139 computes posterior probabilities over a baseline model of no association and six competing association
140 models, each of which assumes a particular mode of inheritance (MOI, dominant or recessive) and
141 consequence class of aetiological variant (high-impact, moderate-impact or 5' UTR). The PPA is
142 obtained by summing the posterior probabilities over all six association models. The association model
143 with the greatest posterior probability (the modal model) determines the inferred MOI and class of
144 aetiological variant. Conditional on an association model, BeviMed models the pathogenicity of each
145 included rare variant. In the model, participants with at least one pathogenic allele (under a dominant
146 MOI), or at least as many pathogenic alleles as the ploidy (under a recessive MOI), have a *pathogenic*
147 *configuration of alleles*, which determines their risk of case status. For each rare disease class, we
148 selected a set of unrelated cases based on pedigree information provided by the 100KGP and compared
149 them to participants not in the case set who belonged to different pedigrees and to an MSUP, also
150 provided by the 100KGP. To account for correlation between case sets, we only recorded the association
151 for each gene having the highest PPA within a given Disease Group. Using a significance threshold of
152 $PPA > 0.95$, we identified 260 significant associations, 241 of which were documented by the PanelApp
153 gene panel database,¹⁴ an expert-curated and annotated resource containing gene lists with high,
154 medium or low levels of prior supporting evidence of causality for rare diseases (**Fig. 1b**). Out of the 241
155 known associations that we identified, 43 (17.8%) were with Disease Sub Groups. For example, within
156 each of the nine known genes associated with the Disease Sub Group 'Posterior segment abnormalities,'
157 the set of cases explained by variants with a posterior probability of pathogenicity > 0.8 comprised a
158 mixture of participants with five different Specific Diseases (**Extended Data Fig. 7**). This demonstrates
159 that participants with different Specific Diseases belonging to the same Disease Sub Group sometimes
160 share defects in the same gene, which confirms that treating Disease Sub Groups, not just Specific
161 Diseases, as case sets, boosts statistical power.

162
163 Out of the 241 associations identified as previously known according to PanelApp, 237 (98.3%) had an
164 inferred MOI that was consistent with the MOIs listed for the relevant gene. Of these, the consistent MOI
165 was found in the matched panel (223 associations), in the notes for the matched panel (five associations)
166 or in the MOIs listed for an alternative relevant panel (nine associations) in PanelApp (**Source data for**
167 **Fig. 1**). This provided independent evidence that the genetic associations we labelled as known (without
168 reference to MOI information) are genuinely supported by evidence in the literature, further
169 demonstrating the accuracy of BeviMed's inference. Of the four known associations with an inferred MOI
170 that was incongruous with PanelApp, two had supporting evidence for the inferred MOI in the literature
171 that was absent from PanelApp: *EDA* with dominant 'Ectodermal dysplasia without a known gene
172 mutation'¹⁵ and *AICDA* with dominant 'Primary immunodeficiency'¹⁶. The two associations with an MOI

173 that was unsupported in the literature were between *UCHL1* and dominant 'Inherited optic neuropathies'
174 and between *SLC39A8* and dominant 'Intellectual disability'.
175

176 Among 5,253 of the probands included in our analysis, the table of clinically reported variants available
177 from the 100KGP Rare Diseases Main Programme at the time of this study comprised 4,907 distinct
178 variants that had been classified as pathogenic or likely pathogenic in 1,863 genes. For 855 of these
179 genes, aetiological variants had been reported for only one family, suggesting that many genes which
180 are aetiological in the 100KGP are not identifiable by statistical association. Nevertheless, across the 260
181 associations identified by BeviMed, 2,536 distinct rare variants had a posterior probability of
182 pathogenicity >0.8 conditional on the modal model and were observed as part of a pathogenic
183 configuration of alleles in a case (**Source data for Fig. 1**). Interestingly, among the subset of 2,485
184 variants contributing to the 241 known associations, only 1,604 featured in the table of clinically reported
185 variants.
186

187 We found 19 previously unidentified genetic associations. To select a shortlist for further investigation,
188 we assigned a plausibility score (range 0–3) based on three sources of additional evidence (**Table 1**).
189 Firstly, we considered evidence of purifying selection from gnomAD v2.1.1. Any dominant associations
190 with high-impact variants in a gene having a probability of loss-of-function intolerance (pLI) >0.9 or with
191 moderate-impact variants in a gene having a Z-score >2 were considered to be supported by population
192 genetic metrics of purifying selection. To avoid disadvantaging recessive associations, which are unlikely
193 to leave a detectable signature of purifying selection in gnomAD even if genuine, they were considered
194 to be supported by default. Secondly, we considered co-segregation data: any association for which
195 variants having a posterior probability of pathogenicity conditional on the modal model >0.8 tracked with
196 case status in at least three additional family members and for which no affected relatives lacked the
197 pertinent variants were considered to be supported by co-segregation. Thirdly, we performed a
198 comprehensive review of the literature for each gene and made a subjective assessment of whether an
199 association was supported by biological function or previously known disease associations for related
200 genes. In total, three genetic associations had a plausibility score of 3 and were therefore investigated
201 further by gathering additional experimental evidence and looking for replication in other sequenced rare
202 disease collections.
203

204 **Variants in *ERG* are responsible for primary lymphoedema**

205 BeviMed identified a dominant genetic association between high-impact variants in *ERG* and the Specific
206 Disease 'Primary lymphoedema,' a group of genetic conditions caused by abnormal development of
207 lymphatic vessels or failure of lymphatic function^{17,18}. Three such variants were responsible for the high
208 PPA, with locations ranging from codon 182 to 463 on the canonical Ensembl transcript
209 ENST00000288319.12. One of the probands had two unaffected parents without the variant allele—one
210 sequenced by the 100KGP and the other by Sanger sequencing—suggesting the truncating
211 heterozygous variant had appeared *de novo*. A participant in a fourth family who had been enrolled to
212 the 100KGP for an unrelated condition also carried a predicted loss-of-function variant in *ERG*. Upon
213 manual chart review, this participant had features associated with this unrelated condition, but additional
214 features consistent with primary lymphoedema, providing internal replication within the discovery cohort
215 (**Fig. 2a**).
216

217 The affected father of the proband with the variant encoding p.S182Afs*22 was called homozygous for
218 the reference allele, initially suggesting a lack of co-segregation of the variant with the disease in that
219 pedigree. However, a review of the GS read alignments for the father revealed that two out of the 48
220 reads overlapping that position supported the alternative allele. Specifically, these reads contained a
221 deletion of a single G within the central poly-G tract of the motif "AGCTGGGGGTGAG." To assess
222 whether this could be the result of erroneous sequencing, we counted the number of such reads in the
223 77,539 genomes in the 100KGP and found that the proband and the father were the only two with more
224 than one such read. This indicated that these reads in the father were unlikely to be erroneous but
225 instead that he was mosaic (**Fig. 2b**), consistent with the observation that his lymphoedema became
226 clinically apparent over two decades later than his daughter, indicating milder disease. A further 130
227 samples collected through the 100KGP had a single read containing the deletion. This number was
228 consistent with observations in the 80 other exonic loci that contain the same 13bp motif (mean: 99.67
229 samples, range: 4 to 149 samples), suggesting that, rather than being mosaic, the 130 samples
230 contained individual sequencing errors. Furthermore, none of the participants who gave these samples
231 had been assigned the Specific Disease 'Primary lymphoedema.'

232
233 *ERG* encodes a critical transcriptional regulator of blood vessel endothelial cell (EC) gene expression¹⁹
234 that is essential for normal vascular development²⁰. However, little is known about the contribution of
235 *ERG* to lymphatic development or how primary lymphoedema could arise from loss-of-function *ERG*
236 variants which affect different parts of the *ERG* protein (**Fig. 2c**). Total cellular expression of *ERG*
237 detected by real-time quantitative polymerase chain reaction (PCR) in purified RNA and by
238 immunoblotting of protein extracts was the same in primary human dermal lymphatic EC (HDLEC) as
239 human umbilical vein EC (HUVEC) (**Fig. 2d** and **Fig. 2e** respectively). Moreover, immunofluorescence
240 microscopy of cultured HDLEC showed that *ERG* expression co-localised with the lymphatic EC nuclear
241 marker PROX1 (**Fig. 2f**) a finding confirmed *in vivo* by immunostaining whole mounts of ear skin from
242 mice at three weeks after birth (**Fig. 2g**). The positions of the p.S182Afs*22 and p.T224Rfs*15 variants
243 suggest nonsense mediated decay and haploinsufficiency as a possible disease mechanism. The other
244 two variants, however, are located in the final exon of *ERG* and may therefore evade nonsense mediated
245 decay. We studied both types of variant in more detail to explore potential disease mechanisms. In
246 HEK293 cells, which do not express endogenous *ERG*, overexpression of wild type *ERG* cDNA
247 recapitulated the nuclear expression pattern observed in the HDLEC and mouse ear skin models.
248 However, overexpression of each of the *ERG* mutant cDNAs resulted in mislocalisation of *ERG* outside
249 of the nucleus, in the cytosol (**Fig. 2h–i, Extended Data Fig. 8**), preventing it from binding to DNA and
250 exerting its function as a transcription factor²¹. Together, these data confirm high levels of *ERG*
251 expression within the nuclei of the lymphatic endothelium consistent with a transcription regulatory
252 function during lymphangiogenesis. They also suggest that in the primary lymphoedema cases, defective
253 lymphangiogenesis may result from reduced *ERG* availability in the nucleus either because of
254 haploinsufficiency resulting from nonsense mediated decay or mislocalisation.

255 256 **Variants in *PMEPA1* result in Loeys-Dietz syndrome**

257 BeviMed identified a dominant genetic association between high-impact variants in *PMEPA1* and the
258 Specific Disease 'Familial Thoracic Aortic Aneurysm Disease' (FTAAD). The variant with the highest
259 conditional probability of pathogenicity was an insertion of one cytosine within a seven-cytosine stretch in
260 the last exon of the canonical Ensembl transcript ENST00000341744.8. This variant, which is predicted

261 to induce a p.S209Qfs*3 frameshift, was observed in three FTAAD pedigrees of European ancestry in
262 the 100KGP discovery cohort. We replicated the association in three additional collections of cases.
263 Firstly, the same variant was identified independently in eight affected members of three pedigrees of
264 Japanese ancestry in a separate Japanese patient group. Secondly, a single-cytosine deletion within the
265 same poly-cytosine stretch as the previous variant, and encoding p.S209Afs*61, was found in an FTAAD
266 case enrolled in a separate collection of 2,793 participants in the 100KGP Pilot Programme. Lastly, we
267 identified a family in Belgium wherein the affected members carried a five base-pair deletion in the same
268 stretch of poly-cytosines inducing a frameshift two residues upstream of the other two variants
269 (p.P207Qfs*3).

270
271 All pedigrees exhibited dominant inheritance of aortic aneurysm disease with incomplete penetrance and
272 skeletal features including pectus deformity, scoliosis and arachnodactyly with complete penetrance,
273 which co-segregated with the respective variants in genotyped participants (**Fig. 3a**). To assess whether
274 *PMEPA1* families affected by FTAAD form a phenotypically distinct subgroup, we analysed the HPO
275 terms assigned to the 593 FTAAD families in both programmes of the 100KGP. Using a permutation-
276 based method^{22,23} based on Resnik's semantic similarity measure²⁴, we found that the four 100KGP
277 *PMEPA1* families were significantly more similar to each other than to other FTAAD families chosen at
278 random ($p=5.7 \times 10^{-3}$). To characterise the *PMEPA1* phenotype in greater detail, we compared the
279 prevalence of each of the HPO terms in the minimal set of terms present in at least three of the four
280 families with the prevalence in the other FTAAD families. We identified four HPO terms related to the
281 musculoskeletal system that were significantly enriched (**Fig. 3b**), echoing the phenotypic characteristics
282 of the syndromic aortopathy Loeys-Dietz syndrome^{25,26}.

283
284 To understand the molecular mechanisms underlying this defect, we examined the protein-protein
285 interactions²⁷ for *PMEPA1* and the complete set of high-confidence genes in the 'Thoracic aortic
286 aneurysm or dissection' PanelApp panel. *PMEPA1* encodes a negative regulator of Transforming Growth
287 Factor β (TGF β) signalling²⁸, a pathway previously implicated in multiple aortopathies, including Loeys-
288 Dietz syndrome²⁹. The genes underlying known forms of Loeys-Dietz syndrome encode part of a tightly
289 interacting subgroup of proteins in the TGF β pathway, in which there is a direct interaction between the
290 proteins encoded by *SMAD2*, *SMAD3* and *PMEPA1* (**Fig. 3c**). As the two candidate variants occur in the
291 last exon of the transcript, they are likely to evade nonsense-mediated decay³⁰. However, their truncating
292 effects are predicted to remove a PPxY interaction motif, while leaving the SMAD interaction motif intact
293 (**Fig. 3d**), possibly affecting binding between *PMEPA1* and *SMAD2/3*, and altering TGF β signalling
294 through a gain-of-function mechanism.

295 296 **Variants in *GPR156* lead to recessive congenital hearing loss**

297 BeviMed identified a recessive genetic association between high-impact variants in *GPR156* and the
298 Specific Disease 'Congenital hearing impairment'. Two high-impact variants in *GPR156* were responsible
299 for the strong evidence of association: a one base pair deletion predicting p.S207Vfs*113 and a one
300 base pair insertion predicting p.P718Lfs*86 with respect to the canonical Ensembl transcript
301 ENST00000464295.6. One family contained two affected siblings who were both homozygous for the
302 p.S207Vfs*113 variant inherited from heterozygous parents. In a second family, there were also two
303 affected siblings, in this case compound heterozygous for the same p.S207Vfs*113 variant that was
304 maternally inherited and a different p.P718Lfs*86 variant that was paternal. Using GeneMatcher³¹, we

305 identified a third pedigree from Saudi Arabia with biallelic truncating variants in *GPR156*. This
306 consanguineous pedigree contained four siblings with hearing impairment, all of whom were
307 homozygous for a variant predicting p.S642Afs*162 (**Fig. 4a**). The eight affected individuals in these
308 three families all had congenital non-syndromic bilateral sensorineural hearing loss (see **Extended Data**
309 **Fig. 9** for illustrative audiograms).

310
311 *GPR156* encodes probable G-protein coupled receptor 156, which has sequence homology to the class
312 C GABAB receptors³². Although previously designated as an orphan receptor, *GPR156* has recently
313 been identified as a critical regulator of stereocilia orientation on hair cells of the auditory epithelium and
314 other mechanosensory tissues³³. Its expression is highly restricted to hair cells in the inner ear³⁴.
315 Disruption of stereocilia is a common pathogenic mechanism underlying many human Mendelian hearing
316 loss disorders³⁵ and the over-expression of *GPR156* in hair cells relative to surrounding cells was
317 commensurate with the over-expression of the 21 genes currently implicated in hearing impairment
318 having a Gene Ontology (GO) term relating to stereocilia (**Fig. 4b**). By immunostaining of the Corti and
319 vestibular system from wild type mice, we found that GPR156 strongly co-localises with actin at the
320 apical surface of the outer and inner hair cells of the organ of Corti (**Fig. 4c**).

321
322 The p.S207Vfs*113 variant is located in the sixth of 10 exons of *GPR156* and therefore predicts absent
323 expression through nonsense mediated decay of the *GPR156* mRNA. In contrast, the p.S642Afs*162
324 and p.P718Lfs*86 variants both occur within the final *GPR156* exon and likely result in expression of
325 abnormal GPR156 with an altered amino acid sequence and premature truncation of the cytoplasmic tail
326 (**Fig. 4d**). To determine the effect of the variants on protein expression, we transfected Cos7 cells, which
327 do not express *GPR156* endogenously, with constructs containing cDNAs for wild type *GPR156* or
328 *GPR156* containing each of the three mutant alleles, tagged with a green fluorescent protein (GFP)
329 reporter. While cells transfected with wild type sequence expressed GPR156-GFP fusion protein
330 robustly, cells transfected with the mutant constructs either did not express the protein appreciably or
331 exhibited markedly reduced expression, suggesting that all three of the truncated proteins are degraded
332 (**Fig. 4e**). These data suggest that the biallelic chain truncating variants in *GPR156* cause a congenital
333 hearing loss by preventing expression of GPR156 protein, thereby disrupting stereocilia formation in the
334 auditory epithelium.

335 336 **Discussion**

337 The standardisation of GS within a healthcare system, together with powerful frameworks for genetic and
338 phenotypic data processing and statistical analysis, promises to advance the resolution of the remaining
339 unknown aetiologies of rare diseases. We have developed a lightweight and easily deployable relational
340 database, the Rareservoir, for genetic analysis of rare diseases using approaches such as BeviMed. In
341 one unified analysis, we identified 260 associations, of which 241 had been published previously in a
342 body of work spanning several decades of genetics research. Our results give an upper bound on the
343 false discovery rate (FDR) of 7.3%. In contrast, a recent analysis of 57,000 samples in the 100KGP
344 reported 249 known and 579 previously unidentified associations³⁶, giving an upper bound on the FDR of
345 70%, which suggests that our analytical approach has a greater specificity for a given sensitivity. The
346 associations spanned 86 disease classes across a wide range of organ systems. Interestingly, only 64%
347 of the variants contributing substantially to the known associations were present in the table of clinically
348 reported variants available at the time of this study. This suggests that, as cohorts grow larger, the

349 results of statistical inference could help guide the clinical reporting process. The case sets we used in
350 our genetic association analysis were based on the formal disease classifications used by the 100KGP.
351 Some of the case sets, such as 'Intellectual disability' (5,529 probands), are particularly large and likely
352 to be highly genetically heterogeneous, potentially limiting the power of our analyses. Careful partitioning
353 of heterogeneous case sets using individual-level HPO terms⁶ has the potential to boost power. Of the
354 19 previously unidentified associations, we shortlisted, replicated and validated three. These three
355 aetiologies involve genes that had not previously been implicated in any of these human diseases. The
356 remaining 16 associations include further plausible hypotheses. For example, *LRRC7*, which we
357 identified to be associated with intellectual disability, encodes a brain-specific protein in post-synaptic
358 densities³⁷, and *Lrrc7*-deficient mice exhibit a neuro-behavioural phenotype³⁸. *USP33*, which we found to
359 be associated with early-onset hypertension, encodes a deubiquitinating enzyme implicated in regulating
360 expression of the β 2-adrenergic receptor regulation³⁹. These and other candidates will require replication
361 and validation before they can be considered causative genes.

362
363 The present study has several limitations. Firstly, approximately 82% of the participants in the 100KGP
364 are of European ancestry. While this percentage is in line with the proportion of residents in England and
365 Wales reporting their ethnic group as White in the 2011 UK census (86%), its large magnitude constrains
366 power to identify causative variants specific to other ancestry groups. Secondly, of the 260 case sets
367 analysed, 28 contained fewer than 5 probands, limiting power to identify the causes of the corresponding
368 disease classes and highlighting the need for continued enrolment of patients with ultra-rare disorders.
369 Thirdly, we have only considered SNVs and indels in coding genes. The exploration of rare variation in
370 non-coding genes and in regulatory elements of the genome may help identify further etiologies. Lastly,
371 we focused our attention on monogenic models of rare disorders, even though the genetic etiologies of
372 certain rare diseases may be polygenic. In addition, important variation in clinical presentation of
373 monogenic disorders may be explained by polygenic effects. These limitations point towards multiple
374 promising avenues of future research to uncover the remaining unknown genetic determinants of rare
375 diseases.

376 **Acknowledgements**

378 This research was made possible through access to the data and findings generated by the 100,000
379 Genomes Project. The 100,000 Genomes Project is managed by Genomics England Limited (a wholly
380 owned company of the Department of Health and Social Care). The 100,000 Genomes Project is funded
381 by the National Institute for Health Research and NHS England. The Wellcome Trust, Cancer Research
382 UK and the Medical Research Council have also funded research infrastructure. The 100,000 Genomes
383 Project uses data provided by patients and collected by the National Health Service as part of their care
384 and support. GS was performed by Illumina at Illumina Laboratory Services and was overseen by
385 Genomics England. We thank all NHS clinicians who have contributed clinical phenotype data to the
386 100,000 Genomes rare diseases programme, and all staff at Genomics England who have contributed to
387 the sequencing, maintenance of the research environment and assembly of the standard bioinformatic
388 files that were required for our analyses. We thank the participants of the rare diseases program who
389 made this research possible. We are grateful to Vaughan Keeley for providing access to paternal DNA
390 (*ERG*), Frances Elmslie for inviting a patient to the clinic (*ERG*), and Thomas Jaworek for technical
391 assistance (*GPR156*). D. Greene was supported by the Cambridge BHF Centre of Research Excellence
392 [RE/18/1/34212] and Wellcome Collaborative Award 219506/Z/19/Z. V. Hartill was supported by

393 MRC/NIHR Clinical Academic Research Partnership MR/V037617/1. G. Birdsey and K. Frudd were
 394 funded by BHF project grant PG/17/33/32990. G. Birdsey and D. Pirri were funded by BHF project grant
 395 PG/20/16/35047. E. Sackey was supported by Swiss Federal National Fund for Scientific Research
 396 n°CRSII5_177191/1. P. Ostergaard and S. Mansour were supported by Medical Research Council grant
 397 MR/P011543/1 and BHF grant RG/17/7/33217. K. Freson was supported by KU Leuven BOF grant
 398 C14/19/096 and FWO grant G072921N. Work at the University of Maryland Baltimore was supported by
 399 the NIDCD/NIH grant R01DC016295 to Z. Ahmed. M. Al-Owain, K. Ramzan and F. Imtiaz were
 400 supported by King Salman Center for Disability Research grant 85722. E. Turro was supported by the
 401 Mindich Child Health and Development Institute, the Charles Bronfman Institute for Personalized
 402 Medicine and the Lowy Foundation USA.

403
 404 **Author contributions**

405 D. Greene developed software, conducted analyses and co-wrote the paper. G.E.R.C. provided genetic
 406 and phenotypic data and access to the Genomics England Research Environment. C. Thys performed
 407 experiments and interpreted results. B. D. Gelb provided biological interpretation and feedback on the
 408 manuscript. K. Freson designed and supervised experiments, provided biological interpretation and
 409 contributed to writing the paper. A. Mumford provided clinical oversight, provided biological interpretation
 410 and contributed to writing the paper. E. Turro oversaw the study and co-wrote the paper. The following
 411 contributions relate to the three gene-specific vignettes. *ERG*: D. Pirri, K. Frudd and E. Sackey
 412 performed experiments and interpreted results. S. Mansour and C. L. S. Turner provided additional
 413 clinical information. P. Ostergaard coordinated validation and contributed to writing the paper. G. Birdsey
 414 designed and supervised experiments and contributed to writing the paper. *PMEPA1*: I. Yamanaka and
 415 N. Boeckx conducted experiments and interpreted results. P. Brennan, V. Hartill, J. Harvengt, T. Kosho,
 416 M. Masuno and T. Ohata provided clinical information. T. Morisaki and B. Loeys oversaw clinical and
 417 experimental studies. H. Morisaki recruited the Japanese cases, conducted experiments, interpreted and
 418 analysed results, and oversaw genetic studies. *GPR156*: H. Stewart provided additional clinical
 419 information for the compound heterozygous family. K. Taibah clinically evaluated and recruited the
 420 p.S642Afs*162 family. A. Giese, K. Ramzan and S. Riaz conducted experiments and interpreted results.
 421 M. Al-Owain assisted with experiments, interpreted results and contributed clinical information. S.
 422 Riazuddin, F. Imtiaz and Z. M. Ahmed designed and supervised experiments, analysed results, and
 423 provided reagents and tools.

424
 425 **Competing interests** No authors have competing interests.
 426
 427

Gene	Case set	Level	Cases	Controls	Variant class	MOI	pLI	Z	Co-segregation evidence	Biological function and existing disease associations	Score	Replication
<i>ERG</i>	Primary lymphoedema	SD	94	55,400	High	Dom	0.96	2.53	Co-segregation in 2 affected and 1 unaffected relatives (mosaicism in one affected parent).	ETS family transcription factor <i>ERG</i> is a critical regulator of endothelial lineage specification, vascular development, angiogenesis, and endothelial homeostasis. ^{40,20} .	3	Internal (case enrolled for a different Specific Disease)

<i>GPR156</i>	Congenital hearing impairment	SD	510	54,739	High	Re c	0	1.0 4	Co-segregation in 2 affected and 4 unaffected relatives.	G protein-coupled receptor that regulates hair cell orientation in mechanosensory epithelia including in murine auditory epithelium ³³ .	3	Riyadh
<i>PMEPA1</i>	Familial Thoracic Aortic Aneurysm Disease	SD	574	54,858	High	Do m	0. 94	1.2 1	Co-segregation in 3 affected relatives and distinctive phenotypic features.	Negative regulator of Transforming Growth Factor β (TGF β) signalling ²⁸ . Aberrant TGF β signalling is implicated in multiple Mendelian aortopathies ²⁹ .	3	100KGP pilot, Antwerp, Tokyo
<i>LRRC7</i>	Intellectual disability	SD	5,52 9	46,401	High	Do m	1	3.6		Brain-specific scaffold protein in post-synaptic densities ³⁷ . LRRC7-inactivated mice have neuro-behavioural phenotype ³⁸	2	
<i>USP33</i>	Extreme early-onset hypertension	SD	182	55,305	High	Do m	0. 86	2.1		Deubiquitinating enzyme implicated in multiple cellular processes, including regulation of expression of the β 2-adrenergic receptor ³⁹ , a critical regulator of circulatory function and blood pressure ⁴¹ .	2	
<i>ARPC3</i>	Charcot-Marie-Tooth disease	SD	549	54,856	Moderate	Do m	0. 22	0.3 9		Component of the Arp2/3 complex that regulates polymerisation of F-actin, abundant in axonal neurofilaments. Multiple Mendelian axonal filamentopathies manifest as Charcot-Marie-Tooth disease ⁴² . ArpC3-inactivation in mice causes axon dysfunction ⁴³ .	1	
<i>FMN1</i>	Congenital hearing impairment	SD	510	54,738	High	Re c	0	- 1.5 3	Co-segregation in 2 unaffected relatives.	Formin family protein involved in linear actin and microtubule polymerisation ⁴⁴ . Pathogenic variants in the formin DIAPH1 cause hearing loss via cytoskeletal disruption in auditory stereocilia ⁴⁵ .	1	
<i>RAB35</i>	Familial Hypercholesterolaemia	SD	469	55,033	High	Do m	0. 98	2.3 6	Co-segregation in 1 affected relative.	Small GTP-binding proteins that are a regulator of endosomal transport and function.	1	
<i>RAB3A</i>	Hereditary ataxia	SD	905	54,504	Moderate	Do m	0. 95	2.3 2		Small GTP-binding proteins that regulate exocytosis and secretion. Although abundant in brain synaptic vesicles, rab3A-inactivated mice have no neuromuscular phenotype ⁴⁶ .	1	
<i>TUFT1</i>	Epidermolysis bullosa	SD	32	55,459	High	Re c	0	0.9	Co-segregation in 1 affected and 4 unaffected relatives.	Acidic protein that mediates dental enamel mineralisation.	1	
<i>FAM222B</i>	Ultra-rare undescribed monogenic disorders	SD	1,20 5	53,681	Moderate	Do m	0. 29	0.4 2		Uncharacterised nucleosomal protein.	0	
<i>INSL4</i>	Rod Dysfunction Syndrome	SD	58	55,425	Moderate	Do m	0	- 1.4 3		Insulin-like growth factor implicated in trophoblast and bone development.	0	
<i>KRT14</i>	Young onset tumour syndromes	DS G	256	55,207	Moderate	Re c	0	0.8 4	Co-segregation in 2 unaffected relatives.	Component of keratin intermediate filaments in epithelial cells. Pathogenic variants cause Epidermolysis	0	

											bullosa simplex IA-D (AD/AR); Dermatopathia pigmentosa reticularis (AD); Naegeli-Franceschetti-Jadassohn syndrome (AD).		
<i>MPPE1</i>	Primary ciliary dyskinesia	SD	105	55,360	High	Re c	0	0.35	Co-segregation in 2 unaffected relatives.		Metallophosphoesterase required for transport of GPI-anchor proteins from the endoplasmic reticulum to the Golgi.	0	
<i>PKMYT1</i>	Single autosomal recessive mutation in rare disease	SD	51	55,429	Moderate	Re c	0.22	0.07	Co-segregation in 2 unaffected relatives.		Serine/threonine protein kinase that is a negative regulator of cell entry into mitosis.	0	
<i>RPL10A</i>	Milroy disease	SD	20	55,470	High	Do m	0.85	2.06			Component of the large ribosomal subunit that mediates protein translation.	0	
<i>SERPIN B3</i>	Autosomal recessive congenital ichthyosis	SD	46	55,437	Moderate	Re c	0	-1.66	Co-segregation in 2 unaffected relatives.		Cysteine endopeptidase inhibitor implicated in autocrine/paracrine signalling and cell protein metabolism	0	
<i>SRP9</i>	Ductal plate malformation	SD	54	55,445	High	Do m	0.42	1.13			Component of the signal recognition particle that targets secretory proteins to the endoplasmic reticulum.	0	
<i>WWOX</i>	Gastrointestinal disorders	DSG	59	55,413	Moderate	Re c	0	-4.44	Co-segregation in 1 unaffected relative.		Short-chain dehydrogenase/reductase that acts as a tumour suppressor and apoptosis regulator. Pathogenic variants cause developmental and epileptic encephalopathy 28 and spinocerebellar ataxia 12.	0	

428 **Table 1 | Plausibility scoring of the 21 genetic associations identified by BeviMed.** Each row
429 corresponds to a genetic association between a gene and a case set in the 100KGP Main Programme
430 without prior supporting evidence in PanelApp. Each column gives additional information for each
431 association. Cells contributing to the final score are shown in bold. Rows are sorted by score in
432 descending order and the genes corresponding to associations with a score of three are underlined. The
433 level of the case set in the disease label hierarchy (DSG: Disease Sub Group, SD: Specific Disease), the
434 class of variants and the MOI corresponding to the model with the greatest posterior probability are
435 shown (Dom: dominant; Rec: recessive). A recessive association contributes one point to the score. A
436 pLI >0.9 contributes one point to the score providing the inferred class of aetiological variants is high-
437 impact variants. A Z-score >2 contributes one point to the score providing the inferred class of
438 aetiological variants is moderate-impact variants. Evidence of co-segregation in ≥3 relatives in the
439 100KGP data contributes one point to the score (including mosaicism supported by ≥2 reads containing
440 the alternate allele). Prior evidence of a relevant biological function or disease association contributes
441 one point to the score. The 'Replication' column specifies cohorts in which additional cases were
442 confirmed.

443

444

FIGURE LEGENDS

445

446

447

448

449

Fig. 1 | BeviMed analysis of the 100KGP. **a**, Bars showing the size of each case set used for the genetic association analyses, grouped by Disease Group and coloured by type (Disease Sub Group or Specific Disease). Case sets smaller than 5 are shown as having size 4 to comply with 100KGP policy on limiting participant identifiability. Below, the names and sizes of the case sets for an exemplar Disease Sub Group, 'Cardiovascular disorders', is shown. **b**, BeviMed PPAs >0.95 arranged by Disease

450 Group. Only the strongest association for each gene within a Disease Group is shown. Associations are
451 coloured by their PanelApp evidence level (green, amber or red). Associations that were mapped to
452 PanelApp by manual review, rather than using our automatic matching algorithm, are marked with an
453 asterisk (**Source data for Fig. 1**). Previously unidentified associations are shown in grey. The shape of
454 the points shows whether the association was with a Disease Sub Group (square), or Specific Disease
455 (circle).

457 **Fig. 2 | Loss-of-function variants in *ERG* are responsible for primary lymphoedema.** **a**, Pedigrees
458 for the four probands with loss-of-function variants in the canonical transcript of *ERG*,
459 ENST00000288319.12. **b**, Truncated barchart showing the distribution of the number of reads supporting
460 the p.S182Afs*22 alternate allele in the 100KGP. The embedded windows show the read pileups at this
461 position in the two affected members of the family with the variant encoding p.S182Afs*22. The reads
462 supporting the reference allele are in blue and those supporting the variant allele are in red. **c**, Schematic
463 showing the effects of each variant at the cDNA and amino acid level, and on the protein product with
464 respect to the canonical transcript. **d**, Reverse transcription-PCR amplification of *ERG* mRNA in HDLEC
465 relative to HUVEC. Data are normalised to GAPDH. Statistical significance was assessed using a two-
466 sided Student's *t*-test, n.s.: not significant ($p=0.39$). **e**, Immunoblot (representing two replicates) of
467 HUVEC and HDLEC protein lysates identified several bands corresponding to *ERG* isoforms expressed
468 at similar intensities in both cell types. **f**, Immunofluorescence microscopy (representing three replicates)
469 of HDLEC shows *ERG* (green) nuclear co-localisation with lymphatic endothelial cell nuclear marker
470 PROX1 (violet) and DAPI (blue). HDLEC junctions are shown using an antibody to VE-cadherin (yellow).
471 Scale bar, 50 μ m. **g**, *En face* immunofluorescence confocal microscopy (representing five replicates) of
472 mouse ear skin. Vessels are stained with antibodies to the lymphatic marker PROX1 (violet) and *ERG*
473 (green). Scale bar, 100 μ m. **h**, Exemplar immunofluorescence microscopy image of HEK293 cells
474 overexpressing wild type *ERG* and the p.S182Afs*22 variant *ERG*. Cells were stained for *ERG* (green)
475 and nuclear marker DAPI (blue). Scale bar, 20 μ m. The brightness was optimised for print. **i**, Dot plot of
476 estimated proportion of *ERG* not overlapping the nuclear marker DAPI in each of a set of
477 immunofluorescence microscopy images of HEK293 cells overexpressing different *ERG* cDNAs (20
478 replicates for wild type (WT), 17 replicates per tested mutant). The estimated proportions were
479 significantly higher in each of the variants compared to wild type: $p=1.52 \times 10^{-11}$, 4.10×10^{-13} and 3.03×10^{-5}
480 for each of p.S182Afs*22, p.T224Rfs*15 and p.A447Cfs*19, respectively (two-sided Student's *t*-test).

481
482 **Fig. 3 | Truncating variants in *PMEPA1* result in Loeys-Dietz syndrome.** **a**, Pedigrees for the three
483 probands in the 100KGP (discovery cohort) heterozygous for the frameshift insertion predicting
484 p.S209Qfs*3 and probands from replication cohorts, including: one from the 100KGP pilot programme
485 heterozygous for the frameshift deletion predicting p.S209Afs*61, three of Japanese ancestry
486 heterozygous for p.S209Qfs*3 and one Belgian pedigree heterozygous for a frameshift deletion encoding
487 p.P207Qfs*3. All variant consequences are shown with respect to the canonical transcript of *PMEPA1*,
488 ENST00000341744.8. **b**, HPO terms present in at least three of the four *PMEPA1* FTAAD families,
489 excluding redundant terms within each level of frequency, alongside their frequency in four *PMEPA1*
490 FTAAD families and the other 589 unexplained FTAAD families. Terms are ordered by *p*-value obtained
491 by a Fisher's exact test of association between the term's presence in an FTAAD family and whether the
492 family is one of the four *PMEPA1* families. Terms were declared significant (indicated by an asterisk), or
493 not significant (n.s.) by comparing their Fisher test *p*-values and rank to a null distribution of equivalent

494 pairs obtained by permutation (10,000 replicates). For each rank, the p -value of the term on the 5th
495 percentile was used as an upper bound for declaring an association significant, provided all terms at
496 higher ranks were also significant. The p -values for each term were as follows. Dolichocephaly:
497 $p=2.9 \times 10^{-4}$, Abnormal axial skeleton morphology: $p=6.7 \times 10^{-3}$, Striae distensae: $p=0.013$, Pes planus:
498 $p=0.014$, Ascending tubular aorta aneurysm: $p=0.62$. **c**, Graph showing *PMEPA1* and genes with high
499 evidence (green) of association with FTAAD in PanelApp. Edges connect genes where the string-db
500 v11.5²⁷ confidence score for physical interactions between corresponding proteins was >0.6 . Genes
501 known to be associated with Loeys-Dietz syndrome are highlighted in blue. *PMEPA1* is highlighted
502 yellow. **d**, Schematic showing the effects of each variant at the cDNA and amino acid level, and on the
503 protein product.

504
505 **Fig. 4 | Loss-of-function variants in *GPR156* give rise to recessive congenital hearing loss. a**,
506 Schematic of the three pedigrees with cases homozygous or compound heterozygous for loss-of-function
507 variants in the canonical transcript of *GPR156*, ENST00000464295.6. Blank symbols indicate individuals
508 with an unknown genotype. **b**, Histograms of expression log fold changes for different sets of genes in
509 mouse hair cells compared to surrounding cells: all genes (left) and genes homologous to the human
510 counterparts in the 'Hearing loss' PanelApp panel with and without a stereocilia-related GO term (i.e. a
511 term whose name contained 'stereocilia' or 'stereocilium', or the descendant of such a term) (right). The
512 log fold change for *Gpr156* is shown as a horizontal line. **c**, Maximum intensity projections of confocal Z-
513 stacks in the organ of Corti and vestibular system of a P20 wild type mouse immunostained with
514 *GPR156* antibody (green) and counterstained with phalloidin (red). Top row: overview of the organ of
515 Corti and vestibular system. Middle and bottom rows: magnified images of outer hair cells (OHC) and
516 inner hair cells (IHC), respectively. No stereociliary bundle staining was observed. The punctate staining
517 observed in the organ of Corti was absent or significantly decreased in the utricle of the vestibular
518 system. Scale bars: 10 μm (each image representative of three replicates). **d**, Schematic showing the
519 effects of each variant at the cDNA and amino acid level, and on the protein product. **e**, Exemplar
520 western blot taken from three replicates of GFP-*GPR156* using anti-*GPR156* antibody in untransfected
521 Cos7 cells (Cos7), Cos7 cells transfected with the wild type construct (W) and Cos7 cells transfected
522 with the constructs containing each of the mutant alleles p.S642Afs*162 (S642), p.P718Lfs*86 (P718)
523 and p.S207Vfs*113 (S207).

524 525 526 **References**

- 528 1. Boycott KM et al. International Cooperation to Enable the Diagnosis of All Rare Genetic Diseases. *Am*
529 *J Hum Genet.* 2017; 100(5):695--705.
- 530 2. Ferreira CR. The burden of rare diseases. *Am J Med Genet A.* 2019; 179(6):885--892.
- 531 3. Turro E et al. Whole-genome sequencing of patients with rare diseases in a national health system.
532 *Nature.* 2020; 583(7814):96--102.
- 533 4. Wang Q et al. Rare variant contribution to human disease in 281,104 UK Biobank exomes. *Nature.*
534 2021; 597(7877):527--532.
- 535 5. Kaplanis J et al. Evidence for 28 genetic disorders discovered by combining healthcare and research
536 data. *Nature.* 2020; 586(7831):757--762.

537 6. Greene D, Richardson S, Turro E. Phenotype Similarity Regression for Identifying the Genetic
538 Determinants of Rare Diseases. *Am J Hum Genet.* 2016; 98(3):490--499.

539 7. Hail Team. (2022). Hail 0.2. <https://github.com/hail-is/hail>.

540 8. Lopez J et al. HGVA: the Human Genome Variation Archive. *Nucleic Acids Res.* 2017; 45(W1):W189-
541 W194.

542 9. Greene D, Richardson S, Turro E. A Fast Association Test for Identifying Pathogenic Variants Involved
543 in Rare Diseases. *Am J Hum Genet.* 2017; 101(1):104--114.

544 10. Karczewski KJ et al. The mutational constraint spectrum quantified from variation in 141,456
545 humans. *Nature.* 2020; 581(7809):434--443.

546 11. Rentzsch P, Schubach M, Shendure J, Kircher M. CADD-Splice-improving genome-wide variant
547 effect prediction using deep learning-derived splice scores. *Genome Med.* 2021; 13(1):31.

548 12. Eilbeck K et al. The Sequence Ontology: a tool for the unification of genome annotations. *Genome*
549 *Biol.* 2005; 6(5):R44.

550 13. Howe KL et al. Ensembl 2021. *Nucleic Acids Res.* 2021; 49(D1):D884-D891.

551 14. Martin AR et al. PanelApp crowdsources expert knowledge to establish consensus diagnostic gene
552 panels. *Nat Genet.* 2019; 51(11):1560--1565.

553 15. Korber L, Schneider H, Fleischer N, Maier-Wohlfart S. No evidence for preferential X-chromosome
554 inactivation as the main cause of divergent phenotypes in sisters with X-linked hypohidrotic ectodermal
555 dysplasia. *Orphanet J Rare Dis.* 2021; 16(1):98.

556 16. Kasahara Y et al. Hyper-IgM syndrome with putative dominant negative mutation in activation-
557 induced cytidine deaminase. *J Allergy Clin Immunol.* 2003; 112(4):755--760.

558 17. Martin-Almedina S, Mortimer PS, Ostergaard P. Development and physiological functions of the
559 lymphatic system: insights from human genetic studies of primary lymphedema. *Physiol Rev.* 2021;
560 101(4):1809--1871.

561 18. Gordon K et al. Update and audit of the St George's classification algorithm of primary lymphatic
562 anomalies: a clinical and molecular approach to diagnosis. *J Med Genet.* 2020; 57(10):653--659.

563 19. Kalna V et al. The Transcription Factor ERG Regulates Super-Enhancers Associated With an
564 Endothelial-Specific Gene Expression Program. *Circ Res.* 2019; 124(9):1337--1349.

565 20. Shah AV, Birdsey GM, Randi AM. Regulation of endothelial homeostasis, vascular development and
566 angiogenesis by the transcription factor ERG. *Vascul Pharmacol.* 2016; 86:3--13.

567 21. Hoesel B et al. Sequence-function correlations and dynamics of ERG isoforms. ERG8 is the black
568 sheep of the family. *Biochim Biophys Acta.* 2016; 1863(2):205--218.

569 22. Westbury SK et al. Human phenotype ontology annotation and cluster analysis to unravel genetic
570 defects in 707 cases with unexplained bleeding and platelet disorders. *Genome Med.* 2015; 7(1):36.

571 23. Greene D, Richardson S, Turro E. ontologyX: a suite of R packages for working with ontological data.
572 *Bioinformatics.* 2017; 33(7):1104--1106.

573 24. Resnik P, Others. Semantic similarity in a taxonomy: An information-based measure and its
574 application to problems of ambiguity in natural language. *J. Artif. Intell. Res.(JAIR).* 1999; 11:95--130.

575 25. Ciurica S et al. Arterial Tortuosity. *Hypertension.* 2019; 73(5):951--960.

576 26. Loeys BL et al. A syndrome of altered cardiovascular, craniofacial, neurocognitive and skeletal
577 development caused by mutations in TGFBR1 or TGFBR2. *Nat Genet.* 2005; 37(3):275--281.

578 27. Szklarczyk D et al. STRING v11: protein-protein association networks with increased coverage,
579 supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res.* 2019;
580 47(D1):D607-D613.

- 581 28. Watanabe Y et al. TMEPAI, a transmembrane TGF-beta-inducible protein, sequesters Smad proteins
582 from active participation in TGF-beta signaling. *Mol Cell*. 2010; 37(1):123--134.
- 583 29. Creamer TJ, Bramel EE, MacFarlane EG. Insights on the Pathogenesis of Aneurysm through the
584 Study of Hereditary Aortopathies. *Genes (Basel)*. 2021; 12(2).
- 585 30. Thermann R et al. Binary specification of nonsense codons by splicing and cytoplasmic translation.
586 *EMBO J*. 1998; 17(12):3484--3494.
- 587 31. Sobreira N, Schiettecatte F, Valle D, Hamosh A. GeneMatcher: a matching tool for connecting
588 investigators with an interest in the same gene. *Hum Mutat*. 2015; 36(10):928--930.
- 589 32. Ellaithy A, Gonzalez-Maeso J, Logothetis DA, Levitz J. Structural and Biophysical Mechanisms of
590 Class C G Protein-Coupled Receptor Function. *Trends Biochem Sci*. 2020; 45(12):1049--1064.
- 591 33. Kindt KS et al. EMX2-GPR156-Gai reverses hair cell orientation in mechanosensory epithelia. *Nat*
592 *Commun*. 2021; 12(1):2861.
- 593 34. Scheffer DI, Shen J, Corey DP, Chen ZY. Gene Expression by Mouse Inner Ear Hair Cells during
594 Development. *J Neurosci*. 2015; 35(16):6366--6380.
- 595 35. Miyoshi T et al. Human deafness-associated variants alter the dynamics of key molecules in hair cell
596 stereocilia F-actin cores. *Hum Genet*. 2022; 141(3-4):363--382.
- 597 36. Smedley D et al. 100,000 Genomes Pilot on Rare-Disease Diagnosis in Health Care - Preliminary
598 Report. *N Engl J Med*. 2021; 385(20):1868--1880.
- 599 37. Thalhammer A, Trinidad JC, Burlingame AL, Schoepfer R. Densin-180: revised membrane topology,
600 domain structure and phosphorylation status. *J Neurochem*. 2009; 109(2):297--302.
- 601 38. Chong CH et al. *Lrrc7* mutant mice model developmental emotional dysregulation that can be
602 alleviated by mGluR5 allosteric modulation. *Transl Psychiatry*. 2019; 9(1):244.
- 603 39. Berthouze M, Venkataramanan V, Li Y, Shenoy SK. The deubiquitinases USP33 and USP20
604 coordinate beta2 adrenergic receptor recycling and resensitization. *EMBO J*. 2009; 28(12):1684--1696.
- 605 40. Birdsey GM et al. The endothelial transcription factor ERG promotes vascular stability and growth
606 through Wnt/Beta-catenin signaling. *Dev Cell*. 2015; 32(1):82--96.
- 607 41. Motiejunaite J, Amar L, Vidal-Petiot E. Adrenergic receptors and cardiovascular effects of
608 catecholamines. *Ann Endocrinol (Paris)*. 2021; 82(3-4):193--197.
- 609 42. Munoz-Lasso DC, Roma-Mateo C, Pallardo FV, Gonzalez-Cabo P. Much More Than a Scaffold:
610 Cytoskeletal Proteins in Neurological Disorders. *Cells*. 2020; 9(2).
- 611 43. Zuchero JB et al. CNS myelin wrapping is driven by actin disassembly. *Dev Cell*. 2015; 34(2):152--
612 167.
- 613 44. DeWard AD, Eisenmann KM, Matheson SF, Alberts AS. The role of formins in human disease.
614 *Biochim Biophys Acta*. 2010; 1803(2):226--233.
- 615 45. Ninoyu Y et al. The integrity of cochlear hair cells is established and maintained through the
616 localization of Dia1 at apical junctional complexes and stereocilia. *Cell Death Dis*. 2020; 11(7):536.
- 617 46. Geppert M et al. The role of Rab3A in neurotransmitter release. *Nature*. 1994; 369(6480):493--497.

618 619 **METHODS**

620 621 **Ethics**

622 The 100,000 Genomes project was approved by East of England–Cambridge Central REC REF
623 20/EE/0035. Only participants who provided written informed consent for their data to be used for
624 research were included in the analyses. The study at the University of Maryland was approved by the

625 Institutional Review Board (RAC#2100001) and written informed consent was obtained by clinicians at
626 King Faisal Hospital in Saudi Arabia from the participating individuals. The study of the Japanese
627 ancestry pedigrees bearing *PMEPA1* truncating alleles was approved by the Institutional Review Board
628 of the National Cerebral and Cardiovascular Centre (M14-020) and Sakakibara Heart Institute (16-035),
629 and written informed consent was obtained from the participating individuals.
630

631 **Motivation for developing a sparse relational database**

632 Computational approaches for discovering the aetiologies of rare diseases typically depend on the
633 analysis of a heterogeneous set of files, each of which can be very large and follow a distinct convention.
634 Genotypes, for example, are ordinarily stored in VCFs containing data for one sample or for multiple
635 samples. In the latter case, the data are usually distributed in files covering many different "chunks" of
636 the reference genome. Variant-level information, such as consequence predictions or pathogenicity
637 scores, are typically encoded in strings that require extensive parsing to decode, either from within the
638 VCFs containing the genotypes, or in separate files. Modifying genotype or annotation files, for example
639 in order to incorporate newly generated data, requires rewriting files in their entirety. Phenotype data,
640 pedigree data and the results of statistical inference are stored in a further set of files. Consequently,
641 analyses are often burdensome to conduct and prone to error. Frameworks such as Hail⁷ and
642 OpenCGA⁸ afford greater flexibility but they depend on the centrally organised deployment of a
643 distributed storage system, hindering usability and portability.
644

645 Relational databases are widely used, mature technologies, well known for their speed, reliability,
646 flexibility, structure and extensibility. In the context of rare diseases, a relational database can in principle
647 render the modification, combination and addition of data on samples, variants, genes and other entities
648 efficient, reliable and straightforward to implement using a single query language. Unfortunately, the
649 performance of relational databases degrades substantially when the number of records in a table
650 reaches several billion, and the number of genotypes in a cohort the size of the 100KGP easily
651 surpasses this threshold. However, the MAFs of pathogenic variants with strong effects on rare disease
652 risk are typically kept below 1/1,000 by negative selection, and the proportion of non-homozygous
653 reference genotypes for variants within that MAF stratum is only about 1% of the total (**Extended Data**
654 **Fig. 1**). Consequently, it is possible to construct a compact relational database that includes virtually all
655 the pathogenic variants even in a large cohort such as the 100KGP. This provides an opportunity for
656 exploiting the benefits of a single unified relational database containing the non-homozygous genotypes
657 of rare variants upon which to conduct the entirety of the aetiological discovery process. Furthermore, it
658 provides a natural foundation for developing web applications for the multidisciplinary review of genetic,
659 phenotypic, statistical and other data.
660

661 **Rareservoir**

662 The Rareservoir is a relational database schema and a complementary software package 'rsrv' for
663 working with rare disease data. The database stores data including rare variant genotypes, variant
664 annotations, phenotypes, sample information and pedigrees (**Extended Data Fig. 1**) but it can be
665 extended arbitrarily. A Rareservoir is built through a series of steps from a set of input data and
666 parameters (**Extended Data Fig. 3**). The 'bcftools' program⁴⁷ extracts ('bcftools view') and normalises
667 ('bcftools norm') variants from either a set of single sample genome VCF files (gVCFs) or from a merged
668 VCF. In all steps of the procedure, variants are encoded as RSVR IDs using the 'rsrv enc' tool (see

669 **Encoding RSVR IDs**). Merged VCFs typically contain cohort-wide variant quality information in the
670 FILTER column, which can be used to select variants for processing. However, this is not readily
671 obtained from single gVCFs. To address this, we developed the 'rsvr depth' tool, which computes variant
672 quality pass rates at all positions in the genome based on a random subsample of gVCFs. If the input is
673 a merged VCF, an internal (i.e. within-VCF) allele frequency threshold is applied with bcftools to filter out
674 internally common variants. If the input is a set of single-sample gVCFs, internally common variants are
675 filtered out in two steps, for computational efficiency. Firstly, a set of variants that are statistically almost
676 certain to be common based on a random sample of gVCFs is identified—by default, the variants for
677 which a one-sided binomial test under the null hypothesis that the MAF=0.01 is rejected at a significance
678 level of 10^{-6} (done using the 'rsvr tabulate' tool). Secondly, all gVCFs are read sequentially, filtering out
679 the variants identified in the previous step (using the 'rsvr mix' tool) and those for which the pass rates
680 identified with 'rsvr depth' do not meet the threshold. Retained genotypes are then loaded into a
681 temporary genotype table in the database in order to apply the final internal allele frequency filter by
682 executing an SQL 'DELETE' statement. These variants are then annotated with gnomAD 'probabilistic
683 minor allele frequency' (PMAF) scores³ using the 'rsvr pmaf' tool. The PMAF score is calculated with
684 respect to a given allele frequency threshold t , by evaluating a binomial test (at a significance threshold
685 of 0.05) on the observed frequency of the variant under the null hypothesis that the variant has an allele
686 frequency of t . If, in any gnomAD population, the null is rejected for $t=0.001$ and the allele count is at
687 least 2, the score is set to 0. If the null is rejected for $t=0.0001$, the score is set to 1. If the null is not
688 rejected, the score is set to 2. Finally, if the variant is absent from gnomAD, the score is set to 3. For the
689 non-pseudo autosomal dominant regions of chromosome X, only allele counts for males are used in
690 calculations. Variants are then additionally annotated with their CADD phred scores using the 'rsvr ann'
691 program, and loaded into the VARIANT table. At this point, variants in the VARIANT and GENOTYPE
692 table which have a PMAF score of 0 may be deleted because they are unlikely to be involved in rare
693 diseases. We then annotate the retained variants with predicted transcript consequences for a given set
694 of transcripts specified in a Gene Transfer Format (GTF) file. The 100KGP Rareservoir, uses Ensembl
695 v104 canonical transcripts with a protein-coding biotype, of which >90% are Matched Annotation from
696 NCBI and EMBL-EBI (MANE)⁴⁸ transcripts. The 'rsvr seqfx' program determines a set of SO terms for
697 each interacting transcript-variant pair and encodes them as a CSQ ID, which is added to the
698 CONSEQUENCE table. This table can also hold LOFTEE scores corresponding to a transcript-variant
699 pair. Note that, as LOFTEE scores on the Genomics England Research Environment correspond to
700 Ensembl v99 transcripts, we mapped Ensembl v104 canonical transcripts to the most similar v99
701 transcripts having an identical CDS in order to obtain the LOFTEE scores for the 100KGP Rareservoir,
702 finding a match for >98% of transcripts. The contents of the GTF file are also imported into the database
703 to create tables of transcript features (FEATURE), transcripts (TX) and genes (GENE). Optionally,
704 VARIANT, GENOTYPE and CONSEQUENCE may be filtered for RSVR IDs that have CSQ IDs meeting
705 particular criteria, for instance, in order to retain only variants with protein-coding consequences. The
706 SAMPLE table of metadata and genetic statistics for each sample represented in the input VCF(s) must
707 then be added to the database, including mandatory columns containing the ID, sex, family, and an
708 indicator of inclusion in the maximal unrelated set of samples in the database. The VARIANT,
709 GENOTYPE and CONSEQUENCE tables are indexed by RSVR ID, in order to support fast lookups by
710 genomic location. The SAMPLE table and GENOTYPE table are indexed by sample ID allowing fast
711 lookups by sample. The CONSEQUENCE, TX and GENE tables are indexed by transcript and gene ID,
712 allowing fast lookups of variants based on gene/transcript specific consequences. If sample phenotypes

713 have been encoded using phenotypic terms (e.g. ICD10 codes or HPO terms), terms from the relevant
714 coding systems can be added to a generic PHENOTYPE table mapping code IDs to descriptions, and
715 codes assigned to samples can be added to the SAMPLE_PHENOTYPE table. Disease labels may be
716 added to the CASE_SET table. The majority of the compute time required for building the database is
717 taken by reading the genotype data from the input VCF, which may be executed in parallel over separate
718 regions against a merged VCF or over single gVCFs. The rsrv tool, implemented in C++, executes
719 rapidly, with 'rsrv seqfx' capable of assigning CSQ IDs for all Ensembl v104 canonical transcripts to all
720 variants (over 685M) in gnomAD v3.0 in under 40 minutes on a single core. The 100KGP Rareservoir,
721 which is stored in a SQLite database, returns complex gene-specific queries in under one second. For
722 example: (1) a table with 628 rows containing the moderate and high-impact variants with a PMAF score
723 ≥ 1 in *TTN*, along with the corresponding SeqFx consequence predictions and CADD scores takes 0.57
724 seconds; (2) a table with 1,498 rows containing, for each variant, the samples and genotypes for
725 individuals who carry an alternate allele takes 0.61 seconds; and (3) a classification for each of the
726 77,539 participants into proband with Dilated Cardiomyopathy, relative of such a proband, unrelated
727 control, or relative of a control takes 0.65 seconds. Specific details on implementation of the workflow,
728 code for encoding data as SQL statements compatible with Rareservoir and the mapping between bits in
729 the 64 bit CSQ ID and each SO term assigned by 'rsrv seqfx' can be found in the rsrv software package
730 (see **Code availability**). Software packages rsrv 1.0, bcftools 1.9 and perl 5 were used to build the
731 100KGP Rareservoir.

732

733 **Encoding RSVR IDs**

734 SNVs and indels may be encoded as 64-bit integers called RSVR IDs. In order to compute an RSVR ID
735 for a given variant, the following expression is evaluated:

$$736 \quad c \times 2^{58} + p \times 2^{30} + |r| \times 2^{24} + |a| \times 2^{18} + \sum_{i=1}^{|A|} A_i \times 4^{i-1},$$

737

738 where c is the chromosome number (using 23, 24 and 25 respectively to represent X, Y and MT), p is the
739 position, and $|r|$ and $|a|$ are the lengths of the reference and alternate alleles, respectively. A is a
740 sequence identical to the alternate allele, a , when its length is less than 10, and otherwise equal to the
741 first five followed by the last four elements of a . In the summation, nucleotides are assigned values: A= 0,
742 C = 1, G = 2 and T = 3. The expression evaluates to integers that can be represented using 63 bits,
743 setting the most significant bit to 0 when encoding as 64-bit integers. The chromosome, position,
744 reference and alternate allele lengths and alternate allele bases are thereby encoded respectively by the
745 subsequent 5, 28, 6, 6 and 18 bits (with two bits per base for the alternate allele). This procedure and its
746 inverse are implemented in the 'rsrv enc' and 'rsrv dec' programs respectively. The reference and
747 alternate alleles of input variants are first normalised by removing any redundant identical sequence from
748 the starts and then the ends. The proportion of variants in gnomAD 3.0 weighted by allele count that can
749 be encoded losslessly is 99.3%, while 99.8% can be represented by a distinct RSVR ID. The full variant
750 information corresponding to any encountered ambiguous RSVR ID may be stored in full in a dedicated
751 table. Structural variants that can be represented by a position and length may also be encoded using
752 distinct 64-bit RSVR IDs alongside SNVs and indels by setting the most significant bit to 1, and
753 subsequently encoding the type of structural variant using 2 bits (Deletion 0, Duplication 1, Inversion 2,
754 Insertion 3), the chromosome using 5 bits (as done for SNVs and indels), and the start and length
755 consecutively using 28 bits.

756

757 **Genetic association analysis of 100KGP data**

758 We constructed a Rareservoir in the Genomics England Research Environment containing the
759 PASSing⁴⁹ variants in the merged VCF of 77,539 consented participants in the 100KGP rare diseases
760 programme. This Rareservoir only included variants with a PMAF >0 according to GnomAD v3.0, an
761 internal MAF <0.002 and at least one predicted consequence on a canonical transcript in Ensembl v104.
762 Variants with a greater MAF are unlikely to be highly penetrant for diseases eligible for inclusion in the
763 100KGP and are likely to have, at most, small effects on risk, making them challenging to validate.
764 Variants with a median genotype quality <35 and SNVs with a CADD Phred score <10 were also
765 excluded from the analyses.
766

767 For each of the 269 rare disease classes (**Extended Data Figs. 5–6**), we applied the BeviMed⁹
768 association test to rare variants extracted from the Rareservoir database in each of the 19,663 canonical
769 transcripts belonging to a gene with a 'protein_coding' biotype. The analysis was carried out using R
770 3.6.2, making use of functionality from packages: Matrix 1.2-18, dplyr 0.8.5, bit64 0.9-7, bit 1.1-14, DBI
771 1.1.0, RSQLite 2.1.4 and BeviMed 5.7. The case set for a given disease class and gene was constructed
772 by selecting one case from each pedigree containing at least one person affected with the disease class.
773 For the purposes of the association analysis, participants were labelled 'explained' by a given gene if
774 they had variants in that gene classified as 'pathogenic_variant' or 'likely_pathogenic_variant' in the
775 'gmc_exit_questionnaire' table in the Genomics England Research Environment. To boost power, we
776 used this information to reassign cases who were explained by variants in a different gene to the control
777 group.
778

779 Using BeviMed, we performed a Bayesian comparison of a baseline model of no association and each of
780 six association models defined by a mode of inheritance and a class of aetiological variant:

- 781 1. No association (prior probability: 0.99),
- 782 2. Dominant association with "high"-impact variants having a PMAF ≥ 2 (i.e., corresponding to a
783 target MAF <0.01%) (prior probability: 0.002475),
- 784 3. Dominant association with "moderate"-impact variants having a PMAF ≥ 2 (prior probability:
785 0.002475),
- 786 4. Dominant association with "5' UTR" variants having a PMAF ≥ 2 (prior probability: 0.00005),
- 787 5. Recessive association with "high"-impact variants having a PMAF ≥ 1 (i.e., corresponding to a
788 target MAF <0.1%) (prior probability: 0.002475),
- 789 6. Recessive association with "moderate"-impact variants having a PMAF ≥ 1 (prior probability:
790 0.002475),
- 791 7. Recessive association with "5' UTR" variants having a PMAF ≥ 1 (prior probability: 0.00005).

792 Thus the overall prior probability of association was 0.01 and there was an equal prior probability of
793 dominant and recessive inheritance. The PPA was the sum of the posterior probabilities of models 2
794 through 7. We imposed a stricter PMAF threshold under a dominant MOI than under a recessive MOI
795 because, ceteris paribus, dominant variants are under stronger negative selection than recessive
796 variants. The three groups of variants were selected as follows:

- 797 ● 5' UTR variants: those with a 5_prime_UTR_variant consequence,
- 798 ● High-impact variants: those with any consequence amongst start_lost, stop_lost,
799 frameshift_variant, stop_gained, splice_donor_variant or splice_acceptor_variant, excluding
800 variants with a "low-confidence" LOFTEE score¹⁰,

- Moderate-impact variants: those with any consequence amongst start_lost, stop_lost, frameshift_variant, stop_gained, splice_donor_variant or splice_acceptor_variant, missense_variant or inframe_deletion.

The rationale for embedding variants from the high-impact class in the moderate-impact class is that both types of variant are capable of inducing a loss of function. The prior on the probability that a modelled rare variant is pathogenic, conditional on either the association model mediated by 5' UTR variants or the association model mediated by moderate-impact variants, was set to Beta(2,8). This encodes a prior conditional expectation that 20% of rare variants are pathogenic, which is well suited to missense and 5' UTR variants. However, we specified a distribution with a greater mean for the high-impact models. Specifically, the prior on the probability that a modelled high-impact variant is pathogenic was set to Beta(3,1), which reflects a prior conditional expectation that 75% of rare variants are pathogenic because loss-of-function variants tend to be functionally equivalent to each other. BeviMed reports the posterior probability that each variant is pathogenic conditional on the mode of inheritance and the class of aetiological variant. The methodology is described in further detail in the original BeviMed publication⁹.

We applied the following post-processing of BeviMed results with a PPA >0.95:

- We re-ran BeviMed including all samples (i.e. with relatives of cases and controls). Associations for which the analysis with all samples caused the PPA to fall below 0.9 were filtered out due to conflicting evidence for the association within families.
- We re-ran BeviMed after removing variants absent from affected relatives of the cases. Associations for which this removal caused the PPA to drop below 0.25 were filtered out because they depended on variants that were not shared by affected cases within families.
- To guard against false positives due to incorrect pedigree data, population structure or cryptic relatedness, we applied the following algorithm. We obtained the distribution of the number of rare variants in the Rareservoir shared by pairs of individuals within each assigned ancestry in the 100KGP. The top percentile in each of these distributions was used to indicate potential relatedness between participants in the same population. We re-ran BeviMed after removing cases so as to ensure that no more than one case from any set of potentially related cases sharing a variant were included in the analysis. . Associations for which this analysis caused the PPA to fall below 0.25 were filtered out.

To account for correlation between case sets, for each gene, we removed all but the most strongly associated disease class within each Disease Group before reporting the 260 associations remaining (**Source data for Fig. 1**). Without the post-processing, the number of reported associations would have been 302. Conditional on the modal model underlying each of the 260 associations, we recorded the variants with a posterior probability of pathogenicity >0.8 accounting for at least one case in the 100KGP (**Source data for Fig. 2**).

PanelApp annotation

Significant associations were coloured according to PanelApp¹⁴ (**Fig. 1b**) evidence levels for panel-gene relations (green for high evidence, amber for moderate evidence, and red for low evidence) for panels of type 'Rare Disease 100K', which are organised hierarchically by Disease Sub Group and Disease Group, or of type 'GMS Rare Disease'. Given an association between a gene and a case set (corresponding either to a Specific Disease or a Disease Sub Group), we searched for panels which contained the gene

845 and had the same name as the case set (ignoring case). If such a match was not found, we searched for
846 panels which contained the gene and which belonged to a Disease Sub Group with the same name as
847 the Disease Sub Group of the case set. When this matching rule generated multiple matches, we
848 selected the panel(s) with the highest evidence. If multiple panels still remained, we selected the panel
849 with the smallest number of genes. Associations for which no matching panel in PanelApp could be
850 found were inspected manually to assess whether PanelApp contained an alternative suitable panel
851 (marked with an asterisk in **Fig. 1b**).

852

853 **Shortlisting previously unidentified genetic associations for validation**

854 Several sources of independent evidence were used to shortlist significant associations for validation.

855 For each source, a score of one was awarded if the evidence was supportive, and zero otherwise.

856 Scores were then added over the different sources and used to rank the associations. Associations for
857 which at least three sources of evidence were supportive were taken forward for further investigation.

858 The sources of evidence and qualifying criteria for being considered supportive are listed below. Note
859 that here we refer to variants which had a probability of pathogenicity >0.8 conditional on the modal
860 model as 'probably pathogenic'.

- 861 ● *Counting co-segregating pedigree members.* The pedigrees harbouring pathogenic configurations
862 of probably pathogenic alleles were checked for co-segregation between genotype and affection
863 status. This evidence counted as supportive for associations for which all such pedigrees
864 demonstrated co-segregation, and there were at least three additional relatives who had not been
865 included in the association analysis but for whom there was co-segregation. Note that BAM files
866 for the affected members of pedigrees who were called homozygous reference for probably
867 pathogenic variants were checked for evidence of mosaicism to guard against the possibility that
868 they were falsely portraying a lack of co-segregation.
- 869 ● *pLI and Z-scores.* pLI and Z-scores for depletion of missense variants were obtained from the
870 gnomAD v2.2.1 browser¹⁰. pLI >0.9 for associations in which high impact variants were most
871 strongly associated were counted as supportive, whilst Z-scores greater than 2 for associations in
872 which moderate impact variants were most strongly associated were counted as supportive.
- 873 ● *Recessive association.* Population genetic metrics of purifying selection (pLI scores and Z-
874 scores) are sensitive to depletion of high-impact variants and missense variants, respectively.
875 They are therefore useful measures to corroborate dominant associations. However, these
876 metrics have low sensitivity to identify the signatures of selection against recessive diseases
877 because isolated pathogenic variants in heterozygous form do not lead to a reduction in
878 reproductive fitness. To avoid disadvantaging recessive associations identified by BeviMed, they
879 were assigned a contribution of one point to the score.
- 880 ● *Literature review.* A comprehensive literature review, assessing the gene's role (if any) in
881 biological processes relevant to the disease, other diseases, and a survey of model organisms
882 was undertaken, and determined to be either supportive or not.

883

884 **ERG: Primary endothelial cell culture**

885 Single donor primary human dermal lymphatic endothelial cells (HDLEC) (Promocell, Heidelberg) were
886 cultured in Endothelial Cell Growth Medium MV2 (Promocell). Pooled donor human umbilical vein
887 endothelial cells (HUVEC) (Lonza, Slough) were grown in Endothelial Cell Growth Media-2 (EGM-2)
888 (Lonza). HUVEC and HDLEC were grown on 1% (v/v) gelatin and used between passages 3-5.

889

890 **ERG: Real-time polymerase chain reaction**

891 HUVEC and HDLEC were grown to confluency in a pre-gelatinised 6-well dish. Total RNA was isolated
892 using the RNeasy Mini Kit (Qiagen) and 1 µg of total RNA was transcribed into cDNA using Superscript
893 III Reverse Transcriptase (Thermo Fisher Scientific). Quantitative real-time PCR was performed using
894 PerfCTa SYBR Green FastMix (Quanta Biosciences) on a Bio-Rad CFX96 System. Gene expression
895 values of ERG in HUVEC and HDLEC were normalised to GAPDH expression and compared using the
896 $\Delta\Delta C_t$ method. The following oligonucleotides were used: ERG, 5'-GGAGTGGGCGGTGAAAGA-3' and
897 5'-AAGGATGTCCGGCGTTGTAGC-3'; GAPDH, 5'-CAAGGTCATCCATGACAACTTTG-3' and 5'-
898 GGGCCATCCACAGTCTTCTG-3'.

899

900 **ERG: Immunoblotting analysis**

901 Immunoblotting was performed according to standard conditions. Proteins were labelled with the
902 following primary antibodies: rabbit anti-human ERG antibody (1:1000; ab133264, Abcam) and mouse
903 anti-human GAPDH (1:10000; MAB374, Millipore). Primary antibodies were detected using fluorescently
904 labelled secondary antibodies: goat anti-rabbit IgG DyLight 680 and goat anti-mouse IgG Dylight 800
905 (Thermo Scientific). Detection of fluorescence intensity was performed using an Odyssey CLx imaging
906 system (Li-COR Biosciences, Lincoln) and Odyssey version 4 software.

907

908 **ERG: Immunofluorescence analysis of endothelial cells and mouse tissues**

909 Confluent cultures of HUVEC and HDLEC were fixed with 4% (w/v) paraformaldehyde for 15 minutes and
910 permeabilised with 0.5% (v/v) Triton-X100, before incubation with 3% BSA (w/v) in PBS containing the
911 following primary antibodies: goat anti-human PROX1 antibody (1:100; AF2727, R&D Systems), rabbit
912 anti-human ERG antibody (1:100; ab92513, Abcam), mouse anti-human VE-cadherin (1:100; 555661,
913 BD Biosciences). Secondary antibody incubation was carried out in 3% BSA (w/v) in PBS, using the
914 following antibodies: donkey anti-goat IgG Alexa Fluor-488 (1:1000; A-11055), donkey anti-rabbit IgG
915 Alexa Fluor-555 (1:1000; A-31572), donkey anti-mouse Alexa Fluor-594 (1:1000; A-21203). All
916 secondary antibodies from Thermo Fisher Scientific. Nuclei were visualised using DAPI (4',6-diamidino-
917 2-phenylindole) (Sigma-Aldrich). Confocal microscopy was carried out on a Carl Zeiss LSM780 confocal
918 laser scanning microscope with Zen 3.2 software. All animal experiments were conducted with ethical
919 approval from Imperial College London under UK Home Office Project Licence number PEDBB1586 in
920 compliance with the UK Animals (Scientific Procedures) Act of 1986. Ear tissue was collected from
921 euthanised 3-week old male and female C57BL/6J mice and fixed in 4% (w/v) paraformaldehyde at room
922 temperature for 2h. Tissue was then washed with PBS followed by a blocking and permeabilization step
923 using 3% (w/v) milk in PBST (containing 0.3% (v/v) Triton X-100 in PBS) for 1h at room temperature. The
924 following primary antibodies were used for immunofluorescence staining: goat anti-human PROX1
925 antibody (1:100; AF2727, R&D Systems) and rabbit anti-human ERG antibody (1:100; ab92513, Abcam).
926 Primary antibodies were incubated at 4°C overnight in 3% (w/v) milk in PBST. The following day, tissues
927 were washed three times with PBST over the course of 2h at room temperature. Tissues were incubated
928 with secondary antibodies at room temperature for 2h in 3% milk (w/v) in PBST. Primary antibodies were
929 detected using fluorescently labelled secondary antibodies: donkey anti-goat IgG Alexa Fluor-488 (1:400;
930 A-11055, Thermo Fisher Scientific) and donkey anti-rabbit IgG Alexa Fluor-555; A-31572, Thermo Fisher
931 Scientific). Stained samples were mounted onto glass slides using Fluoromount G (Thermo Fisher

932 Scientific). Images were acquired using Zeiss LSM-780 confocal laser scanning microscope with Zen 3.2
933 software. All confocal images represent maximum intensity projection of Z-stacks of single tiles.

934
935 **ERG: Subcloning and overexpression in HEK293 cells**

936 We subcloned *ERG* (ENST00000288319.12) from HUVEC into the mammalian expression vector
937 pcDNA3.1 (Thermo Fisher). *ERG* variants were generated by site-directed mutagenesis using the
938 Quikchange Lightning kit (Agilent, Stockport, Cheshire) using the wild type *ERG* cDNA as template.
939 Expression of wild type and mutant *ERG* was carried out using Polyethylenimine (PEI; Sigma-Aldrich)
940 transfection reagent in HEK293 cells grown in DMEM (Thermo Fisher) with 10% (v/v) fetal bovine serum.
941 After 24 hr, cells were fixed with 4% (w/v) paraformaldehyde for 15 minutes and permeabilised with 0.5%
942 (v/v) Triton-X100, before incubation with 3% BSA (w/v) in PBS containing mouse monoclonal anti-ERG
943 antibody (1:100; sc-376293, Santa Cruz Biotechnology). Secondary antibody incubation was carried out
944 in 3% BSA (w/v) in PBS, using donkey anti-mouse Alexa Fluor-488 (1:1000; A-21202, Thermo Fisher).
945 Nuclei were visualised using DAPI (4',6-diamidino-2-phenylindole) (Sigma-Aldrich). Confocal microscopy
946 was carried out on a Carl Zeiss LSM780 confocal laser scanning microscope with Zen 3.2 software.

947
948 **ERG: Estimation of nuclear and non-nuclear ERG in HEK293 cells**

949 Each image was read into a pair of channel-specific 1,024 x 1,024 matrices in R 4.2.1 using the readCzi
950 function from the readCzi R package v0.2.0. A pixel was declared to contain a nuclear region if the
951 intensity in the blue channel exceeded 60% of the 95th percentile of blue intensities across all pixels
952 above background (identified as exceeding 1.35×10^{-2} by visual inspection of bimodal intensity
953 histograms). A pixel was declared to contain ERG if the intensity in the green channel exceeded 30% of
954 the 95th percentile of the green intensities within the pixels previously declared to be nuclear. To fill in
955 intranuclear gaps, any non-nuclear pixels adjacent to at least 5 nuclear pixels were declared nuclear.
956 The estimated proportion of ERG that was cytosolic in an image was set to the number of ERG pixels
957 that did not overlap nuclear pixels divided by the number of ERG pixels.

958
959 **GPR156: Western blots**

960 We subcloned *GPR156* from human brain cDNA, into EGFP-N2 vector. The three mutant *GPR156*
961 constructs were generated by mutagenesis using the QuickChange kit (Stratagene, La Jolla, CA) and a
962 wild type *GPR156*-GFP as a template. For expression analysis, the WT and mutants were transfected in
963 COS7 cells grown in DMEM (Gibco, Gaithersburg, MD, USA) with 10% fetal bovine serum. Transfections
964 were performed with Lipofectamine 2000 reagent (Life Technologies). Cells were harvested 48hr after
965 transfection, lysed in buffer containing 1% CHAPS, 100mM NaCl, and 25mM HEPES, pH 7.4 and
966 clarified by centrifugation at 18,407 x g. Lysates (20µg) were run on a 4–20% SDS-PAGE gel.
967 Membrane was blocked with 5% milk then incubated with anti-GPR156 (1:200) and immunoblots
968 developed with HRP conjugated secondary (sheep anti-rabbit) antibody (1:1,000). Comparable loading
969 was checked by stripping and reprobing the blots with anti-GAPDH (1:500) antibodies (Santa Cruz
970 Biotechnology, Heidelberg, Germany).

971
972 **GPR156: Whole mount immunostaining of GPR156 in mouse inner ears**

973 All the animal work was approved by the University of Maryland, Baltimore Institutional Animal Care and
974 use Committee (IACUC 420002). Inner ears were dissected from C57BL/6J mice with a postnatal age of
975 10 days and fixed in 4% paraformaldehyde (PFA) in phosphate buffered saline (PBS) overnight. For

976 whole mount immunostaining, the cochleae were micro-dissected and were subjected to blocking for 1
977 hour with 10% normal goat serum in PBS containing 0.25% tritonX100, followed by overnight incubation
978 at 4°C with anti-GPR156 antibodies (1:200; Cat#PA5-23857; Thermo Fisher) in 3% normal goat serum
979 with PBS. F-Actin was decorated using Phalloidin (1:300). Confocal images were acquired from Zeiss
980 LSM710 confocal microscope and images were processed using ImageJ 1.53t software.

981

982 **Data Availability**

983 Genetic and phenotypic data for the 100KGP study participants are available through the Genomics
984 England Research Environment via application at [https://www.genomicsengland.co.uk/join-a-gecip-](https://www.genomicsengland.co.uk/join-a-gecip-domain)
985 domain. PanelApp gene panels and evidence of associations were obtained using the PanelApp
986 application programming interface (<https://panelapp.genomicsengland.co.uk/api/docs/>) on the 20th
987 October 2021. CADD v1.5 (<https://cadd.gs.washington.edu/>), gnomAD v3.0
988 (<https://cadd.gs.washington.edu/>) and Ensembl v104 (<http://may2021.archive.ensembl.org/index.html>)
989 were used for variant annotation.

990

991 **Code Availability**

992 The rsrv tool and Rareservoir schema are available from <https://github.com/turrogroup/rsrv>.

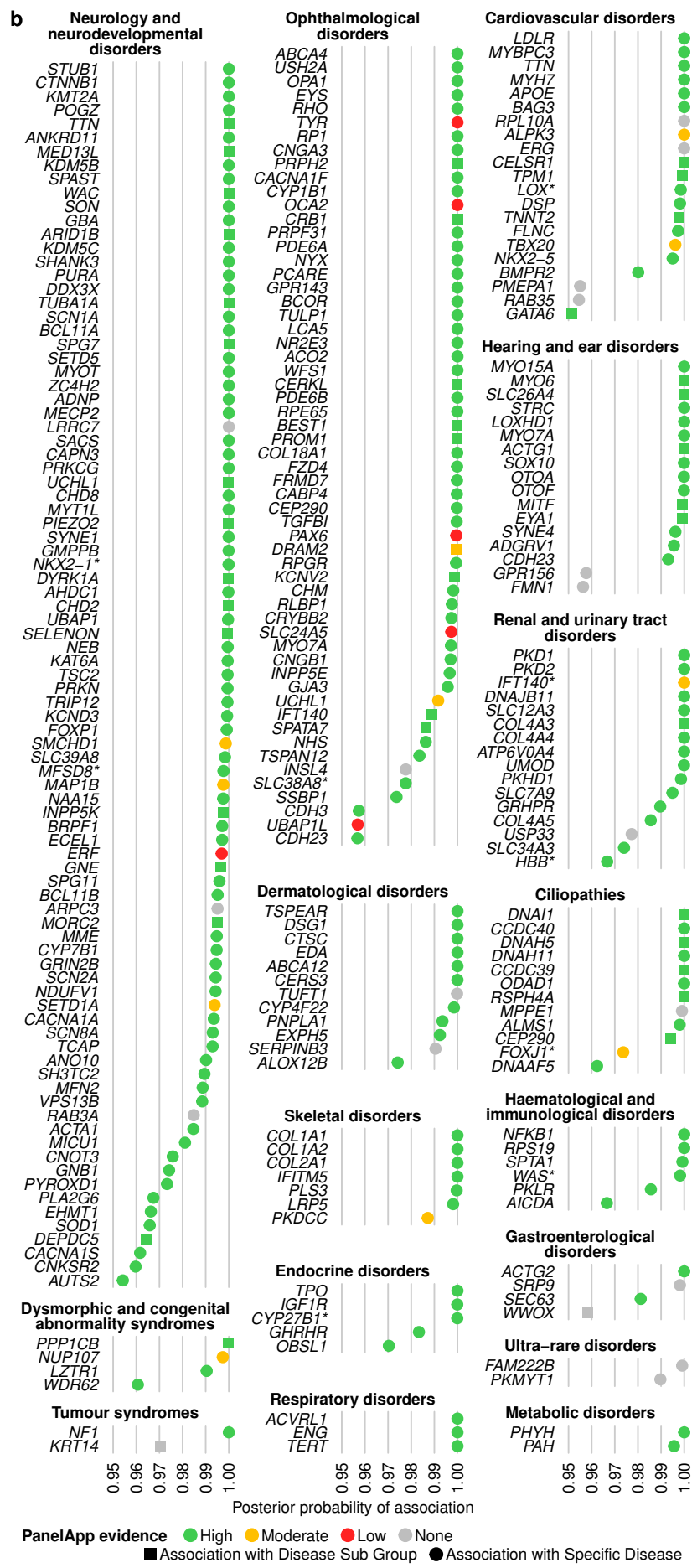
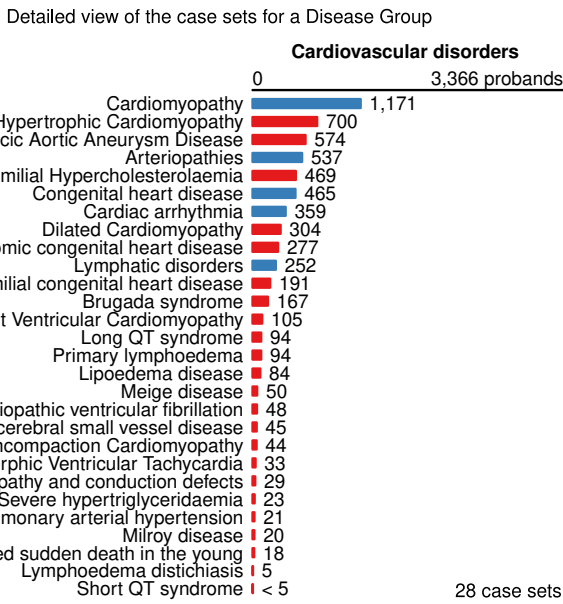
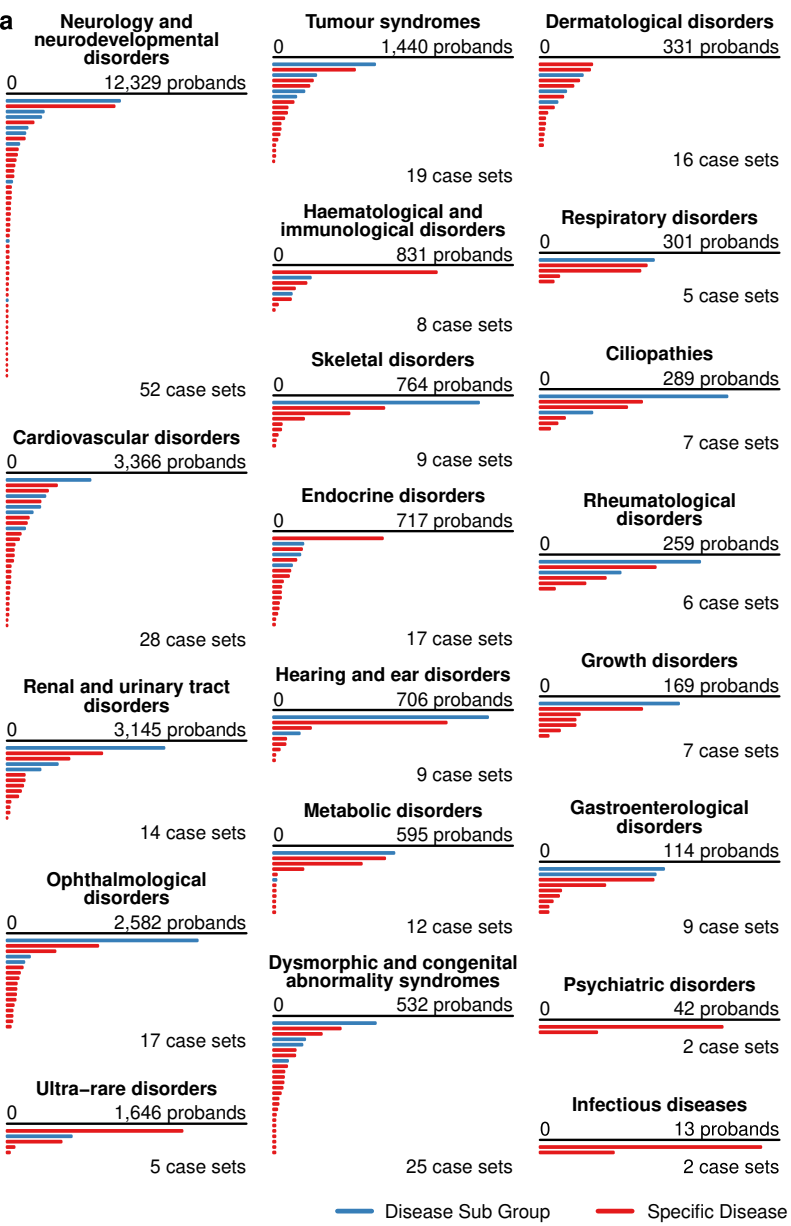
993

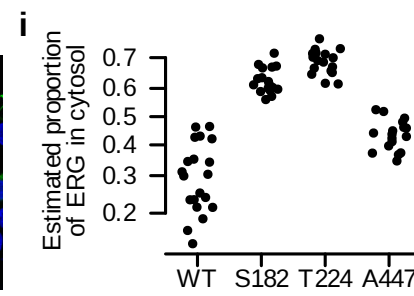
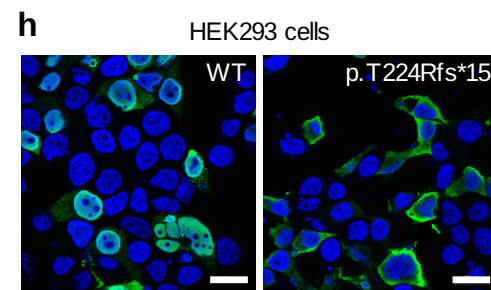
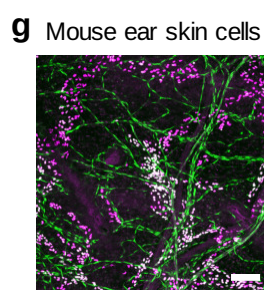
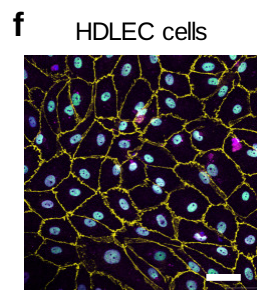
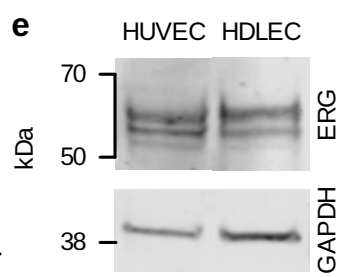
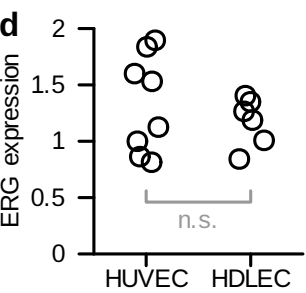
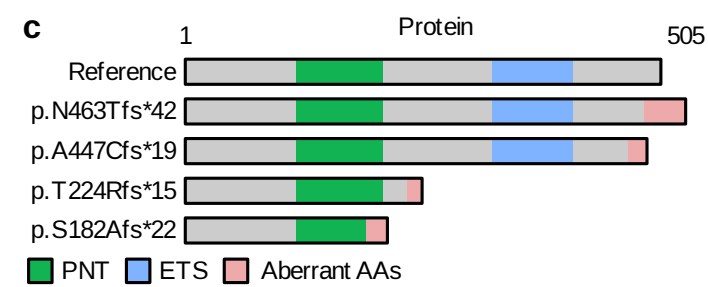
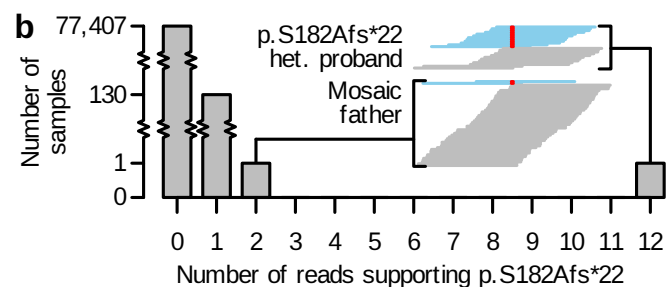
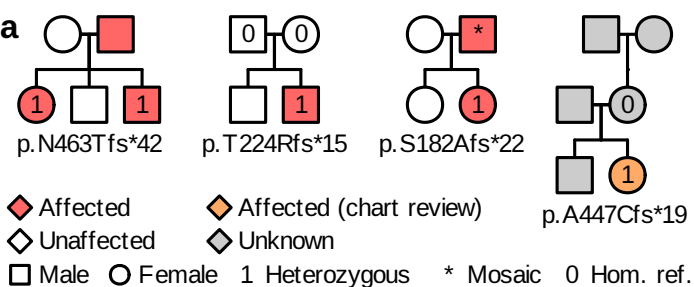
994 **Methods-only references**

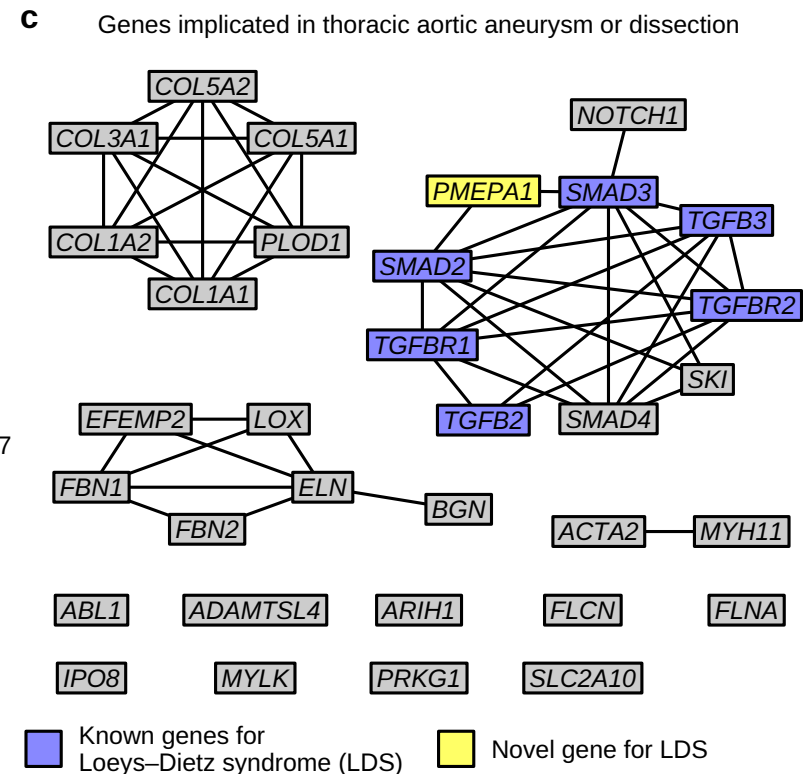
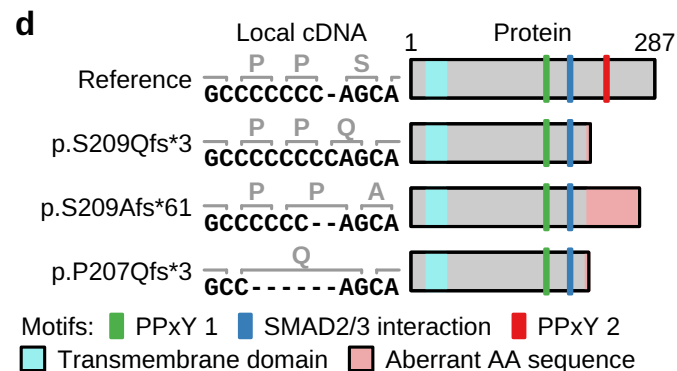
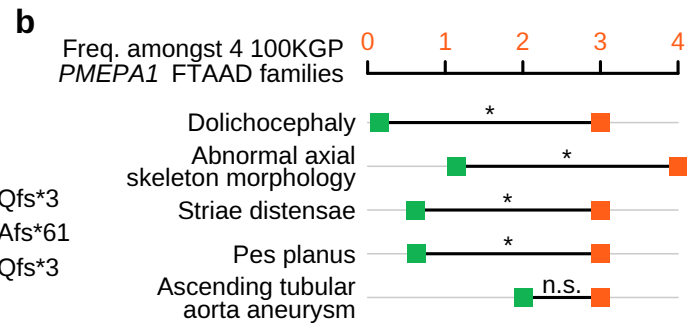
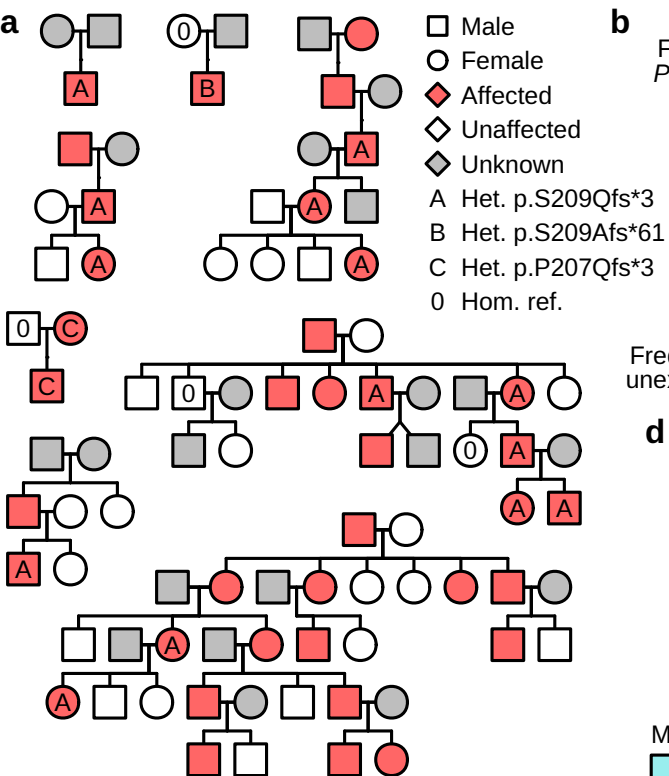
- 995 47. Li H et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*. 2009; 25(16):2078--
996 2079.
- 997 48. Morales J et al. A joint NCBI and EMBL-EBI transcript set for clinical genomics and research. *Nature*.
998 2022; 604(7905):310--315.
- 999 49. Variant QC for 100,000 Genomes Project merged VCF files (2022). [https://research-](https://research-help.genomicsengland.co.uk/display/GERE/Site+QC%2C+FILTER+and+INFO+Fields)
1000 [help.genomicsengland.co.uk/display/GERE/Site+QC%2C+FILTER+and+INFO+Fields](https://research-help.genomicsengland.co.uk/display/GERE/Site+QC%2C+FILTER+and+INFO+Fields)

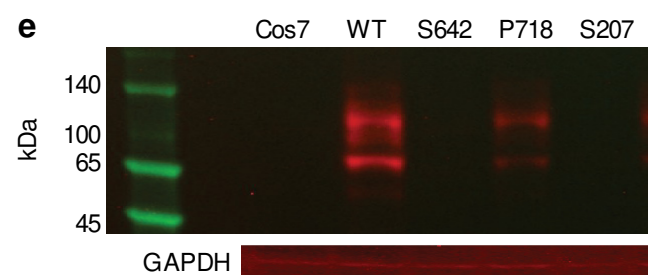
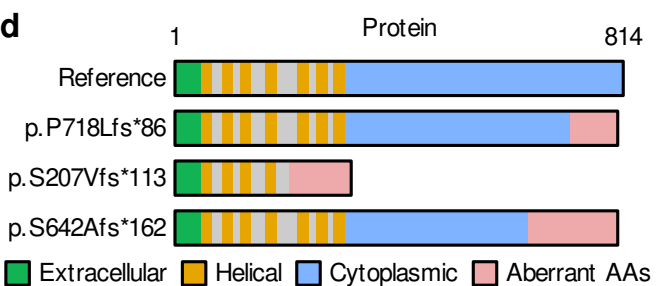
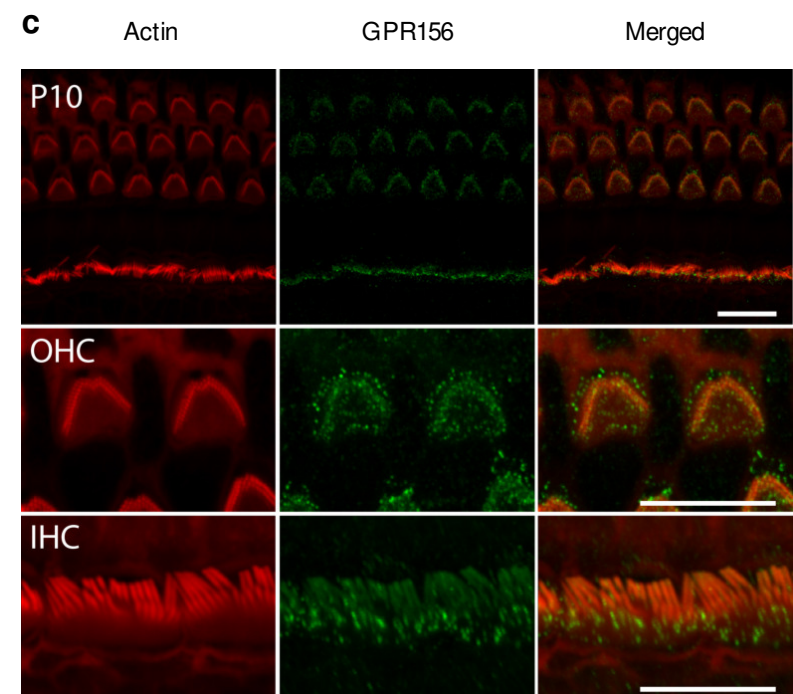
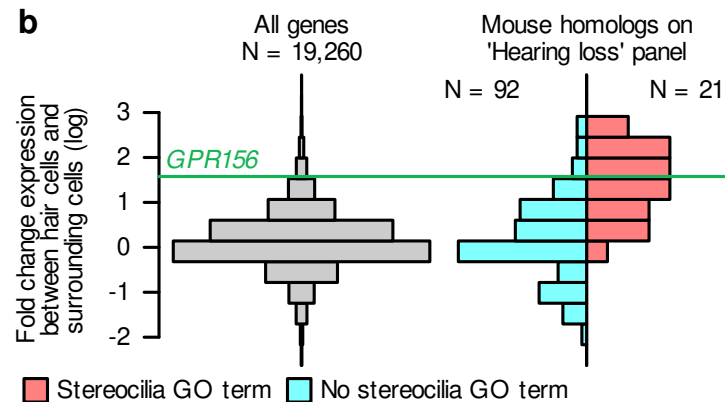
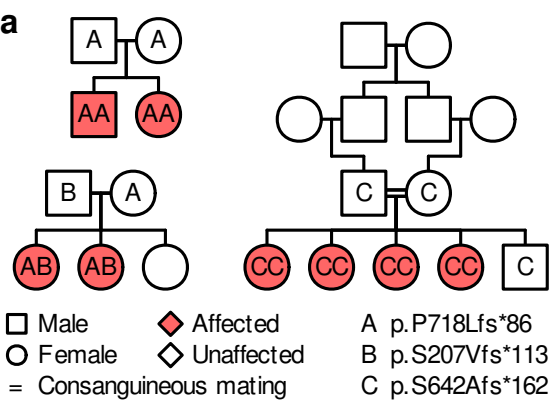
1001

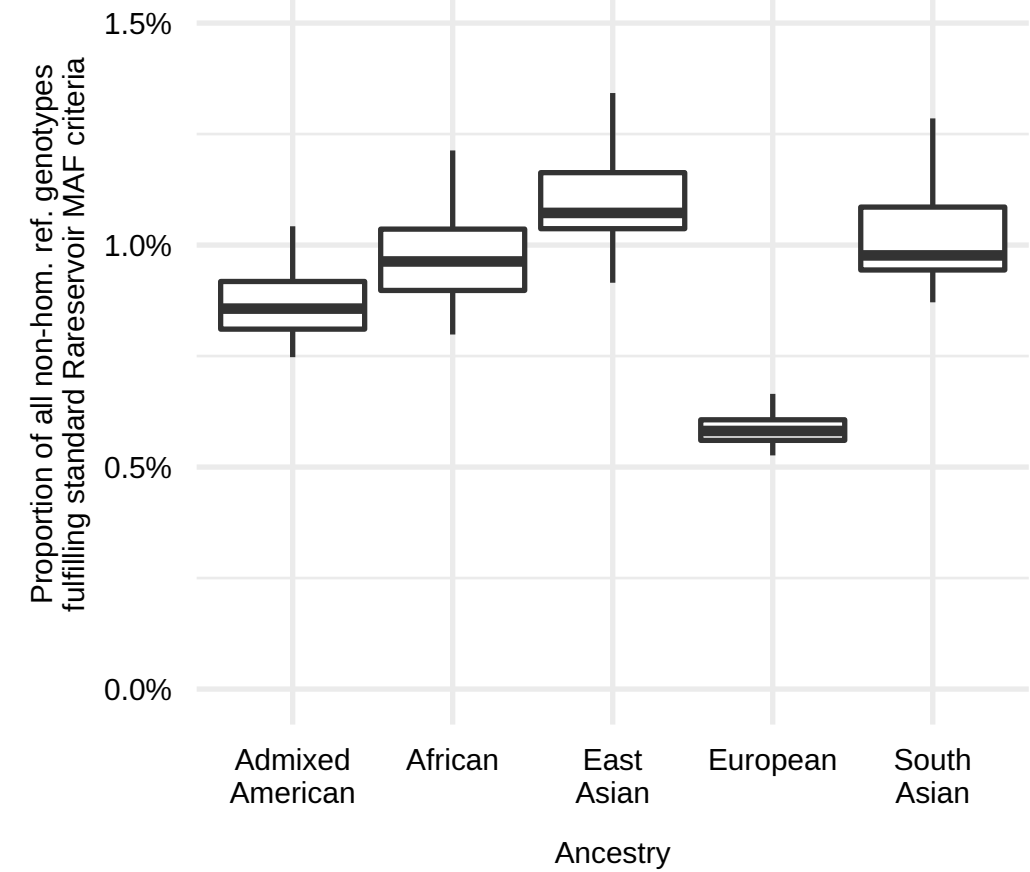
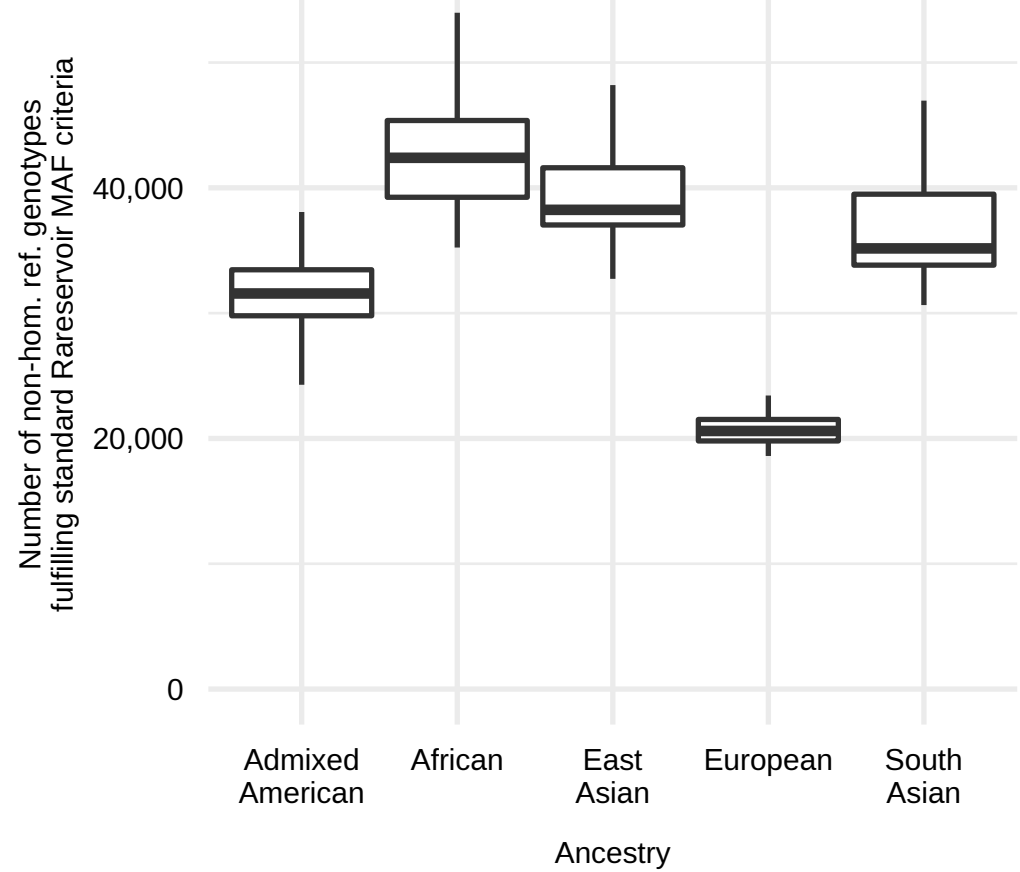
1002

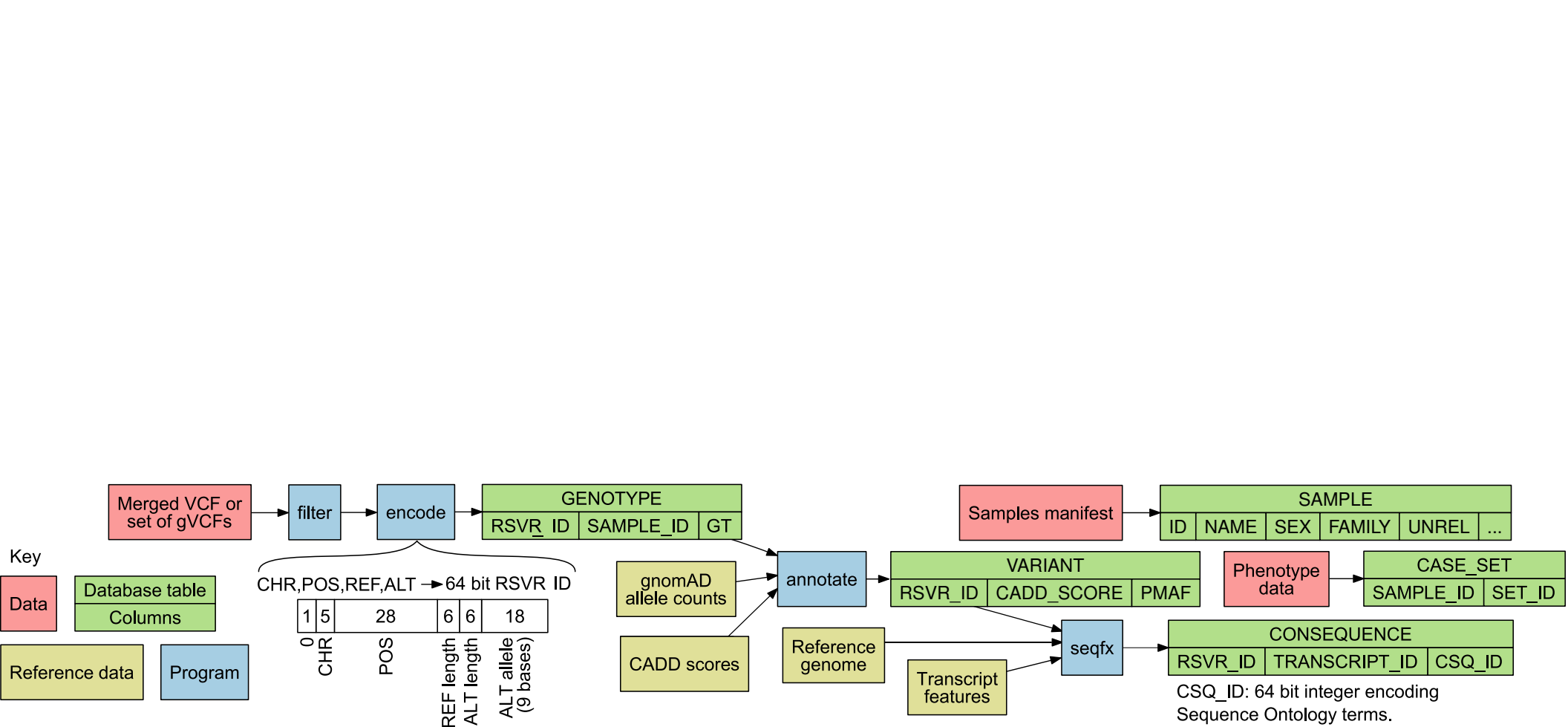


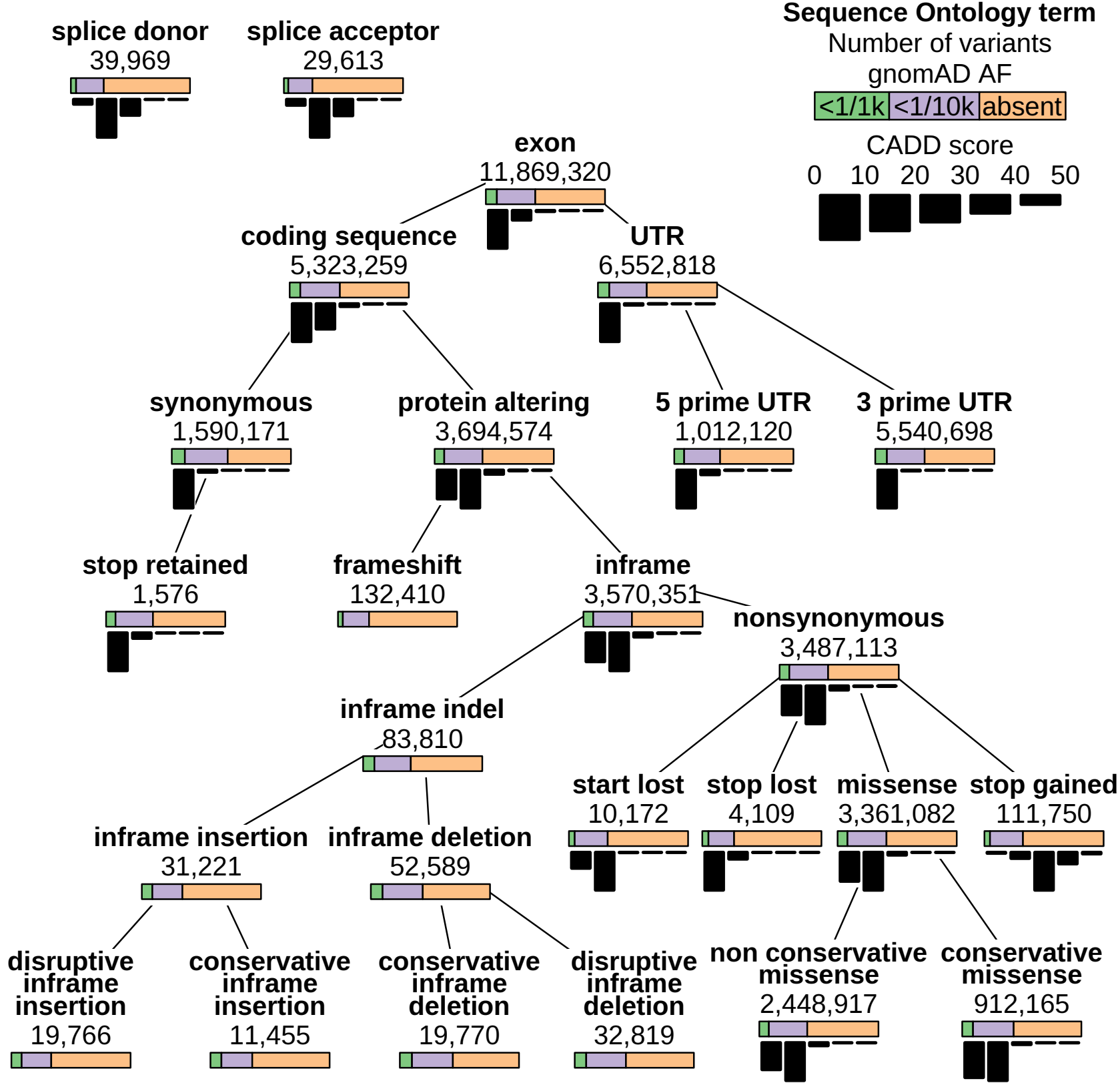




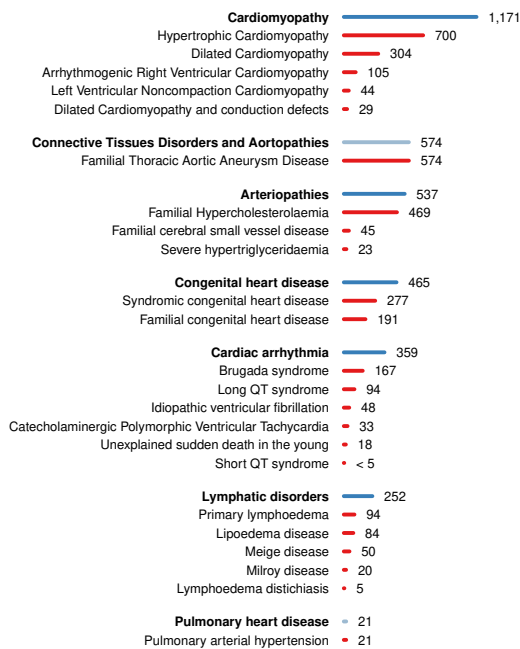


a**b**

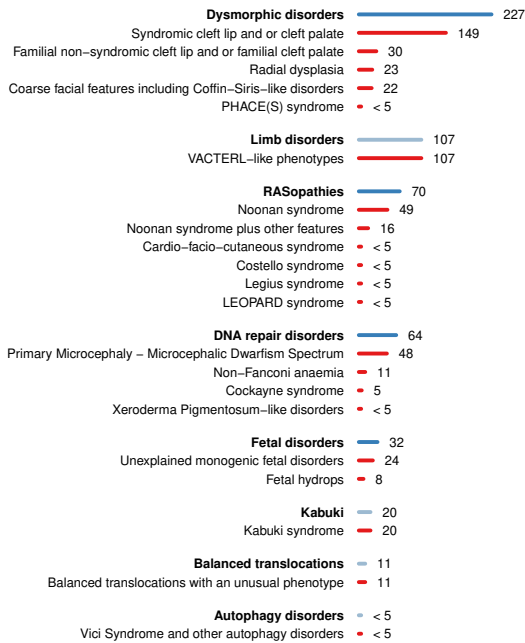




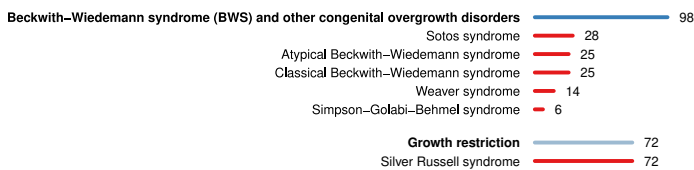
Cardiovascular disorders



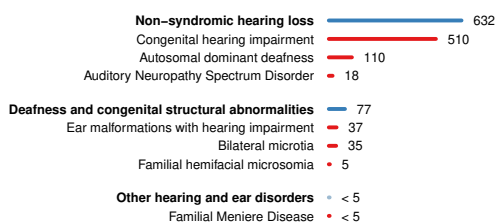
Dysmorphic and congenital abnormality syndromes



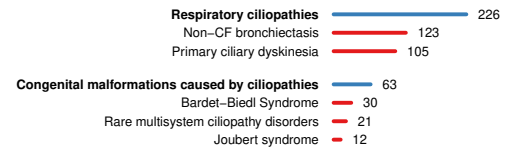
Growth disorders



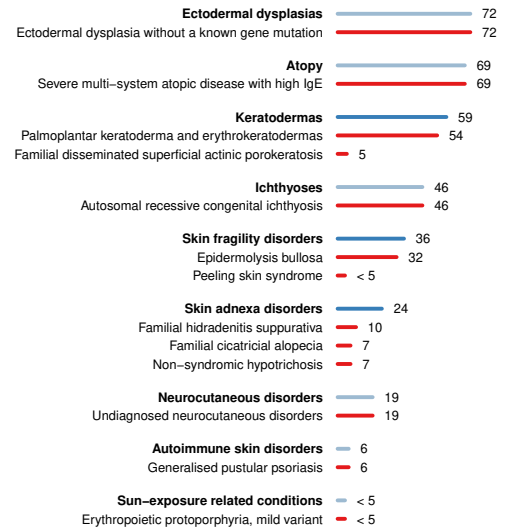
Hearing and ear disorders



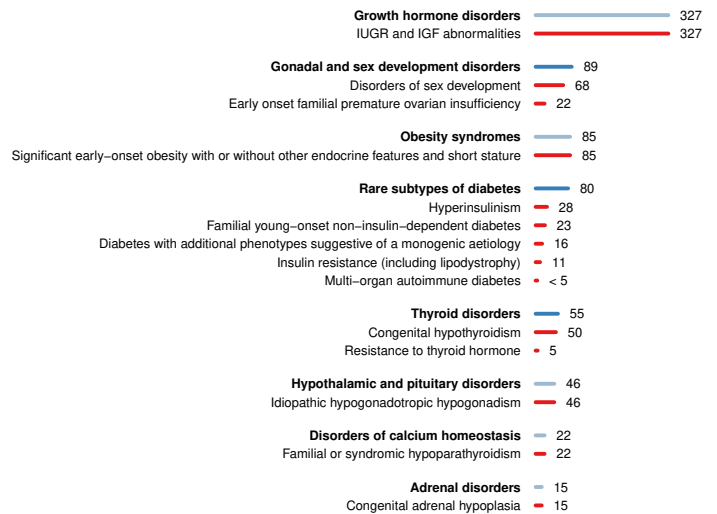
Ciliopathies



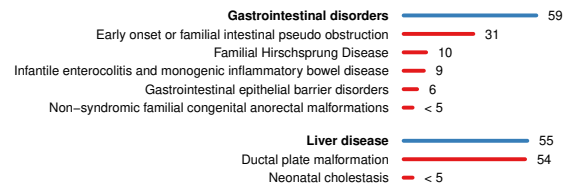
Dermatological disorders



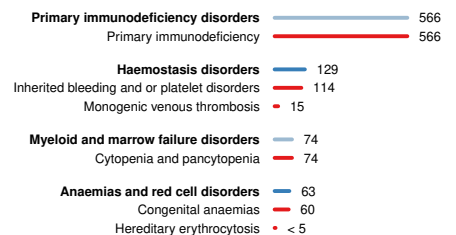
Endocrine disorders



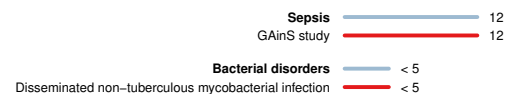
Gastroenterological disorders



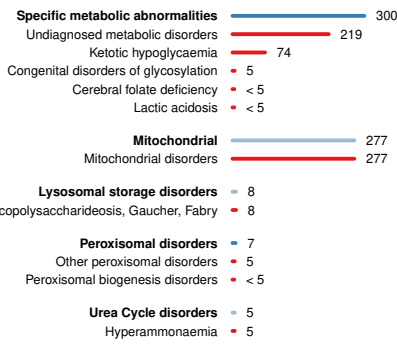
Haematological and immunological disorders



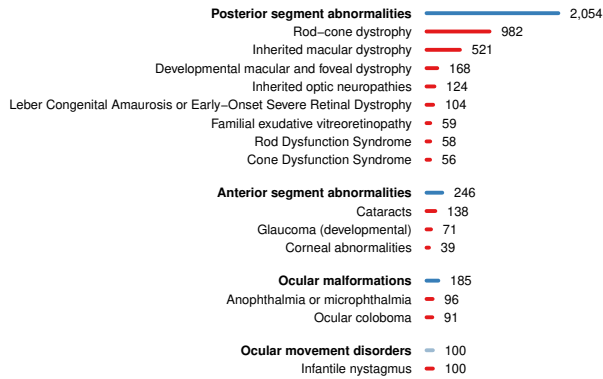
Infectious diseases



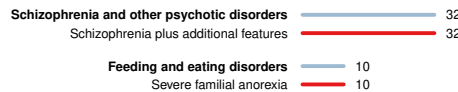
Metabolic disorders



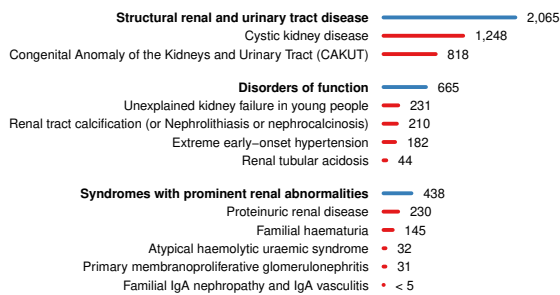
Ophthalmological disorders



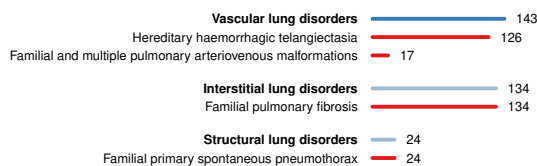
Psychiatric disorders



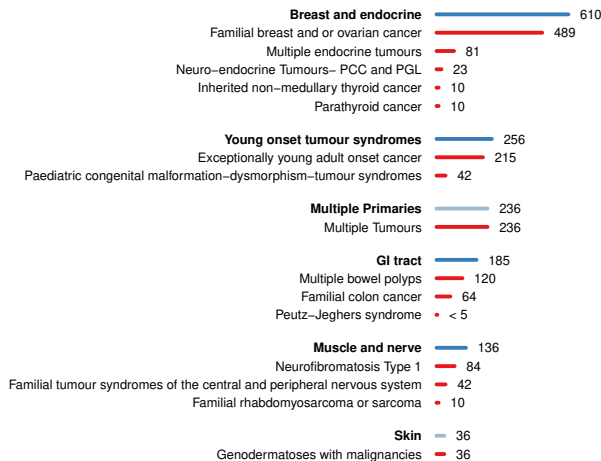
Renal and urinary tract disorders



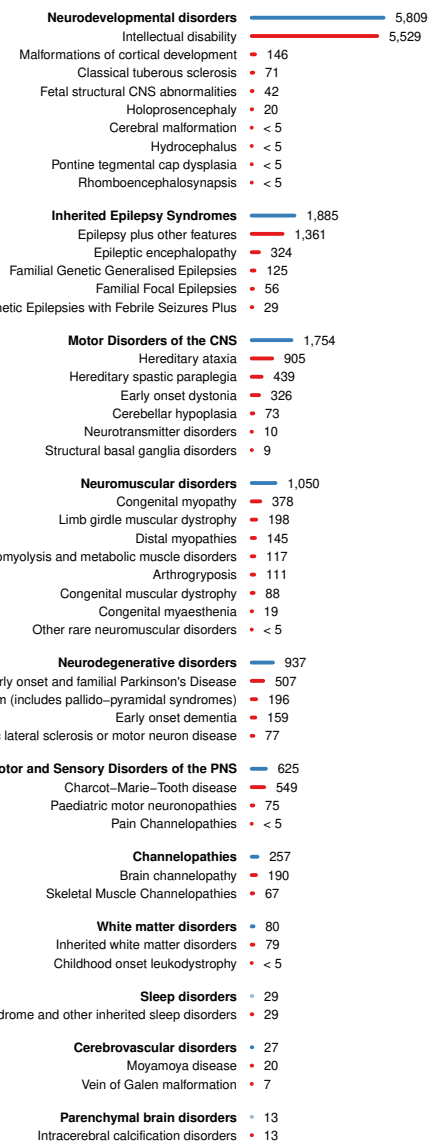
Respiratory disorders



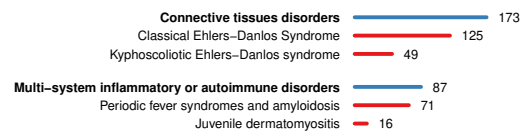
Tumour syndromes



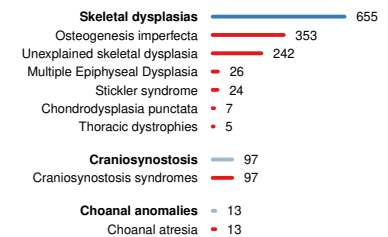
Neurology and neurodevelopmental disorders



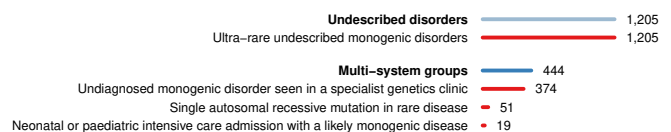
Rheumatological disorders



Skeletal disorders

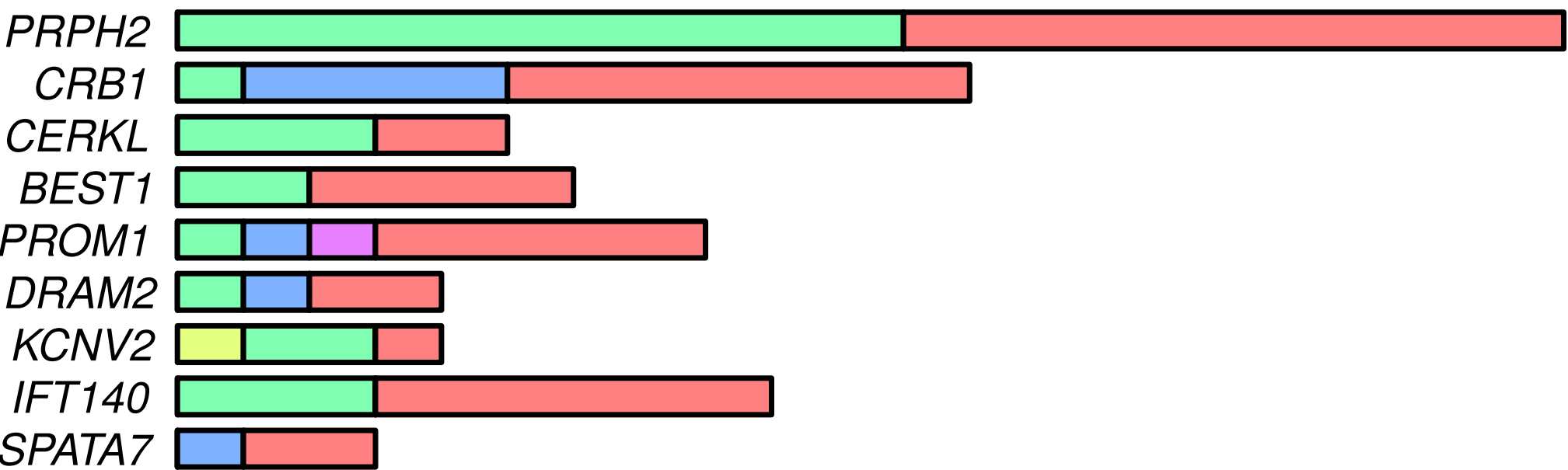


Ultra-rare disorders



Number of PSA cases with an inferred pathogenic configuration of alleles

0 5 10 15 20



Specific Disease:

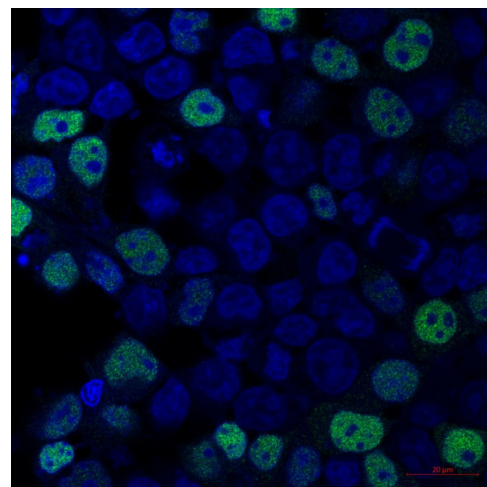
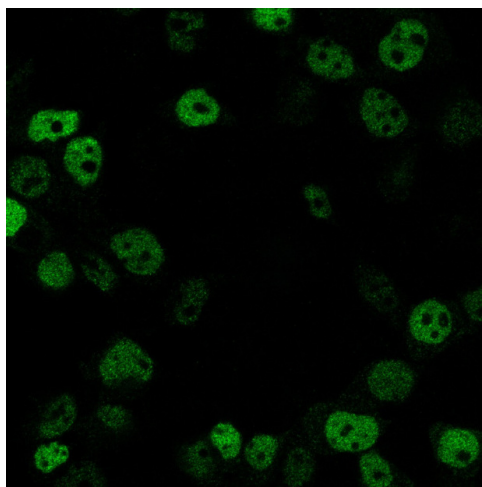
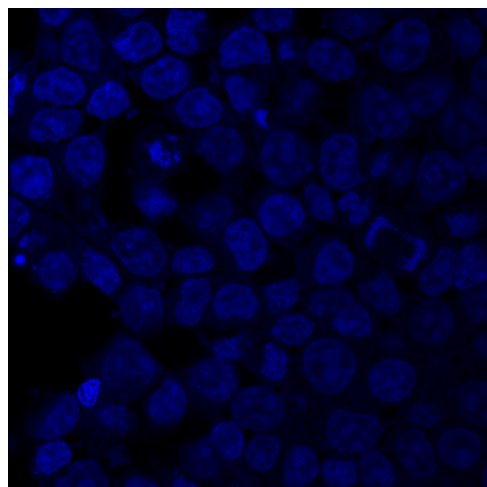
- Rod-cone dystrophy
- Rod Dysfunction Syndrome
- Leber Congenital Amaurosis or Early-Onset Severe Retinal Dystrophy
- Inherited macular dystrophy
- Cone Dysfunction Syndrome

DAPI

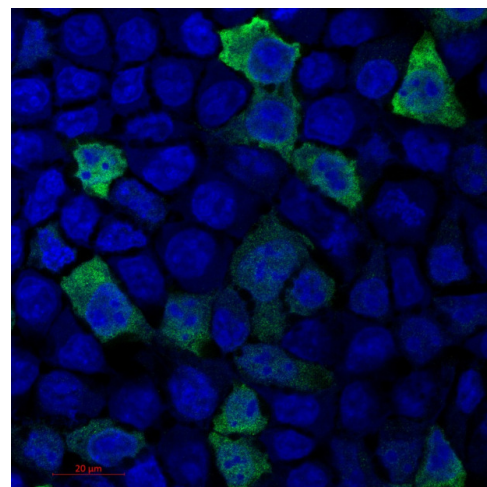
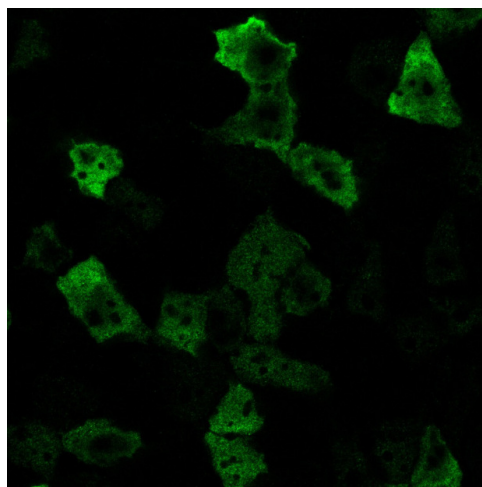
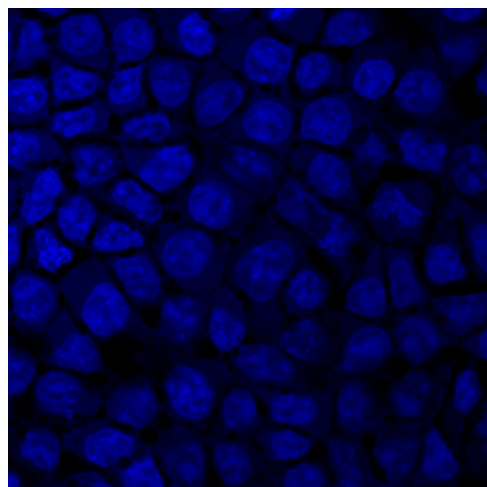
ERG

Merged

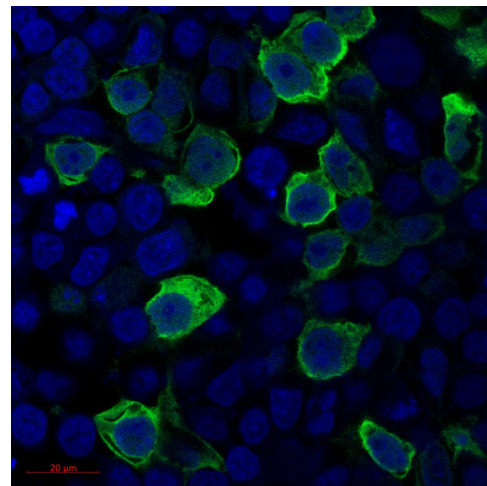
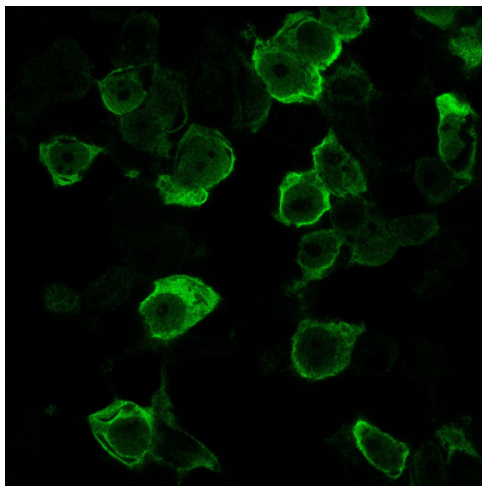
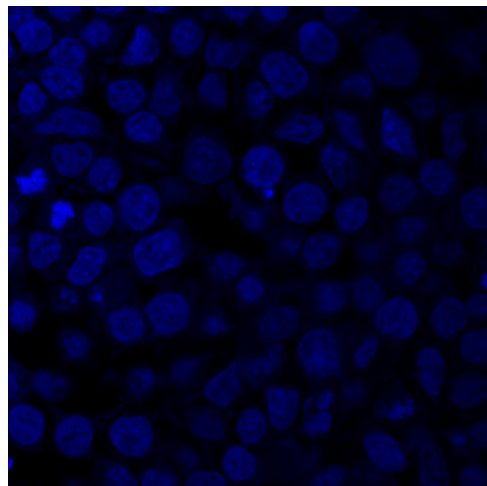
WT



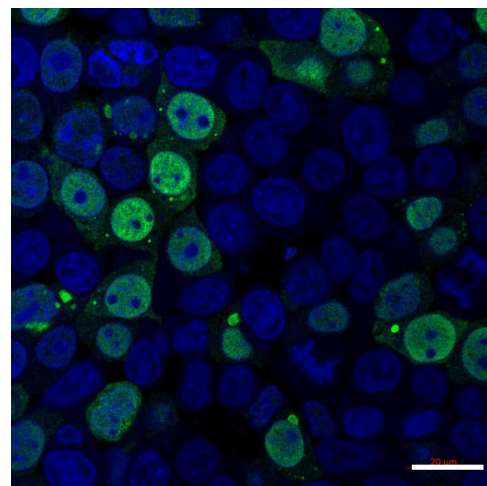
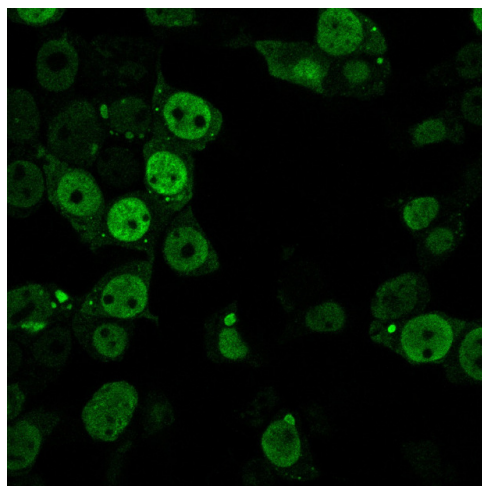
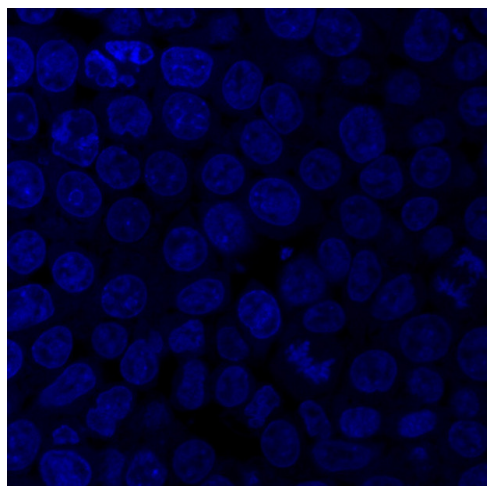
p.S182Afs*22

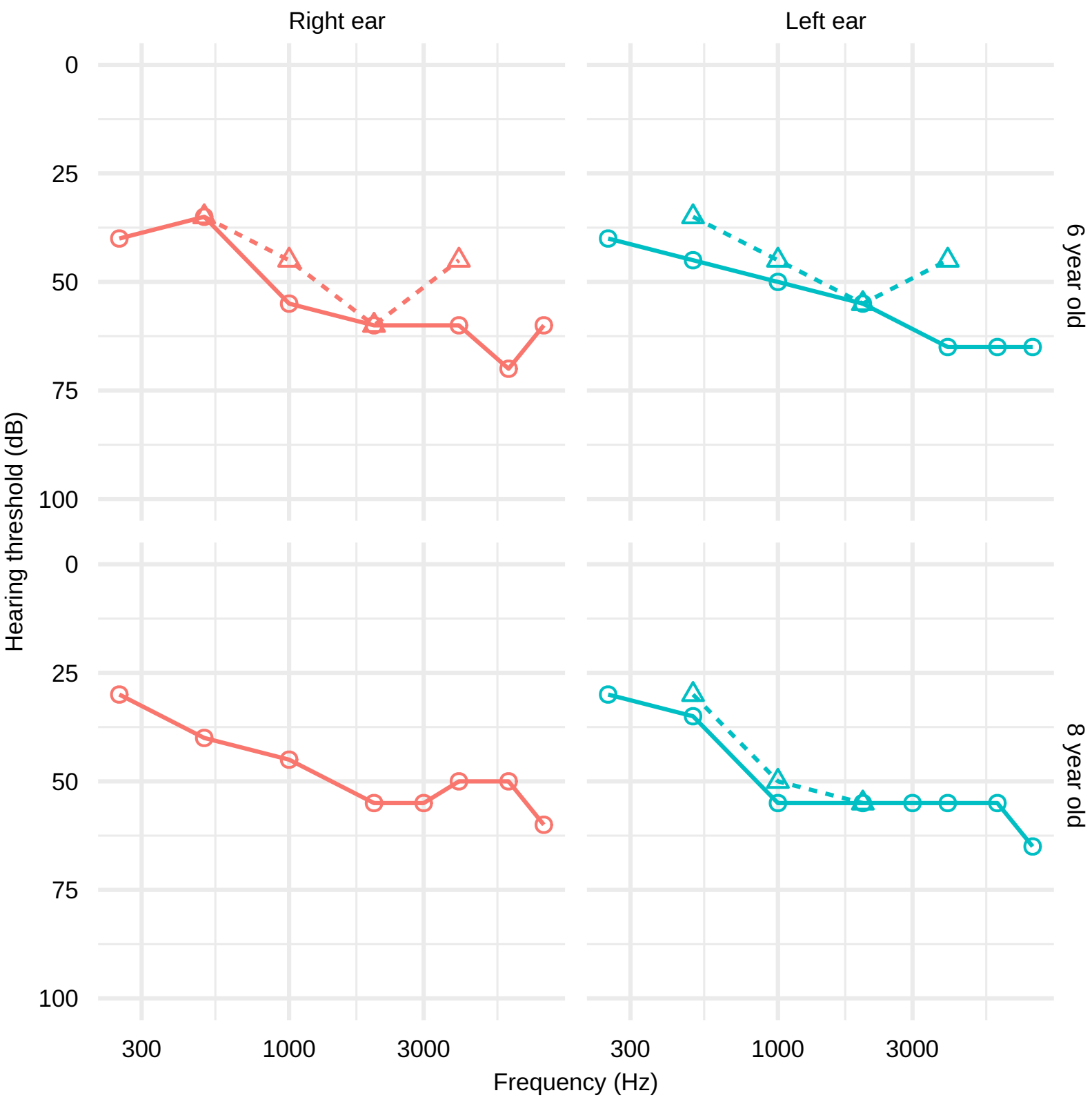


p.T224Rfs*15



p.A447Cfs*19





Type Air Bone