

<sup>1</sup> Cell type-specific contextualisation of the human phenome: towards  
<sup>2</sup> the systematic treatment of all rare diseases

<sup>3</sup> Brian M. Schilder      Kitty B. Murphy      Hiranyamaya Dash      Yichun Zhang  
<sup>4</sup> Robert Gordon-Smith      Jai Chapman      Momoko Otani      Nathan G. Skene

<sup>5</sup> 2025-08-08

## 6 Abstract

7 Rare diseases (RDs) are an extremely heterogeneous and underserved category of medical conditions. While  
8 the majority of RDs are strongly genetic, it remains largely unknown via which physiological mechanisms  
9 genetics cause RD. Therefore, we sought to systematically characterise the cell type-specific mechanisms  
10 underlying all RD phenotypes with a known genetic cause by leveraging the Human Phenotype Ontology  
11 and transcriptomic single-cell atlases of the entire human body from embryonic, foetal, and adult samples.  
12 In total we identified significant associations between 201 cell types and 9,575/11,028 (86.7%) unique pheno-  
13 types across 8,628 RDs, greatly increasing the collective knowledge of RD phenotype-cell type mechanisms.  
14 Next, we sought to systematically identify phenotypes in which the application of these results would have  
15 the greatest clinical impact based on metrics of severity (e.g. lethality, motor/mental impairment) and com-  
16 patibility with gene therapy (e.g. cell type specificity, postnatal treatability). Furthermore, we have made  
17 these results entirely reproducible and freely accessible to the global community to maximise their impact,  
18 including an interactive web portal (<https://neurogenomics-ukdri.dsi.ic.ac.uk/>). To summarise, this work  
19 represents a significant step forward in the mission to treat patients across an extremely diverse spectrum  
20 of serious RDs.

## 21 Introduction

22 While rare diseases (RDs) are individually uncommon, they collectively account for an enormous global  
23 disease burden with over 10,000 recognised RDs affecting at least 300-400 million people globally<sup>1</sup> (1 in  
24 10-20 people)<sup>2</sup>. Over 75% of RDs primarily affect children with a 30% mortality rate by five years of age<sup>3</sup>.  
25 Despite the prevalence and severity of RDs, patients suffering from these conditions are vastly underserved  
26 due to several contributing factors. First, diagnosis is extremely challenging due to the highly variable  
27 clinical presentations of many of these diseases. The diagnostic odyssey can take patients and their families  
28 decades, with an average time to diagnosis of five years<sup>4</sup>. Of those, ~46% receive at least one incorrect  
29 diagnosis and over 75% of all patients never receive any diagnosis<sup>5</sup>. Second, prognosis is also made difficult  
30 by high variability in disease course and outcomes which makes matching patients with effective and timely  
31 treatment plans even more challenging. Finally, even for patients who receive an accurate diagnosis/prognosis,  
32 treatments are currently only available for less than 5% of all RDs<sup>6</sup>. In addition to the scientific challenges of  
33 understanding RDs, there are strong financial disincentives for pharmaceutical and biotechnology companies  
34 to develop expensive therapeutics for exceedingly small RD patient populations with little or no return  
35 on investment<sup>7,8</sup>. Those that have been produced are amongst the world's most expensive drugs, greatly  
36 limiting patients' ability to access it<sup>9,10</sup>. New high-throughput approaches for the development of rare disease  
37 therapeutics could greatly reduce costs (for manufacturers and patients) and accelerate the timeline from  
38 discovery to delivery.

39 A major challenge in both healthcare and scientific research is the lack of standardised medical terminology.

40 Even in the age of electronic healthcare records (EHR) much of the information about an individual's history  
41 is currently fractured across healthcare providers, often with differing nomenclatures for the same conditions.  
42 The Human Phenotype Ontology (HPO) is a hierarchically organised set of controlled clinical terms that  
43 provides a much needed common framework by which clinicians and researchers can precisely communi-  
44 cate patient conditions<sup>14</sup>. The HPO spans all domains of human physiology and currently describes 18,082  
45 phenotypes across 10,300 RDs. Each phenotype and disease is assigned its own unique identifier and organ-  
46 ised as a hierarchical graph, such that higher-level terms describe broad phenotypic categories or *branches*  
47 (e.g. *HP:0033127*: ‘Abnormality of the musculoskeletal system’ which contains 4,495 unique phenotypes)  
48 and lower-level terms describe increasingly precise phenotypes (e.g. *HP:0030675*: ‘Contracture of proximal  
49 interphalangeal joints of 2nd-5th fingers’). It has already been integrated into healthcare systems and clinical  
50 diagnostic tools around the world, with increasing adoption over time<sup>11</sup>. Standardised frameworks like the  
51 HPO also allow us to aggregate relevant knowledge about the molecular mechanisms underlying each RD.  
  
52 Over 80% of RDs have a known genetic cause<sup>15,16</sup>. Since 2008, the HPO has been continuously updated  
53 using curated knowledge from the medical literature, as well as by integrating databases of expert validated  
54 gene-phenotype relationships, such as OMIM<sup>17-19</sup>, Orphanet<sup>20,21</sup>, and DECIPHER<sup>22</sup>. Mappings between  
55 HPO terms to other commonly used medical ontologies (e.g. SNOMED CT<sup>23</sup>, UMLS<sup>24,25</sup>, ICD-9/10/11<sup>26</sup>)  
56 make the HPO even more valuable as a clinical resource (provided in Mappings section of Methods). Many of  
57 these gene annotations are manually or semi-manually curated by expert clinicians from case reports of rare  
58 disease patients in which the causal gene is identified through whole exome or genome sequencing. Currently,  
59 the HPO contains gene annotations for 11,047 phenotypes across 8,631 diseases. Yet genes alone do not tell  
60 the full story of how RDs come to be, as their expression and functional relevance varies drastically across  
61 the multitude of tissues and cell types contained within the human body. Our knowledge of the physiological  
62 mechanisms via which genetics cause pathogenesis is lacking for most RDs, severely hindering our ability to  
63 effectively diagnose, prognose and treat RD patients.  
  
64 Our knowledge of cell type-specific biology has exploded over the course of the last decade and a half,  
65 with numerous applications in both scientific and clinical practices<sup>27-29</sup>. In particular, single-cell RNA-seq  
66 (scRNA-seq) has allowed us to quantify the expression of every gene (i.e. the transcriptome) in individual  
67 cells. More recently, comprehensive single-cell transcriptomic atlases across tissues have also emerged<sup>30,31</sup>.  
68 In particular, the Descartes Human<sup>32</sup> and Human Cell Landscape<sup>33</sup> projects provide comprehensive multi-  
69 system scRNA-seq atlases in embryonic, foetal, and adult human samples from across the human body.  
70 These datasets provide data-driven gene signatures for hundreds of cell subtypes. Given that many disease-  
71 associated genes are expressed in some cell types but not others, we can infer that disruptions to these genes  
72 will have varying impact across cell types. By comparing the aggregated disease gene annotations with  
73 cell type-specific expression profiles, we can therefore uncover the cell types and tissues via which diseases  
74 mediate their effects.

75 Here, we combine and extend several of the most comprehensive genomic and transcriptomic resources  
76 currently available to systematically uncover the cell types underlying granular phenotypes across 8,628  
77 diseases Fig. 1. Conversely, this approach also allows us to better understand the roles of understudied cell  
78 types by observing which phenotypes they tend to associate with. For example, the original authors proposed  
79 that a novel class *AFB+/ALB+* cells may represent hepatoblasts circulating through the bloodstream during  
80 foetal development<sup>34</sup>. Our results support this hypothesis as *AFB+/ALB+* cells were significantly associated  
81 with 12 liver-related phenotypes, as well as 58 blood-related phenotypes.

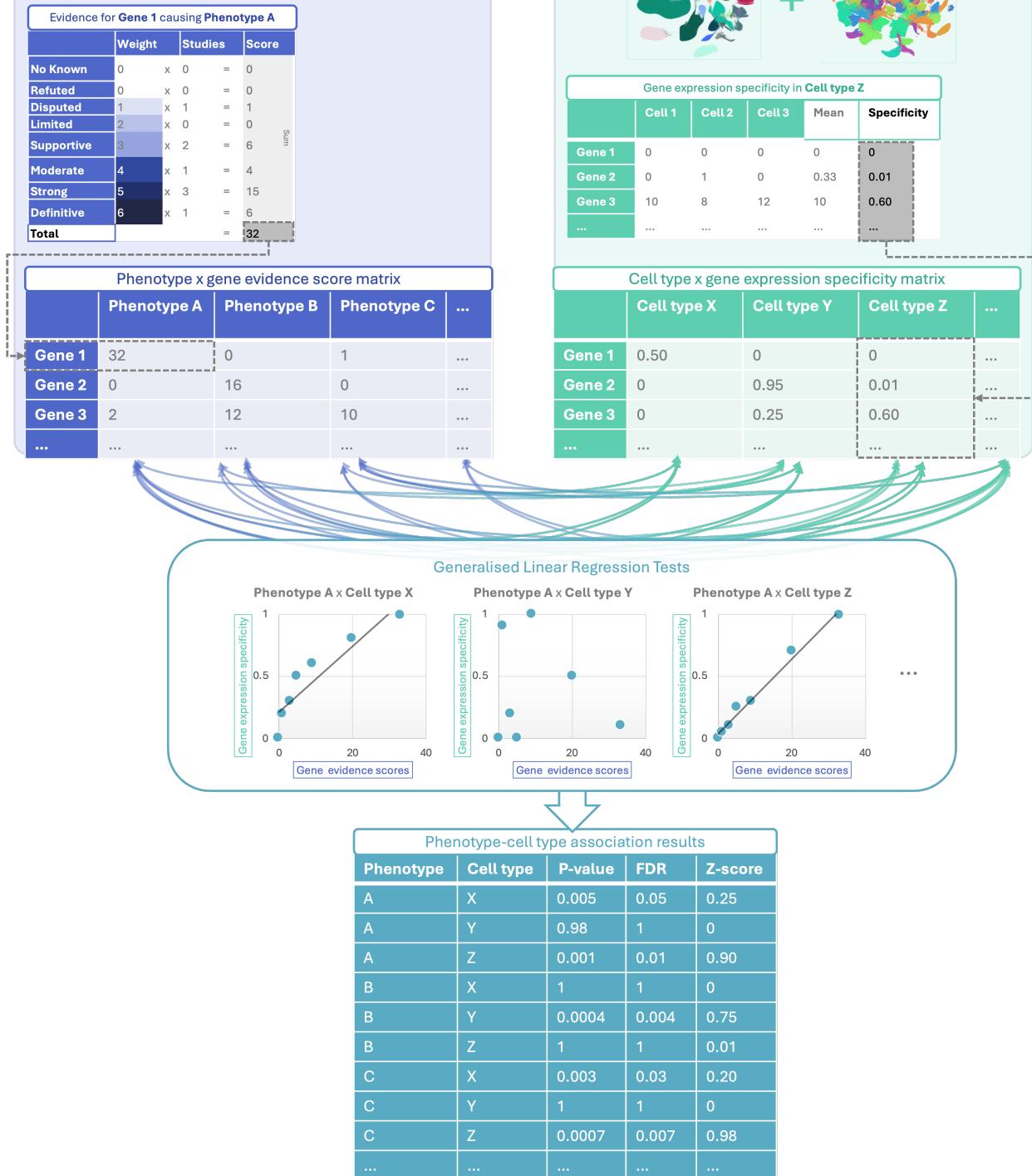
82 Beyond making discoveries in basic science, our phenome-wide cell type associations provide essential context  
83 for the development of novel therapeutics, especially gene therapy modalities such as adeno-associated viral  
84 (AAV) vectors in which advancement have been made in their ability selectively target specific cell types<sup>35,36</sup>.  
85 Precise knowledge of relevant cell types and tissues causing the disease can improve safety by minimising  
86 harmful side effects in off-target cell types and tissues. It can also enhance efficacy by efficiently delivering  
87 expensive therapeutic payloads to on-target cell types and tissues. For example, if a phenotype primarily  
88 effects retinal cells, then the gene therapy would be optimised for delivery to retinal cells of the eye. Using  
89 this information, we developed a high-throughput pipeline for comprehensively nominating cell type-resolved  
90 gene therapy targets across thousands of RD phenotypes. As a prioritisation tool, we sorted these targets  
91 based on the severity of their respective phenotypes, using a generative AI-based approach<sup>37</sup>. Together,  
92 our study dramatically expands the available knowledge of the cell types, organ systems and life stages  
93 underlying RD phenotypes.

## 94 Results

### 95 Phenotype-cell type associations

96 In this study we systematically investigated the cell types underlying phenotypes across the HPO. We hy-  
97 pothesised that genes which are specifically expressed in certain cell types will be most relevant for the proper  
98 functioning of those cell types. Thus, phenotypes caused by disruptions to specific genes will have greater or  
99 lesser effects across different cell types. To test this, we computed associations between the weighted gene  
100 lists for each phenotype with the gene expression specificity for each cell type in our transcriptomic reference  
101 atlases.

102 More precisely, for each phenotype we created a list of associated genes weighted by the strength of the  
103 evidence supporting those associations, imported from the Gene Curation Coalition (GenCC)<sup>38</sup>. Analogously,  
104 we created gene expression profiles for each cell type in our scRNA-seq atlases and then applied normalisation  
105 to compute how specific the expression of each gene is to each cell type. To assess consistency in the  
106 phenotype-cell type associations, we used multiple scRNA-seq atlases: Descartes Human (~4 million single-  
107 nuclei and single-cells from 15 fetal tissues)<sup>32</sup> and Human Cell Landscape (~703,000 single-cells from 49  
108 embryonic, fetal and adult tissues)<sup>33</sup>. We ran a series of linear regression models to test for the relationship



**Figure 1: Multi-modal data fusion reveals the cell types underlying thousands of human phenotypes.** Schematic overview of study design in which we numerically encoded the strength of evidence linking each gene and each phenotype (using the Human Phenotype Ontology and GenCC databases). We then created gene signature profiles for all cell types in the Descartes Human and Human Cell Landscape scRNA-seq atlases. Finally, we iteratively ran generalised linear regression tests between all pairwise combinations of phenotype gene signatures and cell type gene signatures. The resulting associations were then used to nominate cell type-resolved gene therapy targets for thousands of rare diseases.

109 between every unique combination of phenotype and cell type. We applied multiple testing correction to  
110 control the false discovery rate (FDR) across all tests.

111 Within the results using the Descartes Human single-cell atlas, 19,929/ 848,078 (2.35%) tests across 77/  
112 77 (100%) cell types and 7,340/11,047 (66.4%) phenotypes revealed significant phenotype-cell type asso-  
113 ciations after multiple-testing correction (FDR<0.05). Using the Human Cell Landscape single-cell atlas,  
114 26,585/1,358,916 (1.96%) tests across 124/124 (100%) cell types and 9,049/11,047 (81.9%) phenotypes showed  
115 significant phenotype-cell type associations (FDR<0.05). The median number of significantly associated phe-  
116 notypes per cell type was 252 (Descartes Human) and 200 (Human Cell Landscape), respectively. Overall,  
117 using the Human Cell Landscape reference yielded a greater percentage of phenotypes with at least one  
118 significant cell type association than the Descartes Human reference. This is expected at the Human Cell  
119 Landscape contains a greater diversity of cell types across multiple life stages (embryonic, fetal, adult).

120 Across both single-cell references, the median number of significantly associated cell types per phenotype was  
121 3, suggesting reasonable specificity of the testing strategy. Within the HPO, 8,628/8,631 (~100%) of diseases  
122 gene annotations showed significant cell type associations for at least one of their respective phenotypes. A  
123 summary of the phenome-wide results stratified by single-cell atlas can be found in Table 2.

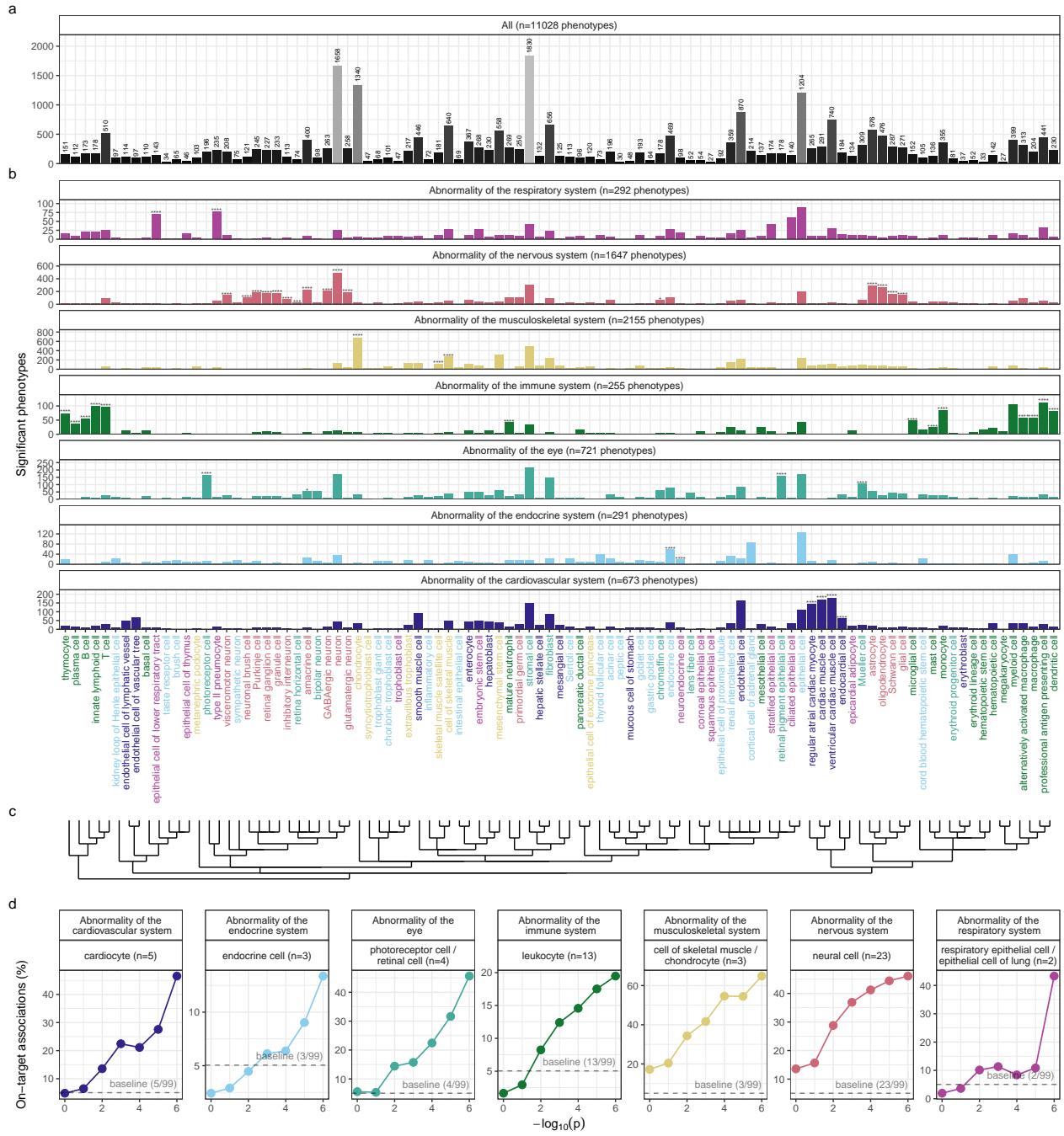
#### 124 Validation of expected phenotype-cell type relationships

125 We intuitively expect that abnormalities of an organ system will often be driven by cell types within that  
126 system. The HPO has broad categories at the higher level of the ontology, enabling us to systematically test  
127 this. For example, phenotypes associated with the heart should generally be caused by cell types of the heart  
128 (i.e. cardiocytes), while abnormalities of the nervous system should largely be caused by neural cells. There  
129 will of course be exceptions to this. For example, some immune disorders can cause intellectual disability  
130 through neurodegeneration. Nevertheless, it is reasonable to expect that abnormalities of the nervous system  
131 will be most often associated with neural cells. All cell types in our single-cell reference atlases were mapped  
132 onto the Cell Ontology (CL); a controlled vocabulary of cell types organised into hierarchical branches  
133 (e.g. neural cell include neurons and glia, which in turn include their respective subtypes).

134 Here, we consider a cell type to be *on-target* relative to a given HPO branch if it belongs to one of the  
135 matched CL branches (see Table 4). Within each high-level branch in the HPO shown in Fig. 2b, we tested  
136 whether each cell type was more often associated with phenotypes in that branch relative to those in all  
137 other branches (including those not shown). We then checked whether each cell type was overrepresented  
138 (at FDR<0.05) within its respective on-target HPO branch, where the number of phenotypes within that  
139 branch. Indeed, we found that all 7 HPO branches were disproportionately associated with on-target cell  
140 types from their respective organ systems.

141 In addition to binary metrics of a cell type being associated with a phenotype or not, we also used association  
142 test p-values as a proxy for the strength of the association. We hypothesized that the more significant the

143 association between a phenotype and a cell type, the more likely it is that the cell type is on-target for its  
144 respective HPO branch. To evaluate whether this, we grouped the association  $-\log_{10}(\text{p-values})$  into 6 bins.  
145 For each HPO-CL branch pairing, we then calculated the proportion of on-target cell types within each bin.  
146 We found that the proportion of on-target cell types increased with increasing significance of the association  
147 ( $\rho = 0.63$ ,  $p = 1.1 \times 10^{-6}$ ). For example, abnormalities of the nervous system with  $-\log_{10}(\text{p-values}) = 1$ ,  
148 only 16% of the associated cell types were neural cells. Whereas for those with  $-\log_{10}(\text{p-values}) = 6$ , 46%  
149 were neural cells despite the fact that this class of cell types only constituted 23% of the total cell types  
150 tested (i.e. the baseline). This shows that the more significant the association, the more likely it is that the  
151 cell type is on-target.



**(a) High-throughput analysis reveals cell types underlying thousands of rare disease phenotypes.** **a**, Some cell types are much more commonly associated with phenotypes than others. Bar height indicates the total number of significant phenotype enrichments per cell type (FDR<0.05) across all branches of the HPO. **b**, Analyses reveal expected and novel cell type associations within high-level HPO branches. Asterisks above each bar indicate whether that cell type was significantly more often enriched in that branch relative to all other HPO branches, including those not shown here, as a proxy for how specifically that cell type is associated with that branch; FDR<0.0001 (\*\*\*\*), FDR<0.001 (\*\*), FDR<0.01 (\*\*), FDR<0.05 (\*). **c**, Ontological relatedness of cell types in the Cell Ontology (CL)<sup>39</sup>. **d**, The proportion of on-target associations (*y*-axis) increases with greater test significance (*x*-axis). Percentage of significant phenotype associations with on-target cell types (second row of facet labels), respective to the HPO branch.

Figure 2

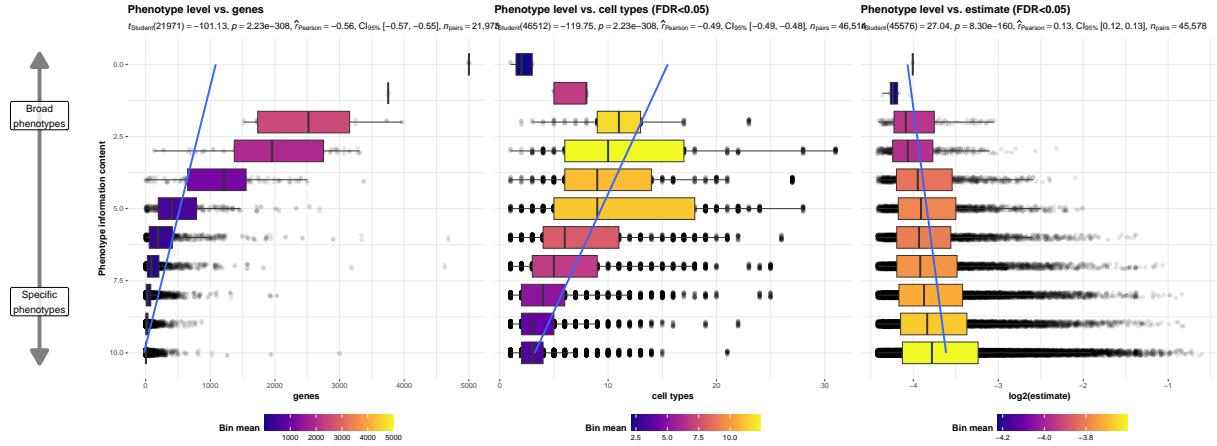
152 **Validation of inter- and intra-dataset consistency**

153 If our methodology works, it should yield consistent phenotype-cell type associations across different datasets.  
154 We therefore tested for the consistency of our results across the two single-cell reference datasets (Descartes  
155 Human vs. Human Cell Landscape) across the subset of overlapping cell types Fig. 11. In total there were  
156 142,285 phenotype-cell type associations to compare across the two datasets (across 10,945 phenotypes and  
157 13 cell types annotated to the exact same CL term. We found that the correlation between p-values of  
158 the two datasets was high ( $\rho=0.91$ ,  $p=5.7 \times 10^{-6}$ ). Within the subset of results that were significant in  
159 both single-cell datasets (FDR<0.05), we found that degree of correlation between the association effect  
160 sizes across datasets was even stronger ( $\rho =0.82$ ,  $p =5.7 \times 10^{-6}$ ). We also checked for the intra-dataset  
161 consistency between the p-values of the foetal and adult samples in the Human Cell Landscape, showing a  
162 very similar degree of correlation as the inter-dataset comparison ( $\rho =0.95$ ,  $p =5.0 \times 10^{-15}$ ). Together,  
163 these results suggest that our approach to identifying phenotype-cell type associations is highly replicable  
164 and generalisable to new datasets.

165 **More specific phenotypes are associated with fewer genes and cell types**

166 Higher levels of the ontology are broad classes of phenotype (e.g. ‘Abnormality of the nervous system’) while  
167 the lower levels can get very detailed (e.g. ‘Spinocerebellar atrophy’). The higher level phenotypes inherit  
168 all genes associated with lower level phenotypes, so naturally they have more genes than the lower level  
169 phenotypes (Fig. 3a;  $\rho =-0.56$ ,  $p =2.2 \times 10^{-308}$ ).

170 Next, we reasoned that the more detailed and specific a phenotype is, the more likely it is to be driven by  
171 one cell type. For example, while ‘Neurodevelopmental abnormality’ could plausibly be driven by any/all  
172 cell types in the brain, it is more likely that ‘Impaired visuospatial constructive cognition’ is driven by fewer  
173 cell types. This was indeed the case, as we observed a strongly significant negative correlation between the  
174 two variables (Fig. 3b;  $\rho =-0.49$ ,  $p =2.2 \times 10^{-308}$ ). We also found that the phenotype-cell type association  
175 effect size increased with greater phenotype specificity, reflecting the decreasing overall number of associated  
176 cell types at each ontological level (Fig. 3c;  $\rho =0.13$ ,  $p =8.3 \times 10^{-160}$ ).



(a) **More specific phenotypes are associated with fewer, more specific genes and cell types.** Information content (IC), is a normalised measure of ontology term specificity. Terms with lower IC represent the broadest HPO terms (e.g. ‘All’), while terms with higher IC indicate progressively more specific HPO terms (e.g. ‘Contracture of proximal interphalangeal joints of 2nd-5th fingers’). Box plots show the relationship between HPO phenotype IC and **a**, the number of genes annotated to each phenotype, **b**, the number of significantly enriched cell types, **c**, the effect sizes (absolute model  $R^2$  estimates after log-transformation) of significant phenotype-cell type association tests. Boxes are coloured by the mean value within each IC bin (after rounding continuous IC values to the nearest integer).

Figure 3

#### 177 Validation of phenotype-cell type associations using biomedical knowledge graphs

178 In order to validate our phenotype-cell type associations without the bias introduced by manually searching  
 179 literature that affirmed our discoveries, we use formalised biomedical knowledge from the scientific community  
 180 stored in a knowledge graph. In particular, the Monarch Knowledge Graph (MKG) is a comprehensive,  
 181 standardised database that aggregates up-to-date knowledge about biomedical concepts and the relationships  
 182 between them. This currently includes 103 well-established phenotype-cell type relationships<sup>40</sup>. We used  
 183 the MKG as a proxy for the field’s current state of knowledge of causal phenotype-cell type associations.  
 184 We evaluated the proportion of MKG associations that were recapitulated by our results Fig. 12. For  
 185 each phenotype-cell type association in the MKG, we computed the percent of cell types recovered in our  
 186 association results at a given ontological distance according to the CL ontology. An ontological distance of 0  
 187 means that our nominated cell type was as close as possible to the MKG cell type after adjusting for the cell  
 188 types available in our single-cell references. Instances of exact overlap of terms between the MKG and our  
 189 results would qualify as an ontological distance of 0 (e.g. ‘monocyte’ vs. ‘monocyte’). Greater ontological  
 190 distances indicate further divergence between the MKG cell type and our nominated cell type. A distance  
 191 of 1 indicating that the MKG cell type was one step away from our nominated cell type in the CL ontology  
 192 graph (e.g. ‘monocyte’ vs. ‘classical monocyte’). The maximum possible percent of recovered terms is capped  
 193 by the percentage of MKG ground-truth phenotypes we were able to find at least one significant cell type  
 194 association for at  $FDR_{pc}$ .

195 In total, our results contained at least one significant cell type associations for 90% of the phenotypes de-

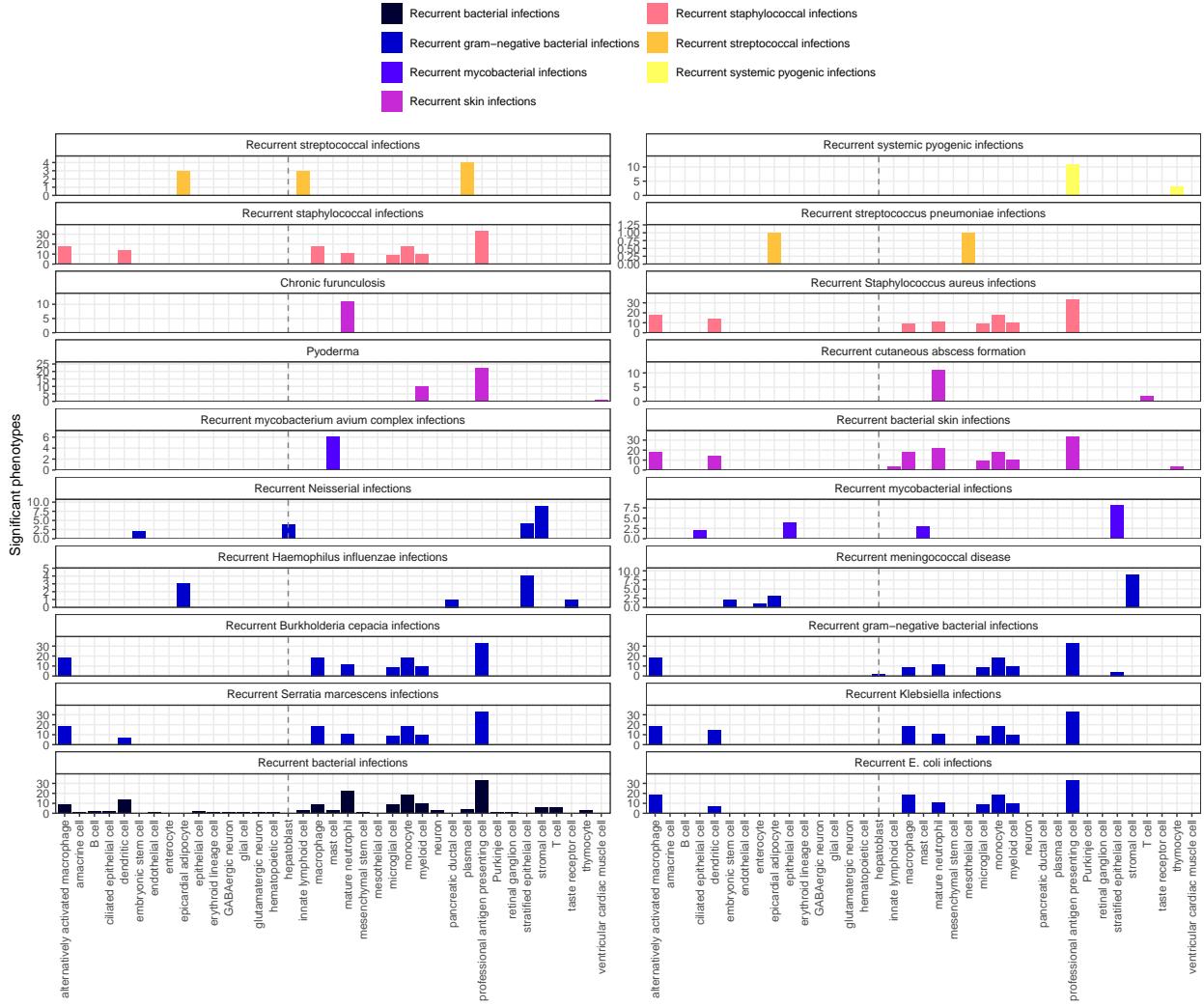
scribed in the MKG. Of these phenotypes, we captured 57% of the MKG phenotype-cell associations at an ontological distance of 0 (i.e. the closest possible Cell Ontology term match). Recall increased with greater flexibility in the matching of cell type annotations. At an ontological distance of 1 (e.g. ‘monocyte’ vs. ‘classical monocyte’), we captured 77% of the MKG phenotype-cell associations. Recall reached a maximum of 90% at a ontological distance of 5. This recall percentage is capped by the proportion of phenotypes for which we were able to find at least one significant cell type association for. It should be noted that we were unable to compute precision as the MKG (and other knowledge databases) only provide true positive associations. Identifying true negatives (e.g. a cell type is definitely never associated with a phenotype) is a fundamentally more difficult task to resolve as it would require proving the null hypothesis. Regardless, these benchmarking tests suggests that our results are able to recover the majority of known phenotype-cell type associations while proposing many new associations.

## Phenome-wide analyses discover novel phenotype-cell type associations

Having established that many of the phenotype-cell type associations align with prior expectations, we then sought to discover novel relationships with undercharacterised phenotypes. We reasoned that recurrent bacterial infections (and all its descendant phenotypes) should primarily be associated with immune cell types. The HPO term ‘Recurrent bacterial infections’ has 19 different descendant phenotypes, e.g. staphylococcal, streptococcal, and Neisserial infections. Each of these phenotypes are associated with partially overlapping subsets of immune cells and other cell types (Fig. 4). As expected, these phenotypes are primarily associated with immune cell types (e.g. macrophages, dendritic cells, T cells, monocytes, neutrophils). Some associations confirm relationships previously suggested in the literature, such as that between ‘Recurrent staphylococcal infections’ and myeloid cells<sup>41–44</sup>. Specifically, our results pinpoint monocytes as the most strongly associated cell subtypes ( $FDR=1.0 \times 10^{-30}$ ,  $\beta=0.18$ ).

Next, we sought to uncover novel, unexpected associations between recurrent bacterial infection phenotypes and cell types. In contrast to all other recurrent infection types, ‘Recurrent Neisserial infections’ highlighted a novel association with hepatoblasts (Descartes Human :  $FDR=1.1 \times 10^{-6}$ ,  $\beta=8.2 \times 10^{-2}$ ). Whilst unexpected, a convincing explanation involves the complement system, a key driver of innate immune response to Neisserial infections. Hepatocytes, which derive from hepatoblasts, produce the majority of complement proteins<sup>45</sup>, and Kupffer cells express complement receptors<sup>46</sup>. In addition, individuals with deficits in complement are at high risk for Neisserial infections<sup>47,48</sup>, and a genome-wide association study in those with a Neisserial infection identified risk variants within complement proteins<sup>49</sup>. While the potential of therapeutically targeting complement in RDs (including Neisserial infections) has been proposed previously<sup>50,51</sup>, performing this in a gene- and cell type-specific manner may help to improve efficacy and reduce toxicity (e.g. due to off-target effects). Importantly, there are over 56 known genes within the complement system<sup>52</sup>, highlighting the need for a systematic, evidence-based approach to identify effective gene targets.

230 Also of note, despite the fact that our datasets contain both hepatoblasts and their mature counterpart, hepa-  
231 tocytes, only the hepatoblasts showed this association. This suggests that the genetic factors that predispose  
232 individuals for risk of Neisserial infections are specifically affecting hepatoblasts before they become fully  
233 differentiated. It is also notable that these phenotypes were the only ones within the ‘Recurrent bacterial  
234 infections’ branch, or even the broader ‘Recurrent infections’ branch, perhaps indicating a unique role for  
235 hepatoblasts in recurrent infectious disease. The only phenotypes within the even broader ‘Abnormality of  
236 the immune system’ HPO branch that significantly associated with mature hepatocytes were ‘Pancreatitis’  
237 ( $FDR=2.1 \times 10^{-2}$ ,  $\beta=5.3 \times 10^{-2}$ ) and ‘Susceptibility to chickenpox’ ( $FDR=1.2 \times 10^{-2}$ ,  $\beta=5.5 \times 10^{-2}$ ) both  
238 of which are well-known to involve the liver<sup>53–55</sup>.



(a) **Association tests reveal that hepatoblasts have a unique role in recurrent Neisserial infections.** Significant phenotype-cell type tests for phenotypes within the branch ‘Recurrent bacterial infections’. Amongst all different kinds of recurrent bacterial infections, hepatoblasts (highlighted by vertical dotted lines) are exclusively enriched in ‘Recurrent gram-negative bacterial infections’. Note that terms from multiple levels of the same ontology branch are shown as separate facets (e.g. ‘Recurrent bacterial infections’ and ‘Recurrent gram-negative bacterial infections’).

Figure 4

239 Phenotypes can be associated with multiple diseases, cell types and genes. In addition to hepatoblasts, ‘Recur-  
 240 rent Neisserial infections’ were also associated with stromal cells ( $FDR=4.6 \times 10^{-6}$ ,  $\beta=7.9 \times 10^{-2}$ ), stratified  
 241 epithelial cells ( $FDR=1.7 \times 10^{-23}$ ,  $\beta=0.15$ ), and embryonic stem cells ( $FDR=5.4 \times 10^{-5}$ ,  $\beta=7.4 \times 10^{-2}$ ).  
 242 ‘Recurrent Neisserial infections’ is a phenotype of 7 different diseases (‘C5 deficiency’, ‘C6 deficiency’, ‘C7  
 243 deficiency’, ‘Complement component 8 deficiency, type II’, ‘Complement factor B deficiency’, ‘Complement  
 244 factor I deficiency’, ‘Mannose-Binding lectin deficiency’). The monogenic nature of these diseases makes it  
 245 very difficult to statistically infer the cell types underlying them. By aggregating these genes to the level of

246 phenotype (the observed symptom) we can better understand the cell types underlying all of these diseases.

247 Having found four distinct cell types associated with RNI, we asked whether the RNI-associated genes were  
248 equally expressed across all of these cell types, or whether they differentially contributed to each of the  
249 associations. RNI provides a convenient case study to investigate this because each of the seven diseases  
250 that have RNI as a phenotype are purely monogenic. This makes it relatively straightforward to demonstrate  
251 how genes can drive associations between cell types, phenotypes and their respective diseases.

252 Diseases that have ‘Recurrent Neisserial infections’ as a phenotype were collected from the HPO annotation  
253 files. Genes that were annotated to a given phenotype (e.g. ‘Recurrent Neisserial infections’) via a particular  
254 disease (e.g. ‘C5 deficiency’) constituted “symptom”-level gene sets. Only diseases whose symptom-level  
255 gene sets had >25% overlap with the driver gene sets for at least one cell type were retained in the network  
256 plot. Using this approach, we were able to construct and refine causal networks tracing multiple scales of  
257 disease biology.

258 This procedure revealed that genetic deficiencies in various complement system genes (e.g. *C5*, *C8*, and  
259 *C7*) are primarily mediated by different cell types (hepatoblasts, stratified epithelial cells, and stromal cells,  
260 respectively). While genes of the complement system are expressed throughout many different tissues and  
261 cell types, these results indicate that different subsets of these genes may mediate their effects through  
262 different cell types. While almost all of these genes show high expression specificity in hepatoblasts, only *C6*,  
263 *C7* and *CFI* meet the threshold for the status of driver genes in stromal cells.

264 Recall that we showed in Fig. 4b that as we approach the leaf nodes of the HPO we tend towards a given  
265 phenotype being associated with a single cell type. Note that mean this in a theoretical sense, as we do  
266 not necessarily demonstrate a single cell type for each phenotype in this particular dataset. However, as  
267 more granular phenotypes are defined over time, we would expect this hypothesis to bear out. The corollary  
268 of this is that we would expect there to be at least four subtypes of the RNI phenotype, as predicted  
269 by the four distinct cell types found to underly this phenotype. This may present as different clinical  
270 courses (e.g. early onset, late onset, relapse-remitting) or biomarkers (e.g. histological) to be reveal in future  
271 examinations of clinical cohorts. Based on this, we predict that forms of RNI caused by genes expressed in  
272 stromal cells would have phenotypic differences from those caused by genes expressed in stratified epithelial  
273 cell. In other words, phenotypic similarity is driven by the underlying causal cell types.

#### 274 **Prioritising phenotypes based on severity**

275 Some phenotypes are more severe than others and thus could be given priority for developing treatments. For  
276 example, ‘Leukonychia’ (white nails) is much less severe than ‘Leukodystrophy’ (white matter degeneration  
277 in the brain). Given the large number of significant phenotype-cell type associations, we needed a way of  
278 prioritising phenotypes for further investigation. We therefore used the large language model GPT-4 to  
279 systematically annotate the severity of all HPO phenotypes<sup>37</sup>.

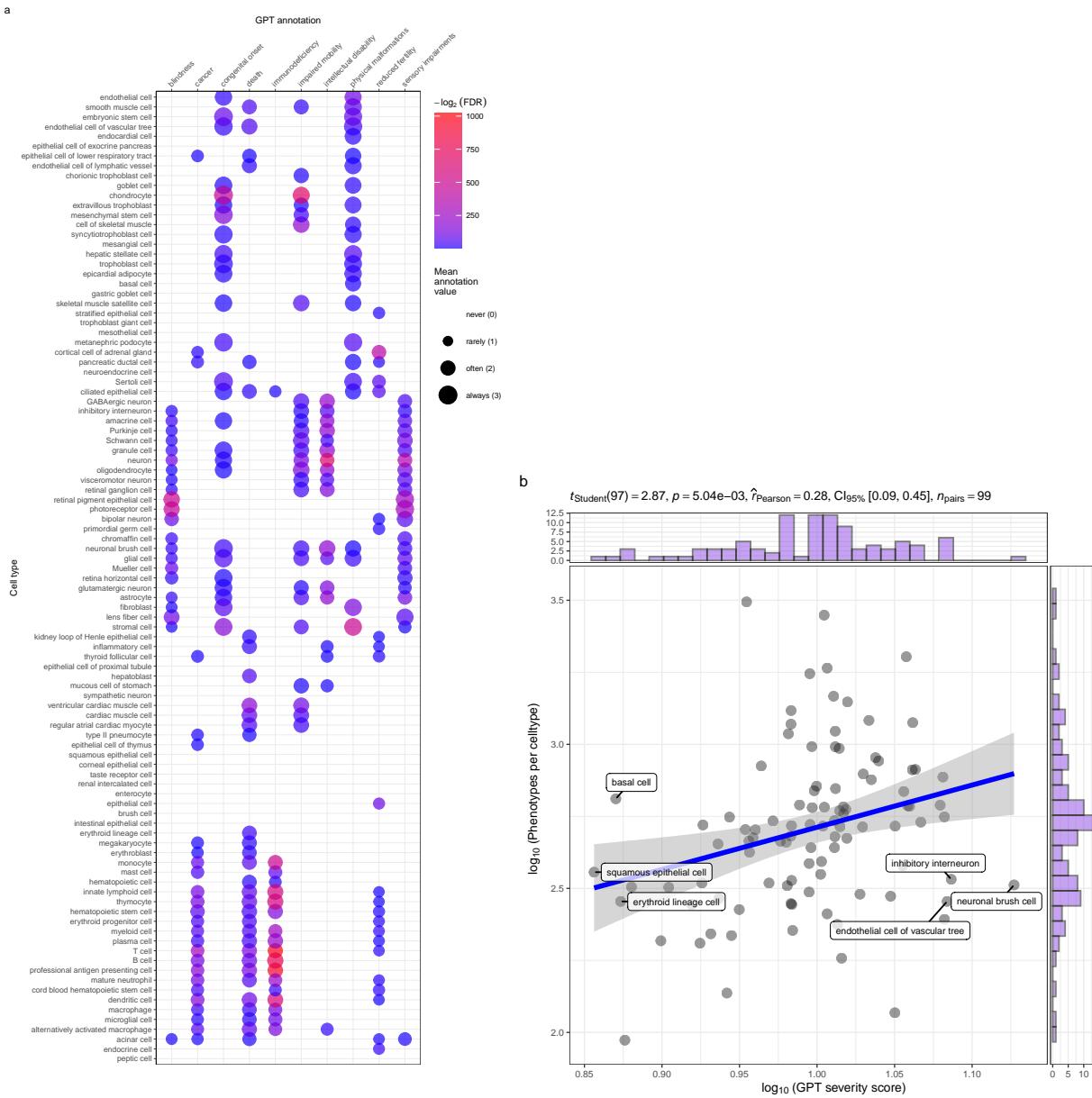
280 Severity annotations were gathered from GPT-4 for 16,982/18,082 (94%) HPO phenotypes in our companion  
281 study<sup>37</sup>. Benchmarking tests of these results using ground-truth HPO branch annotations. For example,  
282 phenotypes within the ‘Blindness’ HPO branch (*HP:0000618*) were correctly annotated as causing blindness  
283 by GPT-4. Across all annotations, the recall rate of GPT-4 annotations was 96% (min=89%, max=100%,  
284 SD=4.5) with a mean consistency score of 91% (min=81%, max=97%, SD=5.7) for phenotypes whose  
285 annotation were collected more than once. This clearly demonstrates the ability of GPT-4 to accurately  
286 annotate phenotypes. This allowed us to begin using these annotations to compute systematically collected  
287 severity scores for all phenotypes in the HPO.

288 From these annotations we computed a weighted severity score metric for each phenotype ranging from 0-100  
289 (100 being the theoretical maximum severity of a phenotype that always causes every annotation). Within  
290 our annotations, the most severe phenotype was ‘Atrophy/Degeneration affecting the central nervous system’  
291 (*HP:0007367*) with a severity score of 47, followed by ‘Anencephaly’ (*HP:0002323*) with a severity score of  
292 45. There were 677 phenotypes with a severity score of 0 (e.g. ‘Thin toenail’). The mean severity score  
293 across all phenotypes was 10 (median=9.4, standard deviation=6.4).

294 We next sought to answer the question “are disruptions to certain cell types more likely to cause severe  
295 phenotypes?”. To address this, we merged the GPT annotations with the significant (FDR<0.05) phenotype-  
296 cell type association results and computed the frequency of each severity annotation per cell type (Fig.  
297 Figure 13). We found that neuronal brush cells were associated with phenotypes that had the highest  
298 average composite severity scores, followed by Mueller cells and glial cells. This suggests that disruptions  
299 to these cell types are more likely to cause generally severe phenotypes. Meanwhile, megakaryocytes were  
300 associated with phenotypes that had the lowest average composite severity scores, suggesting that disruptions  
301 to these cell types can be better tolerated than others.

302 Different aspects of phenotype severity will be more associated with some cell types than others. After  
303 encoding the GPT annotations numerically (0=“never”, 1=“rarely”, 2=“often”, 3=“always”) we computed  
304 the mean encoded value per cell type within each annotation. We then ran a series of one-sided Wilcoxon  
305 rank-sum tests to objectively determine whether some cell types tended to be associated with phenotypes  
306 that more frequently caused certain severity metrics (death, intellectual disability, impaired mobility, etc.)  
307 relative to all other cell types (Fig. 5a). This consistently yielded expected relationships between cell types  
308 (e.g. retinal pigment epithelial cells) and phenotype characteristics (e.g. blindness). Similarly, phenotypes  
309 that more commonly cause death are most commonly associated with ventricular cardiac muscle cells, and  
310 least commonly associated with squamous epithelial cells and bipolar neurons. Analogous patterns of ex-  
311 pected associations are shown consistently across all annotations (e.g. fertility-reducing phenotypes asso-  
312 ciated with cortical cell of adrenal glands, immunodeficiency-causing phenotypes associated with T cells,  
313 mobility-impairing phenotypes associated with chondrocytes, cancer-causing phenotypes associated with T  
314 cells, etc.).

315 We also sought to answer whether the number of phenotypes that a cell type is associated with has a  
316 relationship with the severity of those phenotypes (Fig. 5b). Our working hypothesis is that when a cell type  
317 that affect many different phenotypes is disrupted, the cell type likely performs some critical function that  
318 affect many physiological systems. It also means that the individual phenotypes tend to be more severe than  
319 other phenotypes that involve less critical cell types. Indeed, we found a significant relationship between  
320 number of associated and mean composite phenotype severity ( $p=5.0 \times 10^{-3}$ , Pearson coefficient=0.28).



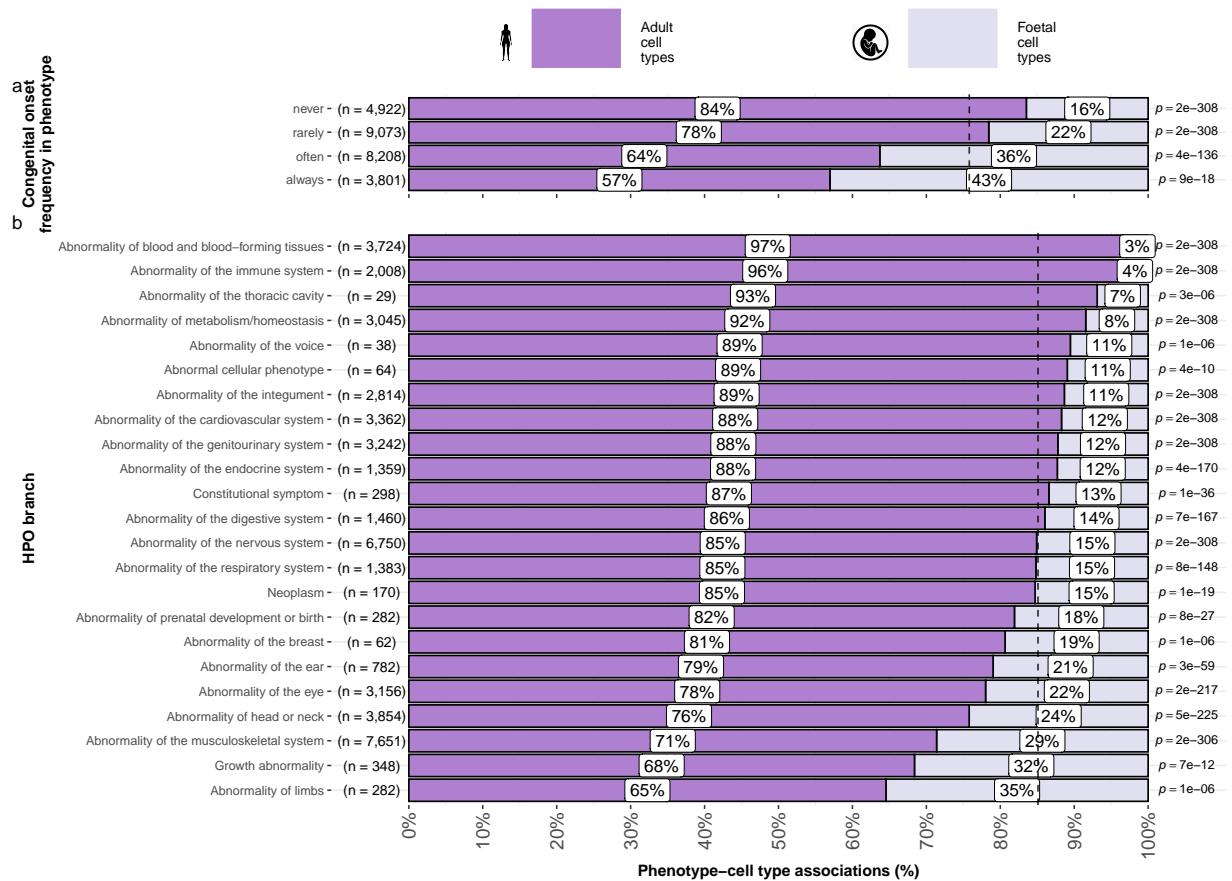
**(a) Genetic disruptions to some cell types cause more clinically severe phenotypes than others.** **a**, Different cell types are associated with different aspects of phenotypic severity. The dot plot shows the mean encoded frequency value for a given severity annotation (0="never", 1="rarely", 2="often", 3="always"; shown as dot size), aggregated by the associated cell type. One-sided Wilcoxon rank-sum tests were performed for each cell type (within each GPT annotation) to determine which cell types more frequently caused severe phenotypes than all other cell types. Dots are colored by  $-\log_2(\text{FDR})$  when Wilcoxon test FDR values were less than 0.05. All dots with non-significant Wilcoxon tests are instead colored grey. Cell types (rows) are clustered according to the p-values of the Wilcoxon tests. **b**, Cell types that affect more phenotypes tend to have more clinically severe consequences. Specifically, the number of phenotypes each cell type is significantly associated with, and the mean composite severity score of each cell type. The cell types with the top/bottom three x/y axis values are labeled to illustrate the cell types that cause the most/least phenotypic disruption when dysfunctional. Side histograms show the density of data points along each axis. Summary statistics for the linear regression are shown in the title ( $t_{\text{Student}}$  = Student t-test statistic,  $p$  = p-value,  $\hat{\rho}_{\text{Pearson}}$  = Pearson correlation coefficient,  $CI_{95\%}$  = confidence intervals,  $n_{\text{pairs}}$  = number of observed data pairs).

Figure 5

321 **Congenital phenotypes are associated with foetal cell types**

322 Which life stage a phenotype affects an individual is clinically important and can have profound implications  
323 for how patients are treated and whether that are treatable with currently available interventions. For  
324 example, beyond a certain point gene therapies may not be an effective means of treating morphological  
325 defects that arise during development. Within the DescartesHuman dataset, 100% of the cells were from  
326 foetal tissues. Meanwhile, the Human Cell Landscape was derived from embryonic, foetal, and adult tissue  
327 samples. Within the Human Cell Landscape, 29% of cell types were found in foetal tissue, and 71% were found  
328 in adult tissues. Many of the cell types in our datasets have both foetal and adult versions (e.g. chondrocytes),  
329 while some only exist in the course of foetal development (e.g. neural crest cells). This presents a unique  
330 opportunity to provide an additional layer of contextualisation in our phenotype-cell type association results  
331 that may provide critical information when determining viable patient treatment options.

332 We reasoned that phenotypes that are most frequently congenital are more likely to be associated with  
333 foetal cell types than adult cell types. As expected, the frequency of congenital onset with each phenotype  
334 (as determined by GPT-4 annotations) was strongly predictive of the proportion of significantly associated  
335 foetal cell types in our results ( $p = 4.7 \times 10^{-261}$ ,  $\chi^2_{Pearson} = 1.2 \times 10^3$ ,  $\hat{V}_{Cramer} = 0.22$ , Fig. 6a). This result is  
336 consistent with the expected role of foetal cell types in development and the aetiology of congenital disorders.



**(a) Foetal vs. adult cell type references provide development context to phenotype aetiology.** **a**, Congenital phenotypes are more often associated with foetal cell types. As a phenotype is more often congenital in nature, the greater proportion of foetal cell types are significantly associated with it. **b**, The proportion of phenotype-cell type association tests that are enriched for foetal cell types within each HPO branch. The p-values to the right of each bar are the results of an additional series of  $\chi^2$  tests to determine whether the proportion of foetal vs. non-foetal cell types significantly different differs from the proportions expected by chance (the dashed vertical line). The foetal silhouette was generated with DALL-E. The adult silhouette is from phylopic.org and is freely available via CC0 1.0 Universal Public Domain Dedication.

Figure 6

- 337 Some branches of the HPO were more commonly enriched in foetal cell types compared to others  
 338 ( $\hat{V}_{Cramer} = 0.22$ ,  $p < 2.2 \times 10^{-308}$ , Fig. 6b). The branch with the greatest proportion of foetal cell type  
 339 enrichments was ‘Abnormality of limbs’ (35%), followed by ‘Growth abnormality’ (32%) and ‘Abnormality  
 340 of the musculoskeletal system’ (29%). Notably, ‘Abnormality of limbs’ branch was most disproportionately  
 341 enriched for foetal cell type associations relative to all other branches (35% cell types). These results align  
 342 well with the fact that physical malformations tend to be developmental in origin.  
 343 Conversely, the HPO branches that were most biased towards adult cell types were ‘Abnormality of blood  
 344 and blood-forming tissues’ (97%), ‘Abnormality of the immune system’ (96%), and ‘Abnormality of the  
 345 thoracic cavity’ (93%).

346 Some phenotypes exclusively involve the foetal version of a cell type, while others exclusively involve the  
347 adult version. We sought to find those phenotypes which had the greatest bias towards either end of this  
348 spectrum. To do so, we designed a metric to identify which phenotypes were more often associated with  
349 foetal cell types than adult cell types. For each phenotype, we calculated the difference in the association  
350 p-values between the foetal and adult version of the equivalent cell type. The resulting metric ranges from 1  
351 (indicating the phenotype is only associated with the foetal version of the cell type) and -1 (indicating the  
352 phenotype is only associated with the adult version of the cell type). To summarise the most foetal-biased  
353 phenotype categories, we ran an ontological enrichment test with the HPO graph Table 7. To identify foetal  
354 cell type-biased phenotype categories, we fed the top 50 phenotypes with the greatest foetal cell type bias  
355 (closer to 1) into the enrichment function Table 8. Conversely, we used the top 50 phenotypes with the  
356 greatest adult cell type bias (closer to -1) to identify adult cell type-biased phenotype categories.

357 The phenotype categories with the greatest bias towards foetal cell types were ‘Abnormal nasal mor-  
358 phology’ ( $p=2.4 \times 10^{-7}$ ,  $\log_2(\text{fold-change})=4.5$ ) and ‘Abnormal external nose morphology’ ( $p=2.5 \times 10^{-6}$ ,  
359  $\log_2(\text{fold-change})=5.4$ ).

360 Specific examples of such phenotypes include ‘Short middle phalanx of the 2nd finger’, ‘Abnormal morphology  
361 of the nasal alae’, and ‘Abnormal labia minora morphology’. Indeed, these phenotypes are morphological  
362 defects apparent at birth caused by abnormal developmental processes.

363 Conversely, the most adult cell type-biased phenotype categories were ‘Abnormal elasticity of skin’  
364 ( $p=3.6 \times 10^{-7}$ ,  $\log_2(\text{fold-change})=6.0$ ) and ‘Abnormally lax or hyperextensible skin’ ( $p=1.3 \times 10^{-5}$ ,  
365  $\log_2(\text{fold-change})=6.0$ ).

366 Specific examples of such phenotypes include ‘Excessive wrinkled skin’ and ‘Paroxysmal supraventricular  
367 tachycardia’ Table 8. It is well known that ageing naturally causes a loss of skin elasticity (due to decreasing  
368 collagen production) and vascular degeneration<sup>56</sup>. Next, we were interested whether some cell types tend to  
369 show strong differences in their phenotype associations between their foetal and adult forms. To test this, we  
370 performed an analogous enrichment procedure as with the phenotypes, except using Cell Ontology terms and  
371 the Cell Ontology graph. This analysis identified the cell type category connective tissue cell ( $p=1.8 \times 10^{-3}$ ,  
372  $\log_2(\text{fold-change})=3.2$ ) as the most foetal-biased cell type. No cell type categories were significantly enriched  
373 for the most adult-biased cell types. This is likely due to the fact that cell types can be disrupted at different  
374 stages of life, resulting in different phenotypes. Thus there the same cell types may be involved in both  
375 the most foetal-biased and adult-biased phenotypes. Together, these findings serve to further validate our  
376 methodology as a tool for identifying the causal cell types underlying a wide range of phenotypes.

### 377 Therapeutic target identification

378 In the above sections, we demonstrated how gene association databases can be used to investigate the cell  
379 types underlying disease phenotypes at scale. While these associations are informative on their own, we

wished to take these results further in order to have a more translational impact. Knowledge of the causal cell types underlying each phenotype can be incredibly informative for scientists and clinicians in their quest to study and treat them. Therapeutic targets with supportive genetic evidence have 2.6x higher success rates in clinical trials<sup>57–59</sup>. Furthermore, knowing which cell types to target with gene therapy can maximise the efficacy of highly expensive payloads, and minimise side effects (e.g. immune reaction to viral vectors). Recent biotechnological advances have greatly enhanced our ability to target specific cell types with gene therapy, making specific and accurate knowledge the correct underlying cell types more pertinent than ever<sup>35,36</sup>.

However, given the sheer number of results, we wished to develop a principled and reproducible approach to filter and rank putative cell type-specific gene targets for diseases where there is the greatest urgent need for improved treatments. We therefore systematically identified putative cell type-specific gene targets for severe phenotypes. First, we transformed our phenotype-cell type association results and merged them with primary data sources (e.g. GenCC gene-disease relationships, scRNA-seq atlas datasets) to create a large table of multi-scale relationships, where each row represented a tetrad of disease-phenotype-cell type-gene relationships. We then filtered non-significant phenotype-cell type relationships (only associations with  $FDR < 0.05$ ) as well as phenotype-gene relationships with strong causal evidence (GenCC score  $> 3$ ). We also removed any phenotypes that were too broad to be clinically useful, as quantified using the information content (IC) ( $IC > 8$ ), which measures the how specific each term is within an ontology (i.e. HPO). Gene-cell type relationships were established by taking genes that had the top 25% expression specificity quantiles within each cell type. When connecting cell types to diseases via phenotypes, we used a symptom intersection threshold of  $>.25$ . Next, we sorted the remaining results in descending order of phenotype severity using the GPT4 composite severity scores described earlier. Finally, to limit the size of the resulting multi-scale networks we took only the top 10 rows, where each row represented a tetrad of disease-phenotype-cell type-gene relationships. This resulted in number of relatively small, high-confidence disease-phenotype-cell type-gene networks that could be reasonably interrogated through manual inspection and network visualisation. For example, if one was interested in the mechanisms causing ‘Recurrent Neisserial infections’, one would need only select all rows that include this phenotype to find all of its most relevant connection to diseases, cell types, and genes.

This yielded putative therapeutic targets for 5,252 phenotypes across 4,819 diseases in 201 cell types and 3,148 genes (Fig. 15). While this constitutes a large number of genes in total, each phenotype was assigned a median of 2.0 gene targets (mean=3.3, min=1, max=10). Relative to the number of genes annotations per phenotype in the HPO overall (median=7.0, mean=62, min=1, max=5,003) this represents a substantial decrease in the number of candidate target genes, even when excluding high-level phenotypes (HPO level $>3.0$ ). It is also important to note that the phenotypes in the prioritised targets list are ranked by their severity, allowing us to distinguish between phenotypes with a high medical urgency (e.g. ‘Hydranencephaly’) from those with lower medical urgency (e.g. ‘Increased mean corpuscular volume’). This can be useful for clinicians, biomedical

415 scientists, and pharmaceutical manufacturers who wish to focus their research efforts on phenotypes with  
416 the greatest need for intervention.

417 Across all phenotypes, epithelial cell were most commonly implicated (838 phenotypes), followed by stromal  
418 cell (626 phenotypes), stromal cell (626 phenotypes), neuron (475 phenotypes), chondrocyte (383 pheno-  
419 types), and endothelial cell (361 phenotypes). Grouped by higher-order ontology category, ‘Abnormality of  
420 the musculoskeletal system’ had the greatest number of enriched phenotypes (959 phenotypes, 857 genes),  
421 followed by ‘Abnormality of the nervous system’ (733 phenotypes, 1,138 genes), ‘Abnormality of head or  
422 neck’ (543 phenotypes, 986 genes), ‘Abnormality of the genitourinary system’ (443 phenotypes, 695 genes),  
423 and ‘Abnormality of the eye’ (377 phenotypes, 545 genes).

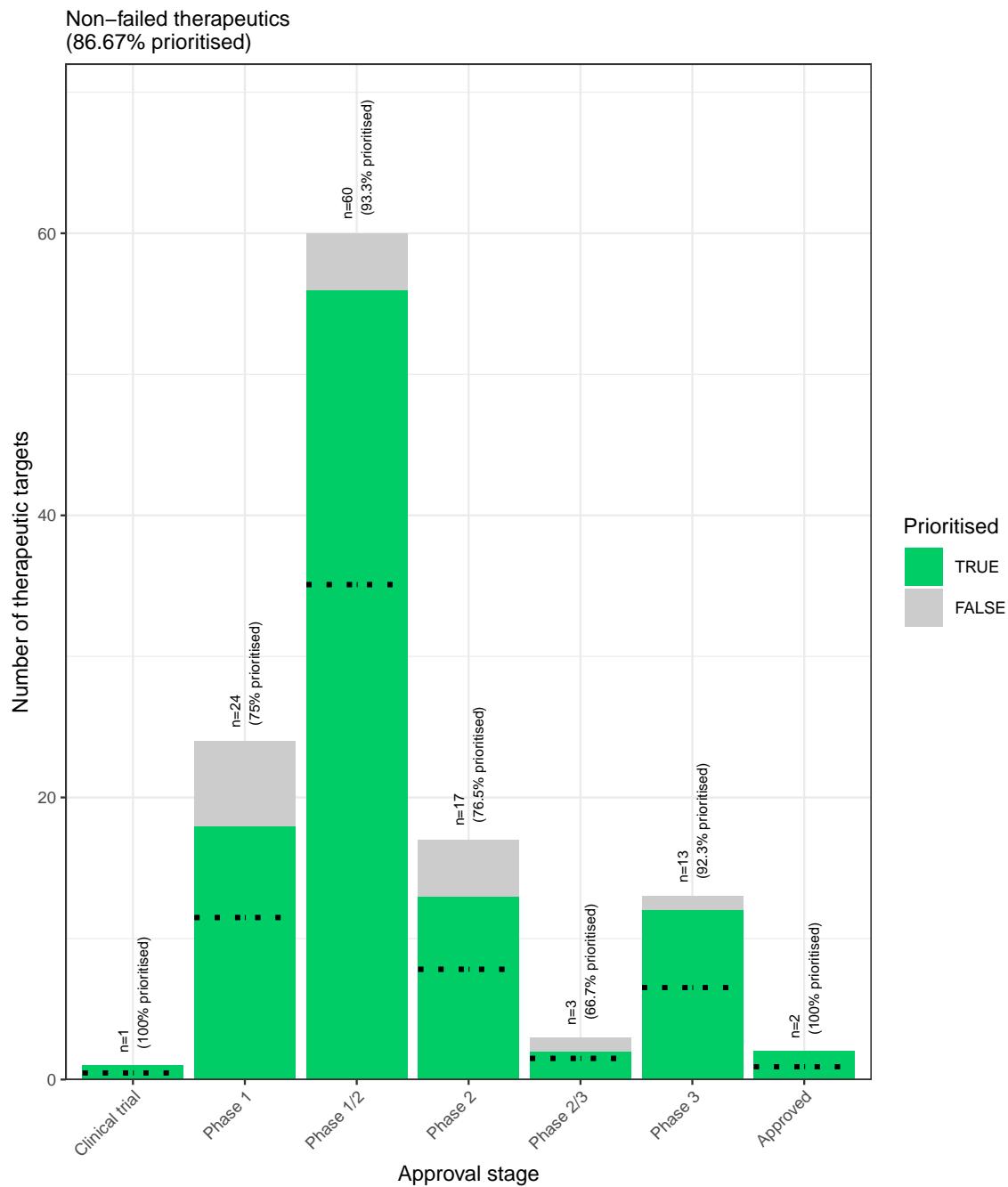
#### 424 Therapeutic target validation

425 To determine whether the genes prioritised by our therapeutic targets pipeline were plausible, we checked  
426 what percentage of gene therapy targets we recapitulated. Data on therapeutic approval status was gathered  
427 from the Therapeutic Target Database (TTD; release 2025-08-08)<sup>60</sup>. Overall, we prioritised 87% (120 total)  
428 of all non-failed existing gene therapy targets (ie. those which are currently approved, investigative, or  
429 undergoing clinical trials). A hypergeometric test confirmed that our prioritised targets were significantly  
430 enriched for non-failed gene therapy targets ( $p = 1.8 \times 10^{-5}$ ). For these hypergeometric tests, the background  
431 gene set was composed of the union of all phenotype-associated genes in the HPO and all gene therapy  
432 targets listed in TTD.

433 Even when considering therapeutics of any kind (Fig. 16), not just gene therapies, we recapitulated 40% of the  
434 non-failed therapeutic targets and 0% of the terminated/withdrawn therapeutic targets (n=1,255). Here we  
435 found that our prioritised targets were highly significantly depleted for failed therapeutics ( $p = 2.2 \times 10^{-142}$ ).  
436 This suggests that our multi-scale evidence-based prioritisation pipeline is capable of selectively identifying  
437 genes that are likely to be effective therapeutic targets.

438 In addition to aggregate enrichment results, we also provide specific examples of successful gene therapies  
439 whose cell type-specific mechanism were recapitulated by our phenotype-cell associations. In particular, our  
440 pipeline nominated the gene *RPE65* within ‘retinal pigment epithelial cells’ as the top target for ‘Fundus  
441 atrophy’ vision-related phenotypes that are hallmarks of ‘Leber congenital amaurosis, type II’ and ‘Se-  
442 vere early-childhood-onset retinal dystrophy’. Indeed, gene therapies targeting *RPE65* within the retina of  
443 patients with these rare genetic conditions are some of the most successful clinical applications of this tech-  
444 nology to date, able to restore vision in many cases<sup>61</sup>. In other cases, a tissue (e.g. liver) may be known to  
445 be causally involved in disease genesis, but the precise causal cell types within that tissue remain unknown  
446 (e.g. hepatocytes, Kupffer cells, Cholangiocytes, Hepatic stellate cells, Natural killer cells, etc.). Tissue-level  
447 investigations (e.g. using bulk transcriptomics or epigenomics) would be dominated by hepatocytes, which  
448 comprise 75% of the liver. Our prioritized gene therapy targets can aid in such scenarios by providing the

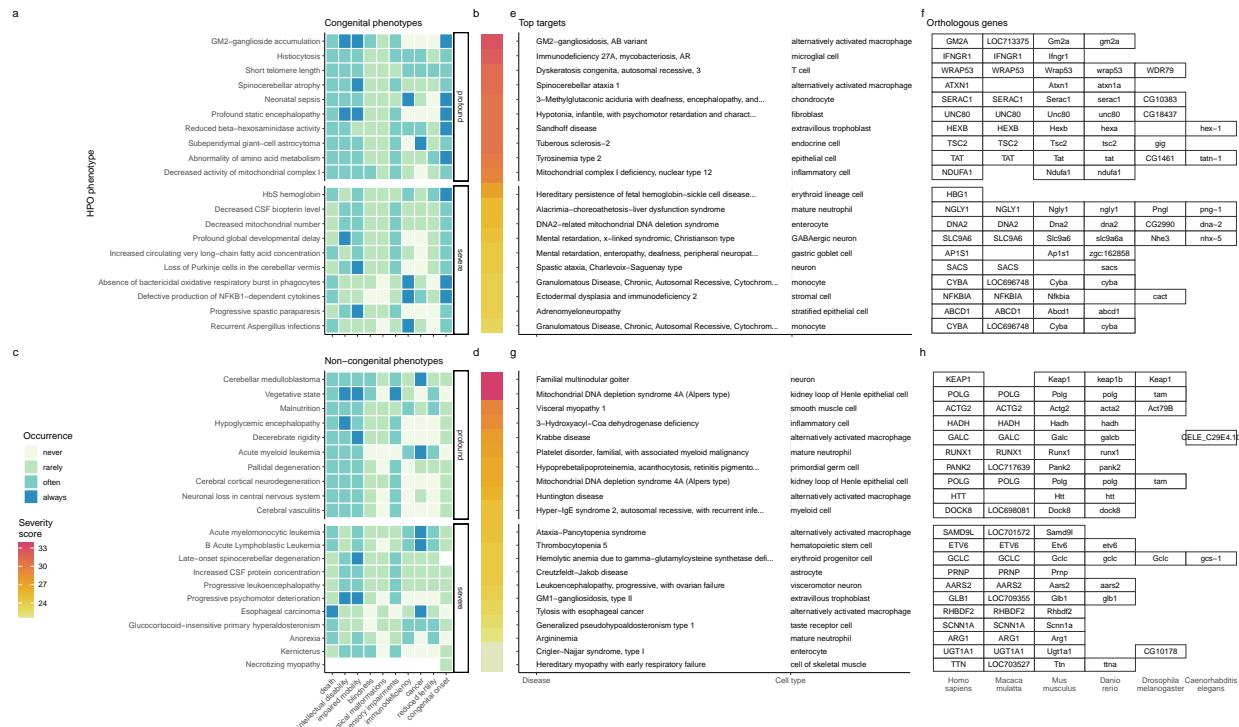
<sup>449</sup> cell type-resolution context most likely to be causal for a given phenotype or set of phenotypes.



(a) **Prioritised targets recapitulate existing gene therapy targets.** The proportion of existing gene therapy targets (documented in the Therapeutic Target Database) recapitulated by our prioritisation pipeline. Therapeutics are stratified by the stage of clinical development they were at during the time of writing. While our prioritized targets did not include any failed ('Terminated') therapies, the fact that only one such therapy exists in the dataset preclude us from making any conclusions about depletion of failed gene therapy targets in our prioritised targets list.

Figure 7

450 **Selected example targets**



(a) **Evidence-based pipeline nominates causal mechanisms to target for gene therapy.** Shown here are the top 40 prioritised gene therapy targets at multiple biological scales, stratified by congenital (top row) vs. non-congenital phenotypes (bottom row) as well as severity class (“profound” or “severe”). In this plot, only the top 10 most severe phenotypes within a given strata/substrata are shown **a,c**. Severity annotation generated by GPT-4. **b,d**, Composite severity scores computed across all severity metrics. **e,g**, Top mediator disease and cell type-specific target for each phenotype. **f,h** top target gene for each phenotype within humans (*Homo sapiens*). We also include the 1:1 ortholog of each human gene in several commonly used animal models, including monkey (*Macaca mulatta*), mouse (*Mus musculus*), zebrafish (*Danio rerio*), fly (*Drosophila melanogaster*) and nematode (*Caenorhabditis elegans*). Boxes are empty where no 1:1 ortholog is known. See supplement Fig. 18 for network plots of cell type-specific gene therapy targets for several severe phenotypes and their associated diseases.

Figure 8

- 451 From our prioritised targets, we selected four phenotype or disease examples: ‘GM2-ganglioside accumulation’,  
452 ‘Spinocerebellar atrophy’, ‘Neuronal loss in central nervous system’. To focus on clinically relevant  
453 phenotypes and reduce overplotting, we limited selection to those with GPT severity scores above 15 Fig. 8a-h.  
454 Selection was based on severity and network simplicity to allow compact visualisation.
- 455 Tay-Sachs disease (TSD) is a fatal neurodegenerative condition caused by *HEXA* deficiency and ganglioside  
456 buildup. We identified alternatively activated macrophages as the cell type most associated with ‘GM2-  
457 ganglioside accumulation’ Fig. 18. This aligns with prior findings of ganglioside accumulation in TSD  
458 macropahges<sup>62,63,64,65</sup>. Our results support macrophages as causal in TSD and the most promising thera-  
459 peutic target.

460 Spinocerebellar atrophy is a progressive neurodegenerative phenotype in disorders like Spinocerebellar ataxia.  
461 Our pipeline implicates M2 macrophages ('Alternatively activated macrophages') as the only causal cell type  
462 Fig. 18. This suggests Purkinje cell loss is downstream of macrophage dysfunction, consistent with microglial  
463 roles in neurodegeneration<sup>66–68</sup>. Our findings provide the first statistically supported link between risk genes  
464 and this cell type, which is supported by relevant mouse models (e.g. *Atxn1*, *Pnpla6*) that replicate cellular  
465 and behavioural disease phenotypes.

466 Despite its broad definition, 'Neuronal loss in central nervous system' was associated with only 3 cell types:  
467 alternatively activated macrophage, macrophage, epithelial cell, specifically M2 macrophages and sinusoidal  
468 endothelial cells Fig. 18.

469 Skeletal dysplasia comprises 450+ disorders affecting bone and cartilage, often leading to lethal outcomes via  
470 organ compression. While surgeries offer partial relief, pharmacological options remain limited. Our analysis  
471 identified chondrocytes as causal Fig. 19, consistent with known gene–cell links (e.g. *SLC26A2*, *COL2A1* in  
472 Achondrogenesis Type 1B and Torrance-type dysplasia). Chondrocyte-targeted therapy may offer long-term  
473 solutions where surgery falls short.

474 Alzheimer's disease (AD), a common neurodegenerative condition, presents with variable symptoms such as  
475 memory loss and proteinopathy. Our analysis shows distinct monogenic AD subtypes associate with different  
476 cell types and phenotypes Fig. 19. For example, AD subtypes 3 and 4 implicate digestive cells ('enterocyte',  
477 'gastric goblet cell'), while AD subtype 2 involves immune cells ('alternatively activated macrophage').  
478 These findings may explain heterogeneity in AD onset and presentation.

479 Parkinson's disease (PD) includes motor and systemic symptoms. PD subtypes 19a and 8 implicate oligo-  
480 dendrocytes and neurons Fig. 19, suggesting *LRRK2* variants act via gliosis in the substantia nigra. Other  
481 PD mechanisms involved chondrocytes (PD 20), amacrine cells (late-onset PD), and respiratory/immune  
482 cells (PD 14). This diversity may underlie PD's multisystem features.

### 483 Experimental model translatability

484 We computed interspecies translatability scores using a combination of both ontological ( $SIM_o$ ) and geno-  
485 typic ( $SIM_g$ ) similarity relative to each homologous human phenotype and its associated genes Fig. 17.  
486 In total, we mapped 1,221 non-human phenotypes (in *Caenorhabditis elegans*, *Danio rerio*, *Mus muscu-*  
*lus*, *Rattus norvegicus*) to 3,319 homologous human phenotypes. Amongst the 5,252 phenotype within our  
488 prioritised therapy targets, 1,788 had viable animal models in at least one non-human species. Per species,  
489 the number of homologous phenotypes was: *Mus musculus* (n=1705) *Danio rerio* (n=244) *Rattus norvegicus*  
490 (n=85) *Caenorhabditis elegans* (n=23). Amongst our prioritised targets with a GPT-4 severity score of >10,  
491 the phenotypes with the greatest animal model similarity were "Rudimentary to absent tibiae" ( $SIM_{og} = 1$ ),  
492 "Hypoglutaminemia" ( $SIM_{og} = 1$ ), "Bilateral ulnar hypoplasia" ( $SIM_{og} = 0.99$ ), "Disproportionate short-  
493 ening of the tibia" ( $SIM_{og} = 0.99$ ), "Acrobrachycephaly" ( $SIM_{og} = 0.98$ ).

494 **Mappings**

495 Mappings from HPO phenotypes and other commonly used medical ontologies were gathered in order to  
496 facilitate use of the results in this study in both clinical and research settings. Direct mappings, with a  
497 cross-ontology distance of 1, are the most precise and reliable. Counts of mappings at each distance are  
498 shown in Table 1. In total, there were 15,105 direct mappings between the HPO and other ontologies, with  
499 the largest number of mappings coming from the UMLS ontology (12,898 UMLS terms).

500 The mappings files can be accessed with the function `HPOExplorer::get_mappings` or directly via the  
501 `HPOExplorer` Releases page on GitHub (<https://github.com/neurogenomics/HPOExplorer/releases/tag/latest>).

503 **Discussion**

504 Investigating RDs at the level of phenotypes offers numerous advantages in both research and clinical  
505 medicine. First, the vast majority of RDs only have one associated gene (7,671/8,631 diseases = 89%).  
506 Aggregating gene sets across diseases into phenotype-centric “buckets” permits sufficiently well-powered  
507 analyses, with an average of ~76 genes per phenotype (median=7) see Fig. 10. Second, we hypothesised  
508 that these phenotype-level gene sets converge on a limited number of molecular and cellular pathways. Per-  
509 turbations to these pathways manifest as one or more phenotypes which, when considered together, tend  
510 to be clinically diagnosed as a certain disease. Third, RDs are often highly heterogeneous in their clinical  
511 presentation across individuals, leading to the creation of an ever increasing number of disease subtypes  
512 (some of which only have a single documented case). In contrast, a phenotype-centric approach enables us  
513 to more accurately describe a particular individual’s version of a disease without relying on the generation  
514 of additional disease subcategories. By characterising an individual’s precise phenotypes over time, we may  
515 better understand the underlying biological mechanisms that have caused their condition. However, in order  
516 to achieve a truly precision-based approach to clinical care, we must first characterise the molecular and  
517 cellular mechanisms that cause the emergence of each phenotype. Here, we provide a highly reproducible  
518 framework that enables this at the scale of the entire genome.

519 Across the 201 cell types and 11,047 RD-associated phenotypes investigated, more than 46,514 significant  
520 phenotype-cell type relationships were discovered. This presents a wealth of opportunities to trace the  
521 mechanisms of rare diseases through multiple biological scales. This in turn enhances our ability to study  
522 and treat causal factors in disease with deeper understanding and greater precision. These results recapitulate  
523 well-known relationships, while providing additional cellular context to many of these known relationships,  
524 and discovering novel relationships.

525 It was paramount to the success of this study to ensure our results were anchored in ground-truth bench-  
526 marks, generated falsifiable hypotheses, and rigorously guarded against false-positive associations. Extensive

validation using multiple approaches demonstrated that our methodology consistently recapitulates expected phenotype-cell type associations (Fig. 2–Fig. 6). This was made possible by the existence of comprehensive, structured ontologies for all phenotypes (the Human Phenotype Ontology) and cell types (the Cell Ontology), which provide an abundance of clear and falsifiable hypotheses for which to test our predictions against. Several key examples include 1) strong enrichment of associations between cell types and phenotypes within the same anatomical systems (Fig. 2b-d), 2) a strong relationship between phenotype-specificity and the strength and number of cell type associations (Fig. 3), 3) identification of the precise cell subtypes involved in susceptibility to various subtypes of recurrent bacterial infections (Fig. 4), 4) a strong positive correlation between the frequency of congenital onset of a phenotype and the proportion of developmental cell types associated with it (Fig. 6)), and 5) consistent phenotype-cell type associations across multiple independent single-cell datasets (Fig. 11).

Unfortunately, there are currently only treatments available for less than 5% of RDs<sup>6</sup>. Novel technologies including CRISPR, prime editing, antisense oligonucleotides, viral vectors, and/or lipid nanoparticles, have been undergone significant advances in the last several years<sup>69–73</sup> and proven remarkable clinical success in an increasing number of clinical applications<sup>74–77</sup>. The U.S. Food and Drug Administration (FDA) recently announced an landmark program aimed towards improving the international regulatory framework to take advantage of the evolving gene/cell therapy technologies<sup>78</sup> with the aim of bringing dozens more therapies to patients in a substantially shorter timeframe than traditional pharmaceutical product development (typically 5–20 years with a median of 8.3 years)<sup>79</sup>. While these technologies have the potential to revolutionise RD medicine, their successful application is dependent on first understanding the mechanisms causing each disease.

To address this critical gap in knowledge, we used our results to create a reproducible and customisable pipeline to nominate cell type-resolved therapeutic targets (Fig. 15–Fig. 8). Targeting cell type-specific mechanisms underlying granular RD phenotypes can improve therapeutic effectiveness by treating the causal root of an individual's conditions<sup>70,80</sup>. A cell type-specific approach also helps to reduce the number of harmful side effects caused by unintentionally delivering the therapeutic to off-target tissues/cell types (which may induce aberrant gene activity), especially when combined with technologies that can target cell surface antigens (e.g. viral vectors)<sup>81</sup>. This has the additional benefit of reducing the minimal effective dose of a therapeutic, which can be both immunogenic and extremely financially costly<sup>9,10,69,72</sup>. Here, we demonstrate the utility of a high-throughput evidence-based approach to RD therapeutics discovery by highlighting several of the most promising therapeutic candidates. Our pipeline takes into account a myriad of factors, including the strength of the phenotype-cell type associations, symptom-cell type associations, cell type-specificity of causal genes, the severity and frequency of the phenotypes, suitability for gene therapy delivery systems (e.g. recombinant adeno-associated viral vectors (rAAV)), as well as a quantitative analysis of phenotypic and genetic animal model translatability (Fig. 17). We validated these candidates by comparing

562 the proportional overlap with gene therapies that are presently in the market or undergoing clinical trials,  
563 in which we recovered 87% of all active gene therapies (Fig. 7, Fig. 16). Despite nominating a large number  
564 of putative targets, hypergeometric tests confirmed that our targets were strongly enriched for targets of  
565 existing therapies that are either approved or currently undergoing clinical trials.

566 From our target prioritisation pipeline results, we highlight cell type-specific mechanisms for ‘GM2-  
567 ganglioside accumulation’ in Tay-Sachs disease, spinocerebellar atrophy in spinocerebellar ataxia, and  
568 ‘Neuronal loss in central nervous system’ in a variety of diseases (Fig. 8). Of interest, all three of these  
569 neurodegenerative phenotypes involved alternatively activated (M2) macrophages. The role of macrophages  
570 in neurodegeneration is complex, with both neuroprotective and neurotoxic functions, including the  
571 clearance of misfolded proteins, the regulation of the blood-brain barrier, and the modulation of the immune  
572 response<sup>82</sup>. We also recapitulated prior evidence that microglia, the resident macrophages of the nervous  
573 system, are causally implicated in Alzheimer’s disease (AD) (Fig. 19)<sup>83</sup>. An important contribution of our  
574 current study is that we were able to pinpoint the specific phenotypes of AD caused by macrophages to  
575 neurofibrillary tangles and long-tract signs (reflexes that indicate the functioning of spinal long fiber tracts).  
576 Other AD-associated phenotypes were caused by other cell types (e.g. gastric goblet cells, enterocytes).

577 It should be noted that our study has several key limitations. First, while our cell type datasets are amongst  
578 the most comprehensive human scRNA-seq references currently available, they are nevertheless missing  
579 certain tissues, cell types (e.g. spermatocytes, oocytes), and life stages (post-natal childhood, senility). It is  
580 also possible that we have not captured certain cell state signatures that only occur in disease (e.g. disease-  
581 associated microglia<sup>84,85</sup>). Though we reasoned that using only control cell type signatures would mitigate  
582 bias towards any particular disease, and avoid degradation of gene signatures due to loss of function mutations.  
583 Second, the collective knowledge of gene-phenotype and gene-disease associations is far from complete and  
584 we fully anticipate that these annotations will continue to expand and change well into the future. It is  
585 for this reason we designed this study to be easily reproduced within a single containerised script so that  
586 we (or others) may rerun it with updated datasets at any point. Finally, causality is notoriously difficult  
587 to prove definitively from associative testing alone, and our study is not exempt from this rule. Despite  
588 this, there are several reasons to believe that our approach is able to better approximate causal relationships  
589 than traditional approaches. First, we did not intentionally preselect any subset of phenotypes or cell types  
590 to investigate here. Along with a scaling prestep during linear modelling, this means that all the results  
591 are internally consistent and can be directly compared to one another (in stark contrast to literature meta-  
592 analyses). Furthermore, for the phenotype gene signatures we used expert-curated GenCC annotations<sup>86,87</sup>  
593 to weight the current strength of evidence supporting a causal relationship between each gene and phenotype.  
594 This is especially important for phenotypes with large genes lists (thousands of annotations) for which some  
595 of the relationships may be tenuous. Within the cell type references, we deliberately chose to use specificity  
596 scores (rather than raw gene expression) as this normalisation procedure has previously been demonstrated

597 to better distinguish between signatures of highly similar cell types/subtypes<sup>88</sup>.

598 Common ontology-controlled frameworks like the HPO open a wealth of new opportunities, especially when  
599 addressing RDs. Services such as the Matchmaker Exchange<sup>89,90</sup> have enabled the discovery of hundreds of  
600 underlying genetic etiologies, and led to the diagnosis of many patients. This also opens the possibility of  
601 gathering cohorts of geographically dispersed patients to run clinical trials, the only viable option for treat-  
602 ment in many individuals. To further increase the number of individuals who qualify for these treatments,  
603 as well as the trial sample size, proposals have been made deviate from the traditional single-disease clinical  
604 trial model and instead perform basket trials on groups of RDs with shared molecular etiologies (SaME)<sup>91</sup>.

605 Moving forward, we are now actively seeking industry and academic partnerships to begin experimentally  
606 validating our multi-scale target predictions and exploring their potential for therapeutic translation. Never-  
607 theless, there are more promising therapeutic targets here than our research group could ever hope to pursue  
608 by ourselves. In the interest of accelerating research and ensuring RD patients are able to benefit from this  
609 work as quickly as possible, we have decided to publicly release all of the results described in this study.  
610 These can be accessed in multiple ways, including through a suite of R packages as well as a web app, the  
611 Rare Disease Celltyping Portal (<https://neurogenomics-ukdri.dsi.ic.ac.uk/>). The latter allows our results to  
612 be easily queried, filtered, visualised, and downloaded without any knowledge of programming. Through  
613 these resources we aim to make our findings useful to a wide variety of RD stakeholders including subdomain  
614 experts, clinicians, advocacy groups, and patients.

## 615 Conclusions

616 In this study we aimed to develop a methodology capable of generating high-throughput phenome-wide  
617 predictions while preserving the accuracy and clinical utility typically associated with more narrowly focused  
618 studies. With the rapid advancement of gene therapy technologies, and a regulatory landscape that is  
619 evolving to better meet the needs of a large and diverse patient population, there is finally momentum to  
620 begin to realise the promise of genomic medicine. This has especially important implications for the global  
621 RD community which has remained relatively neglected. Here, we have provided a scalable, cost-effective,  
622 and fully reproducible means of resolving the multi-scale, cell-type specific mechanisms of virtually all rare  
623 diseases.

## 624 Methods

### 625 Human Phenotype Ontology

626 The latest version of the HPO (release 2024-02-08) was downloaded from the EMBL-EBI Ontology Lookup  
627 Service<sup>92</sup> and imported into R using the `HPOExplorer` package. This R object was used to extract ontolog-  
628 ical relationships between phenotypes as well as to assign absolute and relative ontological levels to each  
629 phenotype. The latest version of the HPO phenotype-to-gene mappings and phenotype annotations were

630 downloaded from the official HPO GitHub repository and imported into R using `HPOExplorer`. This contains  
631 lists of genes associated with phenotypes via particular diseases, formatted as three columns in a table (gene,  
632 phenotype, disease).

633 However, not all genes have equally strong evidence of causality with a disease or phenotype, especially when  
634 considering that the variety of resources used to generate these annotations (OMIM, Orphanet, DECIPHER)  
635 use variable methodologies (e.g. expert-curated review of the medical literature vs. automated text mining  
636 of the literature). Therefore we imported data from the Gene Curation Coalition (GenCC)<sup>86,87</sup>, which (as  
637 of 2025-08-02) 24,112 evidence scores across 7,566 diseases and 5,533 genes. Evidence scores are defined  
638 by GenCC using a standardised ordinal rubric which we then encoded as a semi-quantitative score ranging  
639 from 0 (no evidence of disease-gene relationship) to 6 (strongest evidence of disease-gene relationship) (see  
640 Table 5). As each Disease-Gene pair can have multiple entries (from different studies) with different levels  
641 of evidence, we then summed evidence scores per Disease-Gene pair to generate aggregated Disease-by-Gene  
642 evidence scores. This procedure can be described as follows.

643 Let us denote:

- 644 •  $D$  as diseases.  
645 •  $P$  as phenotypes in the HPO.  
646 •  $G$  as genes  
647 •  $S$  as the evidence scores describing the strength of the relationship between each Disease-Gene pair.  
648 •  $M_{ij}$  as the aggregated Disease-by-Gene evidence score matrix.

$$M_{ij} = \sum_{k=1}^f D_i G_j S_k$$

649 Next, we extracted Disease-Gene-Phenotype relationships from the annotations file distributed by the HPO  
650 (*phenotype\_to\_genes.txt*). This provides a list of genes associated with phenotypes via particular diseases,  
651 but does not include any strength of evidence scores.

652 Here we define: -  $A_{ijk}$  as the Disease-Gene-Phenotype relationships. -  $D_i$  as the  $i$ th disease. -  $G_j$  as the  $j$ th  
653 gene. -  $P_k$  as the  $k$ th phenotype.

$$A_{ijk} = D_i G_j P_k$$

654 In order to assign evidence scores to each Phenotype-Gene relationship, we combined the aforementioned  
655 datasets from GenCC ( $M_{ij}$ ) and HPO ( $A_{ijk}$ ) by merging on the gene and disease ID columns. For each

656 phenotype, we then computed the mean of Disease-Gene scores across all diseases for which that phenotype  
 657 is a symptom. This resulted in a final 2D tensor of Phenotype-by-Gene evidence scores ( $L_{ij}$ ):

658

659

660

Tensor of Phenotype-by-Gene  
evidence scores

$$L_{ij} = \begin{cases} \frac{\sum_{k=1}^f D_i G_j P_k}{f}, & \text{if } D_i G_j \in A, \\ 1, & \text{if } D_i G_j \notin A \end{cases}$$

661

662

663

Tensor of Disease-by-Gene  
evidence scores

Phenotype

Disease-by-Gene-by-Phenotype  
relationships

#### 664 Construction of the tensor of Phenotype-by-Gene evidence scores.

665

666

667 Histograms of evidence score distributions at each step in processing can be found in Fig. 9.

#### 668 Single-cell transcriptomic atlases

669 In this study, the gene by cell type specificity matrix was constructed using the Descartes Human transcriptome  
 670 atlas of foetal gene expression, which contains a mixture of single-nucleus and single-cell RNA-seq  
 671 data (collected with sci-RNA-seq3)<sup>32</sup>. This dataset contains 377,456 cells representing 77 distinct cell types  
 672 across 15 tissues. All 121 human foetal samples ranged from 72 to 129 days in estimated postconceptual age.  
 673 To independently replicate our findings, we also used the Human Cell Landscape which contains single-cell  
 674 transcriptomic data (collected with microwell-seq) from embryonic, foetal, and adult human samples across  
 675 49 tissues<sup>33</sup>.

676 Specificity matrices were generated separately for each transcriptomic atlas using the R package EWCE  
 677 (v1.11.3)<sup>88</sup>. Within each atlas, cell types were defined using the authors' original freeform annotations  
 678 in order to preserve the granularity of cell subtypes as well as incorporate expert-identified rare cell types.  
 679 Cell types were only aligned and aggregated to the level of corresponding Cell Ontology (CL)<sup>39</sup>  
 680 annotations afterwards when generating summary figures and performing cross-atlas analyses. Using the original  
 681 gene-by-cell count matrices from each single-cell atlas, we computed gene-by-cell type expression specificity  
 682 matrices as follows. Genes with very no expression across any cell types were considered to be uninformative  
 683 and were therefore removed from the input gene-by-cell matrix  $F(g, i, c)$ .

684 Next, we calculated the mean expression per cell type and normalised the resulting matrix to transform it

685 into a gene-by-cell type expression specificity matrix ( $S_{gc}$ ). In other words, each gene in each cell type had  
 686 a 0-1 score where 1 indicated the gene was mostly specifically expressed in that particular cell type relative  
 687 to all other cell types. This procedure was repeated separately for each of the single-cell atlases and can be  
 688 summarised as:

689

690 **Compute mean expression of each gene per cell type**

$$S_{gc} = \frac{\sum_{i=1}^{|L|} F_{gic}}{\sum_{r=1}^k \left( \frac{\sum_{i=1}^{|L|} F_{gic}}{N_c} \right)}$$

691 **Compute row sums of  
mean gene-by-cell type matrix**

692

#### 694 **Phenotype-cell type associations**

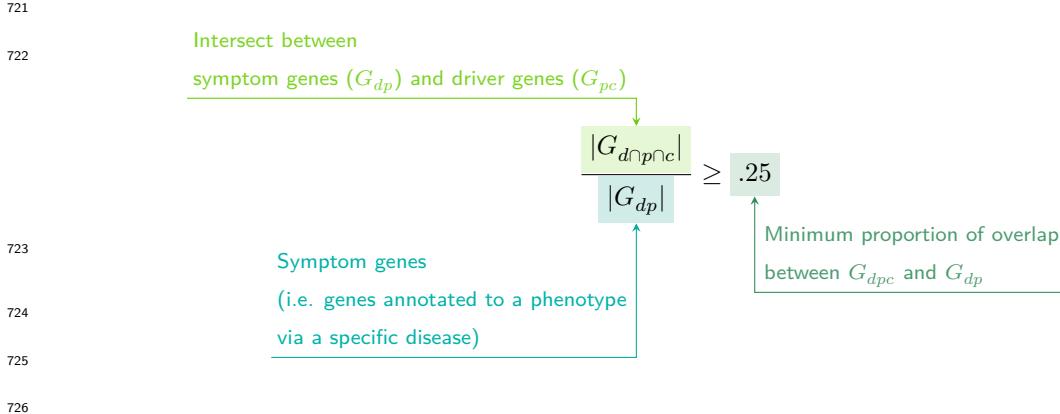
695 To test for relationships between each pairwise combination of phenotype (n=11,047) and cell type (n=201)  
 696 we ran a series of univariate generalised linear models implemented via the `stats::glm` function in R. First,  
 697 we filtered the gene-by-phenotype evidence score matrix ( $L_{ij}$ ) and the gene-by-cell type expression specificity  
 698 matrix ( $S_{gc}$ ) to only include genes present in both matrices (n=4,949 genes in the Descartes Human analyses;  
 699 n=4,653 genes in the Human Cell Landscape analyses). Then, within each matrix any rows or columns with a  
 700 sum of 0 were removed as these were uninformative data points that did not vary. To improve interpretability  
 701 of the results  $\beta$  coefficient estimates across models (i.e. effect size), we performed a scaling prestep on all  
 702 dependent and independent variables. Initial tests showed that this had virtually no impact on the total  
 703 number of significant results or any of the benchmarking metrics based on p-value thresholds Fig. 2. This  
 704 scaling prestep improved our ability to rank cell types by the strength of their association with a given  
 705 phenotype as determined by separate linear models.

706 We repeated the aforementioned procedure separately for each of the single-cell references. Once all results  
 707 were generated using both cell type references (2,206,994 association tests total), we applied Benjamini-  
 708 Hochberg false discovery rate<sup>93</sup> (denoted as  $FDR_{pc}$ ) to account for multiple testing. Of note, we applied  
 709 this correction across all results at once (as opposed to each single-cell reference separately) to ensure the  
 710  $FDR_{pc}$  was stringently controlled for across all tests performed in this study.

#### 711 **Symptom-cell type associations**

712 Here we define a symptom as a phenotype as it presents within the context of the specific disease. The features  
 713 of a given symptom can be described as the subset of genes annotated to phenotype  $p$  via a particular disease

714  $d$ , denoted as  $G_{dp}$  (see Fig. 10). To attribute our phenotype-level cell type enrichment signatures to specific  
 715 diseases, we first identified the gene subset that was most strongly driving the phenotype-cell type association  
 716 by computing the intersect of genes that were both in the phenotype annotation and within the top 25%  
 717 specificity percentile for the associated cell type. We then computed the intersect between symptom genes  
 718 ( $G_{dp}$ ) and driver genes ( $G_{pc}$ ), resulting in the gene subset  $G_{d \cap p \cap c}$ . Only  $G_{d \cap p \cap c}$  gene sets with 25% or greater  
 719 overlap with the symptom gene subset ( $G_{dp}$ ) were kept. This procedure was repeated for all phenotype-cell  
 720 type-disease triads, which can be summarised as follows:



## 727 Validation of expected phenotype-cell type relationships

728 We first sought to confirm that our tests (across both single-cell references) were able to recover expected  
 729 phenotype-cell type relationships across seven high-level branches within the HPO (Fig. 2), including ab-  
 730 normalities of the cardiovascular system, endocrine system, eye, immune system, musculoskeletal system,  
 731 nervous system, and respiratory system. Within each branch the number of significant tests in a given  
 732 cell type were plotted (Fig. 2b). Mappings between freeform annotations (the level at which we performed  
 733 our phenotype-cell type association tests) provided by the original atlas authors and their closest CL term  
 734 equivalents were provided by CellxGene<sup>30</sup>. CL terms along the *x-axis* of Fig. 2b were assigned colours corre-  
 735 sponding to which HPO branch showed the greatest number of enrichments (after normalising within each  
 736 branch to account for differences in scale). The normalised colouring allows readers to quickly assess which  
 737 HPO branch was most often associated with each cell type, while accounting for differences in the number  
 738 of phenotypes across branches. We then ran a series of Analysis of Variance (ANOVA) tests to determine  
 739 whether (within a given branch) a given cell type was more often enriched ( $FDR < 0.05$ ) within that branch  
 740 relative to all of the other HPO branches of an equivalent level in the ontology (including all branches not  
 741 shown in Fig. 2b). After applying Benjamini-Hochberg multiple testing correction<sup>93</sup> (denoted as  $FDR_{b,c}$ ),  
 742 we annotated each respective branch-by-cell type bar according to the significance (\*\*\*\* :  $FDR_{b,c} < 1e-04$ ,  
 743 \*\*\* :  $FDR_{b,c} < 0.001$ , \*\* :  $FDR_{b,c} < 0.01$ , \* :  $FDR_{b,c} < 0.05$ ). Cell types in Fig. 2a-b were ordered along  
 744 the *x-axis* according to a dendrogram derived from the CL ontology (Fig. 2c), which provides ground-truth

745 semantic relationships between all cell types (e.g. different neuronal subtypes are grouped together).  
746 As an additional measure of the accuracy of our phenotype-cell types test results we identified conceptually  
747 matched branches across the HPO and the CL (Fig. 2d and Table 6). For example, ‘Abnormality of the  
748 cardiovascular system’ in the HPO was matched with ‘cardiocytes’ in the CL which includes all cell types  
749 specific to the heart. Analogously, ‘Abnormality of the nervous system’ in the HPO was matched with ‘neural  
750 cell’ in the CL which includes all descendant subtypes of neurons and glia. This cross-ontology matching  
751 was repeated for each HPO branch and can be referred to as on-target cell types. Within each branch, the  
752  $-\log_{10}(FDR_{pc})$  values of on-target cell types were binned by rounding to the nearest integer (*x-axis*) and  
753 the percentage of tests for on-target cell types relative to all cell types were computed at each bin (*y-axis*)  
754 (Fig. 2d). The baseline level (dotted horizontal line) illustrates the percentage of on-target cell types relative  
755 to the total number of observed cell types. Any percentages above this baseline level represent greater than  
756 chance representation of the on-target cell types in the significant tests.

#### 757 Validation of inter- and intra-dataset consistency

758 We tested for inter-dataset consistency of our phenotype-cell type association results across different single-  
759 cell reference datasets (Descartes Human and Human Cell Landscape). For all tests reported here, the  
760 relevant association metrics (p-values or effect size) were first averaged to the level of ancestral HPO terms  
761 (5 levels down the hierarchy) to reduce figure size. For association tests with exactly matching Cell Ontology  
762 ID across the two references, we tested for a relationship between the p-values generated with each of the  
763 references by fitting linear regression model (`stats::lm` via the R function `ggstatsplot::ggscatterstats`).  
764 Next, we performed an additional linear regression between the effect sizes (each GLM model’s  $R^2$  estimates  
765 after applying a  $\log_2$  fold-change transformation) of all significant phenotype-cell type associations ( $FDR <$   
766 0.05) with exactly matching cell types across the two references.

767 We also tested for intra-dataset consistency within the Human Cell Landscape by running additional linear  
768 regressions between the phenotype-cell type association test statistics of the foetal and the adult samples (us-  
769 ing both p-values and model  $R^2$  estimates). While we would not expect the same exact cell type associations  
770 across different developmental stages, we would nevertheless expect there to be some degree of correlation  
771 between the developing and mature versions of the same cell types.

#### 772 More specific phenotypes are associated with fewer genes and cell types

773 To explore the relationship between HPO phenotype specificity and various metrics from our results, we  
774 computed the information content (IC) scores for each term in the HPO. IC is a measure of how much  
775 specific information a term within an ontology contains. In general, terms deeper in an ontology (closer to the  
776 leaves) are more specific, and thus informative, than terms at the very root of the ontology (e.g. ‘Phenotypic  
777 abnormality’). Where  $k$  denotes the number of offspring terms (including the term itself) and  $N$  denotes the

778 total number of terms in the ontology, IC can be calculated as:

$$IC = -\log\left(\frac{k}{N}\right)$$

779 Next, IC scores were quantised into 10 bins using the `ceiling` R function to improve visualisation. We  
780 then performed a series of linear regressions between phenotype binned IC scores and: 1) number of genes  
781 annotated per HPO phenotype, 2) the number of significantly associated cell types per HPO phenotype, and  
782 3) the model estimate of each significant phenotype-cell type associations (at FDR < 0.05) after taking the  
783 log of the absolute value ( $\log_2(|estimate|)$ ).

#### 784 Monarch Knowledge Graph recall

785 Finally, we gathered known phenotype-cell type relationships from the Monarch Knowledge Graph (MKG),  
786 a comprehensive database of links between many aspects of disease biology<sup>40</sup>. This currently includes 103  
787 links between HPO phenotypes (n=103) and CL cell types (n=79). Of these, we only considered the 82  
788 phenotypes that we were able to test given that our ability to generate associations was dependent on  
789 the existence of gene annotations within the HPO. We considered instances where we found a significant  
790 relationship between exactly matching pairs of HPO-CL terms as a hit.

791 However, as the cell types in MKG were not necessarily annotated at the same level as our single-cell refer-  
792 ences, we considered instances where the MKG cell type was an ancestor term of our cell type (e.g. ‘myeloid  
793 cell’ vs. ‘monocyte’), or *vice versa*, as hits. We also adjusted ontological distance by computing the ratio  
794 between the observed ontological distance and the smallest possible ontological distance for that cell type  
795 given the cell type that were available in our references ( $dist_{adjusted} = \left(\frac{dist_{observed}+1}{dist_{minimum}+1}\right) - 1$ ). This provides  
796 a way of accurately measuring how dissimilar our identified cell types were for each phenotype-cell type  
797 association (Fig. 12).

#### 798 Prioritising phenotypes based on severity

799 Only a small fraction of the the phenotypes in HPO (<1%) have metadata annotations containing informa-  
800 tion on their time course, consequences, and severity. This is due to the time-consuming nature of manually  
801 annotating thousands of phenotypes. To generate such annotations at scale, we previously used Generative  
802 Pre-trained Transformer 4 (GPT-4), a large language model (LLM) as implemented within OpenAI’s Appli-  
803 cation Programming Interface (API)<sup>37</sup>. After extensive prompt engineering and ground-truth benchmarking,  
804 we were able to acquire annotations on how often each phenotype directly causes intellectual disability, death,  
805 impaired mobility, physical malformations, blindness, sensory impairments, immunodeficiency, cancer, re-  
806 duced fertility, or is associated with a congenital onset. These criteria were previously defined in surveys  
807 of medical experts as a means of systematically assessing phenotype severity<sup>94</sup>. Responses for each metric

were provided in a consistent one-word format which could be one of: ‘never’, ‘rarely’, ‘often’, ‘always’. This procedure was repeated in batches (to avoid exceeding token limits) until annotations were gathered for 16,982/18,082 HPO phenotypes.

We then encoded these responses into a semi-quantitative scoring system ('never'=0, 'rarely'=1, 'often'=2, 'always'=3), which were then weighted by multiplying a semi-subjective scoring of the relevance of each metric to the concept of severity on a scale from 1.0-6.0, with 6.0 being the most severe ('death'=6, 'intellectual\_disability'=5, 'impaired\_mobility'=4, 'physical\_malformations'=3, 'blindness'=4, 'sensory\_impairments'=3, 'immunodeficiency'=3, 'cancer'=3, 'reduced\_fertility'=1, 'congenital\_onset'=1). Finally, the product of the score was normalised to a quantitative severity score ranging from 0-100, where 100 is the theoretical maximum severity score. This phenotype severity scoring procedure can be expressed as follows.

819 Let us denote:

- $p$  : a phenotype in the HPO.
  - $j$  : the identity of a given annotation metric (i.e. clinical characteristic, such as ‘intellectual disability’ or ‘congenital onset’).
  - $W_j$ : the assigned weight of metric  $j$ .
  - $F_j$ : the maximum possible value for metric  $j$ , equal to 3 (“always”). This value is equivalent across all  $j$  annotations.
  - $F_{pj}$  : the numerically encoded value of annotation metric  $j$  for phenotype  $p$ .
  - $NSS_p$ : the final composite severity score for phenotype  $p$  after applying normalisation to align values to a 0-100 scale and ensure equivalent meaning regardless of which other phenotypes are being analysed in addition to  $p$ . This allows for direct comparability of severity scores across studies with different sets of phenotypes.

$$NSS_p = \frac{\sum_{j=1}^m (F_{pj} \times W_j)}{\sum_{j=1}^m (\max\{F_j\} \times W_j)} \times 100$$

Normalised Severity Score  
for each phenotype

Sum of weighted annotation values  
across all metrics

Numerically encoded annotation value  
of metric  $j$  for phenotype  $p$

Weight for metric  $j$

Theoretical maximum severity score

837 Using the numerically encoded GPT annotations (0=“never”, 1=“rarely”, 2=“often”, 3=“always”) we com-  
838 puted the mean encoded value per cell type within each annotation. One-sided Wilcoxon rank-sum tests  
839 were run using the `rstatix::wilcox_test()` function to test whether each cell type was associated with  
840 more severe phenotypes relative to all other cell types. This procedure was repeated for severity annotation  
841 independently (death, intellectual disability, impaired mobility, etc.) Fig. 5a. Next, we performed a Pear-  
842 son correlation test between the number of phenotypes that a cell type is significantly associated with (at  
843 FDR<0.05) has a relationship with the mean composite GPT severity score of those phenotypes (Fig. 5b).  
844 This was performed using the `ggstatsplot::ggscatterstats()` R function.

845 **Congenital phenotypes are associated with foetal cell types**

846 The GPT-4 annotations also enabled us to assess whether foetal cell types were more often significantly  
847 associated with congenital phenotypes in our Human Cell Landscape results as this single-cell reference  
848 contained both adult and foetal versions of cell types (Fig. 6). To do this, we performed a chi-squared ( $\chi^2$ )  
849 test on the proportion of significantly associated cell types containing any of the substrings ‘fetal’, ‘fetus’,  
850 ‘primordial’, ‘hESC’ or ‘embryonic’ (within cell types annotations from the original Human Cell Landscape  
851 authors<sup>33</sup>) vs. those associated without, stratified by how often the corresponding phenotype had a congenital  
852 onset according to the GPT phenotype annotations (including ‘never’, ‘rarely’, ‘often’, ‘always’). In addition,  
853 a series of  $\chi^2$  tests were performed within each congenital onset frequency strata, to determine whether the  
854 observed proportion of foetal cell types vs. non-foetal cell types significantly deviated from the proportions  
855 expected by chance.

856 We next tested whether the proportion of tests with significant associations with foetal cell types varied  
857 across the major HPO branches using a  $\chi^2$  test. We also performed separate  $\chi^2$  test within each branch to  
858 determine whether the proportion of significant associations with foetal cell types was significantly different  
859 from chance.

860 Next, we aimed to create a continuous metric from -1 to 1 that indicated how biased each phenotype is  
861 towards associations with the foetal or adult form of a cell type. For each phenotype we calculated the  
862 foetal-adult bias score as the difference in the association p-values between the foetal and adult version  
863 of the equivalent cell type (foetal-adult bias :  $p_{adult} - p_{foetal} = \Delta p \in [-1, 1]$ ). A score of 1 indicates the  
864 phenotype is only associated with the foetal version of the cell type and -1 indicates the phenotype is only  
865 associated with the adult version of the cell type.

866 In order to summarise higher-order HPO phenotype categories that were most biased towards foetal  
867 or adult cell types, ontological enrichment tests were run on the phenotypes with the top/bottom  
868 50 greatest/smallest foetal-adult bias scores. The enrichment tests were performed using the  
869 `simona::dag_enrich_on_offsprings` function, which uses a hypergeometric test to determine whether a  
870 list of terms in an ontology are enriched for offspring terms (descendants) of a given ancestor term within

871 the ontology. Phenotypes categories with an HPO ontological enrichment a p-value < 0.05 were considered  
872 significant.

873 We were similarly interested in which higher-order cell type categories tended to be most commonly associated  
874 with these strongly foetal-/adult-biased phenotype s. Another set of ontological enrichment tests were run on  
875 the cell types associated with the top/bottom 50 phenotypes from the previous analysis. The CL ontology-  
876 aligned IDs for each group cell types were fed into the `simona::dag_enrich_on_offsprings` using the CL  
877 ontology. Significantly enriched cell type categories were defined as those with a CL ontological enrichment  
878 p-value < 0.05.

#### 879 Therapeutic target identification

880 We developed a systematic and automated strategy for identifying putative cell type-specific gene targets  
881 for each phenotype based on a series of filters at phenotype, cell type, and gene levels. The entire target  
882 prioritisation procedure can be replicated with a single function: `MSTExplorer::prioritise_targets`. This  
883 function automates all of the reference data gathering (e.g. phenotype metadata, cell type metadata, cell  
884 type signature reference, gene lengths, severity tiers) and takes a variety of arguments at each step for greater  
885 customisability. Each step is described in detail in Table 3. Phenotypes that often or always caused physical  
886 malformations (according to the GPT-4 annotations) were also removed from the final prioritised targets  
887 list, as these were unlikely to be amenable to gene therapy interventions. Finally, phenotypes were sorted  
888 by their composite severity scores such that the most severe phenotypes were ranked the highest.

#### 889 Therapeutic target validation

890 To assess whether our prioritised therapeutic targets were likely to be viable, we computed the overlap  
891 between our gene targets and those of existing gene therapies at various stages of clinical development  
892 (Fig. 7). Gene targets were obtained for each therapy from the Therapeutic Target Database (TTD; release  
893 2025-08-08) and mapped onto standardised HUGO Gene Nomenclature Committee (HGNC) gene symbols  
894 using the `orthogene` R package. We stratified our overlap metrics according to whether the therapies had  
895 failed (unsuccessful clinical trials or withdrawn), or were non-failed (successful or ongoing clinical trials).  
896 We then conducted hypergeometric tests to determine whether the observed overlap between our prioritised  
897 targets and the non-failed therapy targets was significantly greater than expected by chance (i.e. enrichment).  
898 We also conducted a second hypergeometric test to determine whether the observed overlap between our  
899 prioritised targets and the failed therapy targets was significantly less than expected by chance (i.e. depletion).  
900 Finally, we repeated the analysis against all therapeutic targets, not just those of gene therapies, to determine  
901 whether our prioritised targets had relevance to other therapeutic modalities.

902 **Experimental model translatability**

903 To improve the likelihood of successful translation between preclinical animal models and human patients,  
904 we created an interspecies translatability prediction tool for each phenotype nominated by our gene therapy  
905 prioritised pipeline (Fig. 17). First, we extracted ontological similarity scores of homologous phenotypes  
906 across species from the MKG<sup>40</sup>. Briefly, the ontological similarity scores ( $SIM_o$ ) are computed for each  
907 homologous pair of phenotypes across two ontologies by calculating the overlap in homologous phenotypes  
908 that are ancestors or descendants of the target phenotype. Next, we generated genotypic similarity scores  
909 ( $SIM_g$ ) for each homologous phenotype pair by computing the proportion of 1:1 orthologous genes using  
910 gene annotation from their respective ontologies. Interspecies orthologs were also obtained from the MKG.  
911 Finally, both scores are multiplied together to yield a unified ontological-genotypic similarity score ( $SIM_{og}$ ).

912 **Novel R packages**

913 To facilitate all analyses described in this study and to make them more easily reproducible by others, we  
914 created several open-source R packages. [KGExplorer](#) imports and analyses large-scale biomedical knowledge  
915 graphs and ontologies. [HPOExplorer](#) aids in managing and querying the directed acyclic ontology graph  
916 within the HPO. [MSTExplorer](#) facilitates the efficient analysis of many thousands of phenotype-cell type  
917 association tests, and provides a suite of multi-scale therapeutic target prioritisation and visualisation func-  
918 tions. These R packages also include various functions for distributing the post-processed results from this  
919 study in an organised, tabular format. Of note, `MSTExplorer::load_example_results` loads all summary  
920 statistics from our phenotype-cell type tests performed here.

921 **Rare Disease Celltyping Portal**

922 To further increase the ease of access for stakeholders in the RD community without the need for program-  
923 matic experience, we developed a series of web apps to interactively explore, visualise, and download the  
924 results from our study. Collectively, these web apps are called the Rare Disease Celltyping Portal. The  
925 website can be accessed at <https://neurogenomics-ukdri.dsi.ic.ac.uk/>.

926 The Rare Disease Celltyping Portal integrates diverse datasets, including the HPO, cell types, genes, and phe-  
927 notype severity, into a unified platform that allows users to perform flexible, bidirectional queries. Users can  
928 start from any entry point: either phenotype, cell type, genes, or severity, and seamlessly trace relationships  
929 across these dimensions.

930 The portal provides a dynamic and intuitive exploration experience with its real-time interaction capabil-  
931 ities and responsive interface including network graphs, bar charts, and heat maps. It has the ability to  
932 handle large datasets efficiently and offer fast query response by building with FARM stack (FastAPI, React,  
933 MongoDB). The portal is designed for a broad audience, including researchers, clinicians, and biologists, by  
934 offering user-friendly navigation and interactive visual outputs. By enabling users to intuitively explore com-

plex biological relationships, the portal aims to accelerate rare disease research, enhance diagnostic accuracy, and drive therapeutic innovation.

All code used to generate the website can be found at <https://github.com/neurogenomics/Rare-Disease-Web-Portal>.

## Mappings

Mappings from the HPO to other medical ontologies were extracted from the EMBL-EBI Ontology Xref Service (Oxo; <https://www.ebi.ac.uk/spot/oxo/>) by selecting the National Cancer Institute metathesaurus (NCIm) as the target ontology and either “SNOMED CT”, “UMLS”, “ICD-9” or “ICD-10CM” as the data source. HPO terms were then selected as the ID framework with to mediate the cross-ontology mappings. Mappings between each pair of ontologies were then downloaded, stored in a tabular format, and uploaded to the public **HPOExplorer** Releases page (<https://github.com/neurogenomics/HPOExplorer/releases>).

## Data Availability

All data is publicly available through the following resources:

- Human Phenotype Ontology (<https://hpo.jax.org>)
- GenCC (<https://thegencc.org/>)
- Descartes Human scRNA-seq atlas (<https://cellxgene.cziscience.com/collections/c114c20f-1ef4-49a5-9c2e-d965787fb90c>)
- Human Cell Landscape scRNA-seq atlas (<https://cellxgene.cziscience.com/collections/38833785-fac5-48fd-944a-0f62a4c23ed1>)
- Processed Cell Type Datasets (*ctd\_DescartesHuman.rds* and *ctd\_HumanCellLandscape.rds*; <https://github.com/neurogenomics/MSTExplorer/releases>)
- Gene x Phenotype association matrix (*hpo\_matrix.rds*; <https://github.com/neurogenomics/MSTExplorer/releases>)
- GPT-4 phenotype severity annotations ([https://github.com/neurogenomics/rare\\_disease\\_celltyping/releases/download/latest/gpt\\_check\\_annot.csv.gz](https://github.com/neurogenomics/rare_disease_celltyping/releases/download/latest/gpt_check_annot.csv.gz))
- Full phenotype-cell type association test results [https://github.com/neurogenomics/MSTExplorer/releases/download/v0.1.10/phenomix\\_results.tsv.gz](https://github.com/neurogenomics/MSTExplorer/releases/download/v0.1.10/phenomix_results.tsv.gz)
- Rare Disease Celltyping Portal (<https://neurogenomics-ukdri.dsi.ic.ac.uk/>)

## Code Availability

All code is made freely available through the following GitHub repositories:

- **KGExplorer** (<https://github.com/neurogenomics/KGExplorer>)

- HPOExplorer (<https://github.com/neurogenomics/HPOExplorer>)
- MSTExplorer (<https://github.com/neurogenomics/MSTExplorer>)
- Code to replicate analyses ([https://github.com/neurogenomics/rare\\_disease\\_celltyping](https://github.com/neurogenomics/rare_disease_celltyping))
- Cell type-specific gene target prioritisation ([https://neurogenomics.github.io/RareDiseasePrioritisation/reports/prioritise\\_targets](https://neurogenomics.github.io/RareDiseasePrioritisation/reports/prioritise_targets))
- Complement system gene list (<https://www.genenames.org/data/genegroup/#!/group/492>)

## 972 Acknowledgements

973 We would like to thank the following individuals for their insightful feedback and assistance with data  
974 resources: Sarah J. Marzi, Gerton Lunter, Peter Robinson, Melissa Haendel, Ben Coleman, Nico Matentzoglu,  
975 Shawn T. O’Neil, Alan E. Murphy, Sarada Gurung.

## 976 Funding

977 This work was supported by a UK Dementia Research Institute (UK DRI) Future Leaders Fellowship  
978 [MR/T04327X/1] and the UK DRI which receives its funding from UK DRI Ltd, funded by the UK Medical  
979 Research Council, Alzheimer’s Society and Alzheimer’s Research UK.

## 980 References

- 981 1. Ferreira, C. R. The burden of rare diseases. *Am. J. Med. Genet. A* **179**, 885–892 (2019).
- 982 2. Zhu, Q. *et al.* An integrative knowledge graph for rare diseases, derived from the genetic and rare  
diseases information center (GARD). *J. Biomed. Semantics* **11**, 13 (2020).
- 983 3. Rare diseases BioResource.
- 984 4. Marwaha, S., Knowles, J. W. & Ashley, E. A. A guide for the diagnosis of rare and undiagnosed  
disease: Beyond the exome. *Genome Med.* **14**, 23 (2022).
- 985 5. Molster, C. *et al.* Survey of healthcare experiences of australian adults living with rare diseases.  
*Orphanet J. Rare Dis.* **11**, 30 (2016).
- 986 6. Halley, M. C., Smith, H. S., Ashley, E. A., Goldenberg, A. J. & Tabor, H. K. A call for an integrated  
approach to improve efficiency, equity and sustainability in rare disease research in the united states.  
*Nat. Genet.* **54**, 219–222 (2022).
- 987 7. Institute of Medicine (US) Committee on Accelerating Rare Diseases Research and Orphan Product  
Development, Field, M. J. & Boat, T. F. *Coverage and Reimbursement: Incentives and Disincentives  
for Product Development*. (National Academies Press (US), 2010).
- 988 8. Yates, N. & Hinkel, J. The economics of moonshots: Value in rare disease drug development. *Clin.  
Transl. Sci.* **15**, 809–812 (2022).

- 989 9. Nuijten, M. Pricing zolgensma - the world's most expensive drug. *J Mark Access Health Policy* **10**,  
2022353 (2022).
- 990 10. Thielen, F. W., Heine, R. J. S. D., Berg, S. van den, Ham, R. M. T. T. & Groot, C. A. U. Towards  
sustainability and affordability of expensive cell and gene therapies? Applying a cost-based pricing  
model to estimate prices for libmeldy and zolgensma. *Cytotherapy* **24**, 1245–1258 (2022).
- 991 11. Gargano, M. A. *et al.* The human phenotype ontology in 2024: Phenotypes around the world. *Nucleic  
Acids Res.* **52**, D1333–D1346 (2024).
- 992 12. Köhler, S. *et al.* Expansion of the human phenotype ontology (HPO) knowledge base and resources.  
*Nucleic Acids Res.* **47**, D1018–D1027 (2019).
- 993 13. Köhler, S. *et al.* The human phenotype ontology in 2021. *Nucleic Acids Res.* **49**, D1207–D1217  
(2021).
- 994 14. Robinson, P. N. *et al.* The human phenotype ontology: A tool for annotating and analyzing human  
hereditary disease. *Am. J. Hum. Genet.* **83**, 610–615 (2008).
- 995 15. Nguengang Wakap, S. *et al.* Estimating cumulative point prevalence of rare diseases: Analysis of the  
orphanet database. *Eur. J. Hum. Genet.* **28**, 165–173 (2020).
- 996 16. Rare diseases, common challenges. *Nat. Genet.* **54**, 215 (2022).
- 997 17. Amberger, J. S., Bocchini, C. A., Scott, A. F. & Hamosh, A. OMIM.org: Leveraging knowledge across  
phenotype-gene relationships. *Nucleic Acids Res.* **47**, D1038–D1043 (2019).
- 998 18. Amberger, J. S. & Hamosh, A. Searching online mendelian inheritance in man (OMIM): A knowl-  
edgebase of human genes and genetic phenotypes. *Curr. Protoc. Bioinformatics* **58**, 1.2.1–1.2.12  
(2017).
- 999 19. McKusick, V. A. Mendelian inheritance in man and its online version, OMIM. *Am. J. Hum. Genet.*  
**80**, 588–604 (2007).
- 1000 20. Maiella, S., Rath, A., Angin, C., Mousson, F. & Kremp, O. [Orphanet and its consortium: Where to  
find expert-validated information on rare diseases]. *Rev. Neurol.* **169 Suppl 1**, S3–8 (2013).
- 1001 21. Weinreich, S. S., Mangon, R., Sikkens, J. J., Teeuw, M. E. en & Cornel, M. C. [Orphanet: A european  
database for rare diseases]. *Ned. Tijdschr. Geneeskd.* **152**, 518–519 (2008).
- 1002 22. Firth, H. V. *et al.* DECIPHER: Database of chromosomal imbalance and phenotype in humans using  
ensembl resources. *Am. J. Hum. Genet.* **84**, 524–533 (2009).
- 1003 23. Chang, E. & Mostafa, J. [The use of SNOMED CT, 2013-2020: a literature review](#). *Journal of the  
American Medical Informatics Association* **28**, 2017–2026 (2021).
- 1004 24. Kim, M. C., Nam, S., Wang, F. & Zhu, Y. [Mapping scientific landscapes in UMLS research: a  
scientometric review](#). *Journal of the American Medical Informatics Association* **27**, 1612–1624 (2020).

- 1005 25. Humphreys, B. L., Del Fiol, G. & Xu, H. [The UMLS knowledge sources at 30: indispensable to current research and applications in biomedical informatics](#). *Journal of the American Medical Informatics Association* **27**, 1499–1501 (2020).
- 1006 26. Krawczyk, P. & Święcicki, Ł. [ICD-11 vs. ICD-10 – a review of updates and novelties introduced in the latest version of the WHO international classification of diseases](#). *Psychiatria Polska* **54**, 7–20 (2020).
- 1007 27. Baysoy, A., Bai, Z., Satija, R. & Fan, R. The technological landscape and applications of single-cell multi-omics. *Nat. Rev. Mol. Cell Biol.* **24**, 695–713 (2023).
- 1008 28. Haque, A., Engel, J., Teichmann, S. A. & Lönnberg, T. A practical guide to single-cell RNA-sequencing for biomedical research and clinical applications. *Genome Med.* **9**, 75 (2017).
- 1009 29. Qi, R. & Zou, Q. Trends and potential of machine learning and deep learning in drug study at Single-Cell level. *Research* **6**, 0050 (2023).
- 1010 30. CZI Single-Cell Biology Program *et al.* CZ CELL×GENE discover: A single-cell data platform for scalable exploration, analysis and modeling of aggregated data. *bioRxiv* 2023.10.30.563174 (2023).
- 1011 31. Svensson, V., Veiga Beltrame, E. da & Pachter, L. A curated database reveals trends in single-cell transcriptomics. *Database* **2020**, (2020).
- 1012 32. Cao, J. *et al.* A human cell atlas of fetal gene expression. *Science* **370**, (2020).
- 1013 33. Han, X. *et al.* Construction of a human cell landscape at single-cell level. *Nature* **581**, 303–309 (2020).
- 1014 34. Cao, J. *et al.* [A human cell atlas of fetal gene expression](#). *Science* **370**, eaba7721 (2020).
- 1015 35. Kawabata, H. *et al.* Improving cell-specific recombination using AAV vectors in the murine CNS by capsid and expression cassette optimization. *Molecular Therapy Methods & Clinical Development* **32**, (2024).
- 1016 36. O’Carroll, S. J., Cook, W. H. & Young, D. [AAV targeting of glial cell types in the central and peripheral nervous system and relevance to human gene therapy](#). *Frontiers in Molecular Neuroscience* **13**, (2021).
- 1017 37. Murphy, K., Schilder, B. M. & Skene, N. G. Harnessing generative AI to annotate the severity of all phenotypic abnormalities within the Human Phenotype Ontology. doi:[10.1101/2024.06.10.24308475](https://doi.org/10.1101/2024.06.10.24308475).
- 1018 38. DiStefano, M. T. *et al.* [The gene curation coalition: A global effort to harmonize gene–disease evidence resources](#). *Genetics in Medicine* **24**, 1732–1742 (2022).
- 1019 39. Diehl, A. D. *et al.* The cell ontology 2016: Enhanced content, modularization, and ontology interoperability. *J. Biomed. Semantics* **7**, 44 (2016).
- 1020 40. Putman, T. E. *et al.* The monarch initiative in 2024: An analytic platform integrating phenotypes, genes and diseases across species. *Nucleic Acids Res.* **52**, D938–D949 (2024).
- 1021 41. Heim, C. E. *et al.* Myeloid-derived suppressor cells contribute to staphylococcus aureus orthopedic biofilm infection. *J. Immunol.* **192**, 3778–3792 (2014).

- 1022 42. Pidwill, G. R., Gibson, J. F., Cole, J., Renshaw, S. A. & Foster, S. J. The role of macrophages in  
staphylococcus aureus infection. *Front. Immunol.* **11**, 620339 (2020).
- 1023 43. Stoll, H. *et al.* Staphylococcal enterotoxins Dose-Dependently modulate the generation of Myeloid-Derived suppressor cells. *Front. Cell. Infect. Microbiol.* **8**, 321 (2018).
- 1024 44. Tebartz, C. *et al.* A major role for myeloid-derived suppressor cells and a minor role for regulatory T cells in immunosuppression during staphylococcus aureus infection. *J. Immunol.* **194**, 1100–1111 (2015).
- 1025 45. Zhou, Z., Xu, M.-J. & Gao, B. Hepatocytes: A key cell type for innate immunity. *Cell. Mol. Immunol.* **13**, 301–315 (2016).
- 1026 46. Dixon, L. J., Barnes, M., Tang, H., Pritchard, M. T. & Nagy, L. E. Kupffer cells in the liver. *Compr. Physiol.* **3**, 785–797 (2013).
- 1027 47. Ladhami, S. N. *et al.* Invasive meningococcal disease in patients with complement deficiencies: A case series (2008-2017). *BMC Infect. Dis.* **19**, 522 (2019).
- 1028 48. Rosain, J. *et al.* Strains responsible for invasive meningococcal disease in patients with terminal complement pathway deficiencies. *J. Infect. Dis.* **215**, 1331–1338 (2017).
- 1029 49. The International Meningococcal Genetics Consortium. Genome-wide association study identifies variants in the CFH region associated with host susceptibility to meningococcal disease. *Nature Genetics* **42**, 772–776 (2010).
- 1030 50. Lung, T. *et al.* The complement system in liver diseases: Evidence-based approach and therapeutic options. *J Transl Autoimmun* **2**, 100017 (2019).
- 1031 51. Reis, E. S. *et al.* Applying complement therapeutics to rare diseases. *Clin. Immunol.* **161**, 225–240 (2015).
- 1032 52. Seal, R. L. *et al.* Genenames.org: The HGNC resources in 2023. *Nucleic Acids Res.* **51**, D1003–D1009 (2023).
- 1033 53. Al-Hamoudi, W. K. Severe autoimmune hepatitis triggered by varicella zoster infection. *World J. Gastroenterol.* **15**, 1004–1006 (2009).
- 1034 54. Brewer, E. C. & Hunter, L. Acute liver failure due to disseminated varicella zoster infection. *Case Reports Hepatol* **2018**, 1269340 (2018).
- 1035 55. Eshchar, J., Reif, L., Waron, M. & Alkan, W. J. Hepatic lesion in chickenpox. A case report. *Gastroenterology* **64**, 462–466 (1973).
- 1036 56. Li, Z. *et al.* [Aging and age-related diseases: From mechanisms to therapeutic strategies](#). *Biogerontology* **22**, 165–187 (2021).
- 1037 57. Nelson, M. R. *et al.* [The support of human genetic evidence for approved drug indications](#). *Nature Genetics* **47**, 856–860 (2015).

- 1038 58. Ochoa, D. *et al.* Human genetics evidence supports two-thirds of the 2021 FDA-approved drugs. *Nature Reviews Drug Discovery* **21**, 551–551 (2022).
- 1039 59. Minikel, E. V., Painter, J. L., Dong, C. C. & Nelson, M. R. Refining the impact of genetic evidence on clinical success. *Nature* 1–6 (2024) doi:[10.1038/s41586-024-07316-0](https://doi.org/10.1038/s41586-024-07316-0).
- 1040 60. Liu, X. *et al.* The therapeutic target database: An internet resource for the primary targets of approved, clinical trial and experimental drugs. *Expert Opin. Ther. Targets* **15**, 903–912 (2011).
- 1041 61. Chiu, W. *et al.* An update on gene therapy for inherited retinal dystrophy: Experience in leber congenital amaurosis clinical trials. *International Journal of Molecular Sciences* **22**, 4534 (2021).
- 1042 62. Fenderson, B. A. Chapter 6 - developmental and genetic diseases. in *Pathology secrets (third edition)* (ed. Damjanov, I.) 98–119 (Mosby, 2009). doi:[10.1016/B978-0-323-05594-9.00006-4](https://doi.org/10.1016/B978-0-323-05594-9.00006-4).
- 1043 63. Vilcaes, A. A., Garbarino-Pico, E., Torres Demichelis, V. & Daniotti, J. L. Ganglioside synthesis by plasma membrane-associated sialyltransferase in macrophages. *International Journal of Molecular Sciences* **21**, 1063 (2020).
- 1044 64. Yohe, H. C., Coleman, D. L. & Ryan, J. L. Ganglioside alterations in stimulated murine macrophages. *Biochimica et Biophysica Acta (BBA) - Biomembranes* **818**, 81–86 (1985).
- 1045 65. Demir, S. A., Timur, Z. K., Ateş, N., Martínez, L. A. & Seyrantepe, V. GM2 ganglioside accumulation causes neuroinflammation and behavioral alterations in a mouse model of early onset tay-sachs disease. *Journal of Neuroinflammation* **17**, 277 (2020).
- 1046 66. Ferro, A., Sheeler, C., Rosa, J.-G. & Cvetanovic, M. Role of microglia in ataxias. *Journal of molecular biology* **431**, 1792–1804 (2019).
- 1047 67. Hol, E. M. & Pasterkamp, R. J. Microglial transcriptomics meets genetics: New disease leads. *Nature Reviews Neurology* 1–2 (2022) doi:[10.1038/s41582-022-00633-w](https://doi.org/10.1038/s41582-022-00633-w).
- 1048 68. Lopes, K. de P. *et al.* Atlas of genetic effects in human microglia transcriptome across brain regions, aging and disease pathologies. *bioRxiv* 2020.10.27.356113 (2020) doi:[10.1101/2020.10.27.356113](https://doi.org/10.1101/2020.10.27.356113).
- 1049 69. Bueren, J. A. & Auricchio, A. Advances and challenges in the development of gene therapy medicinal products for rare diseases. *Hum. Gene Ther.* **34**, 763–775 (2023).
- 1050 70. Bulaklak, K. & Gersbach, C. A. The once and future gene therapy. *Nat. Commun.* **11**, 5820 (2020).
- 1051 71. Godbout, K. & Tremblay, J. P. Prime editing for human gene therapy: Where are we now? *Cells* **12**, (2023).
- 1052 72. Kohn, D. B., Chen, Y. Y. & Spencer, M. J. Successes and challenges in clinical gene therapy. *Gene Ther.* **30**, 738–746 (2023).
- 1053 73. Zhao, Z., Shang, P., Mohanraju, P. & Geijsen, N. Prime editing: Advances and therapeutic applications. *Trends Biotechnol.* **41**, 1000–1012 (2023).
- 1054 74. Darrow, J. J. Luxturna: FDA documents reveal the value of a costly gene therapy. *Drug Discov. Today* **24**, 949–954 (2019).

- 1055 75. Mendell, J. R. *et al.* Single-Dose Gene-Replacement therapy for spinal muscular atrophy. *N. Engl. J. Med.* **377**, 1713–1722 (2017).
- 1056 76. Mueller, C. *et al.* 5 year expression and neutrophil defect repair after gene therapy in alpha-1 antitrypsin deficiency. *Mol. Ther.* **25**, 1387–1394 (2017).
- 1057 77. Russell, S. *et al.* Efficacy and safety of voretigene neparvovec (AAV2-hRPE65v2) in patients with RPE65-mediated inherited retinal dystrophy: A randomised, controlled, open-label, phase 3 trial. *Lancet* **390**, 849–860 (2017).
- 1058 78. Lu, C.-F. FDA takes first step toward international regulation of gene therapies to treat rare diseases. (2024).
- 1059 79. Brown, D. G., Wobst, H. J., Kapoor, A., Kenna, L. A. & Southall, N. Clinical development times for innovative drugs. *Nat. Rev. Drug Discov.* **21**, 793–794 (2022).
- 1060 80. Moffat, J. G., Vincent, F., Lee, J. A., Eder, J. & Prunotto, M. Opportunities and challenges in phenotypic drug discovery: An industry perspective. *Nat. Rev. Drug Discov.* **16**, 531–543 (2017).
- 1061 81. Zhou, Q. & Buchholz, C. J. Cell type specific gene delivery by lentiviral vectors: New options in immunotherapy. *Oncoimmunology* **2**, e22566 (2013).
- 1062 82. Gao, C., Jiang, J., Tan, Y. & Chen, S. [Microglia in neurodegenerative diseases: mechanism and potential therapeutic targets](#). *Signal Transduction and Targeted Therapy* **8**, 1–37 (2023).
- 1063 83. Mcquade, A. & Blurton-jones, M. Microglia in alzheimer’s disease : Exploring how genetics and phenotype influence risk. *Journal of Molecular Biology* 1–13 (2019) doi:[10.1016/j.jmb.2019.01.045](https://doi.org/10.1016/j.jmb.2019.01.045).
- 1064 84. Keren-shaul, H. *et al.* [A unique microglia type associated with restricting development of alzheimer ’s disease](#). *Cell* **169**, 1276–1290.e17 (2017).
- 1065 85. Deczkowska, A. *et al.* [Disease-associated microglia: A universal immune sensor of neurodegeneration](#). *Cell* **173**, 1073–1081 (2018).
- 1066 86. DiStefano, M. T. *et al.* The gene curation coalition: A global effort to harmonize gene-disease evidence resources. *Genet. Med.* **24**, 1732–1742 (2022).
- 1067 87. DiStefano, M. *et al.* P451: The gene curation coalition works to resolve discrepancies in gene-disease validity assertions. *Genetics in Medicine Open* **1**, 100498 (2023).
- 1068 88. Skene, N. G. & Grant, S. G. N. Identification of vulnerable cell types in major brain disorders using single cell transcriptomes and expression weighted cell type enrichment. *Front. Neurosci.* **10**, 16 (2016).
- 1069 89. Osmond, M. *et al.* Outcome of over 1500 matches through the matchmaker exchange for rare disease gene discovery: The 2-year experience of Care4Rare canada. *Genet. Med.* **24**, 100–108 (2022).
- 1070 90. Philippakis, A. A. *et al.* The matchmaker exchange: A platform for rare disease gene discovery. *Hum. Mutat.* **36**, 915–921 (2015).

- 1071 91. Zanello, G. *et al.* Targeting shared molecular etiologies to accelerate drug development for rare  
diseases. *EMBO Mol. Med.* **15**, e17159 (2023).
- 1072 92. Côté, R. *et al.* The ontology lookup service: Bigger and better. *Nucleic Acids Res.* **38**, W155–60  
(2010).
- 1073 93. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach  
to multiple testing. *J. R. Stat. Soc.* (1995).
- 1074 94. Lazarin, G. A. *et al.* Systematic classification of disease severity for evaluation of expanded carrier  
screening panels. *PLoS One* **9**, e114391 (2014).
- 1075 95. Solovyeva, V. V. *et al.* New approaches to tay-sachs disease therapy. *Frontiers in Physiology* **9**,  
(2018).
- 1076 96. Hoffman, J. D. *et al.* Next-generation DNA sequencing of HEXA: A step in the right direction for  
carrier screening. *Molecular Genetics & Genomic Medicine* **1**, 260–268 (2013).
- 1077 97. Sugiyama, K., Tagawa, S. & Toda, M. Methods for visual understanding of hierarchical system struc-  
tures. *IEEE Trans. Syst. Man Cybern.* **11**, 109–125 (1981).

1078

1079

1080 **Supplementary Materials**

1081 **Supplementary Results**

1082 **Selected example targets**

1083 From our prioritised targets, we selected the following four sets of phenotypes or diseases as examples:  
1084 ‘GM2-ganglioside accumulation’, ‘Spinocerebellar atrophy’, ‘Neuronal loss in central nervous system’. Only  
1085 phenotypes with a GPT severity score greater than 15 were considered to avoid overplotting and to focus on  
1086 the more clinically relevant phenotypes Fig. 8a-h. These examples were then selected partly on the basis of  
1087 severity rankings, and partly for their relatively smaller, simpler networks than lent themselves to compact  
1088 visualisations.

1089 Tay-Sachs disease (TSD) is a devastating hereditary condition in which children are born appearing healthy,  
1090 which gradually degrades leading to death after 3-5 years. The underlying cause is the toxic accumulation  
1091 of gangliosides in the nervous system due to a loss of the enzyme produced by *HEXA*. While this could  
1092 in theory be corrected with gene editing technologies, there remain some outstanding challenges. One of  
1093 which is identifying which cell types should be targeted to ensure the most effective treatments. Here  
1094 we identified alternatively activated macrophages as the cell type most strongly associated with ‘GM2-  
1095 ganglioside accumulation’ Fig. 8i. The role of aberrant macrophage activity in the regulation of ganglioside  
1096 levels is supported by observation that gangliosides accumulate within macrophages in TSD<sup>62</sup>, as well as  
1097 experimental evidence in rodent models<sup>63,64,65</sup>. Our results not only corroborate these findings, but propose  
1098 macrophages as the primary causal cell type in TSD, making it the most promising cell type to target in  
1099 therapies.

1100 Another challenge in TSD is early detection and diagnosis, before irreversible damage has occurred. Our  
1101 pipeline implicated extravillous trophoblasts of the placenta in ‘GM2-ganglioside accumulation’. While not  
1102 necessarily a target for gene therapy (as the child is detached from the placenta after birth), checking these  
1103 cells *in utero* for an absence of *HEXA* may serve as a viable biomarker as these cells normally express  
1104 the gene at high levels. Early detection of TSD may lengthen the window of opportunity for therapeutic  
1105 intervention<sup>95</sup>, especially when genetic sequencing is not available or variants of unknown significance are  
1106 found within *HEXA*<sup>96</sup>.

1107 Spinocerebellar atrophy is a debilitating and lethal phenotype that occurs in diseases such as Spinocerebellar  
1108 ataxia and Boucher-Nenhauser syndrome. These diseases are characterised by progressive degeneration of  
1109 the cerebellum and spinal cord, leading to severe motor and cognitive impairments. Our pipeline identified  
1110 M2 macrophages (labeled as the closest CL term ‘Alternatively activated macrophages’ in Fig. 8j) as the  
1111 only causal cell type associated with ‘Spinocerebellar atrophy’. This strongly suggests that degeneration of  
1112 cerebellar Purkinje cells are in fact downstream consequences of macrophage dysfunction, rather than being  
1113 the primary cause themselves. This is consistent with the known role of macrophages, especially microglia, in

1114 neuroinflammation and other neurodegenerative conditions such as Alzheimer's and Parkinsons' disease<sup>66–68</sup>.  
1115 While experimental and postmortem observational studies have implicated microglia in spinocerebellar atro-  
1116 phy previously<sup>66</sup>, our results provide a statistically-supported and unbiased genetic link between known risk  
1117 genes and this cell type. Therefore, targeting M2 microglia in the treatment of spinocerebellar atrophy may  
1118 therefore represent a promising therapeutic strategy. This is aided by the fact that there are mouse models  
1119 that perturb the ortholog of human spinocerebellar atrophy risk genes (e.g. *Atxn1*, *Pnpla6*) and reliably  
1120 recapitulate the effects of this diseases at the cellular (e.g. loss of Purkinje cells), morphological (e.g. atrophy  
1121 of the cerebellum, spinal cord, and muscles), and functional (e.g. ataxia) levels.

1122 Next, we investigated the phenotype 'Neuronal loss in the central nervous system'. Despite the fact that this  
1123 is a fairly broad phenotype, we found that it was only significantly associated with 3 cell types (alternatively  
1124 activated macrophage, macrophage, epithelial cell), specifically M2 macrophages and sinusoidal endothelial  
1125 cells Fig. 8k.

1126 Skeletal dysplasia is a heterogeneous group of over 450 disorders that affect the growth and development of  
1127 bone and cartilage. This phenotype can be lethal when deficient bone growth leads to the constriction of  
1128 vital organs such as the lungs. Even after surgical interventions, these complications continue to arise as the  
1129 child develops. Pharmacological interventions to treat this condition have largely been ineffective. While  
1130 there are various cell types involved in skeletal system development, our pipeline nominated chondrocytes  
1131 as the causal cell type underlying the lethal form of this condition (Fig. 19). Assuringly, we found that  
1132 the disease 'Achondrogenesis Type 1B' is caused by the genes *SLC26A2* and *COL2A1* via chondrocytes.  
1133 We also found that 'Platyspondylic lethal skeletal dysplasia, Torrance type'. Thus, in cases where surgical  
1134 intervention is insufficient, targeting these genes within chondrocytes may prove a viable long-term solution  
1135 for children suffering from lethal skeletal dysplasia.

1136 Alzheimer's disease (AD) is the most common neurodegenerative condition. It is characterised by a set of  
1137 variably penetrant phenotypes including memory loss, cognitive decline, and cerebral proteinopathy. Inter-  
1138 estingly, we found that different forms of early onset AD (which are defined by the presence of a specific  
1139 disease gene) are each associated with different cell types via different phenotypes (Fig. 19). For example,  
1140 AD 3 and AD 4 are primarily associated with cells of the digestive system ('enterocyte', 'gastric goblet  
1141 cell') and are implied to be responsible for the phenotypes 'Senile plaques', 'Alzheimer disease', 'Parietal  
1142 hypometabolism in FDG PET'. Meanwhile, AD 2 is primarily associated with immune cells ('alternatively  
1143 activated macrophage') and is implied to be responsible for the phenotypes 'Neurofibrillary tangles', 'Long-  
1144 tract signs'. This suggests that different forms of AD may be driven by different cell types and phenotypes,  
1145 which may help to explain its variability in onset and clinical presentation.

1146 Finally, Parkinson's disease (PD) is characterised by motor symptoms such as tremor, rigidity, and bradyki-  
1147 nesia. However there are a number of additional phenotypes associated with the disease that span multiple  
1148 physiological systems. PD 19a and PD 8 seemed to align most closely with the canonical understanding of

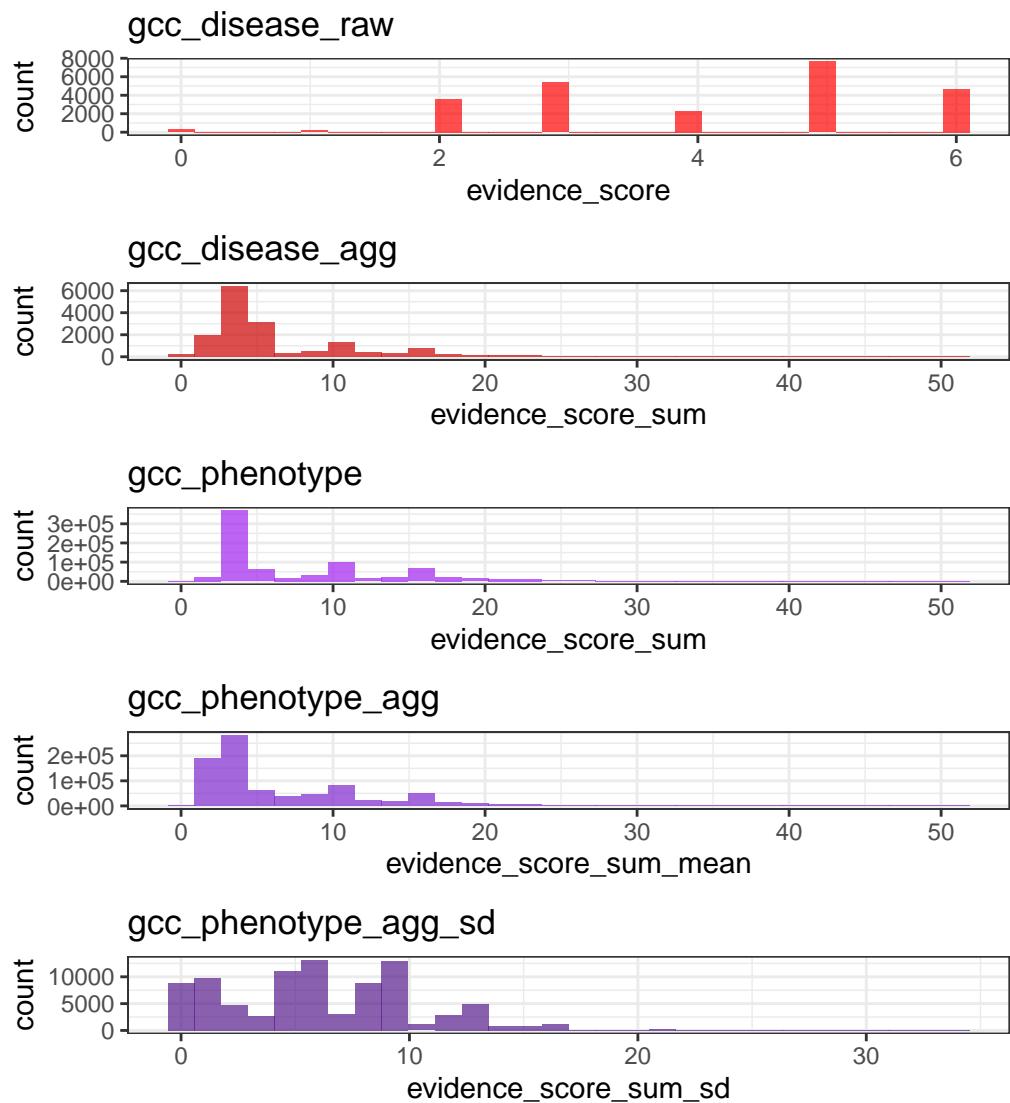
1149 PD as a disease of the central nervous system in that they implicated oligodendrocytes and neurons (Fig. 19).  
1150 Though the reference datasets being used in this study were not annotated at sufficient resolution to distin-  
1151 guish between different subtypes of neurons, in particular dopaminergic neurons. PD 19a/8 also suggested  
1152 that risk variants in *LRRK2* mediate their effects on PD through both myeloid cells and oligodendrocytes  
1153 by causing gliosis of the substantia nigra. The remaining clusters of PD mechanisms revolved around chon-  
1154 drocytes (PD 20), amacrine cells of the eye (hereditary late-onset PD), and the respiratory/immune system  
1155 (PD 14). While the diversity in cell type-specific mechanisms is somewhat surprising, it may help to explain  
1156 the wide variety of cross-system phenotypes frequently observed in PD.

1157 It should be noted that the HPO only includes gene annotations for the monogenic forms of AD and PD.  
1158 However it has previously been shown that there is at least partial overlap in their phenotypic and genetic  
1159 aetiology with respect to their common forms. Thus understanding the monogenic forms of these diseases  
1160 may shed light onto their more common counterparts.

#### 1161 **Phenome-wide analyses discover novel phenotype-cell type associations**

1162 We visualised the putative causal relationships between genes, cell types and diseases associated with RNI as  
1163 a network (Fig. 14). The phenotype ‘Recurrent Neisserial infections’ was connected to cell types through the  
1164 aforementioned association test results ( $FDR < 0.05$ ). Genes that were primarily driving these associations  
1165 (i.e. genes that were both strongly linked with ‘Recurrent Neisserial infections’ and were highly specifically  
1166 expressed in the given cell type) were designated as “driver genes” and retained for plotting. Across all  
1167 phenotypes in the HPO, more specific phenotypes (terms in the HPO with greater IC) are not only more  
1168 specific to certain cell types (Fig. 3b), but are also associated with genes that have greater cell type-specific  
1169 expression within those cell types. Even so, we should note that the choice of which specificity quantiles to  
1170 include is arbitrary. It should also be noted that simply because a gene is not specific to a cell type does not  
1171 mean it is not important for the function of the cell type. Indeed, there are many genes that are ubiquitously  
1172 expressed throughout many tissues in the body and are essential for cell function. Gene expression specificity  
1173 is nevertheless a useful metric to help distinguish many hundreds of cell (sub)types with overlapping gene  
1174 signatures.

1175 Supplementary Figures



(a) **Distribution of GenCC evidence scores at each processing step.** GenCCC (<https://thegencc.org/>) is a database where semi-quantitative scores for the current strength of evidence attributing disruption of a gene as a causal factor in a given disease. “gcc\_disease\_raw” is the distribution of raw GenCC scores before any aggregation. “gcc\_disease\_agg” is the distribution of GenCC scores after aggregating by disease. “gcc\_phenotype” is the distribution of scores after linking each phenotype to one or more disease. “gcc\_phenotype\_agg” is the distribution of scores after aggregating by phenotype, while “gcc\_phenotype\_agg\_sd” is the standard deviation of those aggregated scores.

Figure 9

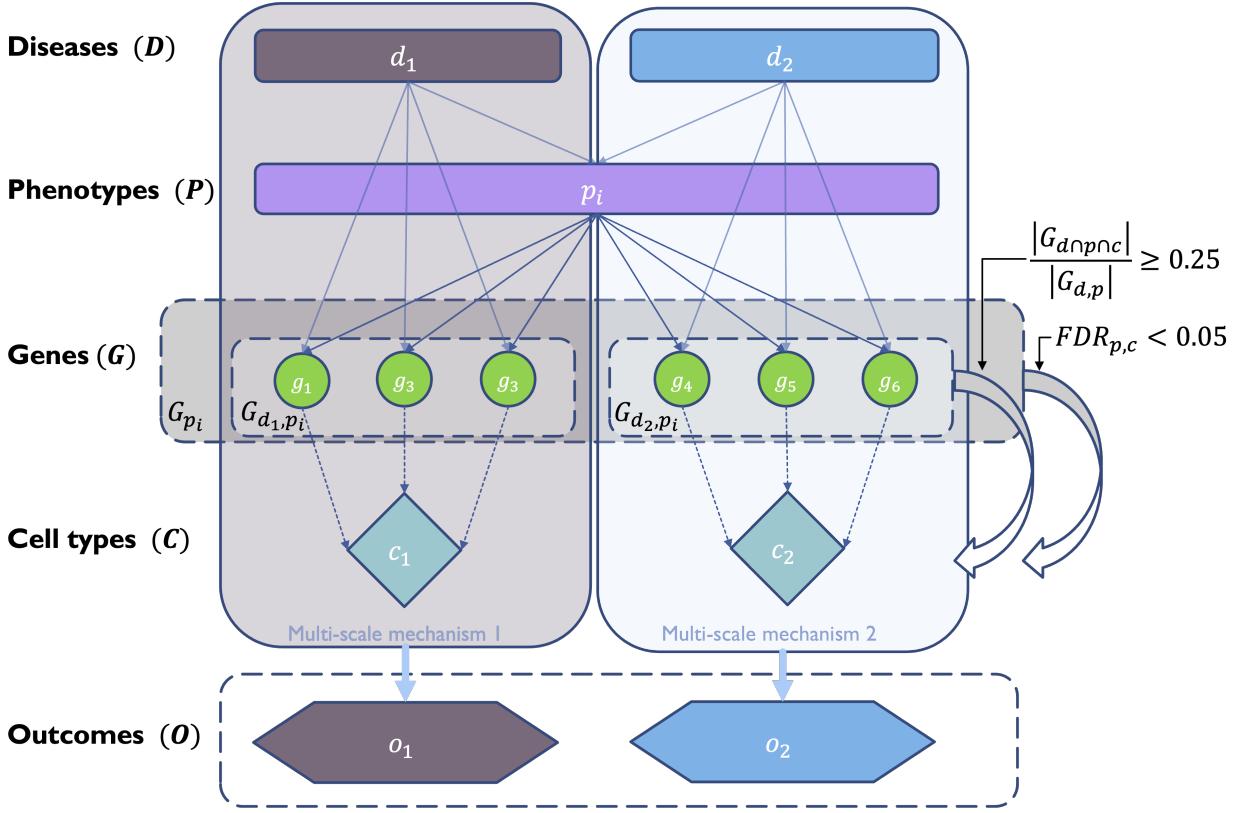
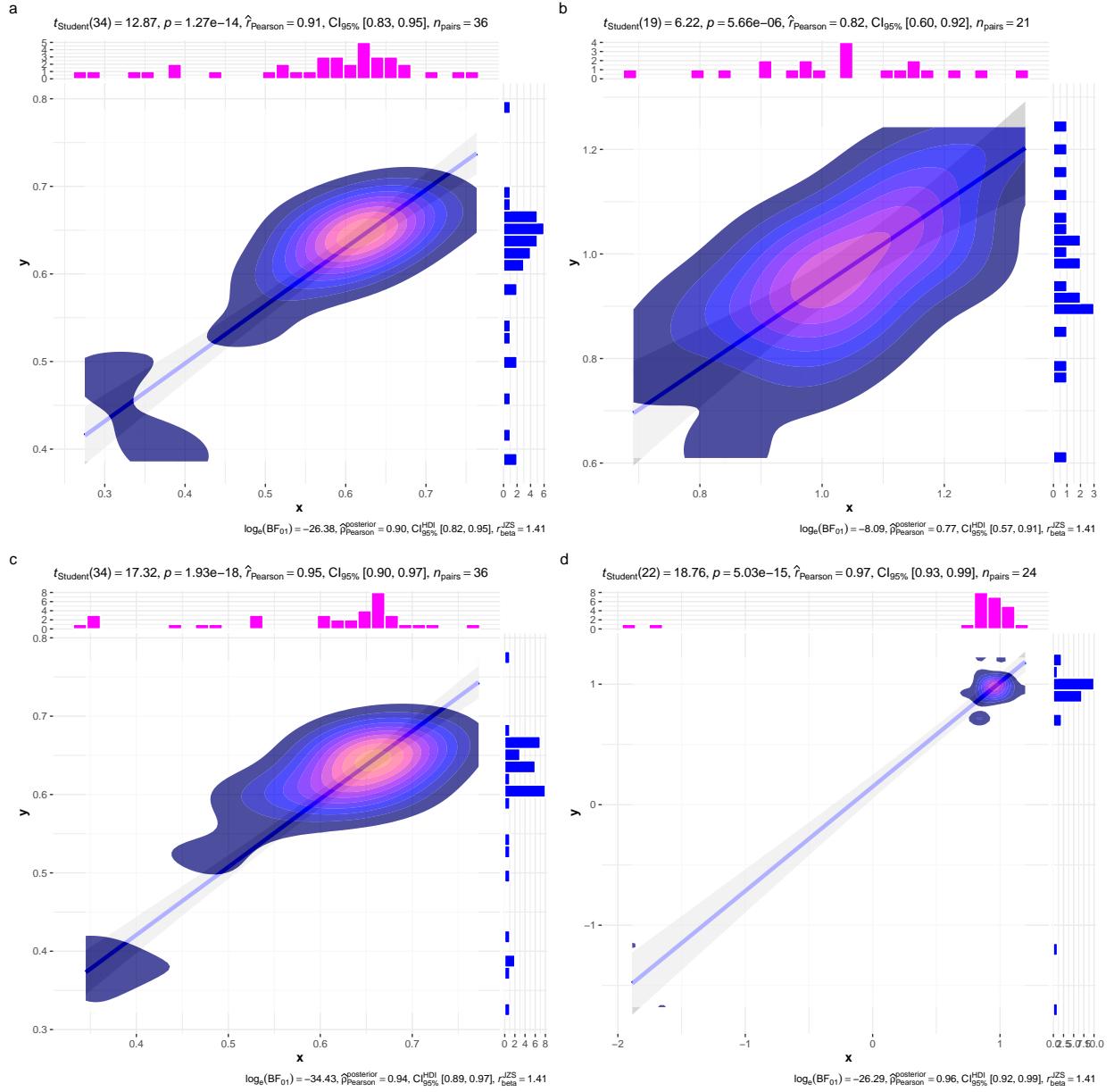
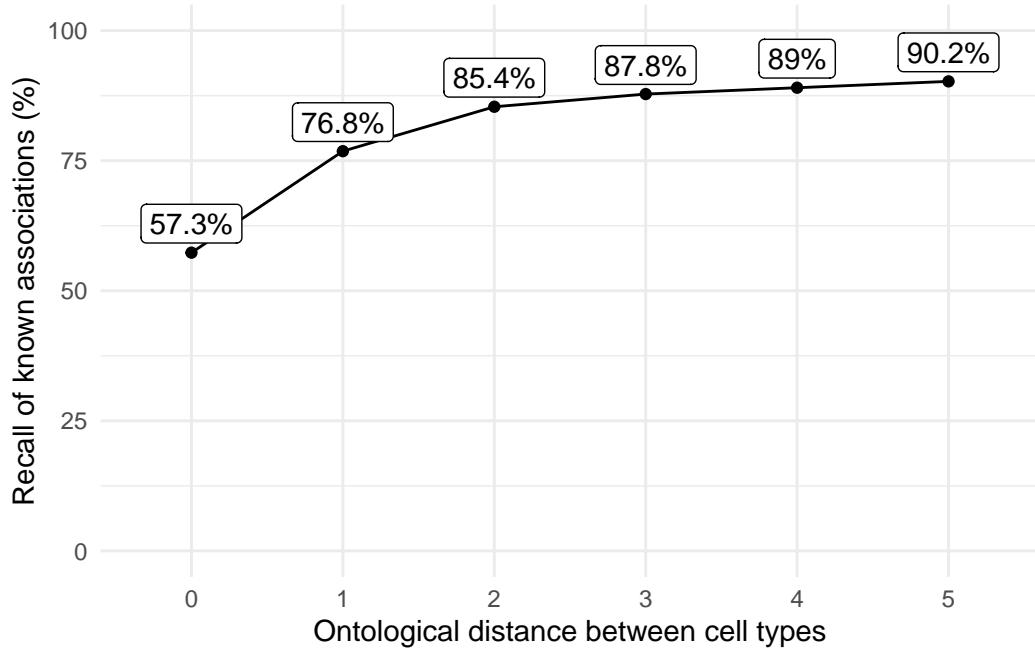


Figure 10: **Diagrammatic overview of multi-scale disease investigation strategy.** Here we provide an abstract example of differential disease aetiology across multiple scales: diseases ( $D$ ), phenotypes ( $P$ ), cell types ( $C$ ), genes ( $G$ ), and clinical outcomes ( $O$ ). In the HPO, genes are assigned to phenotypes via particular diseases ( $G_{dp}$ ). Therefore, the final gene list for each phenotype is aggregated from across multiple diseases ( $G_p$ ). We performed association tests for all pairwise combinations of cell types and phenotypes and filtered results after multiple testing corrections ( $FDR < 0.05$ ). Each phenotype in the context of a given disease is referred to here as a symptom. Links were established between symptoms and cell types through proportional gene set overlap at a minimum threshold of 25%.



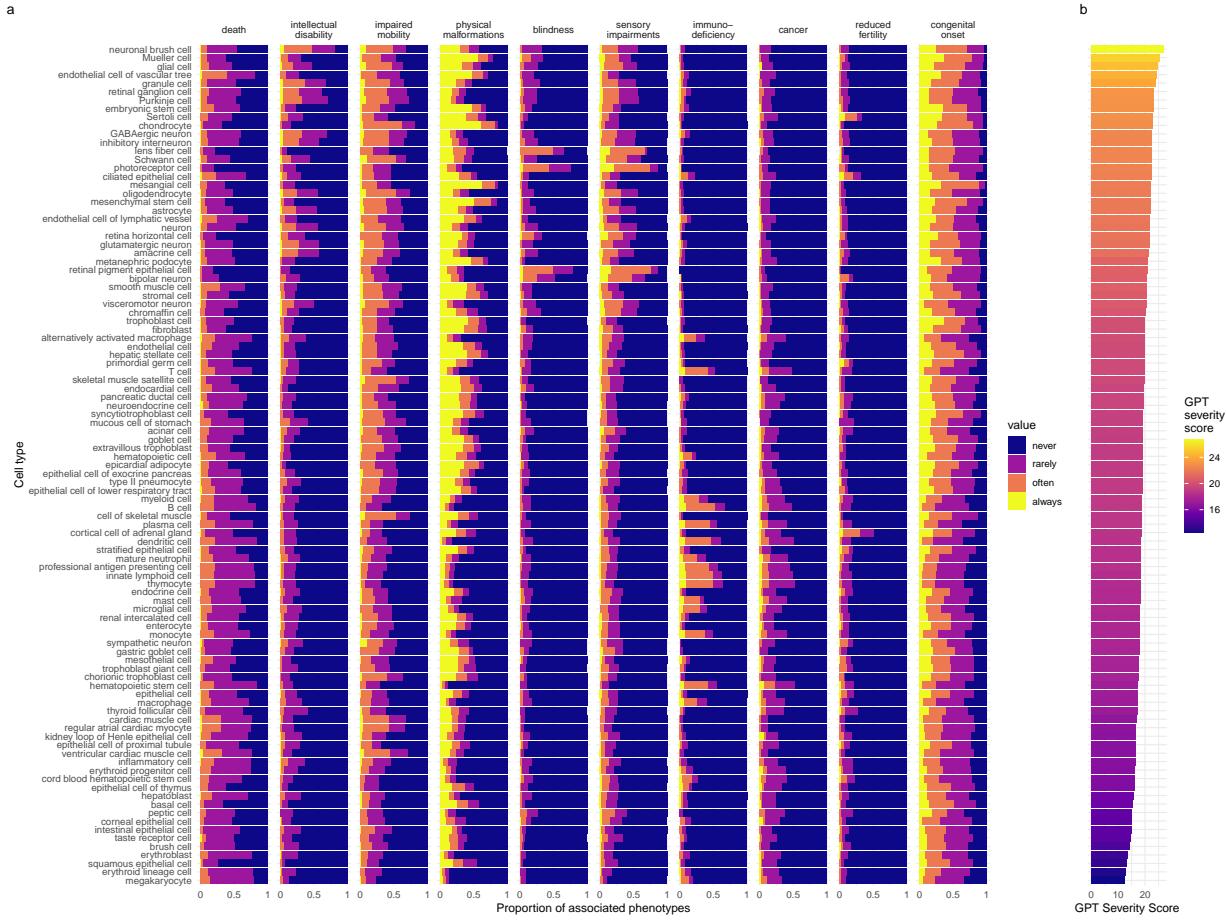
**(a) Inter- and intra-dataset validation across the different CellTypeDataset (CTD) and developmental stages.** Correlations are computed using Pearson correlation coefficient. Point density is plotted using a 2D kernel density estimate. **a** Correlation between the uncorrected p-values from all phenotype-cell type association tests using the Descartes Human vs. Human Cell Landscape CTDs. **b** Correlation between the  $\log_{10}(fold-change)$  from significant phenotype-cell type association tests (FDR<0.05) using the Descartes Human vs. Human Cell Landscape CTDs. **c** Correlation between the uncorrected p-values from all phenotype-cell type association tests using the Human Cell Landscape fetal samples vs. Human Cell Landscape adult samples. **d** Correlation between the  $\log_{10}(fold-change)$  from significant phenotype-cell type association tests (FDR<0.05) using the Human Cell Landscape fetal samples vs. Human Cell Landscape adult samples.

Figure 11



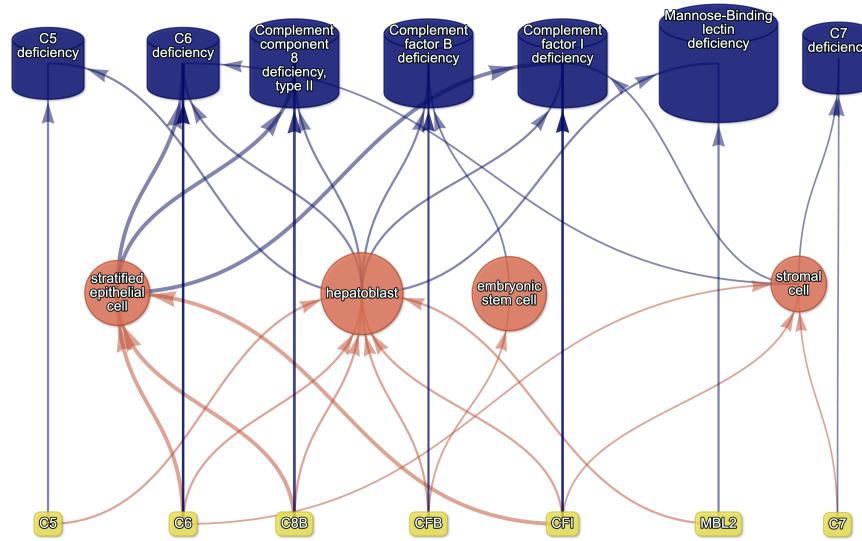
(a) Recall of ground-truth Monarch Knowledge Graph phenotype-cell type relationships at each ontological distance between cell types according to the Cell Ontology.

Figure 12

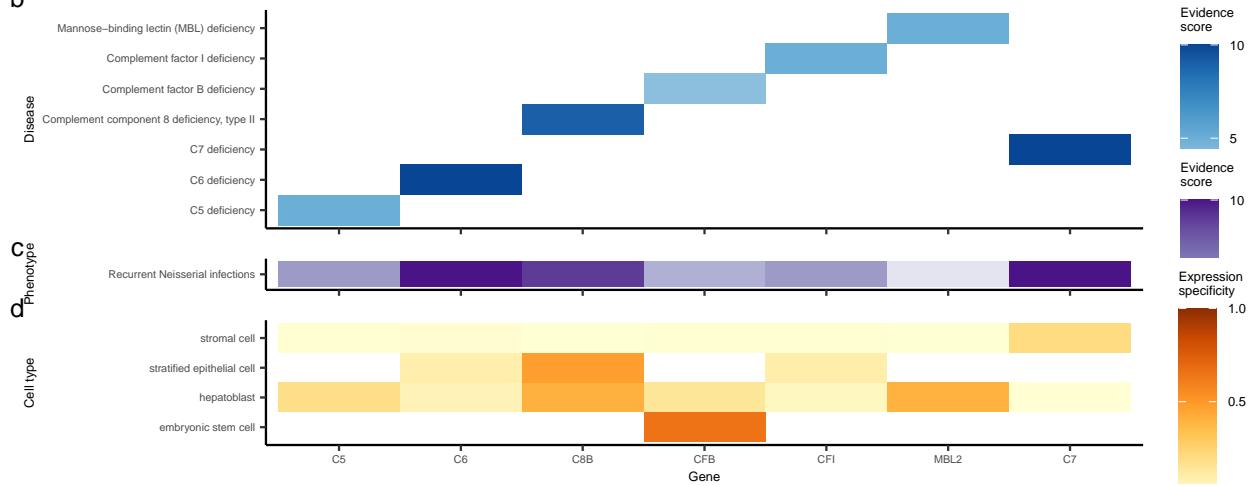


(a) **Cell types ordered by the mean severity of the phenotypes they're associated with.** **a**, The distribution of phenotype severity annotation frequencies aggregated by cell type. **b**, The composite severity score, averaged across all phenotypes associated with each cell type.

a

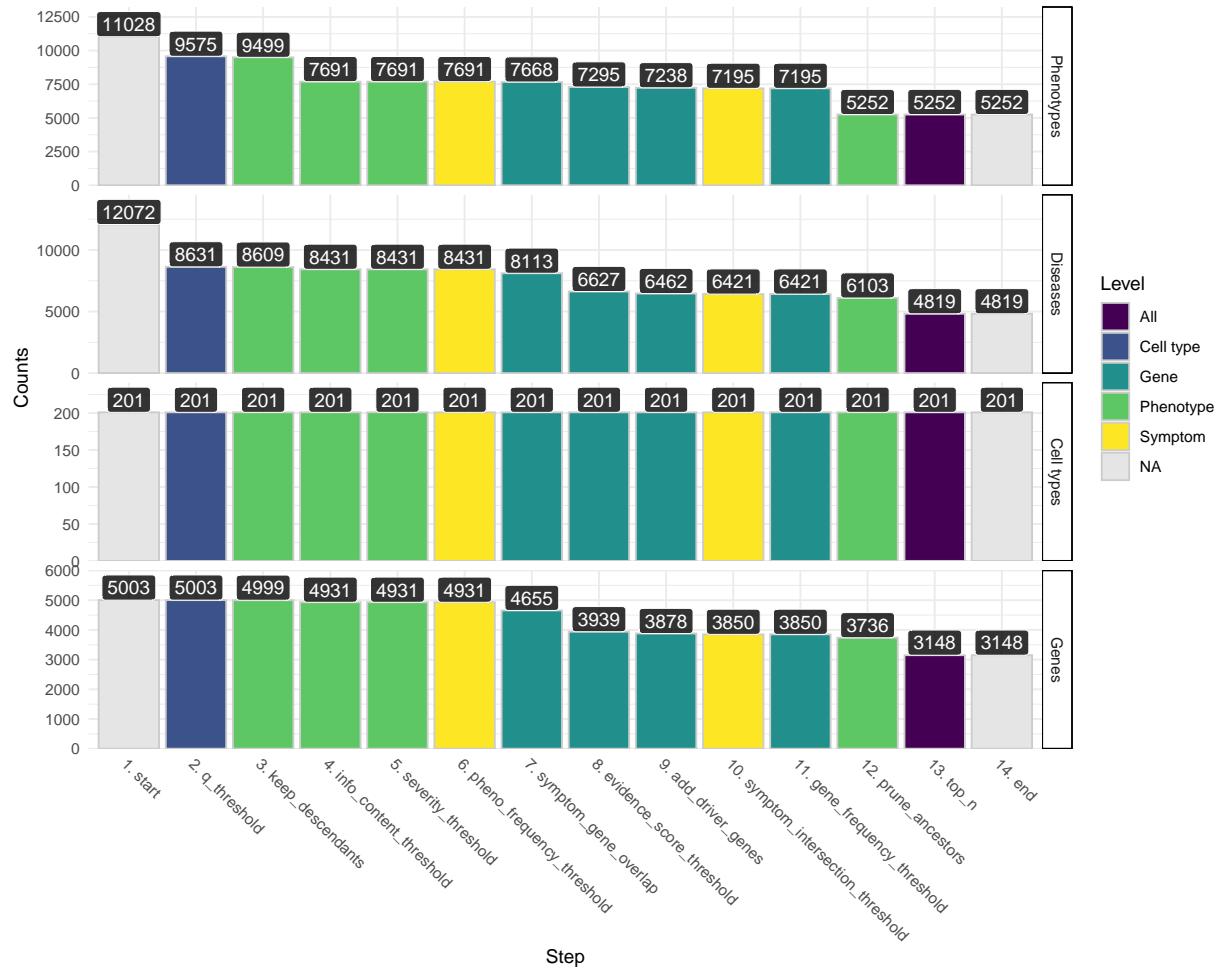


b



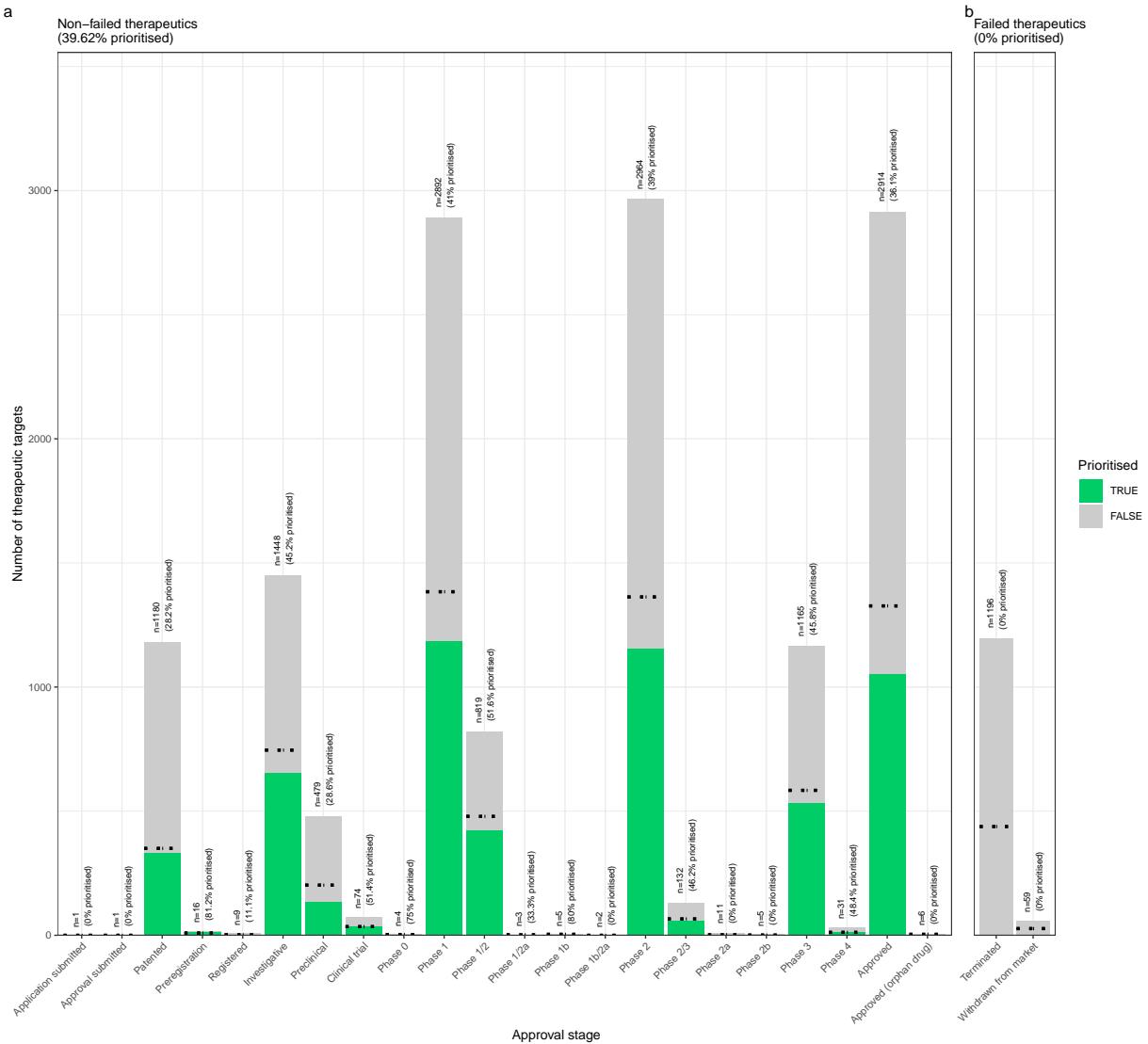
(a) **Causal network of recurrent Neisserial infections (RNI) reveals multi-scale disease aetiology.** RNI is a phenotype in seven different monogenic diseases caused by disruptions to specific complement system genes. Four cell types were significantly associated with RNI. **a**, One can trace how genes causal for RNI (yellow boxes, bottom) mediate their effects through cell types (orange circles, middle) and diseases (blue cylinders, top). Cell types are connected to RNI via association testing ( $FDR < 0.05$ ). Genes shown here have both strong evidence for a causal role in RNI and high expression specificity in the associated cell type. Cell types can be linked to monogenic diseases via the genes specifically expressed in those cell types (i.e. are in the top 25% of cell type specificity expression quantiles). Nodes are arranged using the Sugiyama algorithm<sup>97</sup>. **b** Expression specificity quantiles (1-40 scale) of each driver gene in each cell type (darker = greater specificity). **c** GenCC-derived eevidence scores between the RNI phenotype and each gene. **d** Expression specificity (0 = least specific, 1 = most specific) of each gene in each cell type.

Figure 14



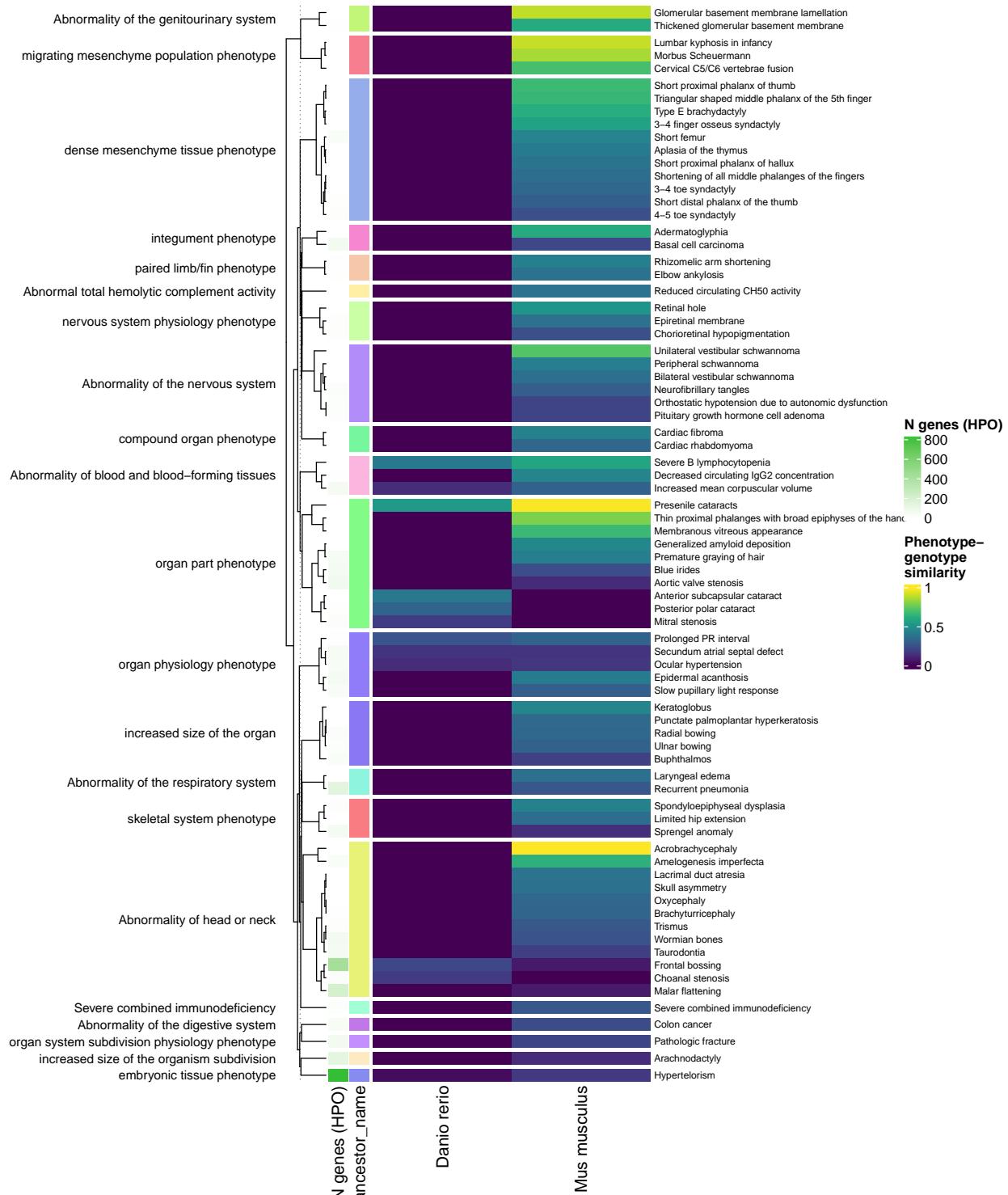
(a) **Prioritised target filtering steps.** This plot visualises the number of unique phenotype-cell type associations, cell types, genes, and phenotypes (*y-axis*) at each filtering step (*x-axis*) within the multi-scale therapeutic target prioritisation pipeline. Each step in the pipeline can be easily adjusted according to user preference and use case. See Table 3 for descriptions and criterion of each filtering step.

Figure 15



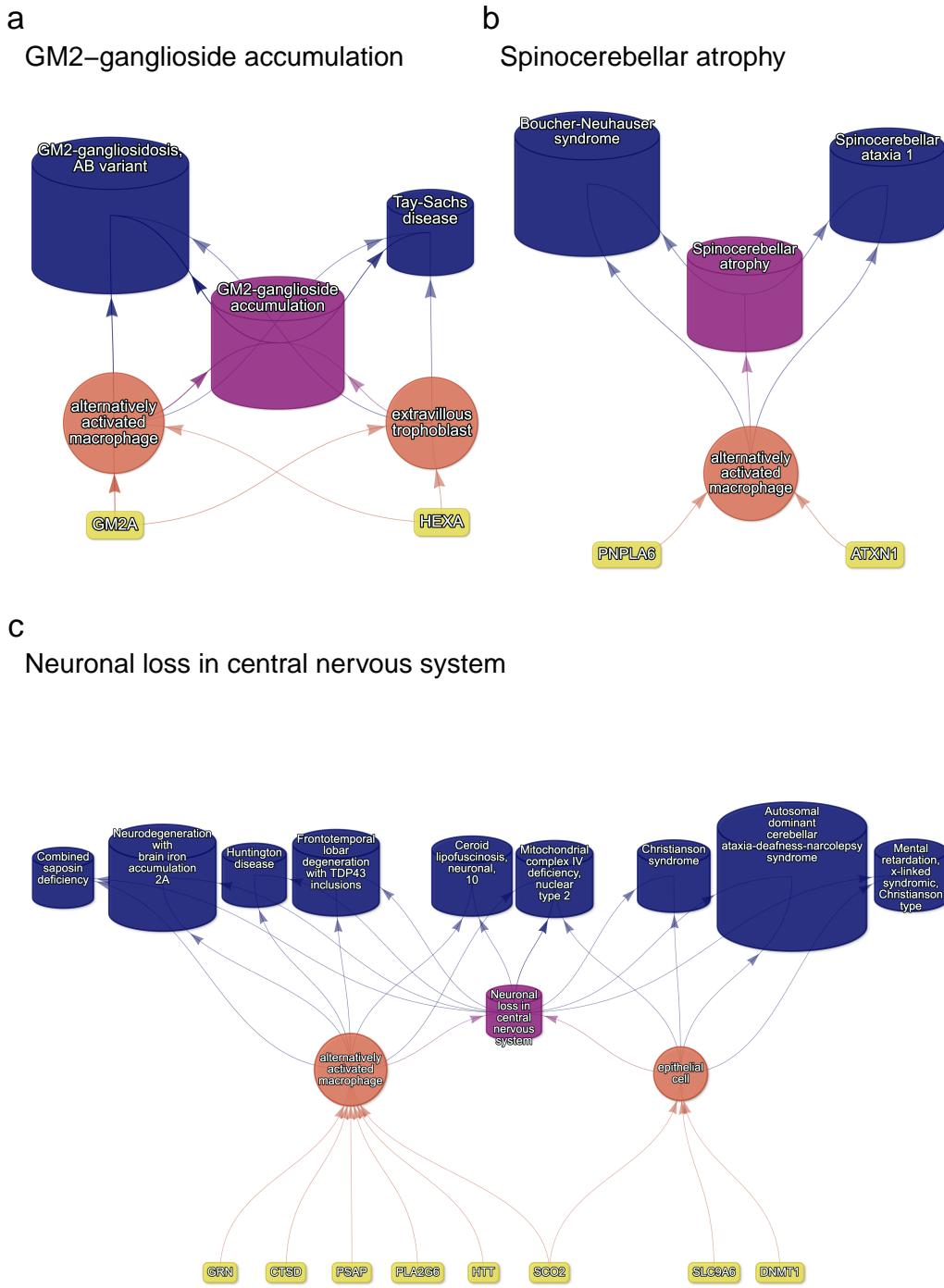
(a) **Validation of prioritised therapeutic targets.** Proportion of existing all therapy targets (documented in the Therapeutic Target Database) recapitulated by our prioritisation pipeline.

Figure 16



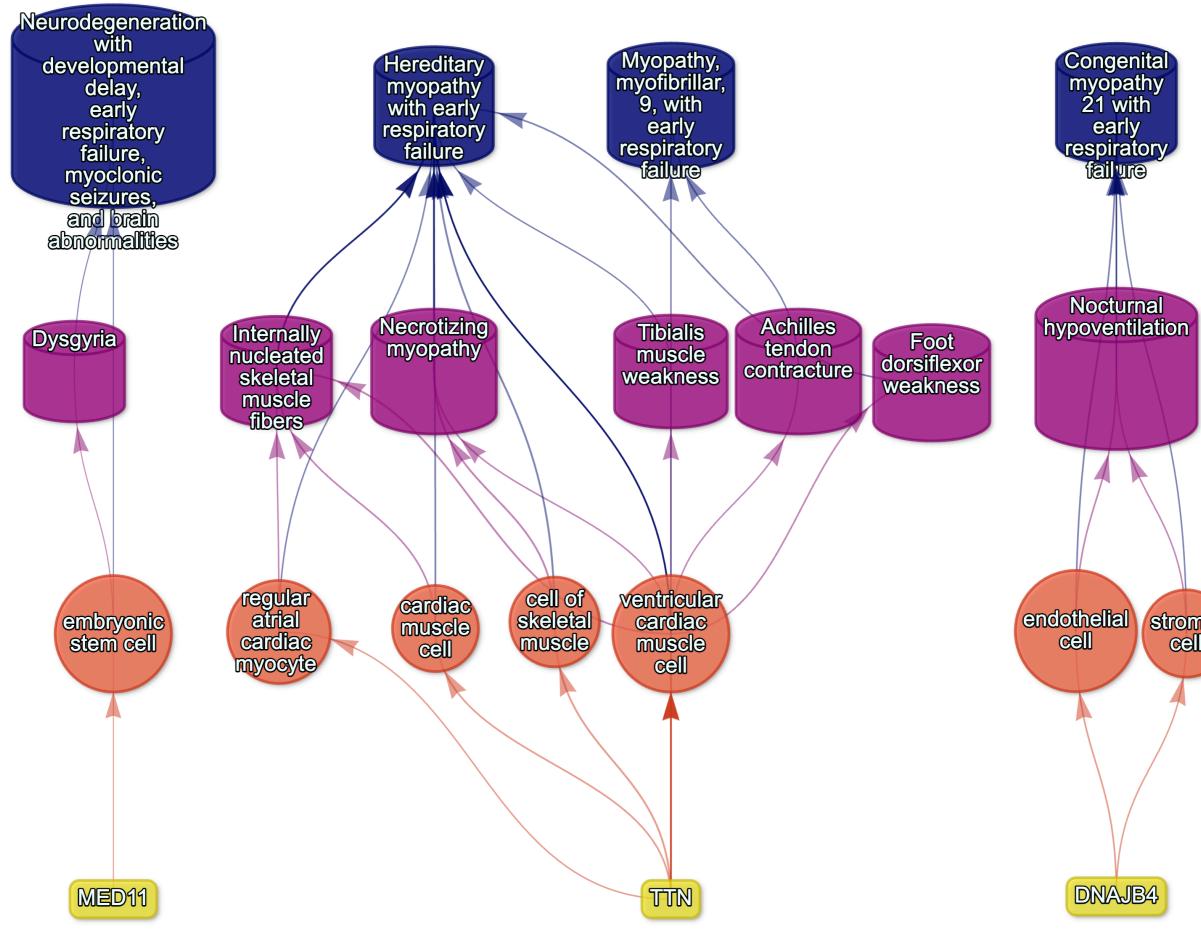
(a) **Identification of translatable experimental models.** Interspecies translatability of the top 200 human phenotypes nominated by the gene therapy prioritised pipeline. Above, the combined ontological-genotypic similarity score ( $SIM_{og}$ ) is displayed as the heatmap fill colour stratified by the model organism (*x-axis*). An additional column (“n\_genes\_db1” on the far left) displays the total number of unique genes annotated to the phenotypic within the HPO. Phenotypes are clustered according to their ontological similarity in the HPO (*y-axis*).

Figure 17



(a) **Causal multi-scale networks reveal cell type-specific therapeutic targets.** Each disease (blue cylinders) is connected to its phenotype (purple cylinders) based on well-established clinical observations recorded within the HPO<sup>11</sup>. Phenotypes are connected to cell types (orange circles) via association testing between weighted gene sets (FDR<0.05). Each cell type is connected to the prioritised gene targets (yellow boxes) based on the driver gene analysis. The thickness of the edges connecting the nodes represent the (mean) fold-change from the bootstrapped enrichment tests. Nodes were spatially arranged using the Sugiyama algorithm<sup>97</sup>.

Figure 18



(a) Respiratory failure

Figure 19: **Example cell type-specific gene therapy targets for phenotypes associated with respiratory failure-related diseases.** Each disease (blue cylinders) is connected to its phenotype (purple cylinders) based on well-established clinical observations recorded within the HPO<sup>11</sup>. Phenotypes are connected to cell types (red circles) via association testing between weighted gene sets ( $FDR < 0.05$ ). Each cell type is connected to the prioritised gene targets (yellow boxes) based on the driver gene analysis. The thickness of the edges connecting the nodes represent the (mean) fold-change from the bootstrapped enrichment tests. Nodes were spatially arranged using the Sugiyama algorithm<sup>97</sup>.

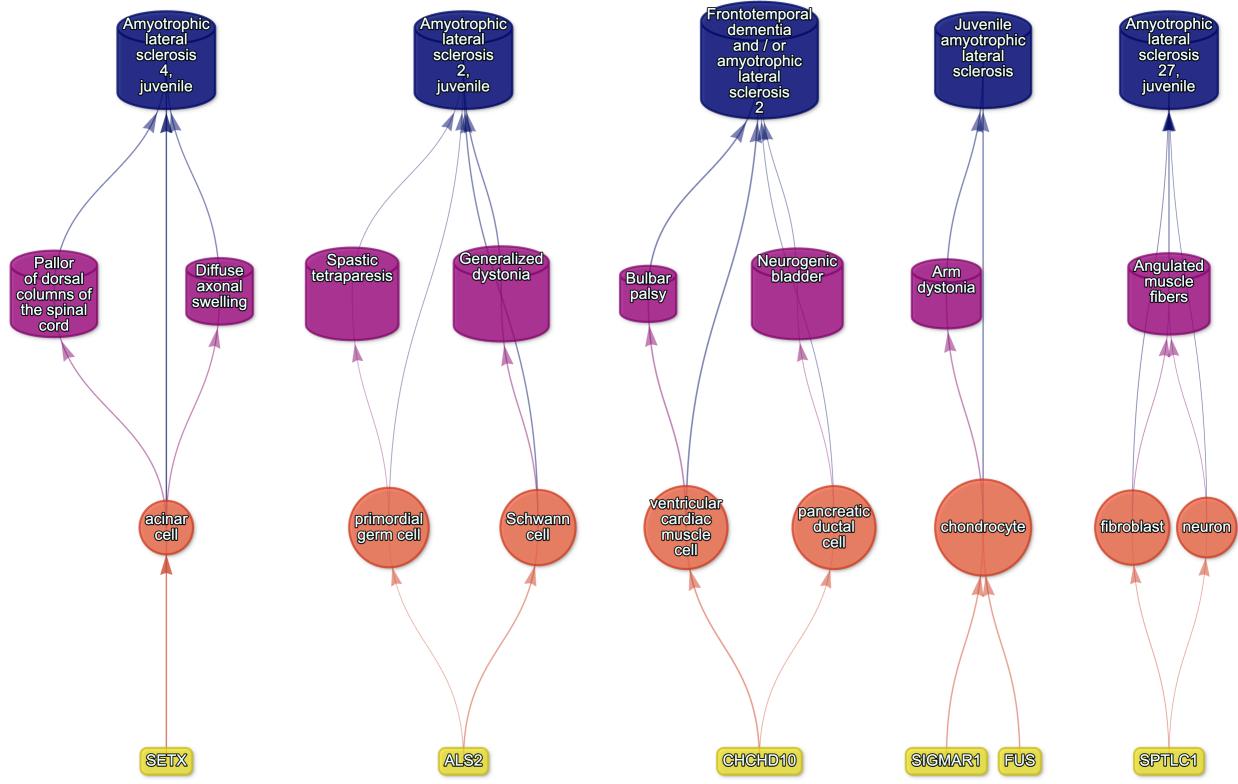


Figure 20: Causal multi-scale network for phenotypes associated with Amyotrophic Lateral Sclerosis (ALS).

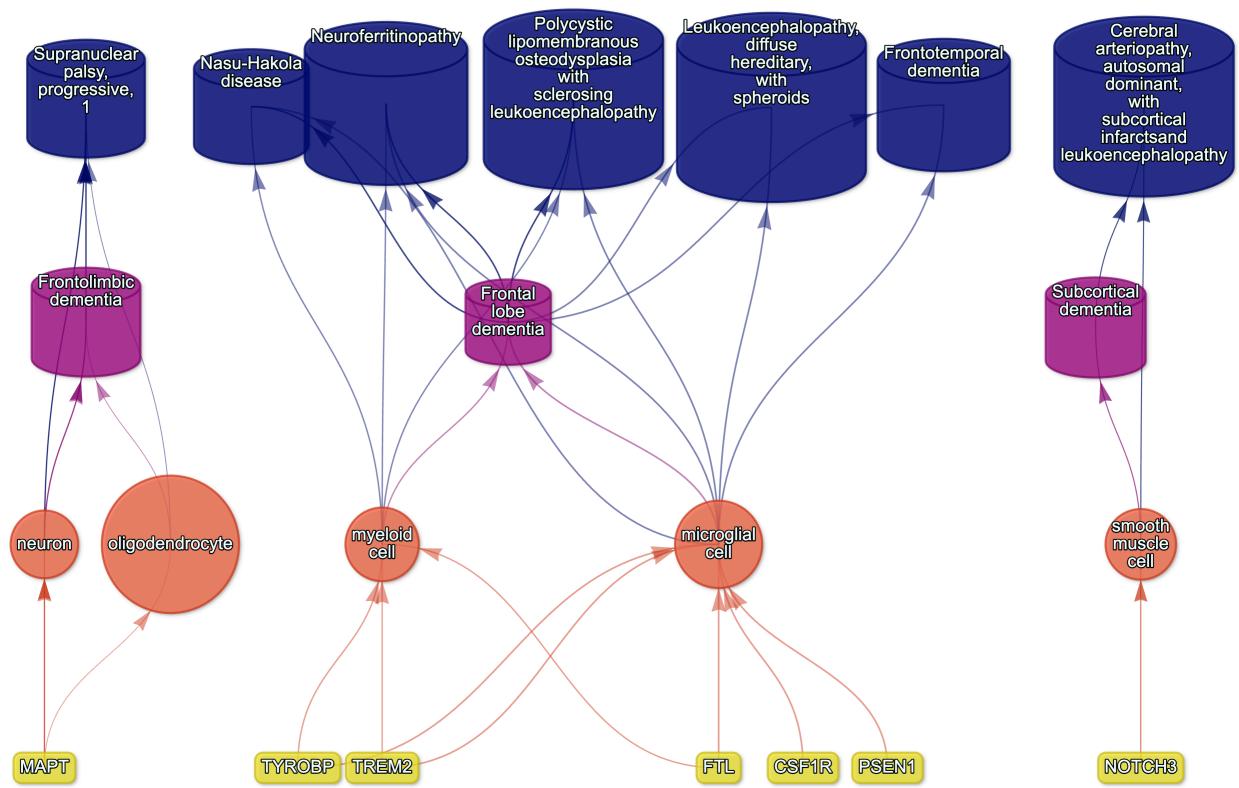


Figure 21: Causal multi-scale network for dementia phenotypes.

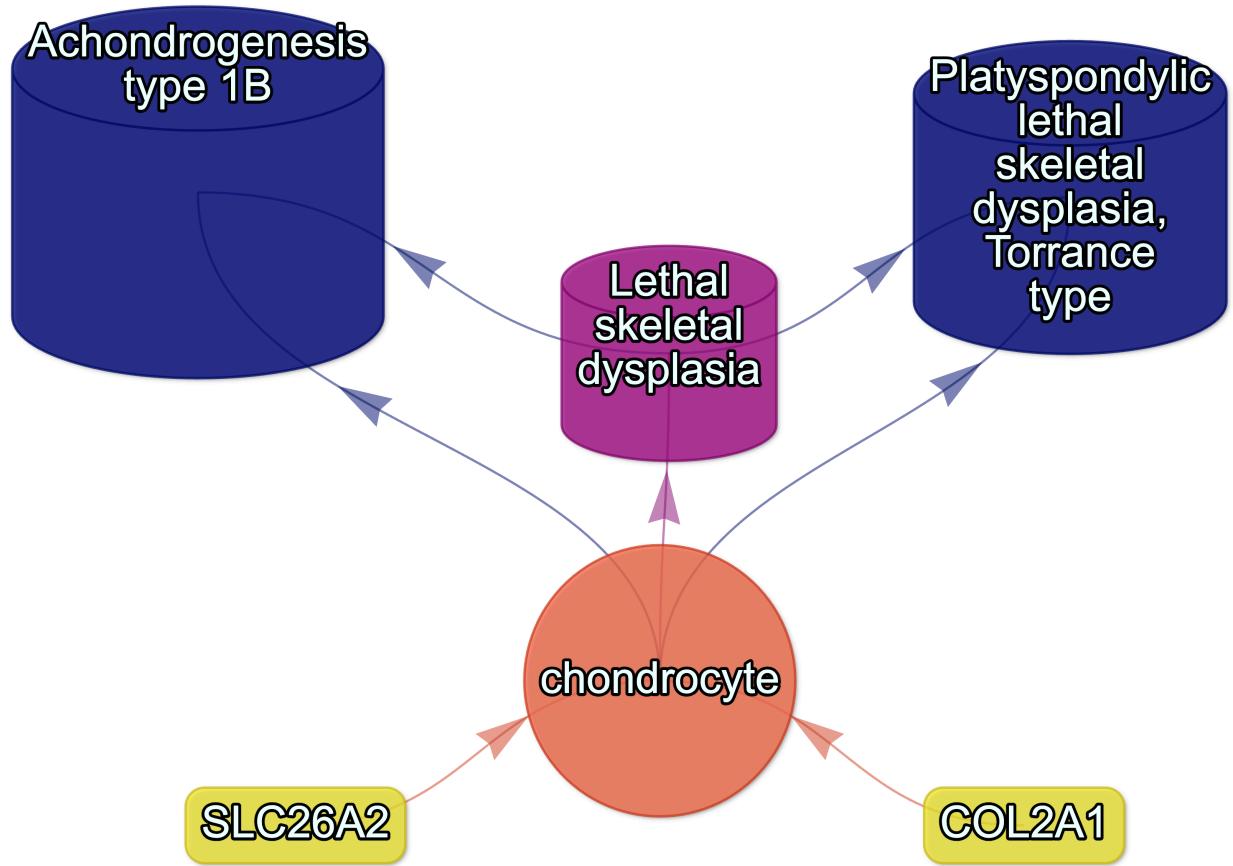


Figure 22: Causal multi-scale network for the phenotype lethal skeletal dysplasia.

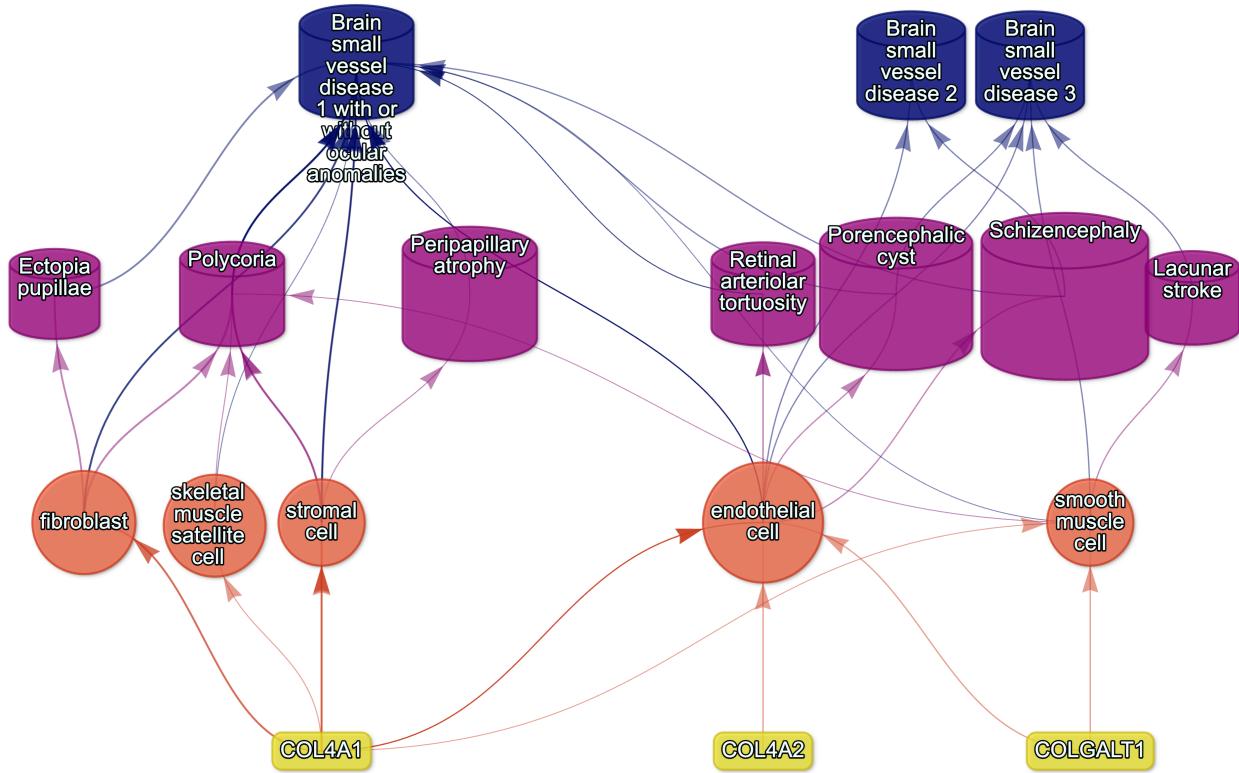


Figure 23: Causal multi-scale network for phenotypes associated with small vessel disease.

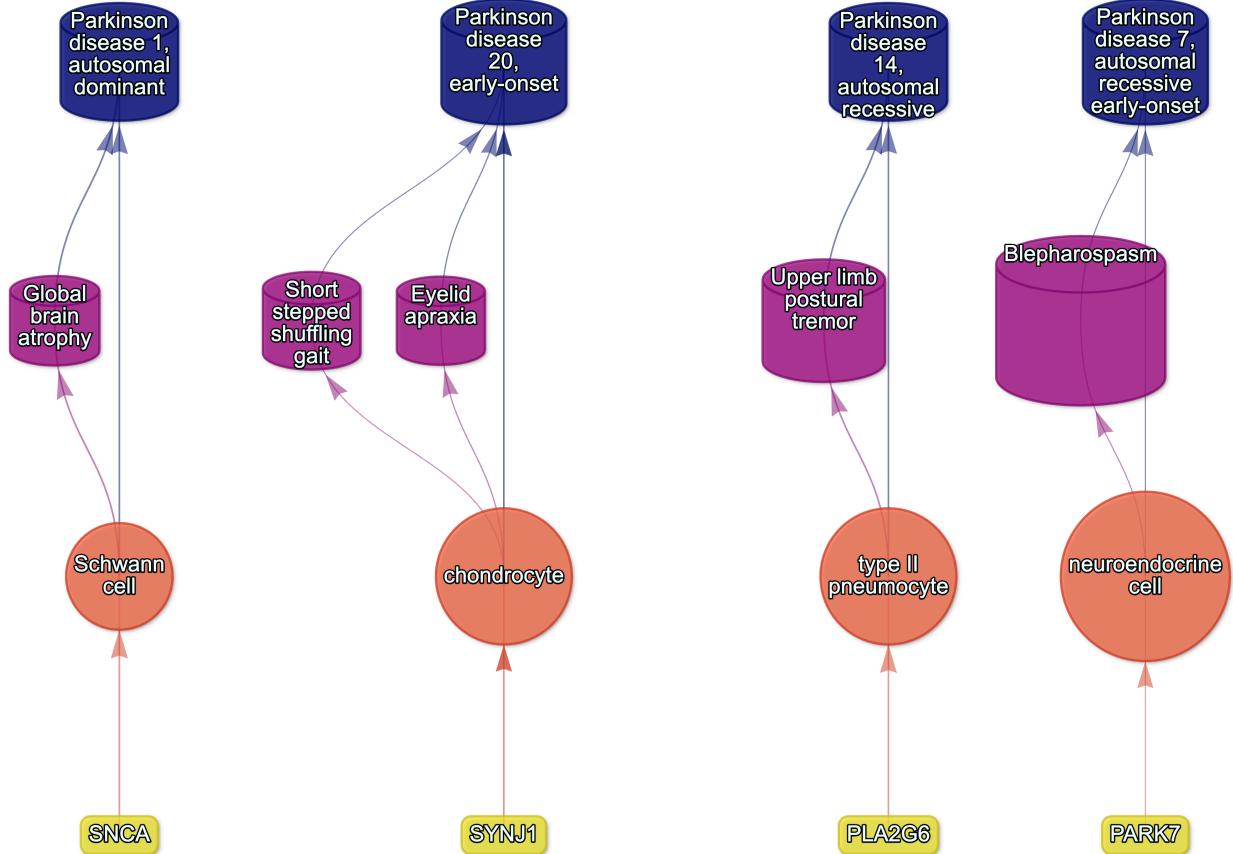


Figure 24: Causal multi-scale network for phenotypes associated with various subtypes of Parkinson's disease.

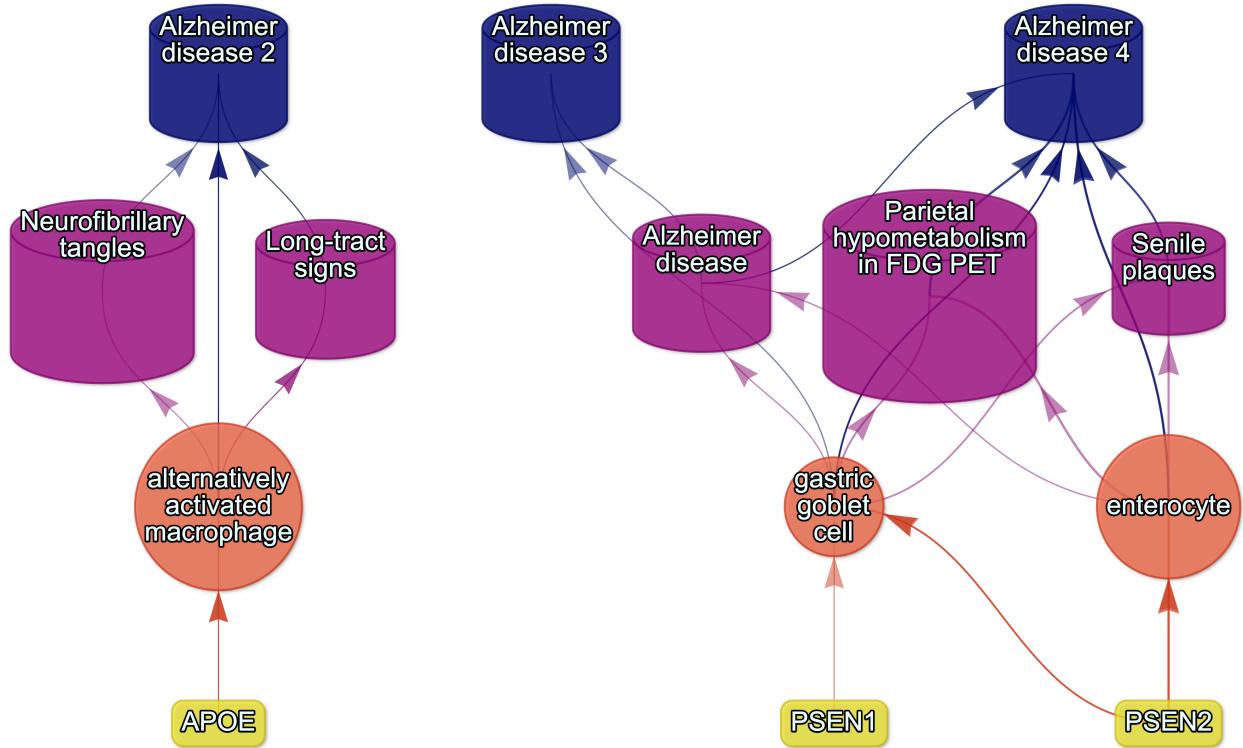


Figure 25: Causal multi-scale network for phenotypes associated with various subtypes of Alzheimer's disease.

1176 Supplementary Tables

Table 1: **Mappings between HPO phenotypes and other medical ontologies.** “source” indicates the medical ontology and “distance” indicates the cross-ontology distance. “source terms” and “HPO terms” indicates the number of unique IDs mapped from the source ontology and HPO respectively. “mappings” is the total number of cross-ontology mappings within a given distance. Some IDs may have more than one mapping for a given source due to many-to-many relationships.

source	distance	source terms	HPO terms	mappings
ICD10	2	25	23	25
ICD10	3	839	876	1170
ICD9	1	21	21	21
ICD9	2	434	306	462
ICD9	3	1052	920	1816
SNOMED	1	4413	3483	4654
SNOMED	2	75	21	78
SNOMED	3	1796	833	9605
UMLS	1	12898	11601	13049
UMLS	2	140	113	142
UMLS	3	1871	1204	11021

Table 3: **Description of each filtering step performed in the multi-scale therapeutic target prioritisation pipeline.** ‘level’ indicates the biological scale at which the step is applied to.

level	step	description
NA	1. start	NA
Cell type	2. q threshold	Keep only cell type-phenotype association results at $q \leq 0.05$ .
Phenotype	3. keep descendants	Remove phenotypes belonging to a certain branch of the HPO, as defined by an ancestor term.
Phenotype	4. info content threshold	Keep only phenotypes with a minimum information criterion score (computed from the HPO).
Phenotype	5. severity threshold	Keep only phenotypes with mean Severity equal to or below the threshold.
Symptom	6. pheno frequency threshold	Keep only phenotypes with mean frequency equal to or above the threshold (i.e. how frequently a phenotype is associated with any diseases in which it occurs).
Gene	7. symptom gene overlap	Ensure that genes nominated at the phenotype-level also appear in the genes overlapping at the cell type-specific symptom-level.
Gene	8. evidence score threshold	Remove genes that are below an aggregate phenotype-gene evidence score threshold.

**Table 3: Description of each filtering step performed in the multi-scale therapeutic target prioritisation pipeline.** ‘level’ indicates the biological scale at which the step is applied to.

level	step	description
Gene	9. add driver genes	Keep only genes that are driving the association with a given phenotype (inferred by the intersection of phenotype-associated genes and gene with high-specificity quantiles in the target cell type).
Symptom	10. symptom intersection threshold	Minimum proportion of genes overlapping between a symptom gene list (phenotype-associated genes in the context of a particular disease) and the phenotype-cell type association driver genes.
Gene	11. gene frequency threshold	Keep only genes at or above a certain mean frequency threshold (i.e. how frequently a gene is associated with a given phenotype when observed within a disease).
Phenotype	12. prune ancestors	Remove redundant ancestral phenotypes when at least one of their descendants already exist.
All	13. top n	Only return the top N targets per variable group (specified with the “group_vars” argument). For example, setting “group_vars” to “hpo_id” and “top_n” to 1 would only return one target (row) per phenotype ID after sorting.
NA	14. end	NA

Table 2: **Summary statistics of enrichment results stratified by single-cell atlas.** Summary statistics at multiple levels (tests, cell types, phenotypes, diseases, cell types per phenotype, phenotypes per cell type) stratified by the single-cell atlas that was used as a cell type signature reference (Descartes Human or Human Cell Landscape).

	DescartesHuman	HumanCellLandscape	all
tests significant	19,929	26,585	46,514
tests	848,078	1,358,916	2,206,994
tests significant (%)	2.35	1.96	2.11
cell types significant	77	124	201
cell types	77	124	201
cell types significant (%)	100	100	100
phenotypes significant	7,340	9,049	9,575
phenotypes tested	11,014	10,959	11,028
phenotypes	11,047	11,047	11,047
phenotypes significant (%)	66.4	81.9	86.7
diseases significant	8,628	8,627	8,628
diseases	8,631	8,631	8,631
diseases significant (%)	100	100	100
cell types per phenotype (mean)	1.81	2.43	4.22
cell types per phenotype (median)	1	2	3
cell types per phenotype (min)	0	0	0
cell types per phenotype (max)	31	28	59
phenotypes per cell type (mean)	259	214	231
phenotypes per cell type (median)	252	200	209
phenotypes per cell type (min)	71	57	57
phenotypes per cell type (max)	696	735	735

Table 4: **Cross-ontology mappings between HPO and CL branches.** The last two columns represent the number of cell types that were overrepresented in the on-target HPO branch and the total number of cell types in that branch. A disaggregated version of this table with all descendant cell type names is available in Table 6.

HPO branch	Phenotypes		Cell types (overrepresented)	Cell types (total)
	(total)	CL branch		
Abnormality of the cardiovascular system	673	cardiocyte	5	6
Abnormality of the endocrine system	291	endocrine cell	3	4
Abnormality of the eye	721	photoreceptor cell/retinal cell	5	5
Abnormality of the immune system	255	leukocyte	14	14
Abnormality of the musculoskeletal system	2155	cell of skeletal muscle/chondrocyte	4	4
Abnormality of the nervous system	1647	neural cell	17	24
Abnormality of the respiratory system	292	respiratory epithelial cell/epithelial cell of lung	3	3

Table 5: **Encodings for GenCC evidence scores.** Assigned numeric values for the GenCC evidence levels.

classification_curie	classification_title	encoding
GENCC:100001	Definitive	6
GENCC:100002	Strong	5
GENCC:100003	Moderate	4
GENCC:100009	Supportive	3
GENCC:100004	Limited	2
GENCC:100005	Disputed Evidence	1
GENCC:100008	No Known Disease Relationship	0
GENCC:100006	Refuted Evidence	0

Table 6: **On-target cell types for each Human Phenotype Ontology (HPO) ancestral branch.** Cell type-phenotype branch pairings were manually curated by comparing high-level HPO terms to terms within the Cell Ontology (CL). Each HPO branch is shown as bolded row dividers. Ancestral CL branch names are shown in the first column, along with the specific CL names and IDs.

CL branch	CL name	CL ID
<b>Abnormality of the cardiovascular system</b>		
cardiocyte	cardiac muscle cell	CL:0000746
cardiocyte	regular atrial cardiac myocyte	CL:0002129
cardiocyte	endocardial cell	CL:0002350
cardiocyte	epicardial adipocyte	CL:1000309
cardiocyte	ventricular cardiac muscle cell	CL:2000046
<b>Abnormality of the endocrine system</b>		
endocrine cell	endocrine cell	CL:0000163
endocrine cell	neuroendocrine cell	CL:0000165
endocrine cell	chromaffin cell	CL:0000166
<b>Abnormality of the eye</b>		
photoreceptor cell / retinal cell	photoreceptor cell	CL:0000210
photoreceptor cell / retinal cell	amacrine cell	CL:0000561
photoreceptor cell / retinal cell	Mueller cell	CL:0000636
photoreceptor cell / retinal cell	retinal pigment epithelial cell	CL:0002586
<b>Abnormality of the immune system</b>		
leukocyte	T cell	CL:0000084
leukocyte	mature neutrophil	CL:0000096
leukocyte	mast cell	CL:0000097
leukocyte	microglial cell	CL:0000129
leukocyte	professional antigen presenting cell	CL:0000145
leukocyte	macrophage	CL:0000235
leukocyte	B cell	CL:0000236
leukocyte	dendritic cell	CL:0000451
leukocyte	monocyte	CL:0000576
leukocyte	plasma cell	CL:0000786
leukocyte	alternatively activated macrophage	CL:0000890
leukocyte	thymocyte	CL:0000893
leukocyte	innate lymphoid cell	CL:0001065
<b>Abnormality of the musculoskeletal system</b>		
cell of skeletal muscle / chondrocyte	chondrocyte	CL:0000138
cell of skeletal muscle / chondrocyte	cell of skeletal muscle	CL:0000188
cell of skeletal muscle / chondrocyte	skeletal muscle satellite cell	CL:0000594
<b>Abnormality of the nervous system</b>		
neural cell	bipolar neuron	CL:0000103
neural cell	granule cell	CL:0000120
neural cell	Purkinje cell	CL:0000121
neural cell	glial cell	CL:0000125
neural cell	astrocyte	CL:0000127
neural cell	oligodendrocyte	CL:0000128
neural cell	microglial cell	CL:0000129
neural cell	neuroendocrine cell	CL:0000165
neural cell	chromaffin cell	CL:0000166
neural cell	photoreceptor cell	CL:0000210
neural cell	inhibitory interneuron	CL:0000498
neural cell	neuron	CL:0000540
neural cell	neuronal brush cell	CL:0000555
neural cell	amacrine cell	CL:0000561
neural cell	GABAergic neuron	CL:0000617
neural cell	Mueller cell	CL:0000636
neural cell	glutamatergic neuron	CL:0000679
neural cell	retinal ganglion cell	CL:0000740
neural cell	retina horizontal cell	CL:0000745
neural cell	Schwann cell	CL:0002573
neural cell	retinal pigment epithelial cell	CL:0002586
neural cell	visceromotor neuron	CL:0005025
neural cell	sympathetic neuron	CL:0011103
<b>Abnormality of the respiratory system</b>		
respiratory epithelial cell / epithelial cell of lung	type II pneumocyte	CL:0002063
respiratory epithelial cell / epithelial cell of lung	epithelial cell of lower respiratory tract	CL:0002632

**Table 7: Some HPO phenotype categories or more biased towards foetal- or adult- versions of the same cell type.** We took the top 50 phenotypes with the greatest bias towards foetal-cell type associations (“Foetal-biased”) and the greatest bias towards adult-cell type associations (“Adult-biased”) and fed each list of terms into ontological enrichment tests to get a summary of the representative HPO branches for each group. The phenotypes most biased towards associations with only the foetal versions of cell type and those biased towards the adult versions of cell types. “FDR” is the False Discovery Rate-adjusted p-value from the enrichment test, “log2-fold enrichment” is the log2 fold-change from the enrichment test, and “depth” is the depth of the enriched HPO term in the ontology.

term	name	FDR	log2-fold enrichment	depth
<b>Foetal-biased</b>				
HP:0005105	Abnormal nasal morphology	0.00	4.5	6
HP:0010938	Abnormal external nose morphology	0.00	5.4	7
HP:0000366	Abnormality of the nose	0.00	3.8	5
HP:0000055	Abnormal female external genitalia morphology	0.00	5.2	6
HP:0000271	Abnormality of the face	0.00	1.9	4
HP:0000234	Abnormality of the head	0.00	1.7	3
HP:0000152	Abnormality of head or neck	0.00	1.6	2
HP:0010460	Abnormality of the female genitalia	0.03	2.8	5
HP:0000811	Abnormal external genitalia	0.03	2.8	5
HP:0000078	Abnormality of the genital system	0.03	1.9	3
<b>Adult-biased</b>				
HP:0010647	Abnormal elasticity of skin	0.00	6.0	5
HP:0008067	Abnormally lax or hyperextensible skin	0.00	6.0	6
HP:0011121	Abnormal skin morphology	0.00	2.4	4
HP:0000951	Abnormality of the skin	0.00	2.1	3
HP:0001574	Abnormality of the integument	0.01	1.6	2
HP:0001626	Abnormality of the cardiovascular system	0.02	1.4	2
HP:0030680	Abnormal cardiovascular system morphology	0.02	1.7	3
HP:0025015	Abnormal vascular morphology	0.04	1.9	4
HP:0030962	Abnormal morphology of the great vessels	0.04	2.7	6

Table 8: **Examples of specific phenotypes that are most biased towards associations with only the foetal versions of cell types (“Foetal-biased”) and those biased towards the adult versions of cell types (“Adult-biased”).** “p-value difference” is the difference in the association p-values between the foetal and adult version of the equivalent cell type (foetal-adult bias :  $p_{adult} - p_{foetal} = \Delta p \in [-1, 1]$ ).

HPO name	HPO ID	CL ID	CL name	p-value difference
<b>Foetal-biased</b>				
Short middle phalanx of the 2nd finger	HP:0009577	CL:0000138	chondrocyte	0.99
Abnormal morphology of the nasal alae	HP:0000429	CL:0000057	fibroblast	0.95
Abnormal labia minora morphology	HP:0012880	CL:0000499	stromal cell	0.94
Acromesomelia	HP:0003086	CL:0000138	chondrocyte	0.93
Left atrial isomerism	HP:0011537	CL:0000163	endocrine cell	0.92
Fixed facial expression	HP:0005329	CL:0000499	stromal cell	0.92
Migraine without aura	HP:0002083	CL:0000163	endocrine cell	0.92
Truncal ataxia	HP:0002078	CL:0000163	endocrine cell	0.92
Anteverted nares	HP:0000463	CL:0000057	fibroblast	0.91
Short 1st metacarpal	HP:0010034	CL:0000138	chondrocyte	0.90
<b>Adult-biased</b>				
Symblepharon	HP:0430007	CL:0000138	chondrocyte	-0.97
Abnormally lax or hyperextensible skin	HP:0008067	CL:0000057	fibroblast	-0.94
Reduced bone mineral density	HP:0004349	CL:0000057	fibroblast	-0.94
Paroxysmal supraventricular tachycardia	HP:0004763	CL:0000138	chondrocyte	-0.93
Lack of skin elasticity	HP:0100679	CL:0000057	fibroblast	-0.92
Excessive wrinkled skin	HP:0007392	CL:0000057	fibroblast	-0.91
Bruising susceptibility	HP:0000978	CL:0000057	fibroblast	-0.91
Corneal opacity	HP:0007957	CL:0000057	fibroblast	-0.90
Broad skull	HP:0002682	CL:0000138	chondrocyte	-0.90
Emphysema	HP:0002097	CL:0000057	fibroblast	-0.89