

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

³ Brian M. Schilder Kitty B. Murphy Yichun Zhang Robert Gordon-Smith
⁴ Jai Chapman Momoko Otani Nathan G. Skene

⁵ 2025-01-05

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 un-
10 derlying all RD phenotypes with a known genetic cause by leveraging the Human Phenotype Ontology and
11 transcriptomic single-cell atlases of the entire human body from embryonic, foetal, and adult samples. In to-
12 tal we identified significant associations between 201 cell types and 9,575/11,028 (86.7%) unique phenotypes
13 across 8,628 RDs. This greatly the collective knowledge of RD phenotype-cell type mechanisms. Next, we
14 sought to systematically identify phenotypes in which the application of these results would have the greatest
15 clinical impact based on metrics of severity (e.g. lethality, motor/mental impairment) and compatibility with
16 gene therapy (e.g. filtering out physical malformations). Furthermore, we have made these results entirely
17 reproducible and freely accessible to the global community to maximise their impact, including an inter-
18 active web portal (https://neurogenomics.github.io/rare_disease_celltyping_apps/home). To summarise,
19 this work represents a significant step forward in the mission to treat patients across an extremely diverse
20 spectrum 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 globally¹ (1 in
24 10-20 people)². Over 75% of RDs primarily affect children with a 30% mortality rate by five years of age³.
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 years⁴. Of those, ~46% receive at least one incorrect
29 diagnosis and over 75% of all patients never receive any diagnosis⁵. 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 RDs⁶. 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^{7,8}. Those that have been produced are amongst the world's most expensive drugs, greatly
36 limiting patients' ability to access it^{9,10}. 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 conditions¹⁴. 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 time¹¹. 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^{15,16}. 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 OMIM¹⁷⁻¹⁹, Orphanet^{20,21}, and DECIPHER²². Mappings between
55 HPO terms to other commonly used medical ontologies (e.g. SNOMED CT²³, UMLS^{24,25}, ICD-9/10/11²⁶)
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 practices²⁷⁻²⁹. 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^{30,31}.
68 In particular, the Descartes Human³² and Human Cell Landscape³³ 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 development³⁴. 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^{35,36}.
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 approach³⁷. 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)³⁸. 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)³² and Human Cell Landscape (~703,000 single-cells from 49
108 embryonic, fetal and adult tissues)³³. We ran a series of linear regression models to test for the relationship

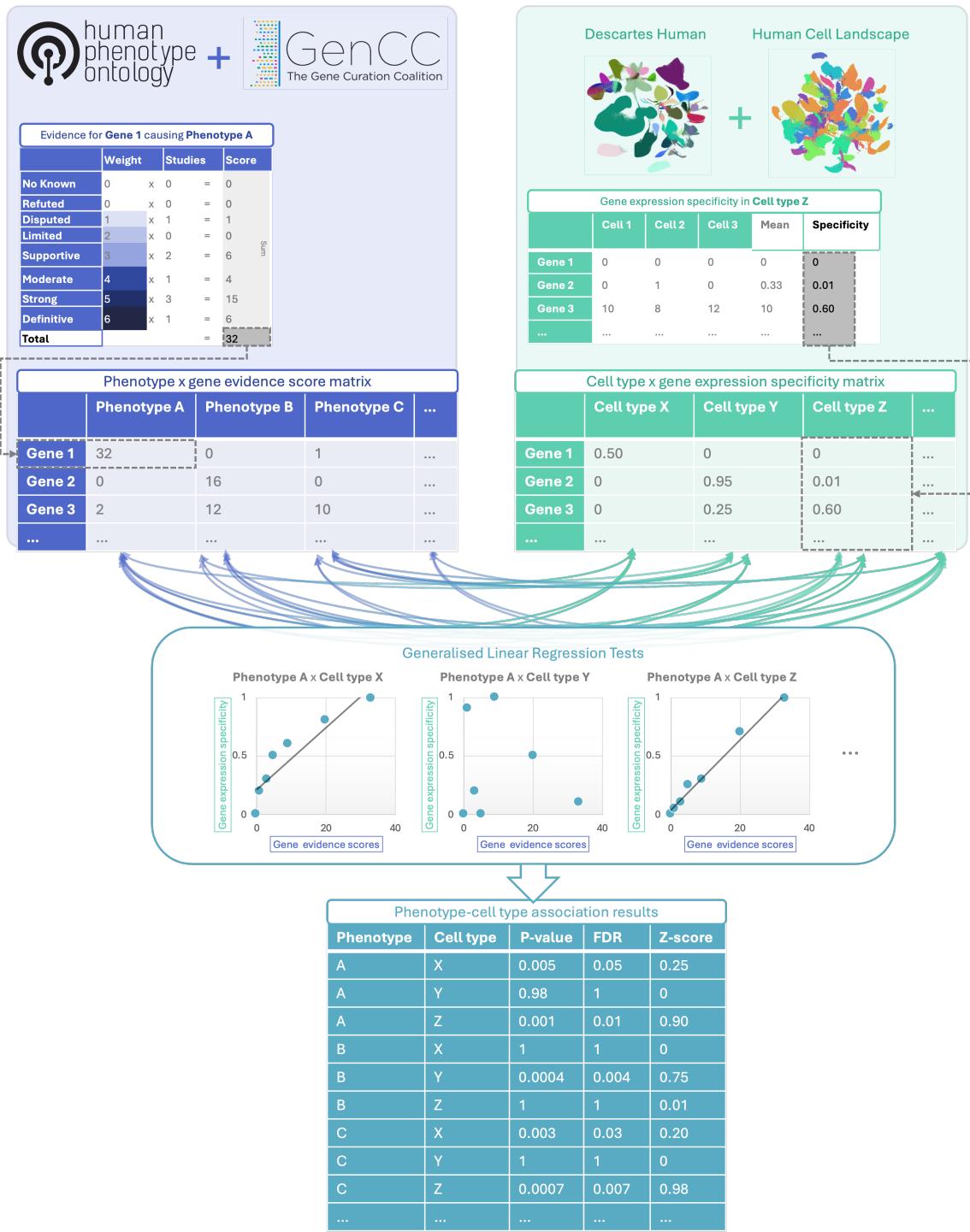


Figure 1: Multi-modal data fusion reveals the cell types underlying thousands of human phenotypes. Schematic overview of study design in which we numerically encoded the strength of evidence linking each gene and each phenotype (using the Human Phenotype Ontology and GenCC databases). We then created 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 3.

124 Validation of expected phenotype-cell type relationships

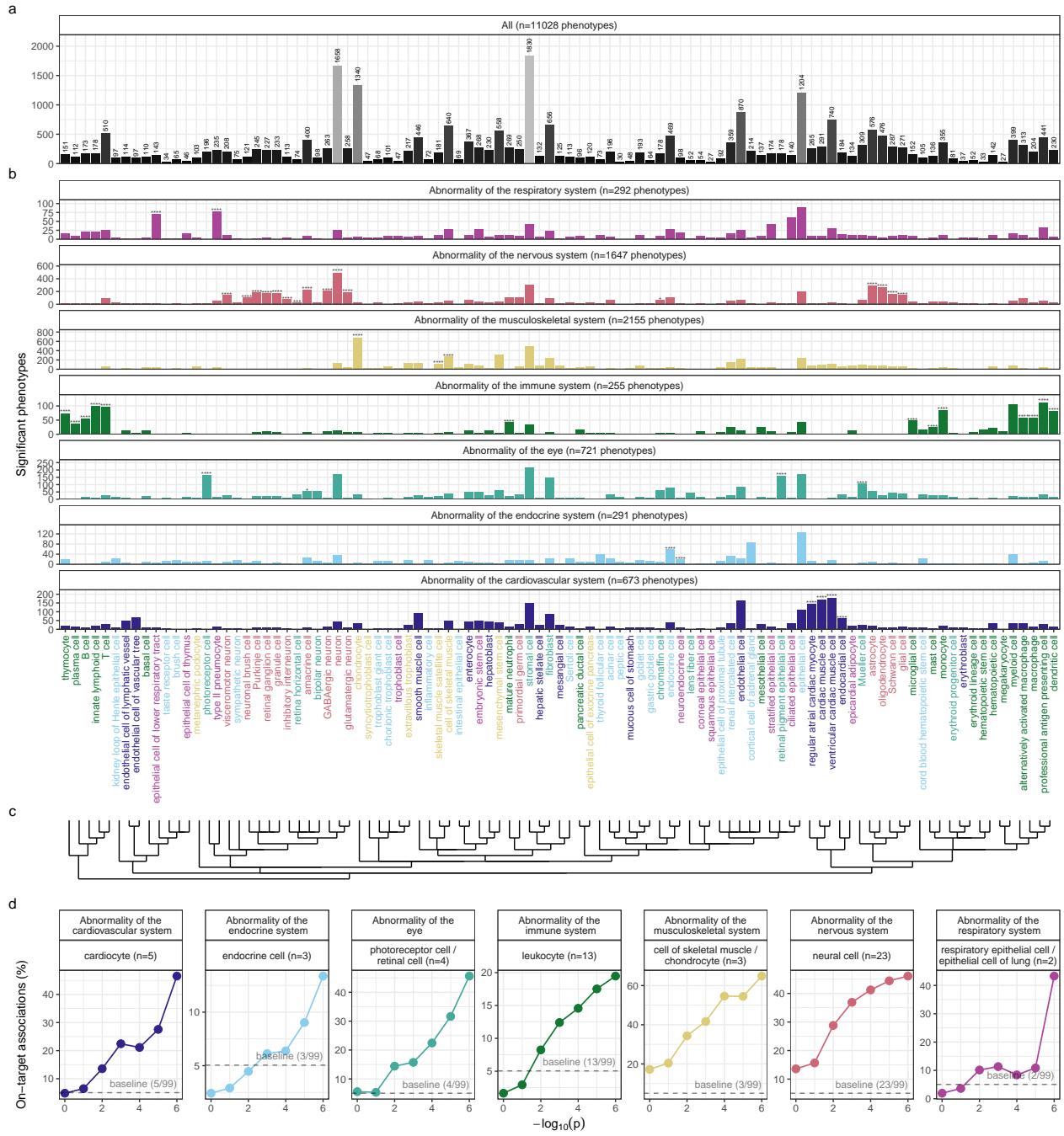
125 We intuitively expect that abnormalities of an organ system will often be driven by cell types within that
126 system. The HPO has broad categories at the higher level of the ontology, enabling us to systematically test
127 this. For example, phenotypes associated with the heart should generally 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 1). 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.

Table 1: 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

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)³⁹. **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

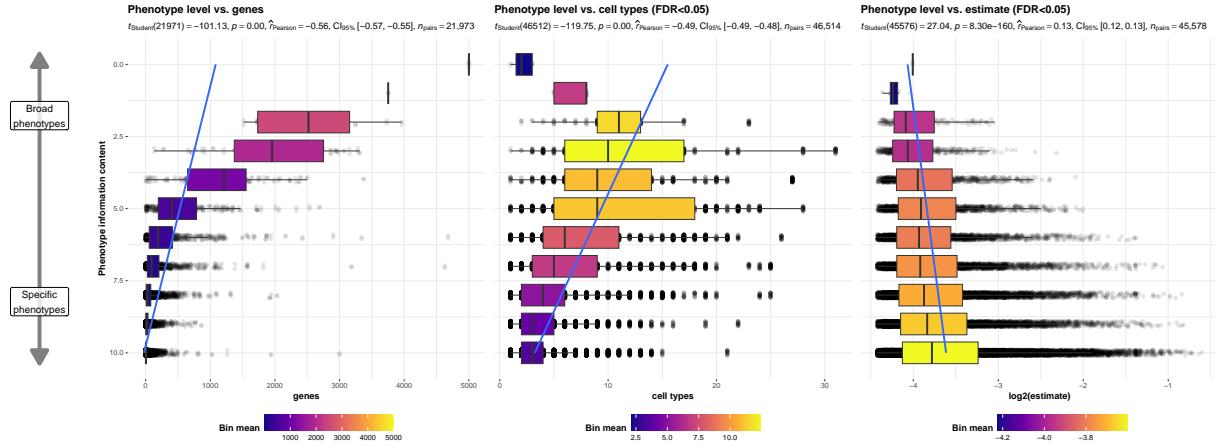
152 **Validation of inter- and intra-dataset consistency**

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. 12. 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 ($\rho=0.49$, $p=1.1 \times 10^{-93}$). Within the subset of results that were significant
159 in both single-cell datasets (FDR<0.05), we found that degree of correlation between the association effect
160 sizes across datasets was even stronger ($\rho=0.72$, $p=1.1 \times 10^{-93}$). 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 ($\rho=0.44$, $p=2.4 \times 10^{-149}$). 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 relationships⁴⁰. 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. 13. 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-

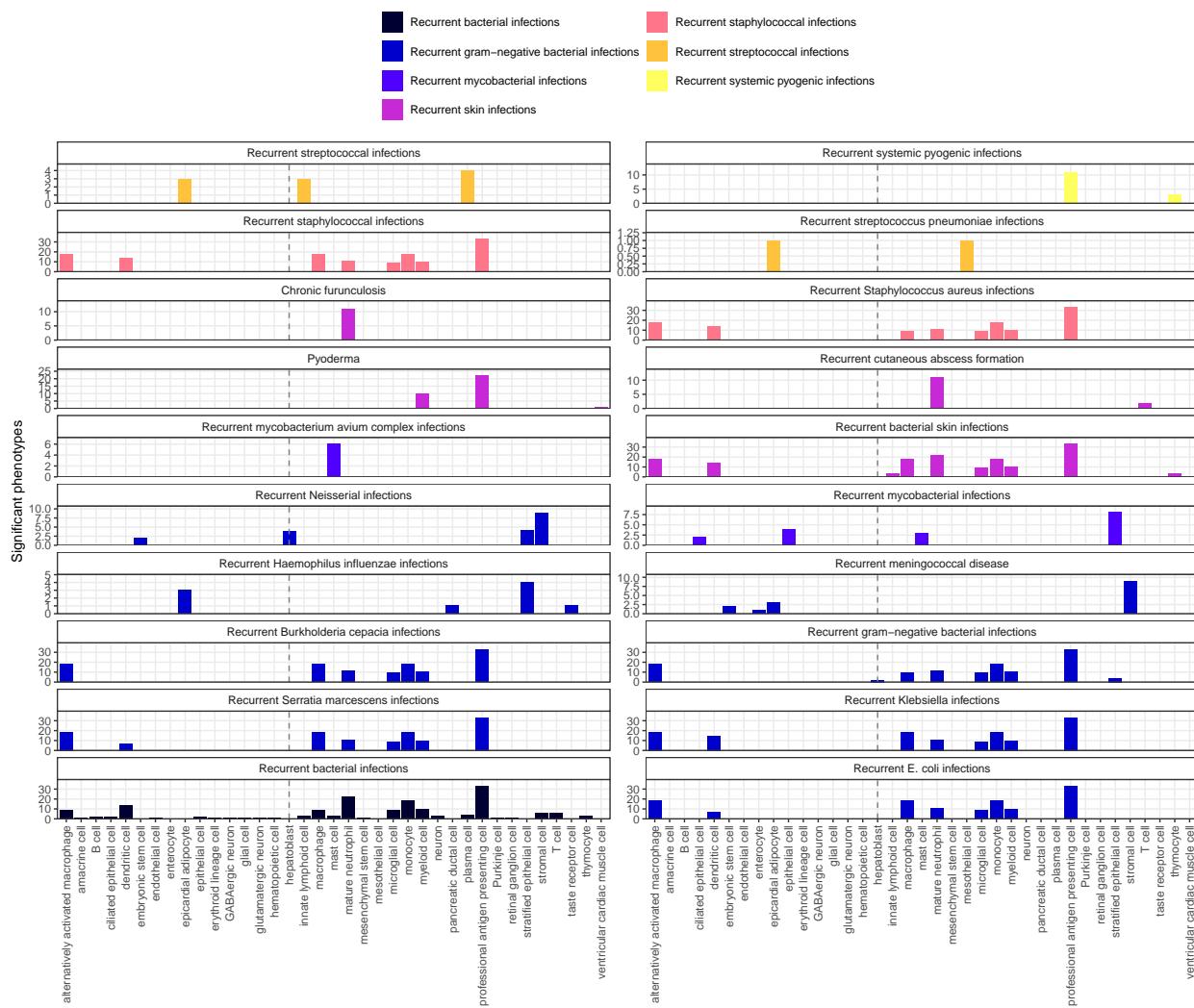
scribed in the MKG. Of these phenotypes, we captured 57% of the MKG phenotype-cell associations at an ontological distance of 0 (i.e. the closest possible Cell Ontology term match). Recall increased with greater flexibility in the matching of cell type annotations. At an ontological distance of 1 (e.g. ‘monocyte’ vs. ‘classical monocyte’), we captured 78% 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^{41–44}. 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 proteins⁴⁵, and Kupffer cells express complement receptors⁴⁶. In addition, individuals with deficits in complement are at high risk for Neisserial infections^{47,48}, and a genome-wide association study in those with a Neisserial infection identified risk variants within complement proteins⁴⁹. While the potential of therapeutically targeting complement in RDs (including Neisserial infections) has been proposed previously^{50,51}, 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 system⁵², highlighting the need for a systematic, evidence-based approach to identify effective gene targets.

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



(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 Next, we visualised the putative causal relationships between genes, cell types and diseases associated with
253 RNI as a network (Fig. 5). The phenotype ‘Recurrent Neisserial infections’ was connected to cell types
254 through the aforementioned association test results ($FDR<0.05$). Genes that were primarily driving these
255 associations (i.e. genes that were both strongly linked with ‘Recurrent Neisserial infections’ and were highly
256 specifically expressed in the given cell type) were designated as “driver genes” and retained for plotting.
257 Across all phenotypes in the HPO, more specific phenotypes (terms in the HPO with greater IC) are not
258 only more specific to certain cell types (Fig. 3b), but are also associated with genes that have greater cell
259 type-specific expression within those cell types. Even so, we should note that the choice of which specificity
260 quantiles to include is arbitrary. It should also be noted that simply because a gene is not specific to a cell
261 type does not mean it is not important for the function of the cell type. Indeed, there are many genes that
262 are ubiquitously expressed throughout many tissues in the body and are essential for cell function. Gene
263 expression specificity is nevertheless a useful metric to help distinguish many hundreds of cell (sub)types
264 with overlapping gene signatures.

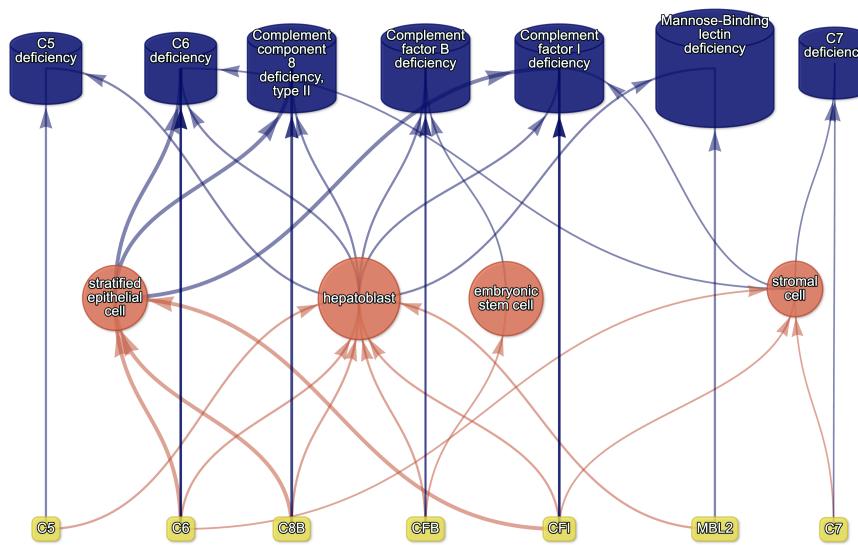
265 Diseases that have ‘Recurrent Neisserial infections’ as a phenotype were collected from the HPO annotation
266 files. Genes that were annotated to a given phenotype (e.g. ‘Recurrent Neisserial infections’) via a particular
267 disease (e.g. ‘C5 deficiency’) constituted “symptom”-level gene sets.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

271 This procedure revealed that genetic deficiencies in various complement system genes (e.g. *C5*, *C8*, and
272 *C7*) are primarily mediated by different cell types (hepatoblasts, stratified epithelial cells, and stromal cells,
273 respectively). While genes of the complement system are expressed throughout many different tissues and

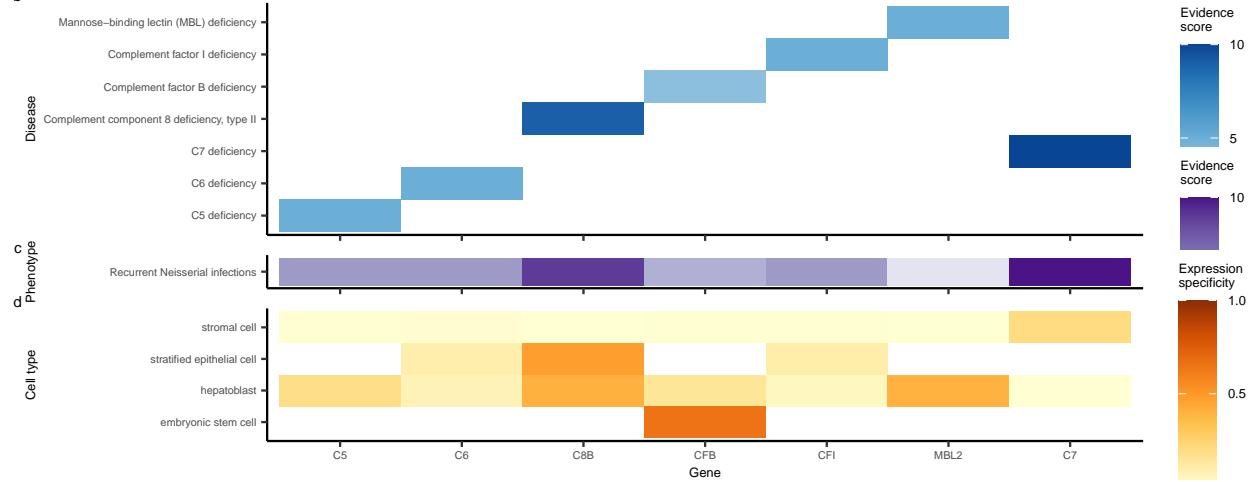
274 cell types, these results indicate that different subsets of these genes may mediate their effects through
275 different cell types. While almost all of these genes show high expression specificity in hepatoblasts, only *C6*,
276 *C7* and *CFI* meet the threshold for the status of driver genes in stromal cells.

277 Recall that we showed in [Fig. 3b] that as we approach the leaf nodes of the HPO we tends towards a given
278 phenotype being associated with a single cell type. Note that mean this in a theoretical sense, as we do
279 not necessarily demonstrate a single cell type for each phenotype in this particular dataset. However, as
280 more granular phenotypes are defined over time, we would expect this hypothesis to bear out. The corollary
281 of this is that we would expect there to be at least four subtypes of the RNI phenotype, as predicted
282 by the four distinct cell types found to underlying this phenotype. This may present as different clinical
283 courses (e.g. early onset, late onset, relapse-remitting) or biomarkers (e.g. histological) to be reveal in future
284 examinations of clinical cohorts. Based on this, we predict that forms of RNI caused by genes expressed in
285 stromal cells would have phenotypic differences from those caused by genes expressed in stratified epithelial
286 cell. In other words, phenotypic similarity is driven by the underlying causal cell types.

a



b



(a) Constructing a multi-scale causal network of disease biology for the phenotype 'Recurrent Neisserial infections' (RNI). RNI is a phenotype in seven different monogenic diseases; these are the only seven genes known to be associated with RNI. Our analyses found four different cell types were significantly associated with RNI. All four of these cell types and all seven diseases are shown in this network plot, allowing us to visualise which cell types each disease-associated gene is acting through. **a**, Starting from the bottom of the plot, one can trace how genes causal for RNI (yellow boxes) mediate their effects through cell types (orange circles) and diseases (blue cylinders). Cell types are connected to RNI via association testing (FDR<0.05). The genes shown here have both strong evidence for a causal role in RNI and high expression specificity in an associated cell type. The cell types can then be linked back to monogenic diseases via the genes specifically expressed in those cell types (i.e. are in the top 1/4 of cell type specificity expression quantiles). Nodes were spatially arranged using the Sugiyama algorithm⁵⁶. **b** Expression specificity quantiles (on a scale from 1-40) of each driver gene in each cell type (darker corresponds to greater specificity).

Figure 5

287 **Prioritising phenotypes based on severity**

288 Some phenotypes are more severe than others and thus could be given priority for developing treatments. For
289 example, ‘Leukonychia’ (white nails) is much less severe than ‘Leukodystrophy’ (white matter degeneration
290 in the brain). Given the large number of significant phenotype-cell type associations, we needed a way of
291 prioritising phenotypes for further investigation. We therefore used the large language model GPT-4 to
292 systematically annotate the severity of all HPO phenotypes³⁷.

293 Severity annotations were gathered from GPT-4 for 16,982/18,082 (94%) HPO phenotypes in our companion
294 study³⁷. Benchmarking tests of these results using ground-truth HPO branch annotations. For example,
295 phenotypes within the ‘Blindness’ HPO branch (*HP:0000618*) were correctly annotated as causing blindness
296 by GPT-4. Across all annotations, the recall rate of GPT-4 annotations was 96% (min=89%, max=100%,
297 SD=4.5) with a mean consistency score of 91% (min=81%, max=97%, SD=5.7) for phenotypes whose
298 annotation were collected more than once. This clearly demonstrates the ability of GPT-4 to accurately
299 annotate phenotypes. This allowed us to begin using these annotations to compute systematically collected
300 severity scores for all phenotypes in the HPO.

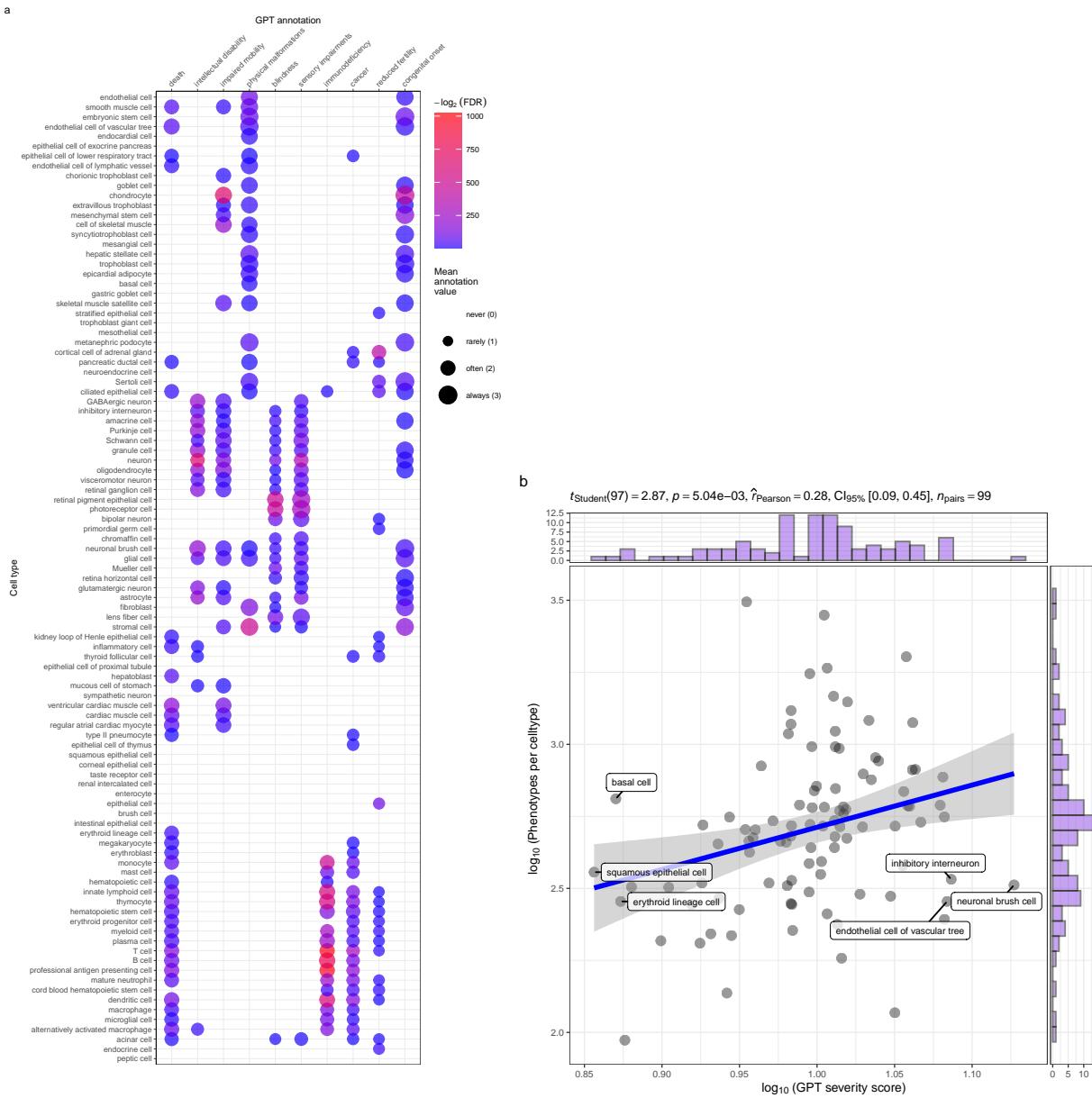
301 From these annotations we computed a weighted severity score metric for each phenotype ranging from 0-100
302 (100 being the theoretical maximum severity of a phenotype that always causes every annotation). Within
303 our annotations, the most severe phenotype was ‘Atrophy/Degeneration affecting the central nervous system’
304 (*HP:0007367*) with a severity score of 47, followed by ‘Anencephaly’ (*HP:0002323*) with a severity score of
305 45. There were 677 phenotypes with a severity score of 0 (e.g. ‘Thin toenail’). The mean severity score
306 across all phenotypes was 10 (median=9.4, standard deviation=6.4).

307 We next sought to answer the question “are disruptions to certain cell types more likely to cause severe
308 phenotypes?”. To address this, we merged the GPT annotations with the significant (FDR<0.05) phenotype-
309 cell type association results and computed the frequency of each severity annotation per cell type (Fig.
310 Figure 14). We found that neuronal brush cells were associated with phenotypes that had the highest
311 average composite severity scores, followed by Mueller cells and glial cells. This suggests that disruptions
312 to these cell types are more likely to cause generally severe phenotypes. Meanwhile, megakaryocytes were
313 associated with phenotypes that had the lowest average composite severity scores, suggesting that disruptions
314 to these cell types can be better tolerated than others.

315 Of course, different aspects of phenotype severity will be more associated with some cells than others. There-
316 fore, after encoding the GPT annotations numerically (0=“never”, 1=“rarely”, 2=“often”, 3=“always”) we
317 computed the mean encoded value per cell type within each annotation. We then ran a series of one-sided
318 Wilcoxon rank-sum tests to objectively determine whether some cell types tended to be associated with
319 phenotypes that more frequently caused certain severity metrics (death, intellectual disability, impaired mo-
320 bility, etc.) relative to all other cell types Fig. 6a. This consistently yielded expected relationships between

321 cell types (e.g. retinal pigment epithelial cell) and phenotype characteristics (e.g. blindness). Similarly, phe-
322 notypes that more commonly cause death are most commonly associated with retinal pigment epithelial
323 cell, and least commonly associated with squamous epithelial cells and bipolar neurons. This pattern is
324 shown consistently across all annotations (e.g. fertility-reducing phenotypes associated with cortical cell of
325 adrenal glands, immunodeficiency-causing phenotypes associated with T cells, mobility-impairing phenotypes
326 associated with chondrocytes, cancer-causing phenotypes associated with T cells, etc.).

327 We also sought to answer whether the number of phenotypes that a cell type is associated with has a
328 relationship with the severity of those phenotypes (Fig. 6b). Our working hypothesis is that when a cell type
329 that affect many different phenotypes is disrupted, the cell type likely performs some critical function that
330 affect many physiological systems. It also means that the individual phenotypes tend to be more severe than
331 other phenotypes that involve less critical cell types. Indeed, we found a significant relationship between
332 number of associated and mean composite phenotype severity ($p=5.0 \times 10^{-3}$, Pearson coefficient=0.28).



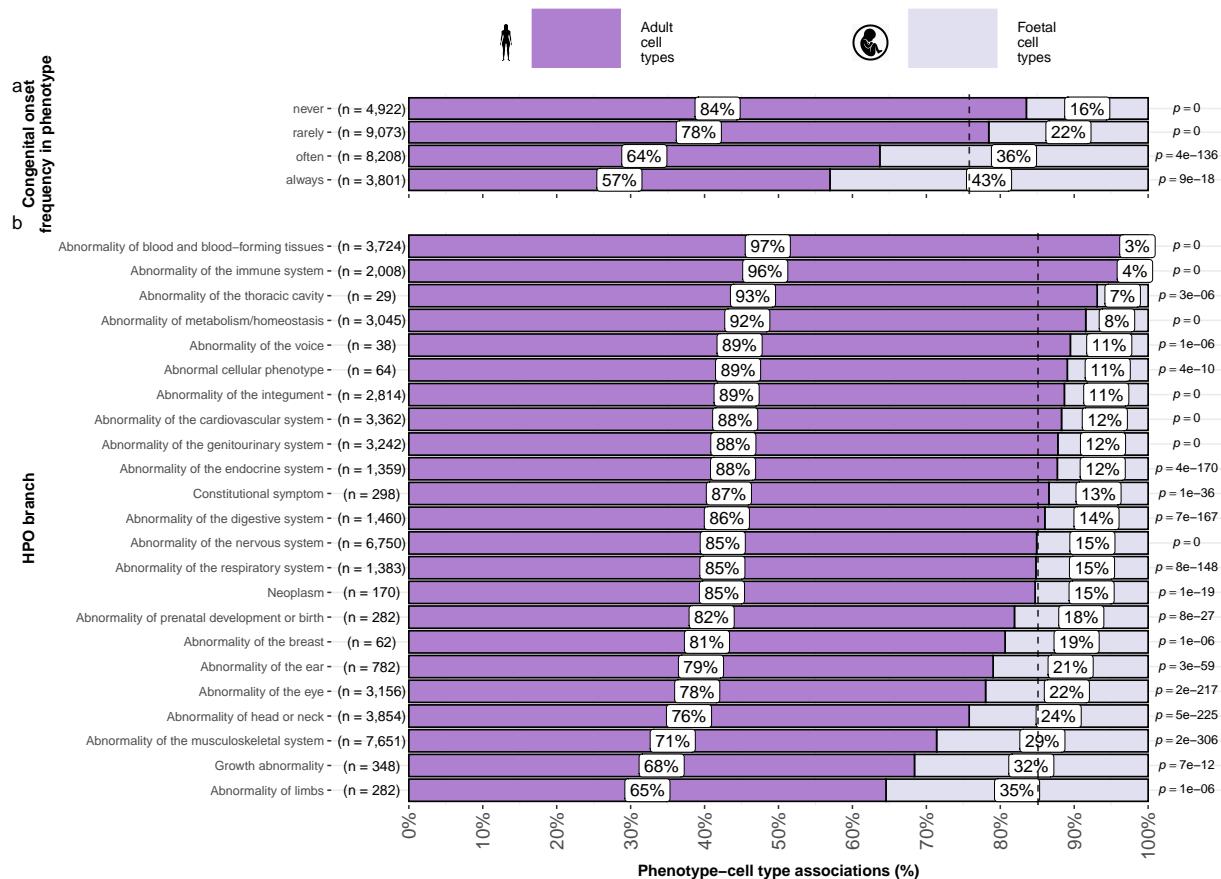
(a) Cell types are differentially associated with phenotypes of varying severity. **a**, The dot plot shows the mean encoded frequency value for a given annotation (0=“never”, 1=“rarely”, 2=“often”, 3=“always”; shown as dot size), aggregated by associated cell type. One-sided Wilcoxon rank-sum tests were performed for each cell type (within each GPT annotation) to determine which cell types had significantly higher frequency values (0-3) 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**, The scatterplot shows a relationship between 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.

Figure 6

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

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

344 We reasoned that phenotypes that are most frequently congenital are more likely to be associated with foetal
345 cell types than adult cell types. As expected, the frequency of congenital onset with each phenotype (as
346 determined by GPT-4 annotations) was strongly predictive of the proportion of significantly associated foetal
347 cell types in our results ($p = 4.7 \times 10^{-261}$, $\chi^2_{Pearson} = 1.2 \times 10^3$, $\hat{V}_{Cramer} = 0.22$). This result is consistent
348 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 χ^2 tests to determine whether the proportion of foetal vs. non-foetal cell types significantly differ 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 7

Upon exploring these results more deeply, we also found that some branches of the HPO were more commonly enriched in foetal cell types compared to others ($\hat{V}_{Cramer}=0.22$, $p<2.2 \times 10^{-308}$). The branch with the greatest proportion of foetal cell type enrichments was ‘Abnormality of limbs’ (35%), followed by ‘Growth abnormality’ (32%) and ‘Abnormality of the musculoskeletal system’ (29%). Most notably, ‘Abnormality of the musculoskeletal system’ was the most enriched branch for foetal cell type associations with a p-value of zzz . These results align well with the fact that physical malformations tend to be developmental in origin. Conversely, the HPO branches that were most biased towards adult cell types were ‘Abnormality of blood and blood-forming tissues’ (97%), ‘Abnormality of the immune system’ (96%), and ‘Abnormality of the thoracic cavity’ (93%).

358 Some phenotypes exclusively involve the foetal version of a cell type, while others exclusively involve the
359 adult version. We sought to find those phenotypes which had the greatest biased towards either end of this
360 spectrum. To do so, we designed a metric to identify which phenotypes were more often associated with
361 foetal cell types than adult cell types. For each phenotype, we calculated the difference in the association
362 p-values between the foetal and adult version of the equivalent cell type. The resulting metric ranges from 1
363 (indicating the phenotype is only associated with the foetal version of the cell type) and -1 (indicating the
364 phenotype is only associated with the adult version of the cell type). To summarise the most foetal-biased
365 phenotype categories, we ran an ontological enrichment test with the HPO graph. To identify foetal cell
366 type-biased phenotype categories, we fed the top 50 phenotypes with the greatest foetal cell type bias (closer
367 to 1) into the enrichment function. Conversely, we used the top 50 phenotypes with the greatest adult cell
368 type bias (closer to -1) to identify adult cell type-biased phenotype categories.

369 The phenotype categories with the greatest bias towards foetal cell types were ‘Abnormal nasal mor-
370 phology’ ($p=2.4 \times 10^{-7}$, $\log_2(\text{fold-change})=4.5$) and ‘Abnormal external nose morphology’ ($p=2.5 \times 10^{-6}$,
371 $\log_2(\text{fold-change})=5.4$). Specific examples of such phenotypes include ‘Short middle phalanx of the 2nd
372 finger’, ‘Abnormal morphology of the nasal alae’, and ‘Abnormal labia minora morphology’. Indeed, these
373 phenotypes are morphological defects apparent at birth caused by abnormal developmental processes.

374 Conversely, the most adult cell type-biased phenotype categories were ‘Abnormal elasticity of skin’
375 ($p=3.6 \times 10^{-7}$, $\log_2(\text{fold-change})=6.0$) and ‘Abnormally lax or hyperextensible skin’ ($p=1.3 \times 10^{-5}$,
376 $\log_2(\text{fold-change})=6.0$). Specific examples of such phenotypes include ‘Excessive wrinkled skin’ and
377 ‘Paroxysmal supraventricular tachycardia’. It is well known that ageing naturally causes a loss of skin
378 elasticity (due to decreasing collagen production) and vascular degeneration⁵⁷. Next, we were interested
379 whether some cell type tend to show strong differences in their phenotype associations between their foetal
380 and adult forms. To test this, we performed an analogous enrichment procedure as with the phenotypes,
381 except using Cell Ontology terms and the Cell Ontology graph. This analysis identified the cell type
382 category connective tissue cell ($p=1.8 \times 10^{-3}$, $\log_2(\text{fold-change})=3.2$) as the most foetal-biased cell type. No
383 cell type categories were significantly enriched for the most adult-biased cell types (Table 10). This is likely
384 due to the fact that cell types can be disrupted at different stages of life, resulting in different phenotypes.
385 Thus there the same cell types may be involved in both the most foetal-biased and adult-biased phenotypes.

386 See Table 8 for the full enrichment results, and Table 9 for specific examples of the most foetal/adult-biased
387 phenotypes. Together, these findings serve to further validate our methodology as a tool for identifying the
388 causal cell types underlying a wide range of phenotypes.

389 Therapeutic target identification

390 In the above sections, we demonstrated how gene association databases can be used to investigate the cell
391 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^{58–60}. 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^{35,36}.

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 both

427 clinicians, biomedical scientists, and pharmaceutical manufacturers who wish to focus their research efforts
428 on phenotypes with the greatest need for intervention.

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

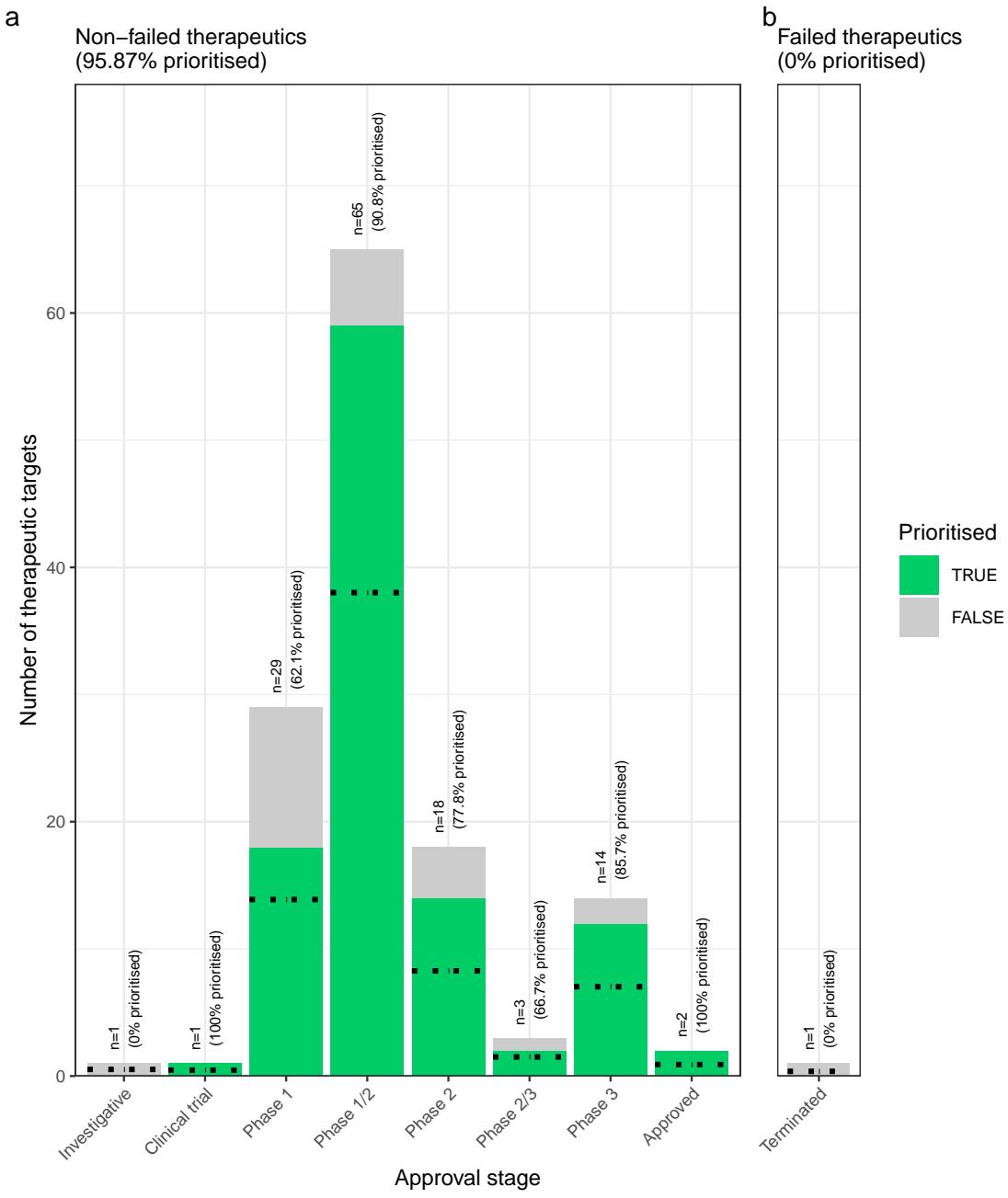
436 Therapeutic target validation

437 To determine whether the genes prioritised by our therapeutic targets pipeline were plausible, we checked
438 what percentage of gene therapy targets we recapitulated. Data on therapeutic approval status was gathered
439 from the Therapeutic Target Database (TTD; release 2025-01-05)⁶¹. Overall, we prioritised 81% of all
440 non-failed existing gene therapy targets. A hypergeometric test confirmed that our prioritised targets were
441 significantly enriched for non-failed gene therapy targets ($p = 1.8 \times 10^{-3}$). Importantly, we did not prioritise
442 any of the failed therapeutics (0%), defined as having been terminated or withdrawn from the market. The
443 hypergeometric test for depletion of failed targets did not reach significance ($p = 0.37$), but this is to be
444 expected as there was only one failed gene therapy target in the TTD database. For these hypergeometric
445 tests, the background gene set was composed of the union of all phenotype-associated genes in the HPO and
446 all gene therapy targets listed in TTD

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

452 In addition to aggregate enrichment results, we also provide specific examples of successful gene therapies
453 whose cell type-specific mechanism were recapitulated by our phenotype-cell associations. In particular,
454 our pipeline nominated the gene *RPE65* within ‘retinal pigment epithelial cell’s as the top target for Fun-
455 dus atrophyUndetectable light- and dark-adapted electroretinogramAbsent foveal reflexEye pokingDelayed
456 early-childhood social milestone developmentPosterior synechiae of the anterior chamberGranular macular
457 appearanceRhegmatogenous retinal detachmentOptic disc drusen vision-related phenotypes that are hall-
458 marks of ‘Leber congenital amaurosis, type II’ and ‘Severe early-childhood-onset retinal dystrophy’. Indeed,
459 gene therapies targeting *RPE65* within the retina of patients with these rare genetic conditions are some of
460 the most successful clinical applications of this technology to date, able to restore vision in many cases⁶².

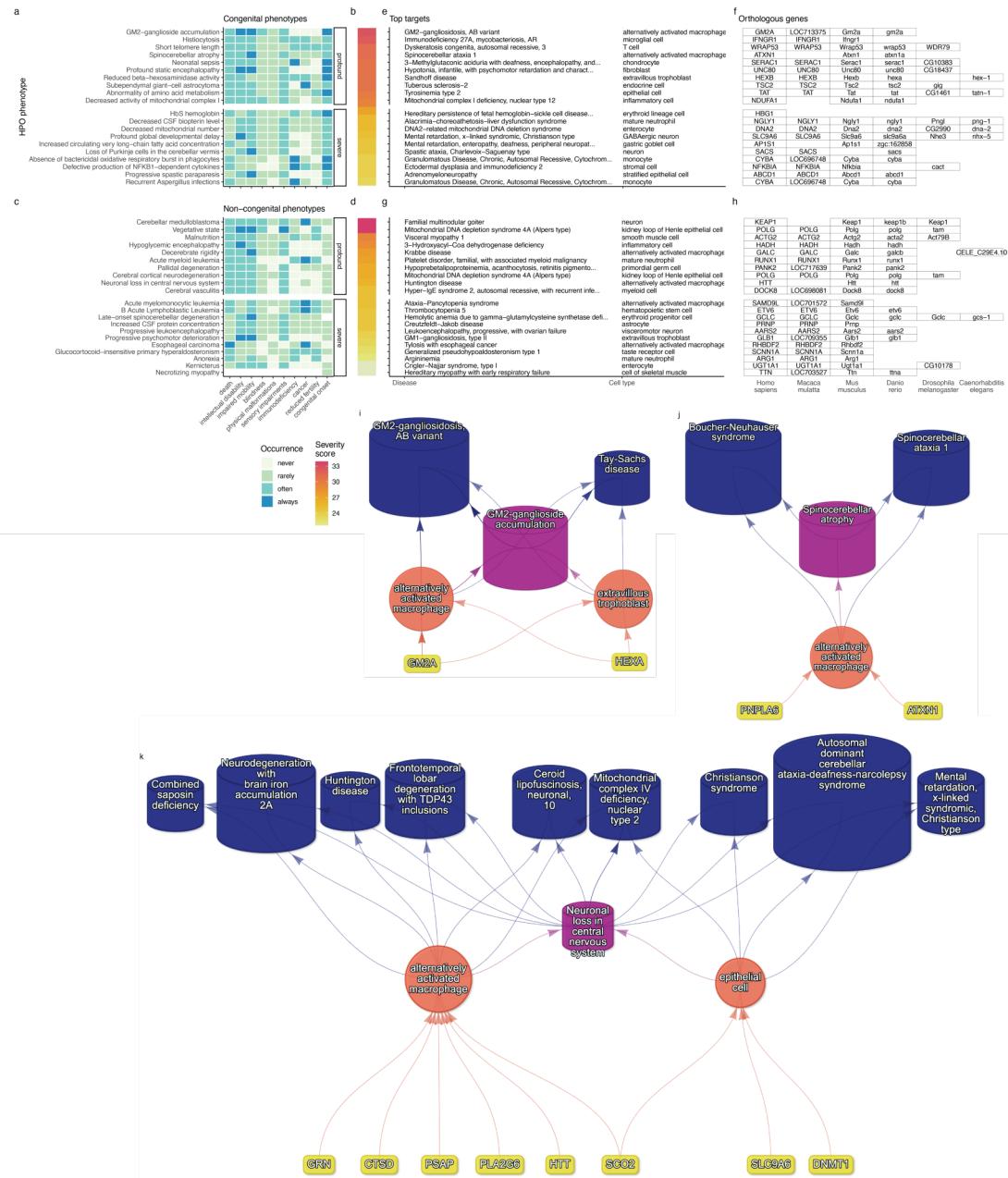
⁴⁶¹ In other cases, a tissue (e.g. liver) may be known to be causally involved in disease genesis, but the precise
⁴⁶² causal cell types within that tissue remain unknown (e.g. hepatocytes, Kupffer cells, Cholangiocytes, Hepatic
⁴⁶³ stellate cells, Natural killer cells, etc.). Tissue-level investigations (e.g. using bulk transcriptomics or epige-
⁴⁶⁴ nomics) would be dominated by hepatocytes, which comprise 75% of the liver. Our prioritized gene therapy
⁴⁶⁵ targets can aid in such scenarios by providing the cell type-resolution context most likely to be causal for a
⁴⁶⁶ 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 8

467 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. **i-k** Example cell type-specific gene therapy targets for several severe phenotypes and their associated diseases. Each disease (blue cylinders) is connected to its phenotype (purple cylinders) based on well-established clinical observations recorded within the HPO¹¹. 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 algorithm⁵⁶.

468 From our prioritised targets, we selected the following four sets of phenotypes or diseases as examples:
469 ‘GM2-ganglioside accumulation’, ‘Spinocerebellar atrophy’, ‘Neuronal loss in central nervous system’. Only
470 phenotypes with a GPT severity score greater than 15 were considered to avoid overplotting and to fo-
471 cus on the more clinically relevant phenotypes. These examples were then selected partly on the basis of
472 severity rankings, and partly for their relatively smaller, simpler networks than lent themselves to compact
473 visualisations.

474 Tay-Sachs disease (TSD) is a devastating hereditary condition in which children are born appearing healthy,
475 which gradually degrades leading to death after 3-5 years. The underlying cause is the toxic accumulation
476 of gangliosides in the nervous system due to a loss of the enzyme produced by *HEXA*. While this could
477 in theory be corrected with gene editing technologies, there remain some outstanding challenges. One of
478 which is identifying which cell types should be targeted to ensure the most effective treatments. Here
479 we identified alternatively activated macrophages as the cell type most strongly associated with ‘GM2-
480 ganglioside accumulation’. The role of aberrant macrophage activity in the regulation of ganglioside levels is
481 supported by observation that gangliosides accumulate within macrophages in TSD⁶³, as well as experimental
482 evidence in rodent models^{64..65,66}. Our results not only corroborate these findings, but propose macrophages
483 as the primary causal cell type in TSD, making it the most promising cell type to target in therapies.

484 Another challenge in TSD is early detection and diagnosis, before irreversible damage has occurred. Our
485 pipeline implicated extravillous trophoblasts of the placenta in ‘GM2-ganglioside accumulation’. While not
486 necessarily a target for gene therapy, checking these cells *in utero* for an absence of *HEXA* may serve as
487 a viable biomarker as these cells normally express the gene at high levels. Early detection of TSD may
488 lengthen the window of opportunity for therapeutic intervention⁶⁷, especially when genetic sequencing is not
489 available or variants of unknown significance are found within *HEXA*⁶⁸.

490 Spinocerebellar atrophy is a debilitating and lethal phenotype that occurs in diseases such as Spinocerebellar
491 ataxia and Boucher-Nenhauser syndrome. These diseases are characterised by progressive degeneration of
492 the cerebellum and spinal cord, leading to severe motor and cognitive impairments. Our pipeline identi-
493 fied M2 macrophages as the only causal cell type associated with ‘Spinocerebellar atrophy’. This strongly
494 suggests that degeneration of cerebellar Purkinje cells are in fact downstream consequences of macrophage
495 dysfunction, rather than being the primary cause themselves. This is consistent with the known role of
496 macrophages, especially microglia, in neuroinflammation and other neurodegenerative conditions such as
497 Alzheimer’s and Parkinsons’ disease⁶⁹⁻⁷¹. While experimental and postmortem observational studies have
498 implicated microglia in spinocerebellar atrophy previously⁶⁹, our results provide a statistically-supported
499 and unbiased genetic link between known risk genes and this cell type. Therefore, targeting M2 microglia in
500 the treatment of spinocerebellar atrophy may therefore represent a promising therapeutic strategy. This is
501 aided by the fact that there are mouse models that perturb the ortholog of human spinocerebellar atrophy
502 risk genes (e.g. *Atxn1*, *Pnpla6*) and reliably recapitulate the effects of this diseases at the cellular (e.g. loss

503 of Purkinje cells), morphological (e.g. atrophy of the cerebellum, spinal cord, and muscles), and functional
504 (e.g. ataxia) levels.

505 Next, we investigated the phenotype ‘Neuronal loss in the central nervous system’. Despite the fact that this
506 is a fairly broad phenotype, we found that it was only significantly associated with 3 cell types (alternatively
507 activated macrophage, macrophage, epithelial cell), specifically M2 macrophages and sinusoidal endothelial
508 cells.

509 Skeletal dysplasia is a heterogeneous group of over 450 disorders that affect the growth and development of
510 bone and cartilage. This phenotype can be lethal when deficient bone growth leads to the constriction of
511 vital organs such as the lungs. Even after surgical interventions, these complications continue to arise as the
512 child develops. Pharmacological interventions to treat this condition have largely been ineffective. While
513 there are various cell types involved in skeletal system development, our pipeline nominated chondrocytes
514 as the causal cell type underlying the lethal form of this condition (Fig. 18). Assuringly, we found that
515 the disease ‘Achondrogenesis Type 1B’ is caused by the genes *SLC26A2* and *COL2A1* via chondrocytes.
516 We also found that ‘Platyspondylic lethal skeletal dysplasia, Torrance type’. Thus, in cases where surgical
517 intervention is insufficient, targeting these genes within chondrocytes may prove a viable long-term solution
518 for children suffering from lethal skeletal dysplasia.

519 Alzheimer’s disease (AD) is the most common neurodegenerative condition. It is characterised by a set of
520 variably penetrant phenotypes including memory loss, cognitive decline, and cerebral proteinopathy. Inter-
521 estingly, we found that different forms of early onset AD (which are defined by the presence of a specific
522 disease gene) are each associated with different cell types via different phenotypes (Fig. 18). For example,
523 AD 3 and AD 4 are primarily associated with cells of the digestive system (‘enterocyte’, ‘gastric goblet
524 cell’) and are implied to be responsible for the phenotypes ‘Senile plaques’, ‘Alzheimer disease’, ‘Parietal
525 hypometabolism in FDG PET’. Meanwhile, AD 2 is primarily associated with immune cells (‘alternatively
526 activated macrophage’) and is implied to be responsible for the phenotypes ‘Neurofibrillary tangles’, ‘Long-
527 tract signs’. This suggests that different forms of AD may be driven by different cell types and phenotypes,
528 which may help to explain its variability in onset and clinical presentation.

529 Finally, Parkinson’s disease (PD) is characterised by motor symptoms such as tremor, rigidity, and bradyki-
530 nesia. However there are a number of additional phenotypes associated with the disease that span multiple
531 physiological systems. PD 19a and PD 8 seemed to align most closely with the canonical understanding of
532 PD as a disease of the central nervous system in that they implicated oligodendrocytes and neurons (Fig. 18).
533 Though the reference datasets being used in this study were not annotated at sufficient resolution to distin-
534 guish between different subtypes of neurons, in particular dopaminergic neurons. PD 19a/8 also suggested
535 that risk variants in *LRRK2* mediate their effects on PD through both myeloid cells and oligodendrocytes
536 by causing gliosis of the substantia nigra. The remaining clusters of PD mechanisms revolved around chon-
537 drocytes (PD 20), amacrine cells of the eye (hereditary late-onset PD), and the respiratory/immune system

538 (PD 14). While the diversity in cell type-specific mechanisms is somewhat surprising, it may help to explain
539 the wide variety of cross-system phenotypes frequently observed in PD.

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

544 **Experimental model translatability**

545 We computed interspecies translatability scores using a combination of both ontological (SIM_o) and geno-
546 typic (SIM_g) similarity relative to each homologous human phenotype and its associated genes Fig. 17.
547 In total, we mapped 278 non-human phenotypes (in *Caenorhabditis elegans*, *Danio rerio*, *Mus musculus*,
548 *Rattus norvegicus*) to 849 homologous human phenotypes. Amongst the 5,252 phenotype within our pri-
549 oritised therapy targets, 354 had viable animal models in at least one non-human species. Per species, the
550 number of homologous phenotypes was: *Danio rerio* (n=214) *Mus musculus* (n=150) *Caenorhabditis elegans*
551 (n=35) *Rattus norvegicus* (n=3). Amongst our prioritised targets with a GPT-4 severity score of >10, the
552 phenotypes with the greatest animal model similarity were ‘Anterior vertebral fusion’ ($SIM_{og} = 0.97$), ‘Disc-
553 like vertebral bodies’ ($SIM_{og} = 0.96$), ‘Metaphyseal enchondromatosis’ ($SIM_{og} = 0.95$), ‘Peripheral retinal
554 avascularization’ ($SIM_{og} = 0.94$), ‘Retinal vascular malformation’ ($SIM_{og} = 0.94$).

555 **Mappings**

Table 2: 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

556 Mappings from HPO phenotypes and other commonly used medical ontologies were gathered in order to
 557 facilitate use of the results in this study in both clinical and research settings. Direct mappings, with a
 558 cross-ontology distance of 1, are the most precise and reliable. Counts of mappings at each distance are
 559 shown in Table 2. In total, there were 15,105 direct mappings between the HPO and other ontologies, with
 560 the largest number of mappings coming from the UMLS ontology (12,898 UMLS terms).

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

564 **Discussion**

565 Investigating RDs at the level of phenotypes offers numerous advantages in both research and clinical
 566 medicine. First, the vast majority of RDs only have one associated gene (7,671/8,631 diseases = 89%).
 567 Aggregating gene sets across diseases into phenotype-centric “buckets” permits sufficiently well-powered
 568 analyses, with an average of ~76 genes per phenotype (median=7) see Fig. 11. Second, we hypothesised
 569 that these phenotype-level gene sets converge on a limited number of molecular and cellular pathways. Per-
 570 turbations to these pathways manifest as one or more phenotypes which, when considered together, tend

571 to be clinically diagnosed as a certain disease. Third, RDs are often highly heterogeneous in their clinical
572 presentation across individuals, leading to the creation of an ever increasing number of disease subtypes
573 (some of which only have a single documented case). In contrast, a phenotype-centric approach enables us
574 to more accurately describe a particular individual's version of a disease without relying on the generation
575 of additional disease subcategories. By characterising an individual's precise phenotypes over time, we may
576 better understand the underlying biological mechanisms that have caused their condition. However, in order
577 to achieve a truly precision-based approach to clinical care, we must first characterise the molecular and
578 cellular mechanisms that cause the emergence of each phenotype. Here, we provide a highly reproducible
579 framework that enables this at the scale of the entire genome.

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

586 It was paramount to the success of this study to ensure our results were anchored in ground-truth bench-
587 marks, generated falsifiable hypotheses, and rigorously guarded against false-positive associations. Extensive
588 validation using multiple approaches demonstrated that our methodology consistently recapitulates expected
589 phenotype-cell type associations (Fig. 2-Fig. 7). This was made possible by the existence of comprehensive,
590 structured ontologies for all phenotypes (the Human Phenotype Ontology) and cell types (the Cell Ontol-
591 ogy), which provide an abundance of clear and falsifiable hypotheses for which to test our predictions against.
592 Several key examples include 1) strong enrichment of associations between cell types and phenotypes within
593 the same anatomical systems (Fig. 2b-d), 2) a strong relationship between phenotype-specificity and the
594 strength and number of cell type associations (Fig. 3), 3) identification of the precise cell subtypes involved
595 in susceptibility to various subtypes of recurrent bacterial infections (Fig. 4), 4) a strong positive correlation
596 between the frequency of congenital onset of a phenotype and the proportion of developmental cell types
597 associated with it (Fig. 7)), and 5) consistent phenotype-cell type associations across multiple independent
598 single-cell datasets (Fig. 12).

599 Unfortunately, there are currently only treatments available for less than 5% of RDs⁶. Novel technologies
600 including CRISPR, prime editing, antisense oligonucleotides, viral vectors, and/or lipid nanoparticles, have
601 been undergone significant advances in the last several years⁷²⁻⁷⁶ and proven remarkable clinical success in
602 an increasing number of clinical applications⁷⁷⁻⁸⁰. The U.S. Food and Drug Administration (FDA) recently
603 announced an landmark program aimed towards improving the international regulatory framework to take
604 advantage of the evolving gene/cell therapy technologies⁸¹ with the aim of bringing dozens more therapies to
605 patients in a substantially shorter timeframe than traditional pharmaceutical product development (typically

606 5-20 years with a median of 8.3 years)⁸². While these technologies have the potential to revolutionise RD
607 medicine, their successful application is dependent on first understanding the mechanisms causing each
608 disease.

609 To address this critical gap in knowledge, we used our results to create a reproducible and customisable
610 pipeline to nominate cell type-resolved therapeutic targets (Fig. 15-Fig. 9). Targeting cell type-specific
611 mechanisms underlying granular RD phenotypes can improve therapeutic effectiveness by treating the causal
612 root of an individual's conditions^{73,83}. A cell type-specific approach also helps to reduce the number of
613 harmful side effects caused by unintentionally delivering the therapeutic to off-target tissues/cell types (which
614 may induce aberrant gene activity), especially when combined with technologies that can target cell surface
615 antigens (e.g viral vectors)⁸⁴. This has the additional benefit of reducing the minimal effective dose of a
616 therapeutic, which can be both immunogenic and extremely financially costly^{9,10,72,75}. Here, we demonstrate
617 the utility of a high-throughput evidence-based approach to RD therapeutics discovery by highlighting several
618 of the most promising therapeutic candidates. Our pipeline takes into account a myriad of factors, including
619 the strength of the phenotype-cell type associations, symptom-cell type associations, cell type-specificity
620 of causal genes, the severity and frequency of the phenotypes, suitability for gene therapy delivery systems
621 (e.g. recombinant adeno-associated viral vectors (rAAV)), as well as a quantitative analysis of phenotypic and
622 genetic animal model translatability (Fig. 17). We validated these candidates by comparing the proportional
623 overlap with gene therapies that are presently in the market or undergoing clinical trials, in which we
624 recovered 81% of all active gene therapies and $NaN \times 10^{-Inf}\%$ of failed gene therapies (Fig. 8, Fig. 16).
625 Despite nominating a large number of putative targets, hypergeometric tests confirmed that our targets were
626 strongly enriched for targets of existing therapies that are either approved or currently undergoing clinical
627 trials.

628 From our target prioritisation pipeline results, we highlight cell type-specific mechanisms for 'GM2-
629 ganglioside accumulation' in Tay-Sachs disease, spinocerebellar atrophy in spinocerebellar ataxia, and
630 'Neuronal loss in central nervous system' in a variety of diseases (Fig. 9). Of interest, all three of these
631 neurodegenerative phenotypes involved alternatively activated (M2) macrophages. The role of macrophages
632 in neurodegeneration is complex, with both neuroprotective and neurotoxic functions, including the
633 clearance of misfolded proteins, the regulation of the blood-brain barrier, and the modulation of the immune
634 response⁸⁵. We also recapitulated prior evidence that microglia, the resident macrophages of the nervous
635 system, are causally implicated in Alzheimer's disease (AD) (Fig. 18)⁸⁶. An important contribution of our
636 current study is that we were able to pinpoint the specific phenotypes of AD caused by macrophages to
637 neurofibrillary tangles and long-tract signs (reflexes that indicate the functioning of spinal long fiber tracts).
638 Other AD-associated phenotypes were caused by other cell types (e.g. gastric goblet cells, enterocytes).

639 It should be noted that our study has several key limitations. First, while our cell type datasets are amongst
640 the most comprehensive human scRNA-seq references currently available, they are nevertheless missing

641 certain tissues, cell types (e.g. spermatocytes, oocytes), and life stages (post-natal childhood, senility). It is
642 also possible that we have not captured certain cell state signatures that only occur in disease (e.g. disease-
643 associated microglia^{87,88}). Though we reasoned that using only control cell type signatures would mitigate
644 bias towards any particular disease, and avoid degradation of gene signatures due to loss of function mutations.
645 Second, the collective knowledge of gene-phenotype and gene-disease associations is far from complete and
646 we fully anticipate that these annotations will continue to expand and change well into the future. It is
647 for this reason we designed this study to be easily reproduced within a single containerised script so that
648 we (or others) may rerun it with updated datasets at any point. Finally, causality is notoriously difficult
649 to prove definitively from associative testing alone, and our study is not exempt from this rule. Despite
650 this, there are several reasons to believe that our approach is able to better approximate causal relationships
651 than traditional approaches. First, we did not intentionally preselect any subset of phenotypes or cell types
652 to investigate here. Along with a scaling prestep during linear modelling, this means that all the results
653 are internally consistent and can be directly compared to one another (in stark contrast to literature meta-
654 analyses). Furthermore, for the phenotype gene signatures we used expert-curated GenCC annotations^{89,90}
655 to weight the current strength of evidence supporting a causal relationship between each gene and phenotype.
656 This is especially important for phenotypes with large genes lists (thousands of annotations) for which some
657 of the relationships may be tenuous. Within the cell type references, we deliberately chose to use specificity
658 scores (rather than raw gene expression) as this normalisation procedure has previously been demonstrated
659 to better distinguish between signatures of highly similar cell types/subtypes⁹¹.

660 Common ontology-controlled frameworks like the HPO open a wealth of new opportunities, especially when
661 addressing RDs. Services such as the Matchmaker Exchange^{92,93} have enabled the discovery of hundreds of
662 underlying genetic etiologies, and led to the diagnosis of many patients. This also opens the possibility of
663 gathering cohorts of geographically dispersed patients to run clinical trials, the only viable option for treat-
664 ment in many individuals. To further increase the number of individuals who qualify for these treatments,
665 as well as the trial sample size, proposals have been made deviate from the traditional single-disease clinical
666 trial model and instead perform basket trials on groups of RDs with shared molecular etiologies (SaME)⁹⁴.

667 Moving forward, we are now actively seeking industry and academic partnerships to begin experimentally
668 validating our multi-scale target predictions and exploring their potential for therapeutic translation. Never-
669 theless, there are more promising therapeutic targets here than our research group could ever hope to pursue
670 by ourselves. In the interest of accelerating research and ensuring RD patients are able to benefit from this
671 work as quickly as possible, we have decided to publicly release all of the results described in this study. These
672 can be accessed in multiple ways, including through a suite of R packages as well as a web app, the Rare Dis-
673 ease Celltyping Portal (https://neurogenomics.github.io/rare_disease_celltyping_apps/home/). The latter
674 allows our results to be easily queried, filtered, visualised, and downloaded without any knowledge of pro-
675 gramming. Through these resources we aim to make our findings useful to a wide variety of RD stakeholders

676 including subdomain experts, clinicians, advocacy groups, and patients.

677 Conclusions

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

686 Methods

687 Human Phenotype Ontology

688 The latest version of the HPO (release releases) was downloaded from the EMBL-EBI Ontology Lookup
689 Service⁹⁵ and imported into R using the `HPOExplorer` package. This R object was used to extract ontolog-
690 ical relationships between phenotypes as well as to assign absolute and relative ontological levels to each
691 phenotype. The latest version of the HPO phenotype-to-gene mappings and phenotype annotations were
692 downloaded from the official HPO GitHub repository and imported into R using `HPOExplorer`. This contains
693 lists of genes associated with phenotypes via particular diseases, formatted as three columns in a table (gene,
694 phenotype, disease).

695 However, not all genes have equally strong evidence of causality with a disease or phenotype, especially when
696 considering that the variety of resources used to generate these annotations (OMIM, Orphanet, DECIPHER)
697 use variable methodologies (e.g. expert-curated review of the medical literature vs. automated text mining
698 of the literature). Therefore we imported data from the Gene Curation Coalition (GenCC)^{89,90}, which (as
699 of 2024-11-02) 22,060 evidence scores across 7,259 diseases and 5,165 genes. Evidence scores are defined
700 by GenCC using a standardised ordinal rubric which we then encoded as a semi-quantitative score ranging
701 from 0 (no evidence of disease-gene relationship) to 6 (strongest evidence of disease-gene relationship) (see
702 Table 5). As each Disease-Gene pair can have multiple entries (from different studies) with different levels
703 of evidence, we then summed evidence scores per Disease-Gene pair to generate aggregated Disease-by-Gene
704 evidence scores. This procedure can be described as follows.

705 Let us denote:

706 • D as diseases.

707 • P as phenotypes in the HPO.

- 708 • G as genes
- 709 • S as the evidence scores describing the strength of the relationship between each Disease-Gene pair.
- 710 • M_{ij} as the aggregated Disease-by-Gene evidence score matrix.

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

711 Next, we extracted Disease-Gene-Phenotype relationships from the annotations file distributed by the HPO
 712 (*phenotype_to_genes.txt*). This provides a list of genes associated with phenotypes via particular diseases,
 713 but does not include any strength of evidence scores.

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

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

716 In order to assign evidence scores to each Phenotype-Gene relationship, we combined the aforementioned
 717 datasets from GenCC (M_{ij}) and HPO (A_{ijk}) by merging on the gene and disease ID columns. For each
 718 phenotype, we then computed the mean of Disease-Gene scores across all diseases for which that phenotype
 719 is a symptom. This resulted in a final 2D tensor of Phenotype-by-Gene evidence scores (L_{ij}):

720

721

722 Tensor of Disease-by-Gene
 evidence scores

723 Tensor of Phenotype-by-Gene
 evidence scores

724 $L_{ij} = \begin{cases} \frac{\sum_{k=1}^f D_i G_j P_k}{f}, & \text{if } D_i G_j \in A, \\ 1, & \text{if } D_i G_j \notin A \end{cases}$

725

726 Construction of the tensor of Phenotype-by-Gene evidence scores.

727

728 Histograms of evidence score distributions at each step in processing can be found in Fig. 10.

730 **Single-cell transcriptomic atlases**

731 In this study, the gene by cell type specificity matrix was constructed using the Descartes Human transcriptome
732 atlas of foetal gene expression, which contains a mixture of single-nucleus and single-cell RNA-seq
733 data (collected with sci-RNA-seq3)³². This dataset contains 377,456 cells representing 77 distinct cell types
734 across 15 tissues. All 121 human foetal samples ranged from 72 to 129 days in estimated postconceptual age.
735 To independently replicate our findings, we also used the Human Cell Landscape which contains single-cell
736 transcriptomic data (collected with microwell-seq) from embryonic, foetal, and adult human samples across
737 49 tissues³³.

738 Specificity matrices were generated separately for each transcriptomic atlas using the R package EWCE
739 (v1.11.3)⁹¹. Within each atlas, cell types were defined using the authors' original freeform annotations
740 in order to preserve the granularity of cell subtypes as well as incorporate expert-identified rare cell types.
741 Cell types were only aligned and aggregated to the level of corresponding Cell Ontology (CL)³⁹
742 annotations afterwards when generating summary figures and performing cross-atlas analyses. Using the original
743 gene-by-cell count matrices from each single-cell atlas, we computed gene-by-cell type expression specificity
744 matrices as follows. Genes with very no expression across any cell types were considered to be uninformative
745 and were therefore removed from the input gene-by-cell matrix $F(g, i, c)$.

746 Next, we calculated the mean expression per cell type and normalised the resulting matrix to transform it
747 into a gene-by-cell type expression specificity matrix ($S_{g,c}$). In other words, each gene in each cell type had
748 a 0-1 score where 1 indicated the gene was mostly specifically expressed in that particular cell type relative
749 to all other cell types. This procedure was repeated separately for each of the single-cell atlases and can be
750 summarised as:

751

752

Compute mean expression of each gene per cell type

Gene-by-cell type specificity matrix

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

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

753

754

755

756 **Phenotype-cell type associations**

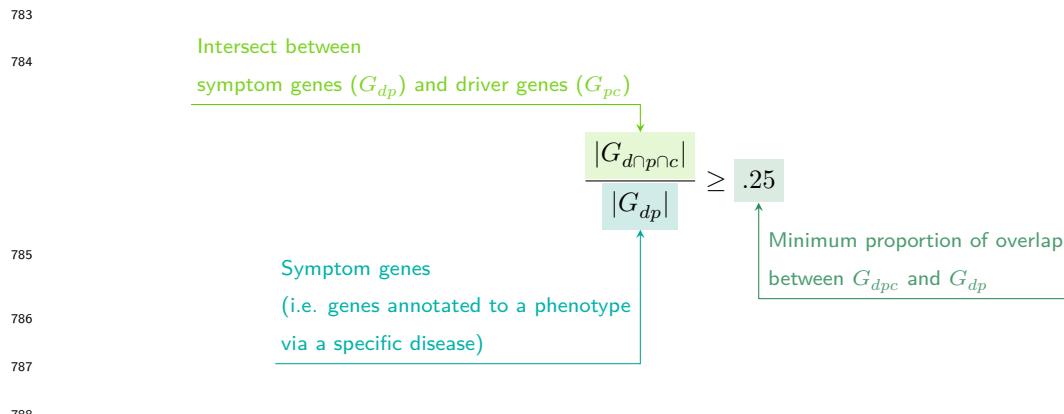
757 To test for relationships between each pairwise combination of phenotype (n=11,047) and cell type (n=201)
758 we ran a series of univariate generalised linear models implemented via the `stats:::glm` function in R. First,

759 we filtered the gene-by-phenotype evidence score matrix (L_{ij}) and the gene-by-cell type expression specificity
 760 matrix (S_{gc}) to only include genes present in both matrices (n=4,949 genes in the Descartes Human analyses;
 761 n=4,653 genes in the Human Cell Landscape analyses). Then, within each matrix any rows or columns with a
 762 sum of 0 were removed as these were uninformative data points that did not vary. To improve interpretability
 763 of the results β coefficient estimates across models (i.e. effect size), we performed a scaling prestep on all
 764 dependent and independent variables. Initial tests showed that this had virtually no impact on the total
 765 number of significant results or any of the benchmarking metrics based on p-value thresholds Fig. 2. This
 766 scaling prestep improved our ability to rank cell types by the strength of their association with a given
 767 phenotype as determined by separate linear models.

768 We repeated the aforementioned procedure separately for each of the single-cell references. Once all results
 769 were generated using both cell type references (2,206,994 association tests total), we applied Benjamini-
 770 Hochberg false discovery rate⁹⁶ (denoted as FDR_{pc}) to account for multiple testing. Of note, we applied
 771 this correction across all results at once (as opposed to each single-cell reference separately) to ensure the
 772 FDR_{pc} was stringently controlled for across all tests performed in this study.

773 Symptom-cell type associations

774 Here we define a symptom as a phenotype as it presents within the context of the specific disease. The features
 775 of a given symptom can be described as the subset of genes annotated to phenotype p via a particular disease
 776 d , denoted as G_{dp} (see Fig. 11). To attribute our phenotype-level cell type enrichment signatures to specific
 777 diseases, we first identified the gene subset that was most strongly driving the phenotype-cell type association
 778 by computing the intersect of genes that were both in the phenotype annotation and within the top 25%
 779 specificity percentile for the associated cell type. We then computed the intersect between symptom genes
 780 (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
 781 overlap with the symptom gene subset (G_{dp}) were kept. This procedure was repeated for all phenotype-cell
 782 type-disease triads, which can be summarised as follows:



789 **Validation of expected phenotype-cell type relationships**

790 We first sought to confirm that our tests (across both single-cell references) were able to recover expected
791 phenotype-cell type relationships across seven high-level branches within the HPO (Fig. 2), including ab-
792 normalities of the cardiovascular system, endocrine system, eye, immune system, musculoskeletal system,
793 nervous system, and respiratory system. Within each branch the number of significant tests in a given
794 cell type were plotted (Fig. 2b). Mappings between freeform annotations (the level at which we performed
795 our phenotype- cell type association tests) provided by the original atlas authors and their closest CL term
796 equivalents were provided by CellxGene³⁰. CL terms along the *x-axis* of Fig. 2b were assigned colours corre-
797 sponding to which HPO branch showed the greatest number of enrichments (after normalising within each
798 branch to account for differences in scale). The normalised colouring allows readers to quickly assess which
799 HPO branch was most often associated with each cell type, while accounting for differences in the number
800 of phenotypes across branches. We then ran a series of Analysis of Variance (ANOVA) tests to determine
801 whether (within a given branch) a given cell type was more often enriched ($FDR < 0.05$) within that branch
802 relative to all of the other HPO branches of an equivalent level in the ontology (including all branches not
803 shown in Fig. 2b). After applying Benjamini-Hochberg multiple testing correction⁹⁶ (denoted as $FDR_{b,c}$),
804 we annotated each respective branch-by-cell type bar according to the significance (**** : $FDR_{b,c} < 1e-04$,
805 *** : $FDR_{b,c} < 0.001$, ** : $FDR_{b,c} < 0.01$, * : $FDR_{b,c} < 0.05$). Cell types in Fig. 2a-b were ordered along
806 the *x-axis* according to a dendrogram derived from the CL ontology (Fig. 2c), which provides ground-truth
807 semantic relationships between all cell types (e.g. different neuronal subtypes are grouped together).

808 As an additional measure of the accuracy of our phenotype-cell types test results we identified conceptually
809 matched branches across the HPO and the CL (Fig. 2d and Table 6). For example, ‘Abnormality of the
810 cardiovascular system’ in the HPO was matched with ‘cardiocytes’ in the CL which includes all cell types
811 specific to the heart. Analogously, ‘Abnormality of the nervous system’ in the HPO was matched with ‘neural
812 cell’ in the CL which includes all descendant subtypes of neurons and glia. This cross-ontology matching
813 was repeated for each HPO branch and can be referred to as on-target cell types. Within each branch, the
814 $-\log_{10}(FDR_{pc})$ values of on-target cell types were binned by rounding to the nearest integer (*x-axis*) and
815 the percentage of tests for on-target cell types relative to all cell types were computed at each bin (*y-axis*)
816 (Fig. 2d). The baseline level (dotted horizontal line) illustrates the percentage of on-target cell types relative
817 to the total number of observed cell types. Any percentages above this baseline level represent greater than
818 chance representation of the on-target cell types in the significant tests.

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

820 We tested for inter-dataset consistency of our phenotype-cell type association results across different single-
821 cell reference datasets (Descartes Human and Human Cell Landscape). First, for association tests with
822 exactly matching Cell Ontology ID across the two references, we tested for a relationship between the p-

823 values generated with each of the references by fitting linear regression model (`stats::lm` via the R function
824 `ggstatsplot::ggscatterstats`). Next, we performed an additional linear regression between the model R^2
825 estimates of all significant phenotype-cell type associations (FDR < 0.05) with exactly matching cell types
826 across the two references.

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

832 More specific phenotypes are associated with fewer genes and cell types

833 To explore the relationship between HPO phenotype specificity and various metrics from our results, we
834 computed the information content (IC) scores for each term in the HPO. IC is a measure of how much
835 specific information a term within an ontology contains. In general, terms deeper in an ontology (closer to the
836 leaves) are more specific, and thus informative, than terms at the very root of the ontology (e.g. ‘Phenotypic
837 abnormality’). Where k denotes the number of offspring terms (including the term itself) and N denotes the
838 total number of terms in the ontology, IC can be calculated as:

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

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

844 Monarch Knowledge Graph recall

845 Finally, we gathered known phenotype-cell type relationships from the Monarch Knowledge Graph (MKG),
846 a comprehensive database of links between many aspects of disease biology⁴⁰. This currently includes 103
847 links between HPO phenotypes (n=103) and CL cell types (n=79). Of these, we only considered the 82
848 phenotypes that we were able to test given that our ability to generate associations was dependent on
849 the existence of gene annotations within the HPO. We considered instances where we found a significant
850 relationship between exactly matching pairs of HPO-CL terms as a hit.

851 However, as the cell types in MKG were not necessarily annotated at the same level as our single-cell refer-
852 ences, we considered instances where the MKG cell type was an ancestor term of our cell type (e.g. ‘myeloid

853 cell' vs. 'monocyte'), or *vice versa*, as hits. We also adjusted ontological distance by computing the ratio
854 between the observed ontological distance and the smallest possible ontological distance for that cell type
855 given the cell type that were available in our references ($dist_{adjusted} = (\frac{dist_{observed}+1}{dist_{minimum}+1}) - 1$). This provides
856 a way of accurately measuring how dissimilar our identified cell types were for each phenotype-cell type
857 association (Fig. 13).

858 **Prioritising phenotypes based on severity**

859 Only a small fraction of the the phenotypes in HPO (<1%) have metadata annotations containing informa-
860 tion on their time course, consequences, and severity. This is due to the time-consuming nature of manually
861 annotating thousands of phenotypes. To generate such annotations at scale, we previously used Generative
862 Pre-trained Transformer 4 (GPT-4), a large language model (LLM) as implemented within OpenAI's Appli-
863 cation Programming Interface (API)³⁷. After extensive prompt engineering and ground-truth benchmarking,
864 we were able to acquire annotations on how often each phenotype directly causes intellectual disability, death,
865 impaired mobility, physical malformations, blindness, sensory impairments, immunodeficiency, cancer, re-
866 duced fertility, or is associated with a congenital onset. These criteria were previously defined in surveys
867 of medical experts as a means of systematically assessing phenotype severity⁹⁷. Responses for each metric
868 were provided in a consistent one-word format which could be one of: 'never', 'rarely', 'often', 'always'. This
869 procedure was repeated in batches (to avoid exceeding token limits) until annotations were gathered for
870 16,982/18,082 HPO phenotypes.

871 We then encoded these responses into a semi-quantitative scoring system ('never'=0, 'rarely'=1, 'often'=2,
872 'always'=3), which were then weighted by multiplying a semi-subjective scoring of the relevance of each
873 metric to the concept of severity on a scale from 1.0-6.0, with 6.0 being the most severe ('death'=6,
874 'intellectual_disability'=5, 'impaired_mobility'=4, 'physical_malformations'=3, 'blindness'=4, 'sen-
875 sory_impairments'=3, 'immunodeficiency'=3, 'cancer'=3, 'reduced_fertility'=1, 'congenital_onset'=1).
876 Finally, the product of the score was normalised to a quantitative severity score ranging from 0-100, where
877 100 is the theoretical maximum severity score. This phenotype severity scoring procedure can be expressed
878 as follows.

879 Let us denote:

- 880 • p : a phenotype in the HPO.
- 881 • j : the identity of a given annotation metric (i.e. clinical characteristic, such as 'intellectual disability'
882 or 'congenital onset').
- 883 • W_j : the assigned weight of metric j .
- 884 • F_j : the maximum possible value for metric j , equal to 3 ("always"). This value is equivalent across all
885 j annotations.

- 886 • F_{pj} : the numerically encoded value of annotation metric j for phenotype p .
- 887 • NSS_p : the final composite severity score for phenotype p after applying normalisation to align values
888 to a 0-100 scale and ensure equivalent meaning regardless of which other phenotypes are being analysed
889 in addition to p . This allows for direct comparability of severity scores across studies with different
890 sets of phenotypes.

891 Sum of weighted annotation values
892 across all metrics
893 Normalised Severity Score
894 for each phenotype
895 Numerically encoded annotation value
896 of metric j for phenotype p
897 Weight for metric j
898 Theoretical maximum severity score
899
$$NSS_p = \frac{\sum_{j=1}^m (F_{pj} \times W_j)}{\sum_{j=1}^m (\max\{F_j\} \times W_j)} \times 100$$

897 Using the numerically encoded GPT annotations (0=“never”, 1=“rarely”, 2=“often”, 3=“always”) we com-
898 puted the mean encoded value per cell type within each annotation. One-sided Wilcoxon rank-sum tests
899 were run using the `rstatix::wilcox_test()` function to test whether each cell type was associated with
900 more severe phenotypes relative to all other cell types. This procedure was repeated for severity annotation
901 independently (death, intellectual disability, impaired mobility, etc.) Fig. 6a. Next, we performed a Pear-
902 son correlation test between the number of phenotypes that a cell type is significantly associated with (at
903 FDR<0.05) has a relationship with the mean composite GPT severity score of those phenotypes (Fig. 6b).
904 This was performed using the `ggstatsplot::ggscatterstats()` R function.

905 Congenital phenotypes are associated with foetal cell types

906 The GPT-4 annotations also enabled us to assess whether foetal cell types were more often significantly
907 associated with congenital phenotypes in our Human Cell Landscape results as this single-cell reference
908 contained both adult and foetal versions of cell types (Fig. 7). To do this, we performed a chi-squared (χ^2)
909 test on the proportion of significantly associated cell types containing any of the substrings ‘fetal’, ‘fetus’,
910 ‘primordial’, ‘hESC’ or ‘embryonic’ (within cell types annotations from the original Human Cell Landscape
911 authors³³) vs. those associated without, stratified by how often the corresponding phenotype had a congenital
912 onset according to the GPT phenotype annotations (including ‘never’, ‘rarely’, ‘often’, ‘always’). In addition,
913 a series of χ^2 tests were performed within each congenital onset frequency strata, to determine whether the
914 observed proportion of foetal cell types vs. non-foetal cell types significantly deviated from the proportions
915 expected by chance.

916 We next tested whether the proportion of tests with significant associations with foetal cell types varied
917 across the major HPO branches using a χ^2 test. We also performed separate χ^2 test within each branch to
918 determine whether the proportion of significant associations with foetal cell types was significantly different
919 from chance.

920 Next, we aimed to create a continuous metric from -1 to 1 that indicated how biased each phenotype
921 is towards associations with the foetal or adult form of a cell type. For each phenotype we calculated the
922 foetal-adult bias score as the difference in the association p-values between the foetal and adult version of the
923 equivalent cell type (foetal-adult bias : $p_{adult} - p_{foetal} = \Delta p \in [-1, 1]$). A score of 1 indicates the phenotype
924 is only associated with the foetal version of the cell type and -1 indicates the phenotype is only associated
925 with the adult version of the cell type. Ontological enrichment tests were then run on the top 50 and bottom
926 50 phenotypes with the greatest foetal-adult bias scores to identify the most foetal-biased and adult-biased
927 phenotype categories, respectively. This was performed using the `simona::dag_enrich_on_offsprings`
928 function, which uses a hypergeometric test to determine whether a list of terms in an ontology are enriched
929 for offspring terms (descendants) of a given ancestor term within the ontology.

930 Therapeutic target identification

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

940 Therapeutic target validation

941 To assess whether our prioritised therapeutic targets were likely to be viable, we computed the overlap
942 between our gene targets and those of existing gene therapies at various stages of clinical development
943 (Fig. 8). Gene targets were obtained for each therapy from the Therapeutic Target Database (TTD; release
944 2025-01-05) and mapped onto standardised HUGO Gene Nomenclature Committee (HGNC) gene symbols
945 using the `orthogene` R package. We stratified our overlap metrics according to whether the therapies had
946 failed (unsuccessful clinical trials or withdrawn), or were non-failed (successful or ongoing clinical trials).
947 We then conducted hypergeometric tests to determine whether the observed overlap between our prioritised
948 targets and the non-failed therapy targets was significantly greater than expected by chance (i.e. enrichment).
949 We also conducted a second hypergeometric test to determine whether the observed overlap between our

950 prioritised targets and the failed therapy targets was significantly less than expected by chance (i.e. depletion).
951 Finally, we repeated the analysis against all therapeutic targets, not just those of gene therapies, to determine
952 whether our prioritised targets had relevance to other therapeutic modalities.

953 Experimental model translatability

954 To improve the likelihood of successful translation between preclinical animal models and human patients,
955 we created an interspecies translatability prediction tool for each phenotype nominated by our gene therapy
956 prioritised pipeline (Fig. 17). First, we extracted ontological similarity scores of homologous phenotypes
957 across species from the MKG⁴⁰. Briefly, the ontological similarity scores (SIM_o) are computed for each
958 homologous pair of phenotypes across two ontologies by calculating the overlap in homologous phenotypes
959 that are ancestors or descendants of the target phenotype. Next, we generated genotypic similarity scores
960 (SIM_g) for each homologous phenotype pair by computing the proportion of 1:1 orthologous genes using
961 gene annotation from their respective ontologies. Interspecies orthologs were also obtained from the MKG.
962 Finally, both scores are multiplied together to yield a unified ontological-genotypic similarity score (SIM_{og}).

963 Novel R packages

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

972 Rare Disease Celltyping Portal

973 To further increase the ease of access for stakeholders in the RD community without the need for program-
974 matic experience, we developed a series of web apps to interactively explore, visualise, and download the
975 results from our study. Collectively, these web apps are called the Rare Disease Celltyping Portal. The
976 landing page for the website was made using HTML, CSS, and javascript and the web apps were created
977 using the Shiny Web application framework for R and deployed on the shinyapps.io server. The website
978 can be accessed at https://neurogenomics.github.io/rare_disease_celltyping_apps/home. All code used to
979 generate the website can be found at https://github.com/neurogenomics/rare_disease_celltyping_apps.

980 **Mappings**

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

987 **Tables**

Table 3: 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: 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.
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.

Table 4: 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
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

988 **Data Availability**

989 All data is publicly available through the following resources:

- 990 • Human Phenotype Ontology (<https://hpo.jax.org>)
- 991 • GenCC (<https://thegencc.org/>)
- 992 • Descartes Human scRNA-seq atlas (<https://cellxgene.cziscience.com/collections/c114c20f-1ef4-49a5-9c2e-d965787fb90c>)
- 993 • Human Cell Landscape scRNA-seq atlas (<https://cellxgene.cziscience.com/collections/38833785-fac5-48fd-944a-0f62a4c23ed1>)
- 994 • Processed Cell Type Datasets (*ctd_DescartesHuman.rds* and *ctd_HumanCellLandscape.rds*; <https://github.com/neurogenomics/MSTExplorer/releases>)
- 995 • Gene x Phenotype association matrix (*hpo_matrix.rds*; <https://github.com/neurogenomics/MSTExplorer/releases>)
- 996 • GPT-4 phenotype severity annotations (https://github.com/neurogenomics/rare_disease_celltyping/releases/download/latest/gpt_check_annot.csv.gz)
- 997 • Full phenotype-cell type association test results https://github.com/neurogenomics/MSTExplorer/releases/download/v0.1.10/phenomix_results.tsv.gz
- 998 • Rare Disease Celltyping Portal (https://neurogenomics.github.io/rare_disease_celltyping_apps/home)
- 999
- 1000
- 1001
- 1002
- 1003
- 1004
- 1005

1006 **Code Availability**

1007 All code is made freely available through the following GitHub repositories:

- 1008 • KGExplorer (<https://github.com/neurogenomics/KGExplorer>)
- 1009 • HPOExplorer (<https://github.com/neurogenomics/HPOExplorer>)
- 1010 • MSTExplorer (<https://github.com/neurogenomics/MSTExplorer>)
- 1011 • Code to replicate analyses (https://github.com/neurogenomics/rare_disease_celltyping)
- 1012 • Cell type-specific gene target prioritisation (https://neurogenomics.github.io/RareDiseasePrioritisation/reports/prioritise_targets)
- 1013 • Complement system gene list (<https://www.genenames.org/data/genegroup/#!/group/492>)
- 1014

1015 **Acknowledgements**

1016 We would like to thank the following individuals for their insightful feedback and assistance with data
1017 resources: Sarah J. Marzi, Gerton Lunter, Peter Robinson, Melissa Haendel, Ben Coleman, Nico Matentzoglu,
1018 Shawn T. O'Neil, Alan E. Murphy, Sarada Gurung.

1019 **Funding**

1020 This work was supported by a UK Dementia Research Institute (UK DRI) Future Leaders Fellowship
1021 [MR/T04327X/1] and the UK DRI which receives its funding from UK DRI Ltd, funded by the UK Medical
1022 Research Council, Alzheimer's Society and Alzheimer's Research UK.

1023 **References**

- 1024 1. Ferreira, C. R. The burden of rare diseases. *Am. J. Med. Genet. A* **179**, 885–892 (2019).
- 1025 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).
- 1026 3. Rare diseases BioResource.
- 1027 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).
- 1028 5. Molster, C. *et al.* Survey of healthcare experiences of australian adults living with rare diseases.
Orphanet J. Rare Dis. **11**, 30 (2016).
- 1029 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).
- 1030 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).
- 1031 8. Yates, N. & Hinkel, J. The economics of moonshots: Value in rare disease drug development. *Clin.
Transl. Sci.* **15**, 809–812 (2022).
- 1032 9. Nuijten, M. Pricing zolgensma - the world's most expensive drug. *J Mark Access Health Policy* **10**,
2022353 (2022).
- 1033 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).
- 1034 11. Gargano, M. A. *et al.* The human phenotype ontology in 2024: Phenotypes around the world. *Nucleic
Acids Res.* **52**, D1333–D1346 (2024).
- 1035 12. Köhler, S. *et al.* Expansion of the human phenotype ontology (HPO) knowledge base and resources.
Nucleic Acids Res. **47**, D1018–D1027 (2019).
- 1036 13. Köhler, S. *et al.* The human phenotype ontology in 2021. *Nucleic Acids Res.* **49**, D1207–D1217
(2021).
- 1037 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).

- 1038 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).
- 1039 16. Rare diseases, common challenges. *Nat. Genet.* **54**, 215 (2022).
- 1040 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).
- 1041 18. Amberger, J. S. & Hamosh, A. Searching online mendelian inheritance in man (OMIM): A knowledgebase of human genes and genetic phenotypes. *Curr. Protoc. Bioinformatics* **58**, 1.2.1–1.2.12 (2017).
- 1042 19. McKusick, V. A. Mendelian inheritance in man and its online version, OMIM. *Am. J. Hum. Genet.* **80**, 588–604 (2007).
- 1043 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).
- 1044 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. Geneeskde.* **152**, 518–519 (2008).
- 1045 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).
- 1046 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).
- 1047 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).
- 1048 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).
- 1049 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).
- 1050 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).
- 1051 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).
- 1052 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).
- 1053 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).

- 1054 31. Svensson, V., Veiga Beltrame, E. da & Pachter, L. A curated database reveals trends in single-cell
transcriptomics. *Database* **2020**, (2020).
- 1055 32. Cao, J. *et al.* A human cell atlas of fetal gene expression. *Science* **370**, (2020).
- 1056 33. Han, X. *et al.* Construction of a human cell landscape at single-cell level. *Nature* **581**, 303–309 (2020).
- 1057 34. Cao, J. *et al.* [A human cell atlas of fetal gene expression](#). *Science* **370**, eaba7721 (2020).
- 1058 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).
- 1059 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).
- 1060 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).
- 1061 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).
- 1062 39. Diehl, A. D. *et al.* The cell ontology 2016: Enhanced content, modularization, and ontology interoperability. *J. Biomed. Semantics* **7**, 44 (2016).
- 1063 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).
- 1064 41. Heim, C. E. *et al.* Myeloid-derived suppressor cells contribute to staphylococcus aureus orthopedic biofilm infection. *J. Immunol.* **192**, 3778–3792 (2014).
- 1065 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).
- 1066 43. Stoll, H. *et al.* Staphylococcal enterotoxins Dose-Dependently modulate the generation of Myeloid-Derived suppressor cells. *Front. Cell. Infect. Microbiol.* **8**, 321 (2018).
- 1067 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).
- 1068 45. Zhou, Z., Xu, M.-J. & Gao, B. Hepatocytes: A key cell type for innate immunity. *Cell. Mol. Immunol.* **13**, 301–315 (2016).
- 1069 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).
- 1070 47. Ladhani, S. N. *et al.* Invasive meningococcal disease in patients with complement deficiencies: A case series (2008–2017). *BMC Infect. Dis.* **19**, 522 (2019).

- 1071 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).
- 1072 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).
- 1073 50. Lung, T. *et al.* The complement system in liver diseases: Evidence-based approach and therapeutic options. *J Transl Autoimmun* **2**, 100017 (2019).
- 1074 51. Reis, E. S. *et al.* Applying complement therapeutics to rare diseases. *Clin. Immunol.* **161**, 225–240 (2015).
- 1075 52. Seal, R. L. *et al.* Genenames.org: The HGNC resources in 2023. *Nucleic Acids Res.* **51**, D1003–D1009 (2023).
- 1076 53. Al-Hamoudi, W. K. Severe autoimmune hepatitis triggered by varicella zoster infection. *World J. Gastroenterol.* **15**, 1004–1006 (2009).
- 1077 54. Brewer, E. C. & Hunter, L. Acute liver failure due to disseminated varicella zoster infection. *Case Reports Hepatol* **2018**, 1269340 (2018).
- 1078 55. Eshchar, J., Reif, L., Waron, M. & Alkan, W. J. Hepatic lesion in chickenpox. A case report. *Gastroenterology* **64**, 462–466 (1973).
- 1079 56. Sugiyama, K., Tagawa, S. & Toda, M. Methods for visual understanding of hierarchical system structures. *IEEE Trans. Syst. Man Cybern.* **11**, 109–125 (1981).
- 1080 57. Li, Z. *et al.* Aging and age-related diseases: From mechanisms to therapeutic strategies. *Biogerontology* **22**, 165–187 (2021).
- 1081 58. Nelson, M. R. *et al.* The support of human genetic evidence for approved drug indications. *Nature Genetics* **47**, 856–860 (2015).
- 1082 59. Ochoa, D. *et al.* Human genetics evidence supports two-thirds of the 2021 FDA-approved drugs. *Nature Reviews Drug Discovery* **21**, 551–551 (2022).
- 1083 60. 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).
- 1084 61. 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).
- 1085 62. 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).
- 1086 63. 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).

- 1087 64. 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).
- 1088 65. 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).
- 1089 66. 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).
- 1090 67. Solovyeva, V. V. *et al.* [New approaches to tay-sachs disease therapy](#). *Frontiers in Physiology* **9**, (2018).
- 1091 68. 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).
- 1092 69. Ferro, A., Sheeler, C., Rosa, J.-G. & Cvetanovic, M. [Role of microglia in ataxias](#). *Journal of molecular biology* **431**, 1792–1804 (2019).
- 1093 70. 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).
- 1094 71. 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).
- 1095 72. 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).
- 1096 73. Bulaklak, K. & Gersbach, C. A. [The once and future gene therapy](#). *Nat. Commun.* **11**, 5820 (2020).
- 1097 74. Godbout, K. & Tremblay, J. P. Prime editing for human gene therapy: Where are we now? *Cells* **12**, (2023).
- 1098 75. Kohn, D. B., Chen, Y. Y. & Spencer, M. J. Successes and challenges in clinical gene therapy. *Gene Ther.* **30**, 738–746 (2023).
- 1099 76. Zhao, Z., Shang, P., Mohanraju, P. & Geijsen, N. Prime editing: Advances and therapeutic applications. *Trends Biotechnol.* **41**, 1000–1012 (2023).
- 1100 77. Darrow, J. J. Luxturna: FDA documents reveal the value of a costly gene therapy. *Drug Discov. Today* **24**, 949–954 (2019).
- 1101 78. Mendell, J. R. *et al.* Single-Dose Gene-Replacement therapy for spinal muscular atrophy. *N. Engl. J. Med.* **377**, 1713–1722 (2017).
- 1102 79. 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).

- 1103 80. 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).
- 1104 81. Lu, C.-F. FDA takes first step toward international regulation of gene therapies to treat rare diseases. (2024).
- 1105 82. 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).
- 1106 83. 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).
- 1107 84. Zhou, Q. & Buchholz, C. J. Cell type specific gene delivery by lentiviral vectors: New options in immunotherapy. *Oncoimmunology* **2**, e22566 (2013).
- 1108 85. 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).
- 1109 86. 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).
- 1110 87. Keren-shaul, H. *et al.* [A unique microglia type associated with restricting development of alzheimer ’s disease](#). *Cell* **169**, 1276–1290.e17 (2017).
- 1111 88. Deczkowska, A. *et al.* [Disease-associated microglia: A universal immune sensor of neurodegeneration](#). *Cell* **173**, 1073–1081 (2018).
- 1112 89. DiStefano, M. T. *et al.* The gene curation coalition: A global effort to harmonize gene-disease evidence resources. *Genet. Med.* **24**, 1732–1742 (2022).
- 1113 90. 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).
- 1114 91. 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).
- 1115 92. 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).
- 1116 93. Philippakis, A. A. *et al.* The matchmaker exchange: A platform for rare disease gene discovery. *Hum. Mutat.* **36**, 915–921 (2015).
- 1117 94. Zanello, G. *et al.* Targeting shared molecular etiologies to accelerate drug development for rare diseases. *EMBO Mol. Med.* **15**, e17159 (2023).
- 1118 95. Côté, R. *et al.* The ontology lookup service: Bigger and better. *Nucleic Acids Res.* **38**, W155–60 (2010).

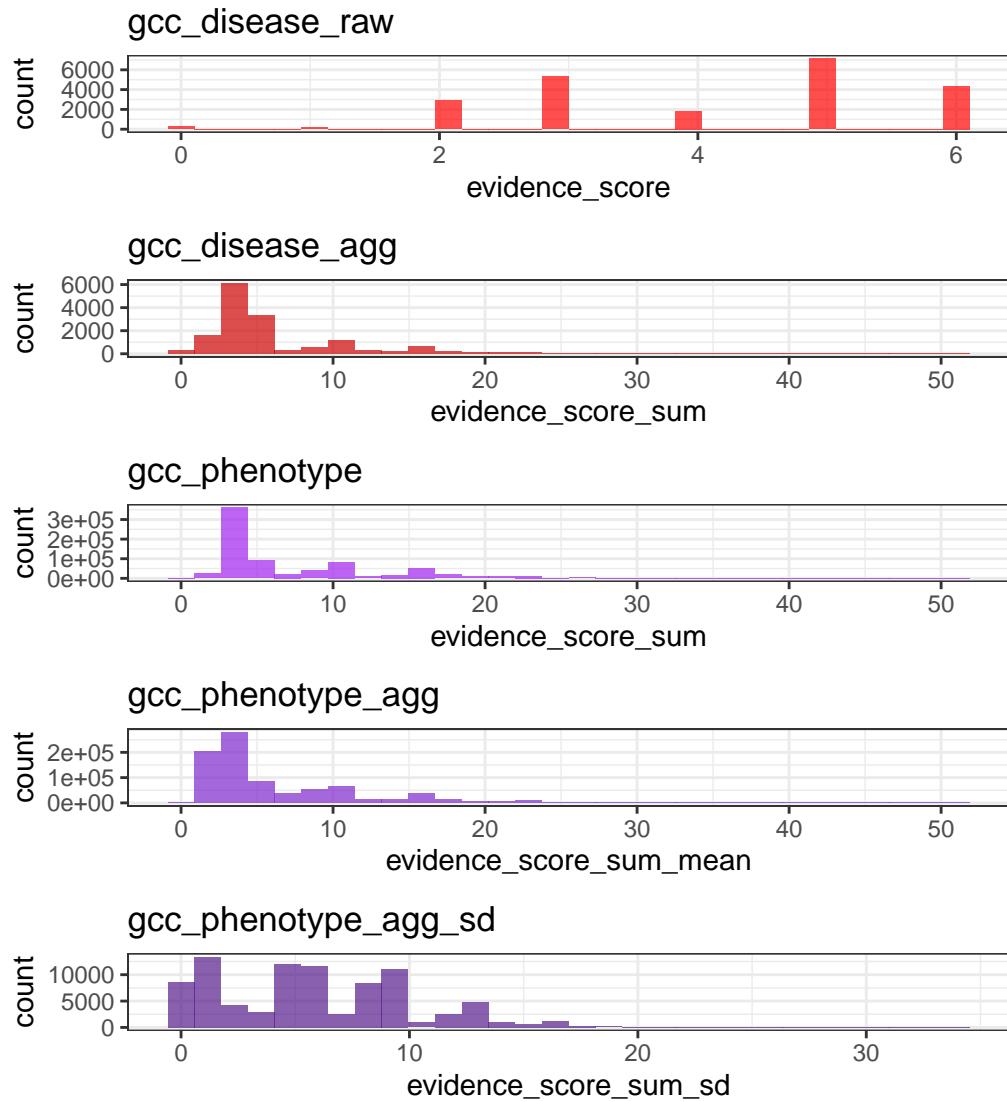
- 1119 96. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach
to multiple testing. *J. R. Stat. Soc.* (1995).
- 1120 97. Lazarin, G. A. *et al.* Systematic classification of disease severity for evaluation of expanded carrier
screening panels. *PLoS One* **9**, e114391 (2014).

1121

1122

₁₁₂₃ **Supplementary Materials**

₁₁₂₄ **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 10

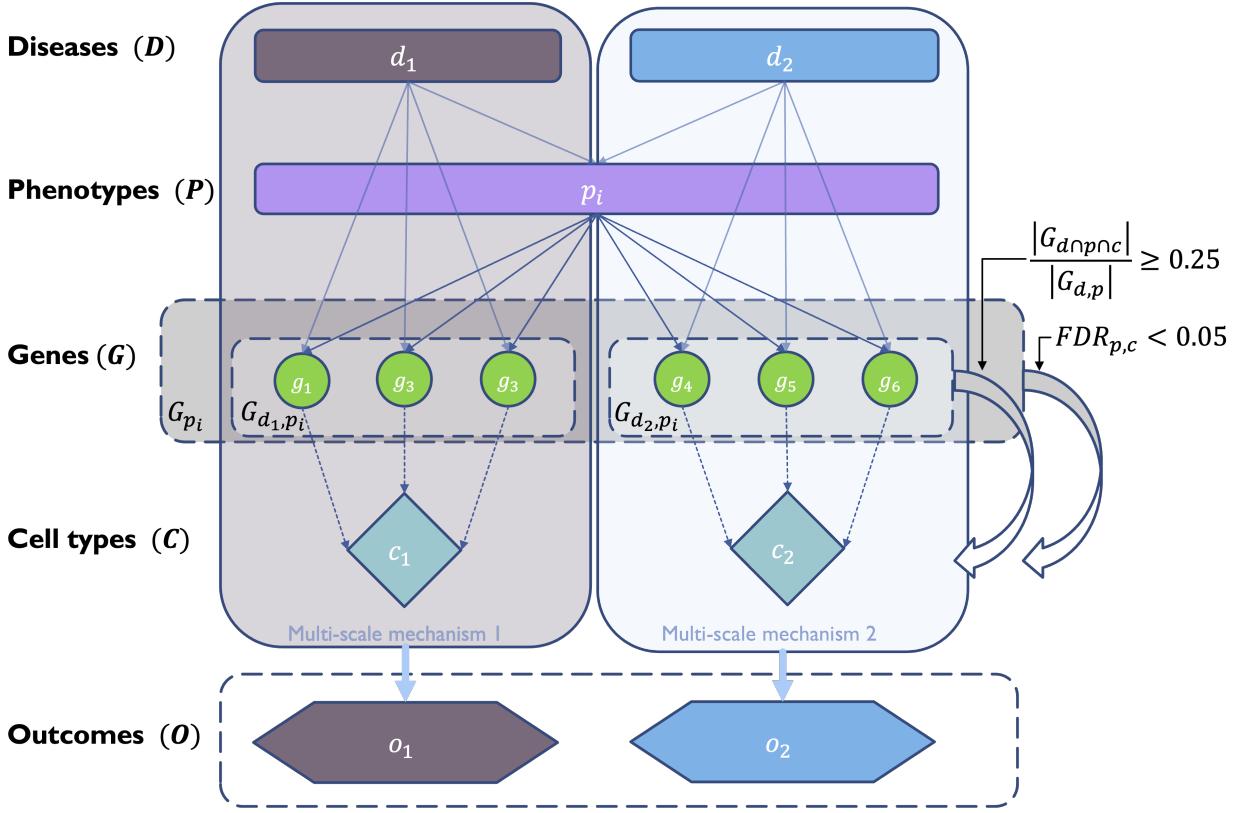
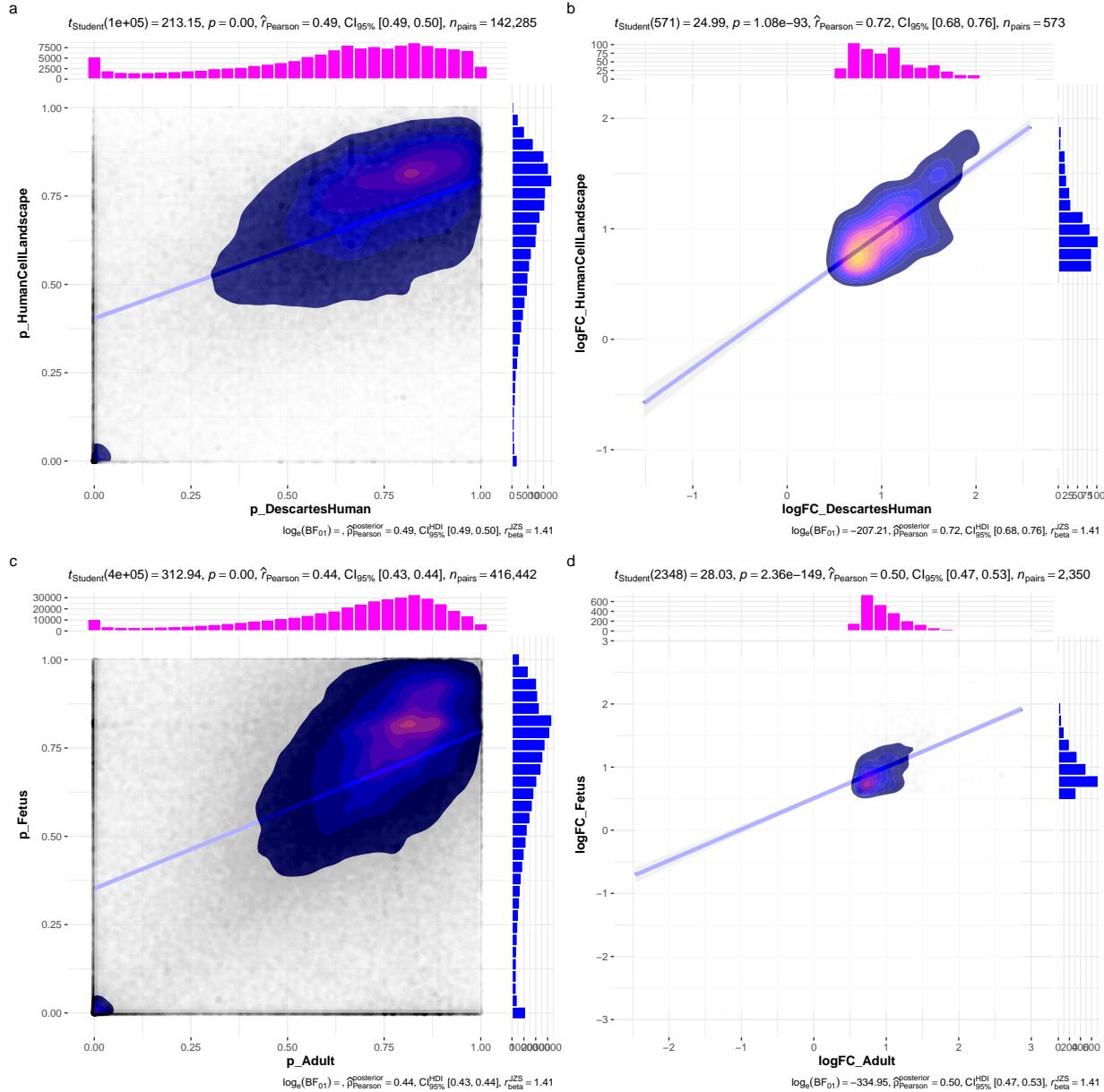
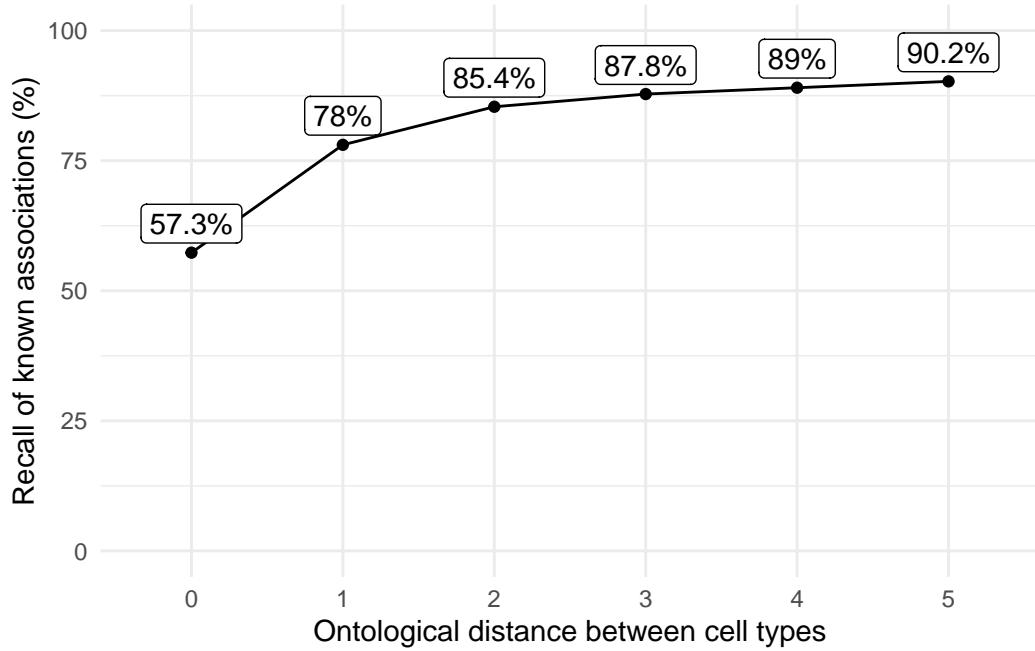


Figure 11: Diagrammatic overview of multi-scale disease investigation strategy. Here we provide an abstract example of differential disease aetiology across multiple scales: diseases (D), phenotypes (P), cell types (C), genes (G), and clinical outcomes (O). In the HPO, genes are assigned to phenotypes 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%.



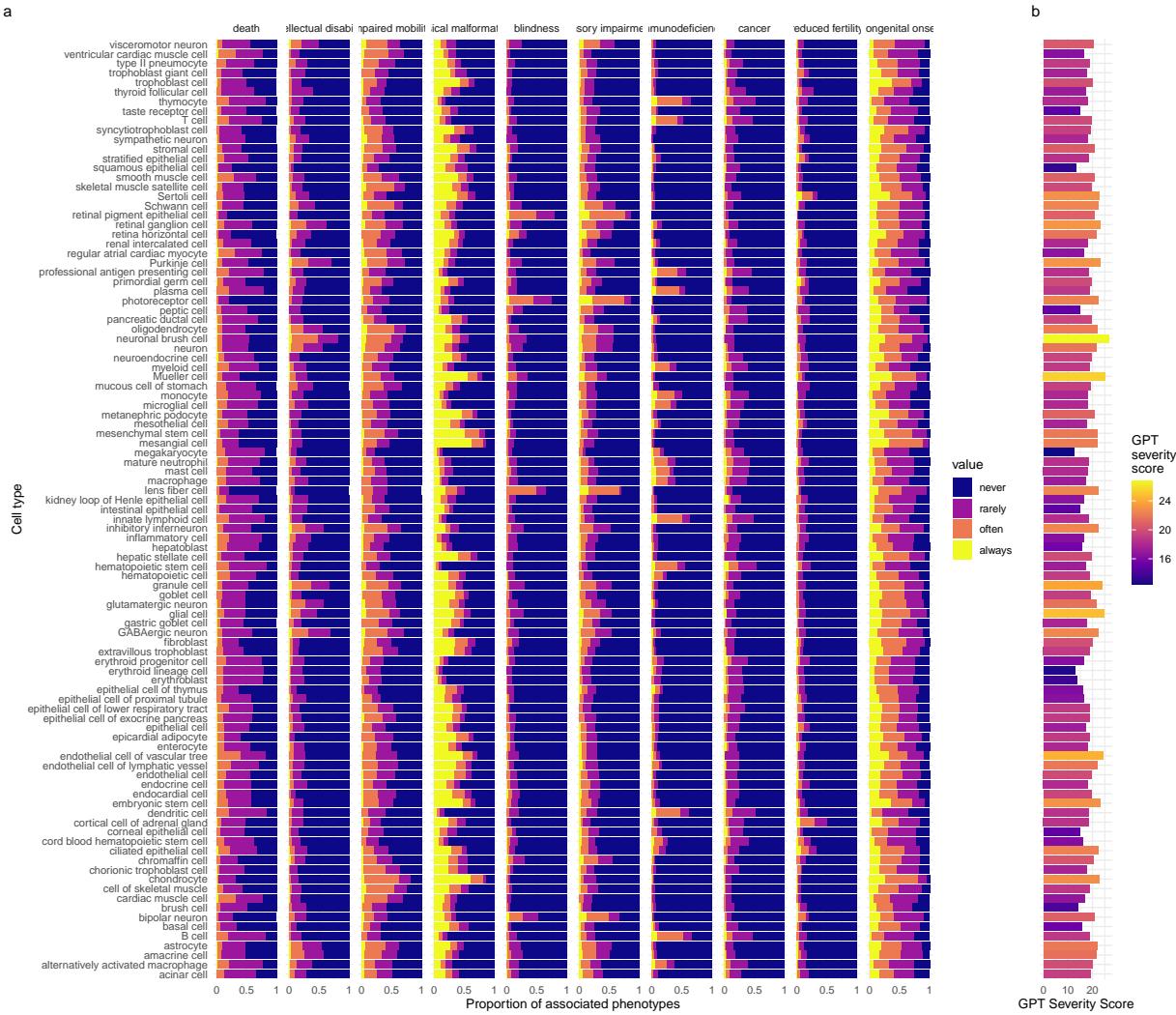
(a) Inter- and intra-dataset validation across the different CellTypeDataset (CTD) and developmental stages. Correlations are computed using Pearson correlation coefficient. Point density is plotted using a 2D kernel density estimate. **a** Correlation between the uncorrected p-values from all phenotype-cell 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 12



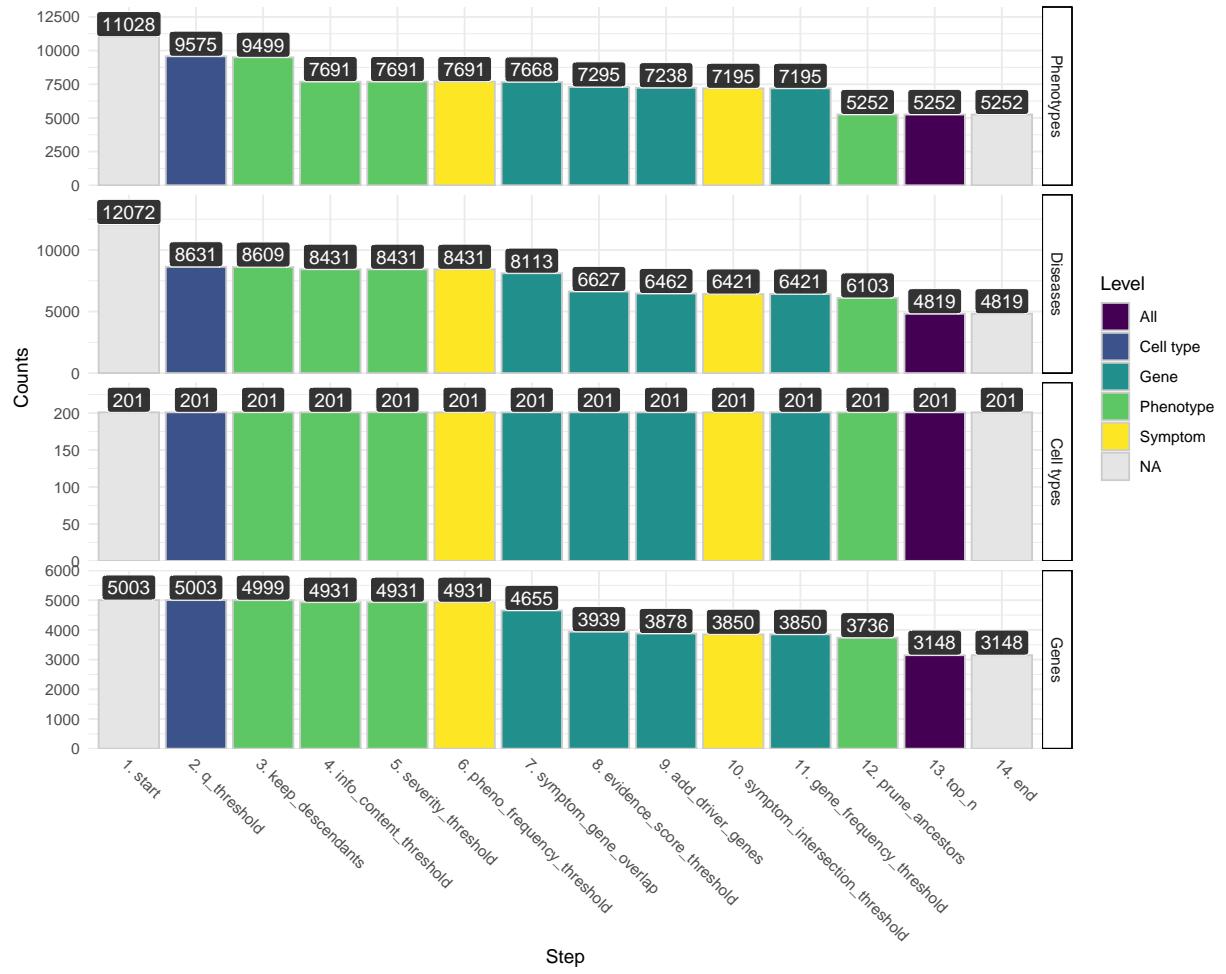
(a) Recall of ground-truth Monarch Knowledge Graph phenotype-cell type relationships at each ontological distance between cell types according to the Cell Ontology.

Figure 13



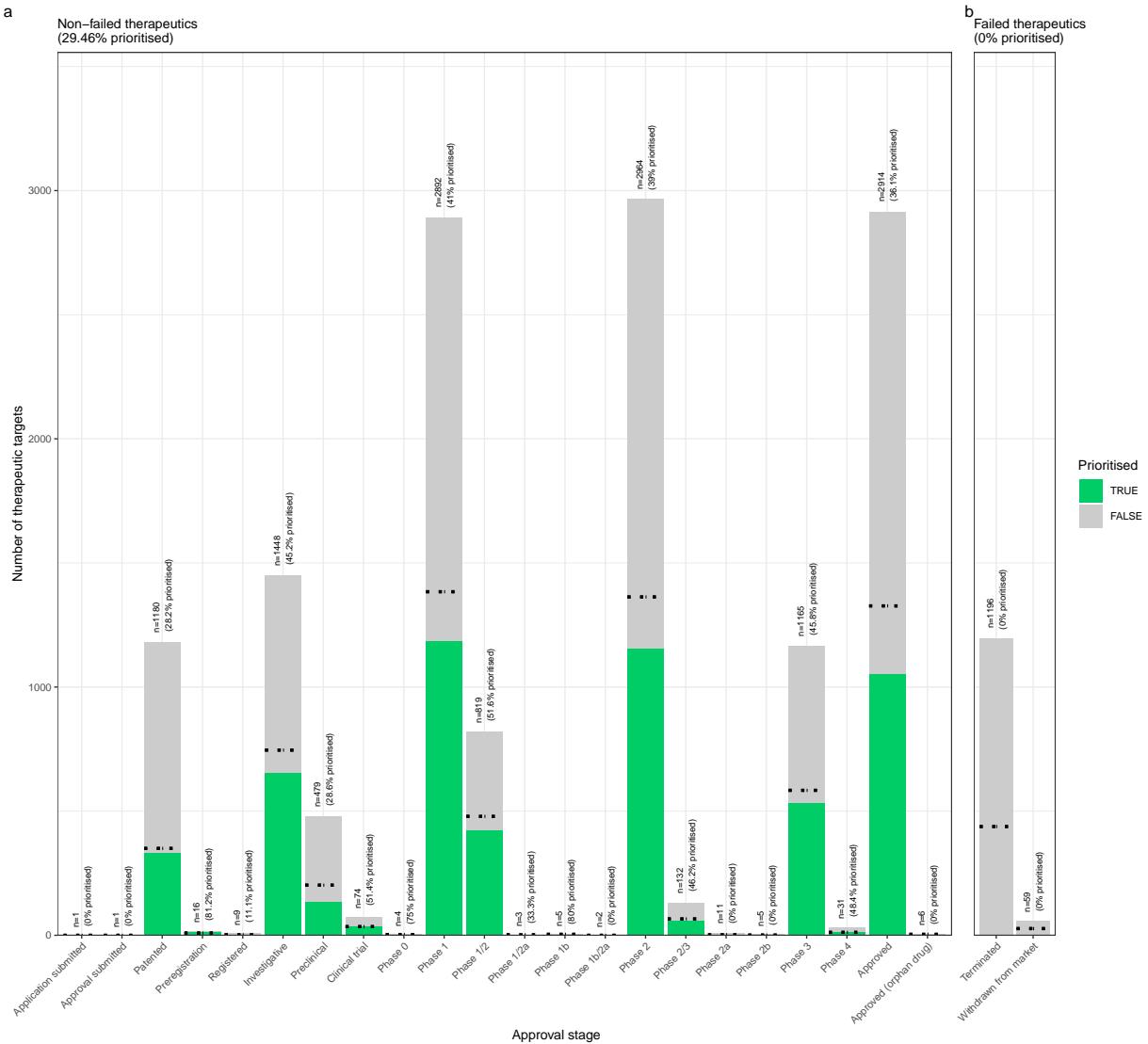
(a) Cell types ordered by the mean severity of the phenotypes they're associated with. **a**, The distribution of phenotype severity annotation frequencies aggregated by cell type. **b**, The composite severity score, averaged across all phenotypes associated with each cell type.

Figure 14



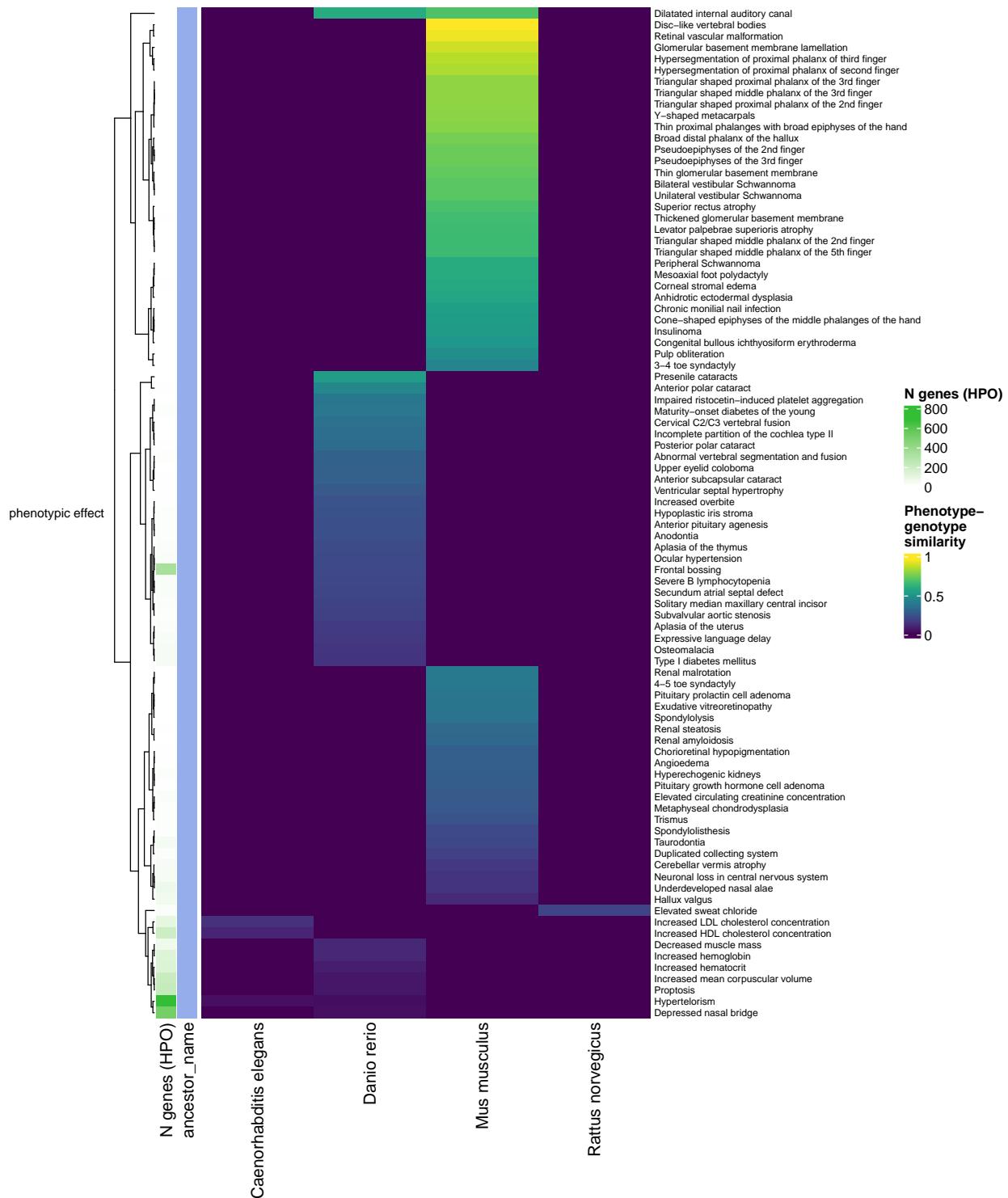
(a) Prioritised target filtering steps. This plot visualises the number of unique phenotype-cell type associations, cell types, genes, and phenotypes (*y-axis*) at each filtering step (*x-axis*) within the multi-scale therapeutic target prioritisation pipeline. Each step in the pipeline can be easily adjusted according to user preference and use case. See Table 4 for descriptions and criterion of each filtering step.

Figure 15



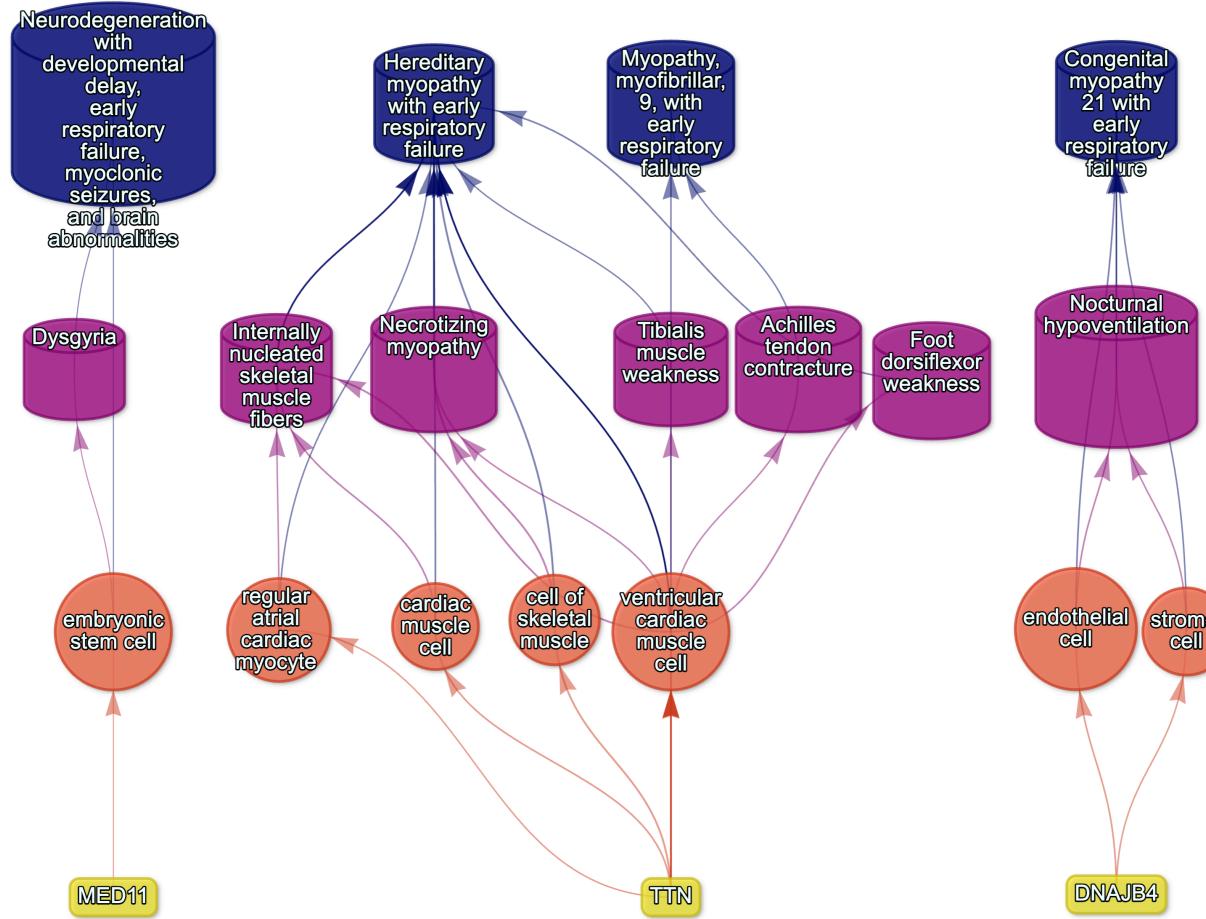
(a) Therapeutics - Validation of prioritised therapeutic targets. Proportion of existing all therapy targets (documented in the Therapeutic Target Database) recapitulated by our prioritisation pipeline.

Figure 16



(a) Identification of translatable experimental models. Interspecies translatability of 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) Respiratory failure

Figure 18: Example cell type-specific gene therapy targets for several severe phenotypes and their associated diseases. Each disease (blue cylinders) is connected to its phenotype (purple cylinders) based on well-established clinical observations recorded within the HPO¹¹. 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 algorithm⁵⁶.

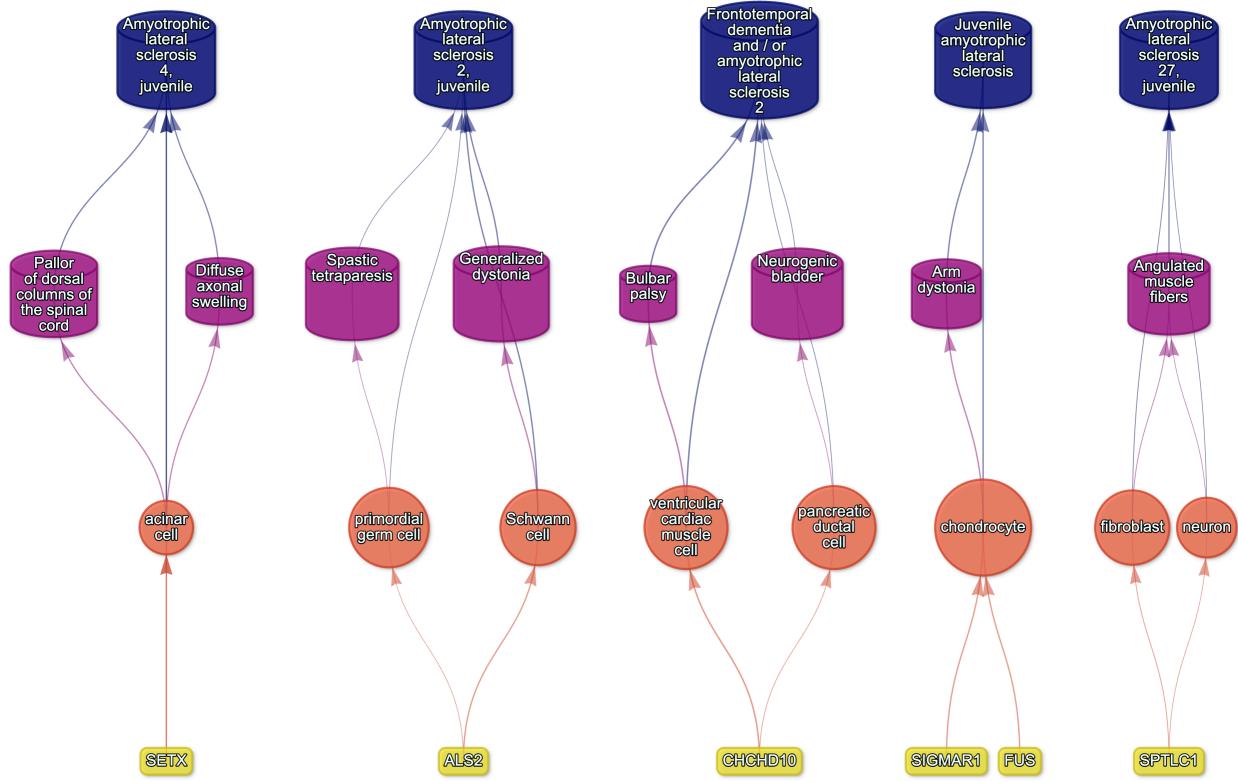


Figure 19: Amyotrophic lateral sclerosis

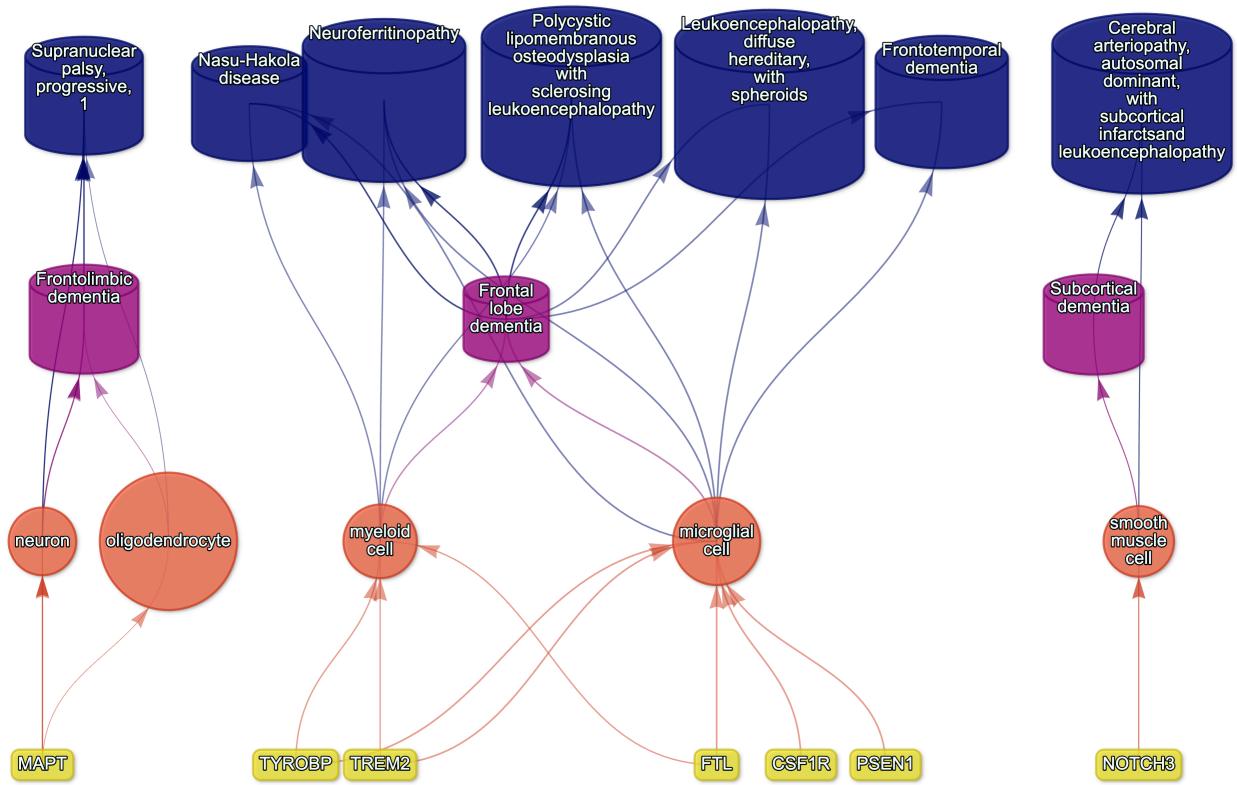


Figure 20: Dementia

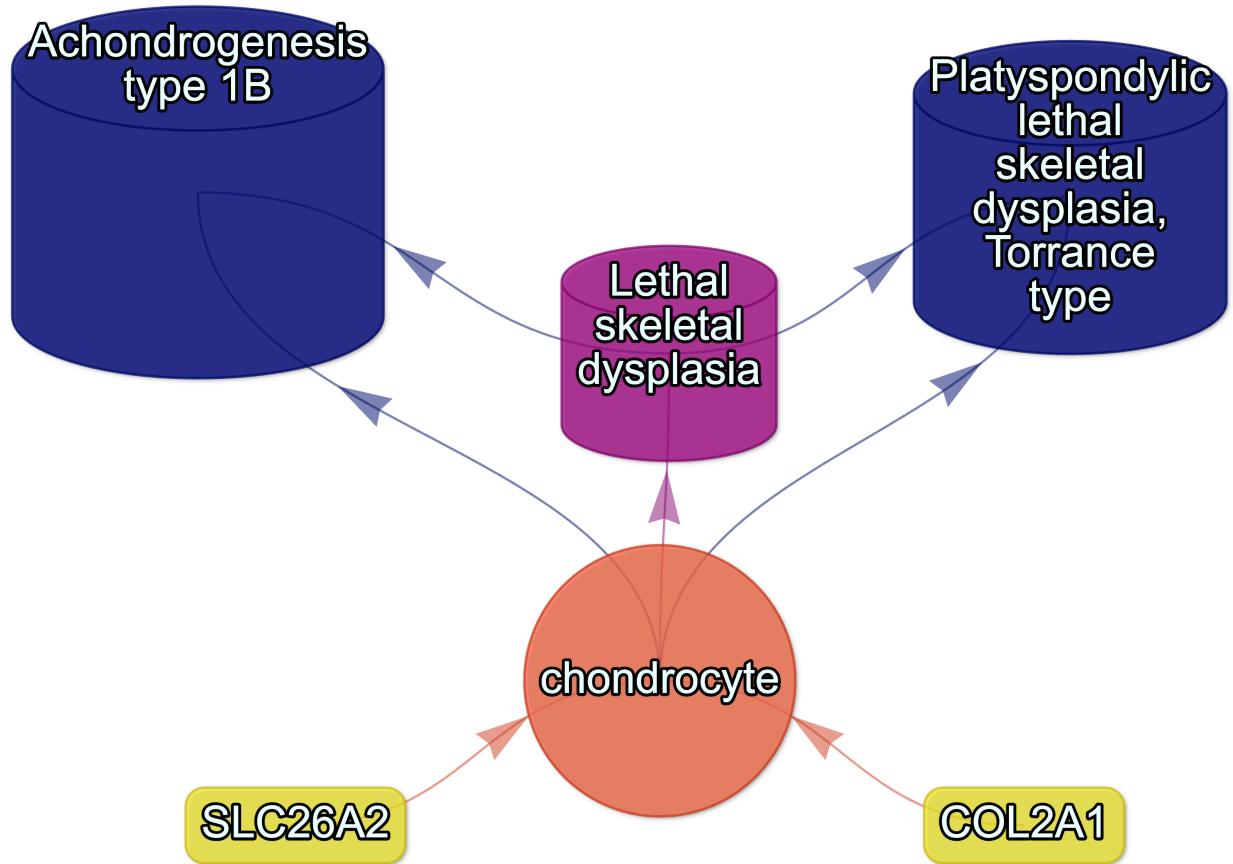


Figure 21: Lethal skeletal dysplasia

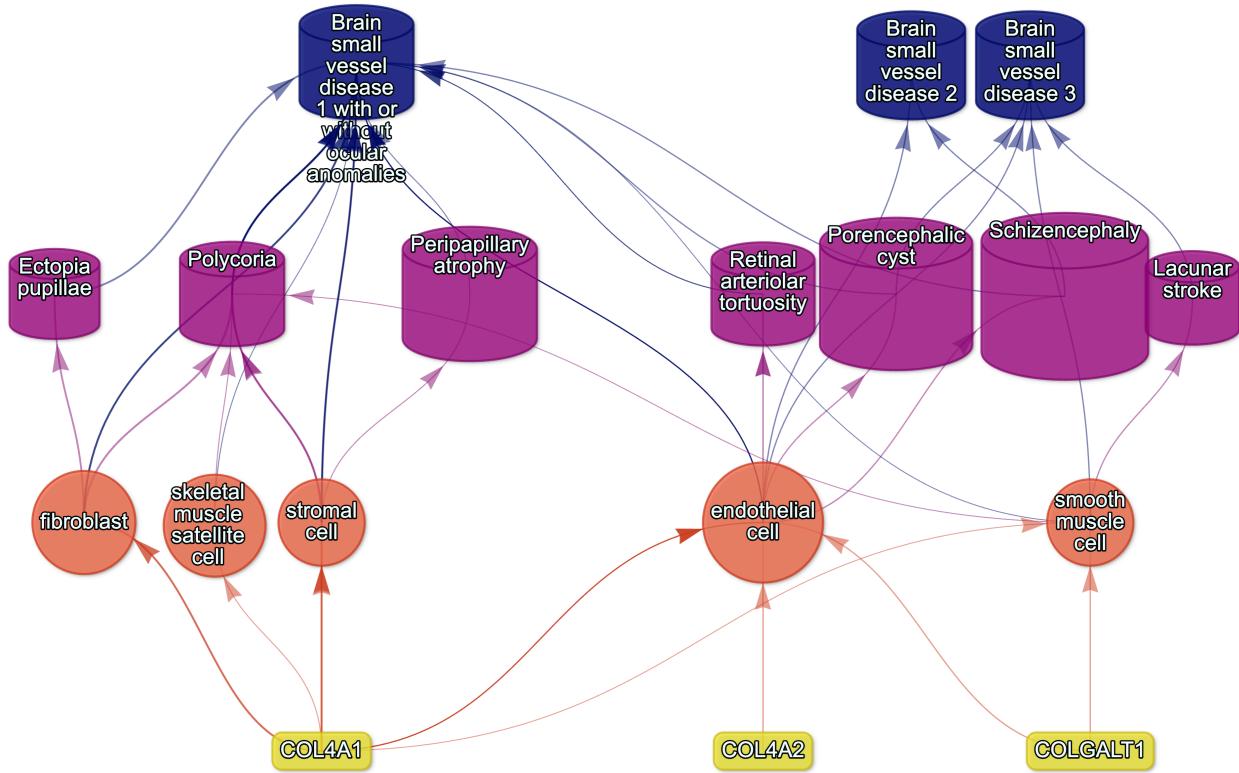


Figure 22: Small vessel disease

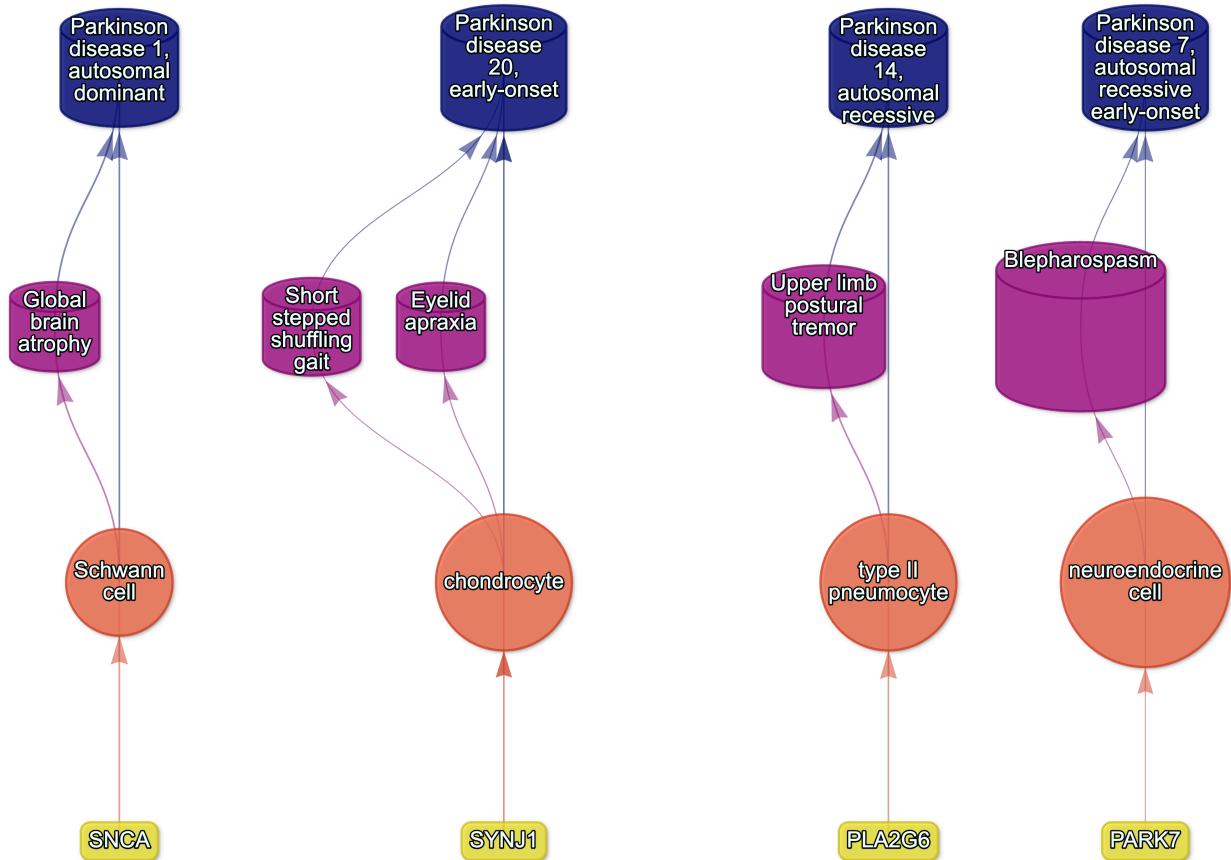


Figure 23: Parkinson's disease

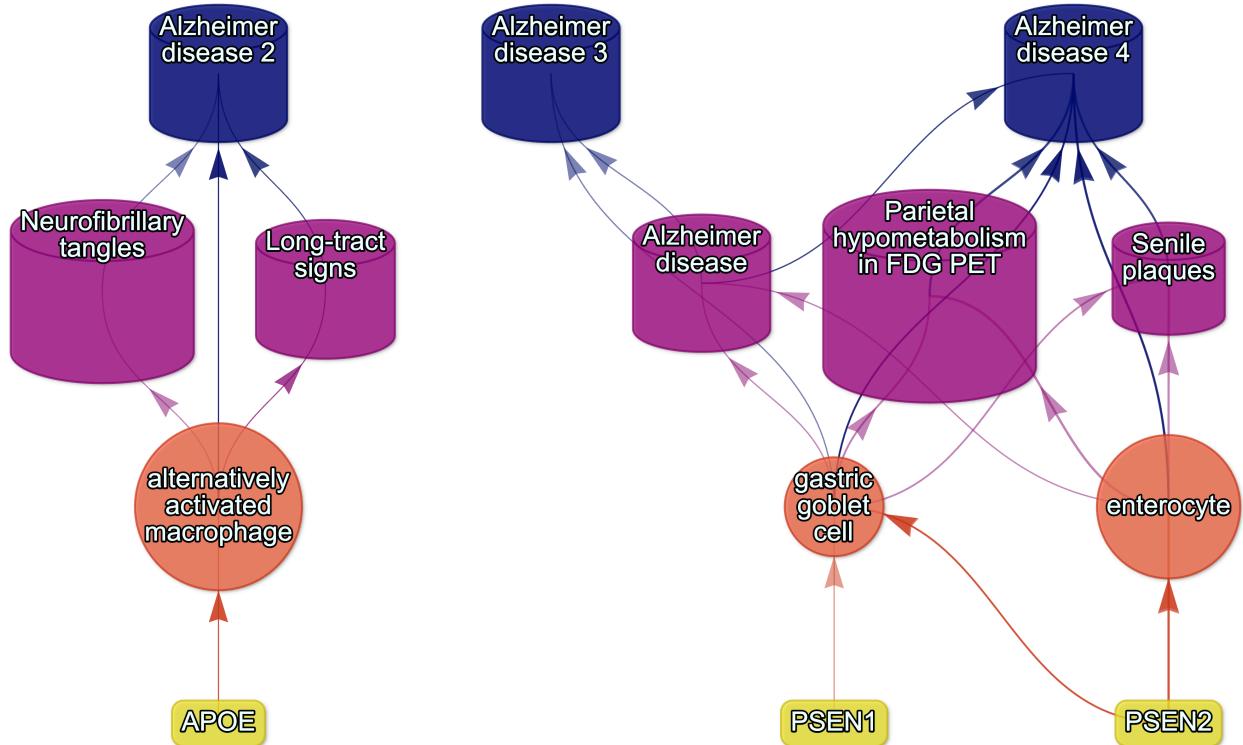


Figure 24: Alzheimer's disease

1125 **Supplementary Tables**

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 HPO ancestral branch.

hpo_branch	cl_branch	cl_name	cl_id
Abnormality of the cardiovascular system	cardiocyte	cardiac muscle cell	CL:0000746
Abnormality of the cardiovascular system	cardiocyte	regular atrial cardiac myocyte	CL:0002129
Abnormality of the cardiovascular system	cardiocyte	endocardial cell	CL:0002350
Abnormality of the cardiovascular system	cardiocyte	epicardial adipocyte	CL:1000309
Abnormality of the cardiovascular system	cardiocyte	ventricular cardiac muscle cell	CL:2000046
Abnormality of the endocrine system	endocrine cell	endocrine cell	CL:0000163
Abnormality of the endocrine system	endocrine cell	neuroendocrine cell	CL:0000165
Abnormality of the endocrine system	endocrine cell	chromaffin cell	CL:0000166
Abnormality of the eye	photoreceptor cell / retinal cell	photoreceptor cell	CL:0000210
Abnormality of the eye	photoreceptor cell / retinal cell	amacrine cell	CL:0000561
Abnormality of the eye	photoreceptor cell / retinal cell	Mueller cell	CL:0000636
Abnormality of the eye	photoreceptor cell / retinal cell	retinal pigment epithelial cell	CL:0002586
Abnormality of the immune system	leukocyte	T cell	CL:0000084
Abnormality of the immune system	leukocyte	mature neutrophil	CL:0000096
Abnormality of the immune system	leukocyte	mast cell	CL:0000097
Abnormality of the immune system	leukocyte	microglial cell	CL:0000129
Abnormality of the immune system	leukocyte	professional antigen presenting cell	CL:0000145
Abnormality of the immune system	leukocyte	macrophage	CL:0000235

Table 6: On-target cell types for each HPO ancestral branch.

hpo_branch	cl_branch	cl_name	cl_id
Abnormality of the immune system	leukocyte	B cell	CL:0000236
Abnormality of the immune system	leukocyte	dendritic cell	CL:0000451
Abnormality of the immune system	leukocyte	monocyte	CL:0000576
Abnormality of the immune system	leukocyte	plasma cell	CL:0000786
Abnormality of the immune system	leukocyte	alternatively activated macrophage	CL:0000890
Abnormality of the immune system	leukocyte	thymocyte	CL:0000893
Abnormality of the immune system	leukocyte	innate lymphoid cell	CL:0001065
Abnormality of the musculoskeletal system	cell of skeletal muscle / chondrocyte	chondrocyte	CL:0000138
Abnormality of the musculoskeletal system	cell of skeletal muscle / chondrocyte	cell of skeletal muscle	CL:0000188
Abnormality of the musculoskeletal system	cell of skeletal muscle / chondrocyte	skeletal muscle satellite cell	CL:0000594
Abnormality of the nervous system	neural cell	bipolar neuron	CL:0000103
Abnormality of the nervous system	neural cell	granule cell	CL:0000120
Abnormality of the nervous system	neural cell	Purkinje cell	CL:0000121
Abnormality of the nervous system	neural cell	glial cell	CL:0000125
Abnormality of the nervous system	neural cell	astrocyte	CL:0000127
Abnormality of the nervous system	neural cell	oligodendrocyte	CL:0000128

Table 6: On-target cell types for each HPO ancestral branch.

hpo_branch	cl_branch	cl_name	cl_id
Abnormality of the nervous system	neural cell	microglial cell	CL:0000129
Abnormality of the nervous system	neural cell	neuroendocrine cell	CL:0000165
Abnormality of the nervous system	neural cell	chromaffin cell	CL:0000166
Abnormality of the nervous system	neural cell	photoreceptor cell	CL:0000210
Abnormality of the nervous system	neural cell	inhibitory interneuron	CL:0000498
Abnormality of the nervous system	neural cell	neuron	CL:0000540
Abnormality of the nervous system	neural cell	neuronal brush cell	CL:0000555
Abnormality of the nervous system	neural cell	amacrine cell	CL:0000561
Abnormality of the nervous system	neural cell	GABAergic neuron	CL:0000617
Abnormality of the nervous system	neural cell	Mueller cell	CL:0000636
Abnormality of the nervous system	neural cell	glutamatergic neuron	CL:0000679
Abnormality of the nervous system	neural cell	retinal ganglion cell	CL:0000740
Abnormality of the nervous system	neural cell	retina horizontal cell	CL:0000745
Abnormality of the nervous system	neural cell	Schwann cell	CL:0002573
Abnormality of the nervous system	neural cell	retinal pigment epithelial cell	CL:0002586
Abnormality of the nervous system	neural cell	visceromotor neuron	CL:0005025

Table 6: On-target cell types for each HPO ancestral branch.

hpo_branch	cl_branch	cl_name	cl_id
Abnormality of the nervous system	neural cell	sympathetic neuron	CL:0011103
Abnormality of the respiratory system	respiratory epithelial cell / epithelial cell of lung	type II pneumocyte	CL:0002063
Abnormality of the respiratory system	respiratory epithelial cell / epithelial cell of lung	epithelial cell of lower respiratory tract	CL:0002632

Table 7: Encodings for Age of Death scores. Assigned numeric values for the Age of Death scores within the HPO annotations.

hpo_id	hpo_name	encoding
HP:0003826	Stillbirth	1
HP:0005268	Miscarriage	1
HP:0034241	Prenatal death	1
HP:0003811	Neonatal death	2
HP:0001522	Death in infancy	3
HP:0003819	Death in childhood	4
HP:0011421	Death in adolescence	5
HP:0100613	Death in early adulthood	6
HP:0033763	Death in adulthood	7
HP:0033764	Death in middle age	7
HP:0033765	Death in late adulthood	8

Table 8: Phenotype enrichment results. The phenotypes most biased towards associations with only the foetal versions of cell types (group=top) and those biased towards the adult versions of cell types (group=bottom). “p_adjust” is the 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 HPO term in the ontology.

group	term	name	p_adjust	log2_fold_enrichment	depth
top	HP:0005105	Abnormal nasal morphology	0.00	4.5	6
top	HP:0010938	Abnormal external nose morphology	0.00	5.4	7
top	HP:0000366	Abnormality of the nose	0.00	3.8	5
top	HP:0000055	Abnormal female external genitalia morphology	0.00	5.2	6
top	HP:0000271	Abnormality of the face	0.00	1.9	4
top	HP:0000234	Abnormality of the head	0.00	1.7	3
top	HP:0000152	Abnormality of head or neck	0.00	1.6	2
top	HP:0010460	Abnormality of the female genitalia	0.03	2.8	5
top	HP:0000811	Abnormal external genitalia	0.03	2.8	5
top	HP:0000078	Abnormality of the genital system	0.03	1.9	3
bottom	HP:0010647	Abnormal elasticity of skin	0.00	6.0	5
bottom	HP:0008067	Abnormally lax or hyperextensible skin	0.00	6.0	6
bottom	HP:0011121	Abnormal skin morphology	0.00	2.4	4
bottom	HP:0000951	Abnormality of the skin	0.00	2.1	3
bottom	HP:0001574	Abnormality of the integument	0.01	1.6	2
bottom	HP:0001626	Abnormality of the cardiovascular system	0.02	1.4	2
bottom	HP:0030680	Abnormal cardiovascular system morphology	0.02	1.7	3
bottom	HP:0025015	Abnormal vascular morphology	0.04	1.9	4
bottom	HP:0030962	Abnormal morphology of the great vessels	0.04	2.7	6

Table 9: Examples of specific phenotypes that are most biased towards associations with only the foetal versions of cell types (group=top) and those biased towards the adult versions of cell types (group=bottom).

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

Table 10: Cell type enrichment results. The cell types that most strongly show a difference between their foetal and adult versions in their phenotype associations. “p_adjust” is the 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 CL term in the ontology.

group	term	name	p_adjust	log2_fold_enrichment	depth
top	CL:0002320	connective tissue cell	0	3.2	1