

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 Hiranyamaya Dash Yichun Zhang
⁴ Robert Gordon-Smith Jai Chapman Momoko Otani Nathan G. Skene

⁵ 2025-08-07

6 Abstract

7 Rare diseases (RDs) are an extremely heterogeneous and underserved category of medical conditions. While
8 the majority of RDs are strongly genetic, it remains largely unknown via which physiological mechanisms
9 genetics cause RD. Therefore, we sought to systematically characterise the cell type-specific mechanisms
10 underlying all RD phenotypes with a known genetic cause by leveraging the Human Phenotype Ontology
11 and transcriptomic single-cell atlases of the entire human body from embryonic, foetal, and adult samples. In
12 total 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 interactive
18 web portal (<https://neurogenomics-ukdri.dsi.ic.ac.uk/>). To summarise, this work represents a significant step
19 forward in the mission to treat patients across an extremely diverse spectrum of serious RDs.

20 Introduction

21 While rare diseases (RDs) are individually uncommon, they collectively account for an enormous global
22 disease burden with over 10,000 recognised RDs affecting at least 300-400 million people globally¹ (1 in
23 10-20 people)². Over 75% of RDs primarily affect children with a 30% mortality rate by five years of age³.
24 Despite the prevalence and severity of RDs, patients suffering from these conditions are vastly underserved
25 due to several contributing factors. First, diagnosis is extremely challenging due to the highly variable
26 clinical presentations of many of these diseases. The diagnostic odyssey can take patients and their families
27 decades, with an average time to diagnosis of five years⁴. Of those, ~46% receive at least one incorrect
28 diagnosis and over 75% of all patients never receive any diagnosis⁵. Second, prognosis is also made difficult
29 by high variability in disease course and outcomes which makes matching patients with effective and timely
30 treatment plans even more challenging. Finally, even for patients who receive an accurate diagnosis/prognosis,
31 treatments are currently only available for less than 5% of all RDs⁶. In addition to the scientific challenges of
32 understanding RDs, there are strong financial disincentives for pharmaceutical and biotechnology companies
33 to develop expensive therapeutics for exceedingly small RD patient populations with little or no return
34 on investment^{7,8}. Those that have been produced are amongst the world's most expensive drugs, greatly
35 limiting patients' ability to access it^{9,10}. New high-throughput approaches for the development of rare disease
36 therapeutics could greatly reduce costs (for manufacturers and patients) and accelerate the timeline from
37 discovery to delivery.

38 A major challenge in both healthcare and scientific research is the lack of standardised medical terminology.
39 Even in the age of electronic healthcare records (EHR) much of the information about an individual's history

40 is currently fractured across healthcare providers, often with differing nomenclatures for the same conditions.
41 The Human Phenotype Ontology (HPO) is a hierarchically organised set of controlled clinical terms that
42 provides a much needed common framework by which clinicians and researchers can precisely communi-
43 cate patient conditions¹⁴. The HPO spans all domains of human physiology and currently describes 18,082
44 phenotypes across 10,300 RDs. Each phenotype and disease is assigned its own unique identifier and organ-
45 ised as a hierarchical graph, such that higher-level terms describe broad phenotypic categories or *branches*
46 (e.g. *HP:0033127*: ‘Abnormality of the musculoskeletal system’ which contains 4,495 unique phenotypes)
47 and lower-level terms describe increasingly precise phenotypes (e.g. *HP:0030675*: ‘Contracture of proximal
48 interphalangeal joints of 2nd-5th fingers’). It has already been integrated into healthcare systems and clinical
49 diagnostic tools around the world, with increasing adoption over time¹¹. Standardised frameworks like the
50 HPO also allow us to aggregate relevant knowledge about the molecular mechanisms underlying each RD.

51 Over 80% of RDs have a known genetic cause^{15,16}. Since 2008, the HPO has been continuously updated
52 using curated knowledge from the medical literature, as well as by integrating databases of expert validated
53 gene-phenotype relationships, such as OMIM^{17–19}, Orphanet^{20,21}, and DECIPHER²². Mappings between
54 HPO terms to other commonly used medical ontologies (e.g. SNOMED CT²³, UMLS^{24,25}, ICD-9/10/11²⁶)
55 make the HPO even more valuable as a clinical resource (provided in Mappings section of Methods). Many of
56 these gene annotations are manually or semi-manually curated by expert clinicians from case reports of rare
57 disease patients in which the causal gene is identified through whole exome or genome sequencing. Currently,
58 the HPO contains gene annotations for 11,047 phenotypes across 8,631 diseases. Yet genes alone do not tell
59 the full story of how RDs come to be, as their expression and functional relevance varies drastically across
60 the multitude of tissues and cell types contained within the human body. Our knowledge of the physiological
61 mechanisms via which genetics cause pathogenesis is lacking for most RDs, severely hindering our ability to
62 effectively diagnose, prognose and treat RD patients.

63 Our knowledge of cell type-specific biology has exploded over the course of the last decade and a half,
64 with numerous applications in both scientific and clinical practices^{27–29}. In particular, single-cell RNA-seq
65 (scRNA-seq) has allowed us to quantify the expression of every gene (i.e. the transcriptome) in individual
66 cells. More recently, comprehensive single-cell transcriptomic atlases across tissues have also emerged^{30,31}.
67 In particular, the Descartes Human³² and Human Cell Landscape³³ projects provide comprehensive multi-
68 system scRNA-seq atlases in embryonic, foetal, and adult human samples from across the human body.
69 These datasets provide data-driven gene signatures for hundreds of cell subtypes. Given that many disease-
70 associated genes are expressed in some cell types but not others, we can infer that disruptions to these genes
71 will have varying impact across cell types. By comparing the aggregated disease gene annotations with
72 cell type-specific expression profiles, we can therefore uncover the cell types and tissues via which diseases
73 mediate their effects.

74 Here, we combine and extend several of the most comprehensive genomic and transcriptomic resources

75 currently available to systematically uncover the cell types underlying granular phenotypes across 8,628
76 diseases Fig. 1. Conversely, this approach also allows us to better understand the roles of understudied cell
77 types by observing which phenotypes they tend to associate with. For example, the original authors proposed
78 that a novel class *AFB+/ALB+* cells may represent hepatoblasts circulating through the bloodstream during
79 foetal development³⁴. Our results support this hypothesis as *AFB+/ALB+* cells were significantly associated
80 with 12 liver-related phenotypes, as well as 58 blood-related phenotypes.

81 Beyond making discoveries in basic science, our phenome-wide cell type associations provide essential context
82 for the development of novel therapeutics, especially gene therapy modalities such as adeno-associated viral
83 (AAV) vectors in which advancement have been made in their ability selectively target specific cell types^{35,36}.
84 Precise knowledge of relevant cell types and tissues causing the disease can improve safety by minimising
85 harmful side effects in off-target cell types and tissues. It can also enhance efficacy by efficiently delivering
86 expensive therapeutic payloads to on-target cell types and tissues. For example, if a phenotype primarily
87 effects retinal cells, then the gene therapy would be optimised for delivery to retinal cells of the eye. Using
88 this information, we developed a high-throughput pipeline for comprehensively nominating cell type-resolved
89 gene therapy targets across thousands of RD phenotypes. As a prioritisation tool, we sorted these targets
90 based on the severity of their respective phenotypes, using a generative AI-based approach³⁷. Together,
91 our study dramatically expands the available knowledge of the cell types, organ systems and life stages
92 underlying RD phenotypes.

93 Results

94 Phenotype-cell type associations

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

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



Evidence for Gene 1 causing Phenotype A

	Weight	Studies	=	Score
No Known	0	x 0	=	0
Refuted	0	x 0	=	0
Disputed	1	x 1	=	1
Limited	2	x 0	=	0
Supportive	3	x 2	=	6
Moderate	4	x 1	=	4
Strong	5	x 3	=	15
Definitive	6	x 1	=	6
Total			=	32

Phenotype x gene evidence score matrix

	Phenotype A	Phenotype B	Phenotype C	...
Gene 1	32	0	1	...
Gene 2	0	16	0	...
Gene 3	2	12	10	...
...

Descartes Human



Human Cell Landscape



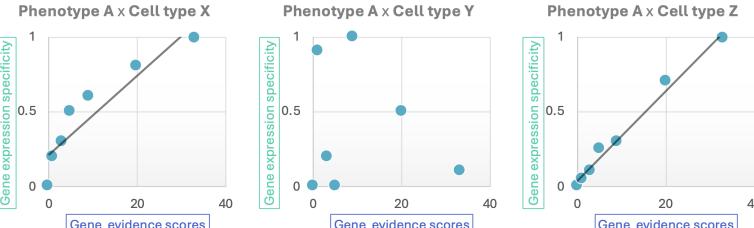
Gene expression specificity in Cell type Z

	Cell 1	Cell 2	Cell 3	Mean	Specificity
Gene 1	0	0	0	0	0
Gene 2	0	1	0	0.33	0.01
Gene 3	10	8	12	10	0.60
...

Cell type x gene expression specificity matrix

	Cell type X	Cell type Y	Cell type Z	...
Gene 1	0.50	0	0	...
Gene 2	0	0.95	0.01	...
Gene 3	0	0.25	0.60	...
...

Generalised Linear Regression Tests



Phenotype-cell type association results

Phenotype	Cell type	P-value	FDR	Z-score
A	X	0.005	0.05	0.25
A	Y	0.98	1	0
A	Z	0.001	0.01	0.90
B	X	1	1	0
B	Y	0.0004	0.004	0.75
B	Z	1	1	0.01
C	X	0.003	0.03	0.20
C	Y	1	1	0
C	Z	0.0007	0.007	0.98
...

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 control the false discovery rate (FDR) across all tests.

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

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

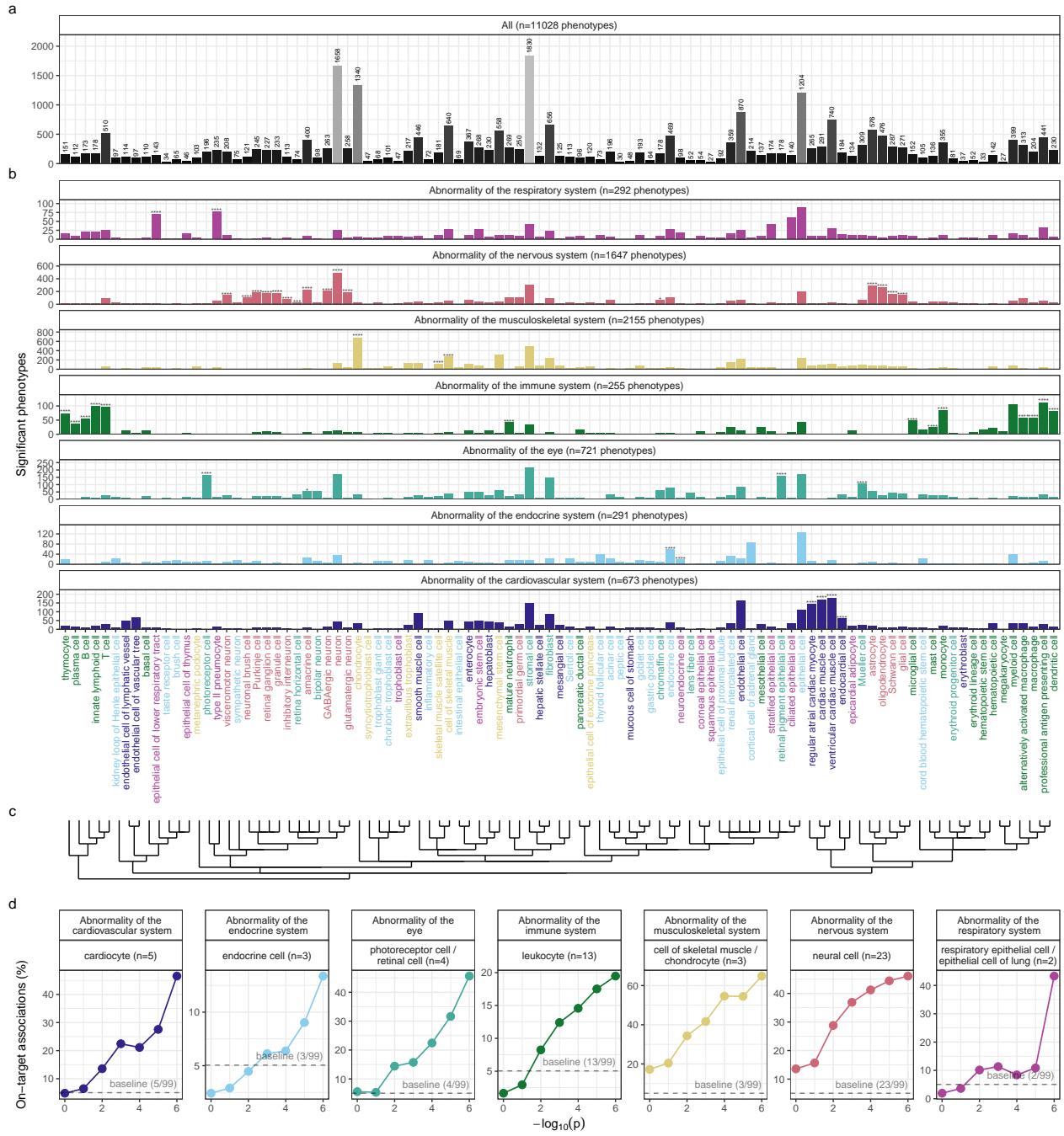
123 Validation of expected phenotype-cell type relationships

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

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

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

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

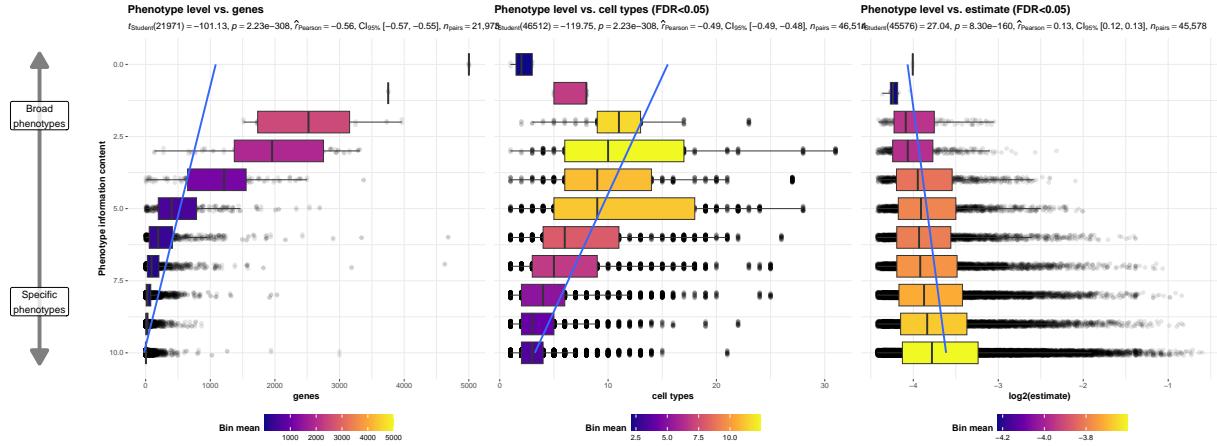
151 **Validation of inter- and intra-dataset consistency**

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

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

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

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

176 Validation of phenotype-cell type associations using biomedical knowledge graphs

177 In order to validate our phenotype-cell type associations without the bias introduced by manually searching
 178 literature that affirmed our discoveries, we use formalised biomedical knowledge from the scientific community
 179 stored in a knowledge graph. In particular, the Monarch Knowledge Graph (MKG) is a comprehensive,
 180 standardised database that aggregates up-to-date knowledge about biomedical concepts and the relationships
 181 between them. This currently includes 103 well-established phenotype-cell type relationships⁴⁰. We used
 182 the MKG as a proxy for the field’s current state of knowledge of causal phenotype-cell type associations.
 183 We evaluated the proportion of MKG associations that were recapitulated by our results Fig. 12. For
 184 each phenotype-cell type association in the MKG, we computed the percent of cell types recovered in our
 185 association results at a given ontological distance according to the CL ontology. An ontological distance of 0
 186 means that our nominated cell type was as close as possible to the MKG cell type after adjusting for the cell
 187 types available in our single-cell references. Instances of exact overlap of terms between the MKG and our
 188 results would qualify as an ontological distance of 0 (e.g. ‘monocyte’ vs. ‘monocyte’). Greater ontological
 189 distances indicate further divergence between the MKG cell type and our nominated cell type. A distance
 190 of 1 indicating that the MKG cell type was one step away from our nominated cell type in the CL ontology
 191 graph (e.g. ‘monocyte’ vs. ‘classical monocyte’). The maximum possible percent of recovered terms is capped
 192 by the percentage of MKG ground-truth phenotypes we were able to find at least one significant cell type
 193 association for at FDR_{pc} .

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

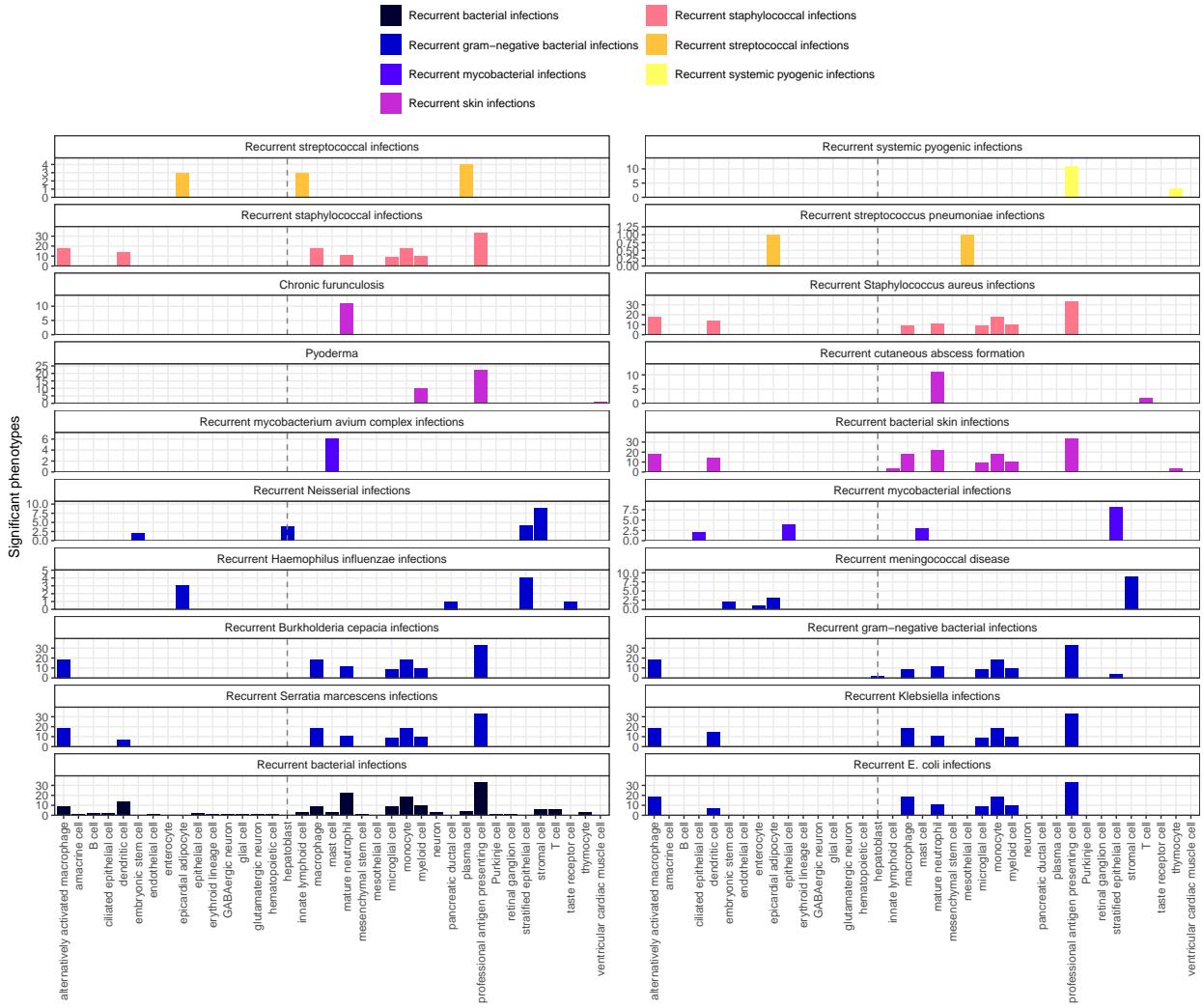
scribed in the MKG. Of these phenotypes, we captured 57% of the MKG phenotype-cell associations at an ontological distance of 0 (i.e. the closest possible Cell Ontology term match). Recall increased with greater flexibility in the matching of cell type annotations. At an ontological distance of 1 (e.g.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

Phenome-wide analyses discover novel phenotype-cell type associations

Having established that many of the phenotype-cell type associations align with prior expectations, we then sought to discover novel relationships with undercharacterised phenotypes. We reasoned that recurrent bacterial infections (and all its descendant phenotypes) should primarily be associated with immune cell types. The HPO term ‘Recurrent bacterial infections’ has 19 different descendant phenotypes, e.g. staphylococcal, streptococcal, and Neisserial infections. Each of these phenotypes are associated with partially overlapping subsets of immune cells and other cell types (Fig. 4). As expected, these phenotypes are primarily associated with immune cell types (e.g. macrophages, dendritic cells, T cells, monocytes, neutrophils). Some associations confirm relationships previously suggested in the literature, such as that between ‘Recurrent staphylococcal infections’ and myeloid cells^{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.

²²⁹ Also of note, despite the fact that our datasets contain both hepatoblasts and their mature counterpart, hepa-
²³⁰ tocytes, only the hepatoblasts showed this association. This suggests that the genetic factors that predispose
²³¹ individuals for risk of Neisserial infections are specifically affecting hepatoblasts before they become fully
²³² differentiated. It is also notable that these phenotypes were the only ones within the ‘Recurrent bacterial
²³³ infections’ branch, or even the broader ‘Recurrent infections’ branch, perhaps indicating a unique role for
²³⁴ hepatoblasts in recurrent infectious disease. The only phenotypes within the even broader ‘Abnormality of
²³⁵ the immune system’ HPO branch that significantly associated with mature hepatocytes were ‘Pancreatitis’
²³⁶ ($FDR=2.1 \times 10^{-2}$, $\beta=5.3 \times 10^{-2}$) and ‘Susceptibility to chickenpox’ ($FDR=1.2 \times 10^{-2}$, $\beta=5.5 \times 10^{-2}$) both
²³⁷ 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

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

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

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

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

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

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

273 **Prioritising phenotypes based on severity**

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

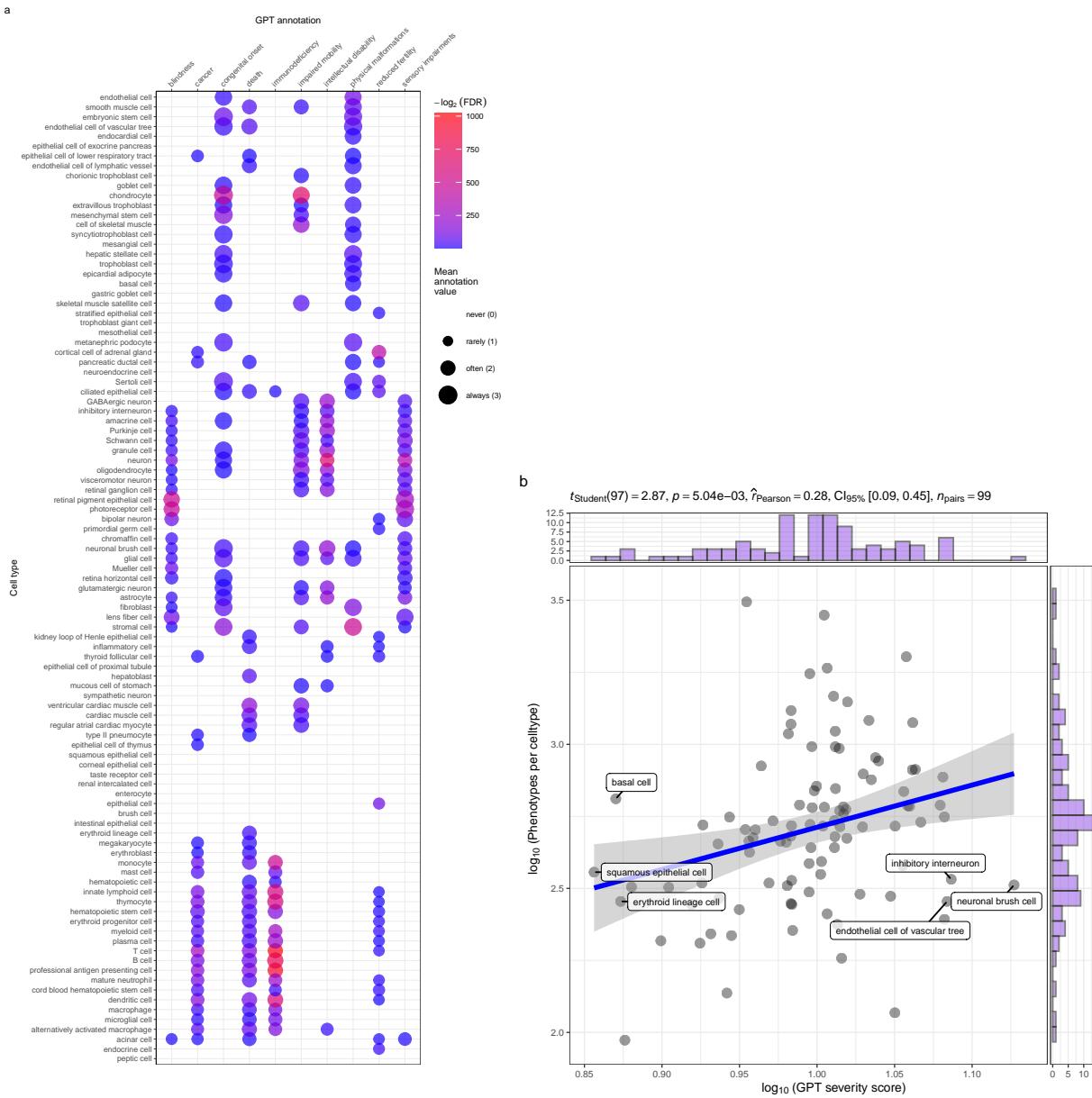
279 Severity annotations were gathered from GPT-4 for 16,982/18,082 (94%) HPO phenotypes in our companion
280 study³⁷. Benchmarking tests of these results using ground-truth HPO branch annotations. For example,
281 phenotypes within the ‘Blindness’ HPO branch (*HP:0000618*) were correctly annotated as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

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

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

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

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



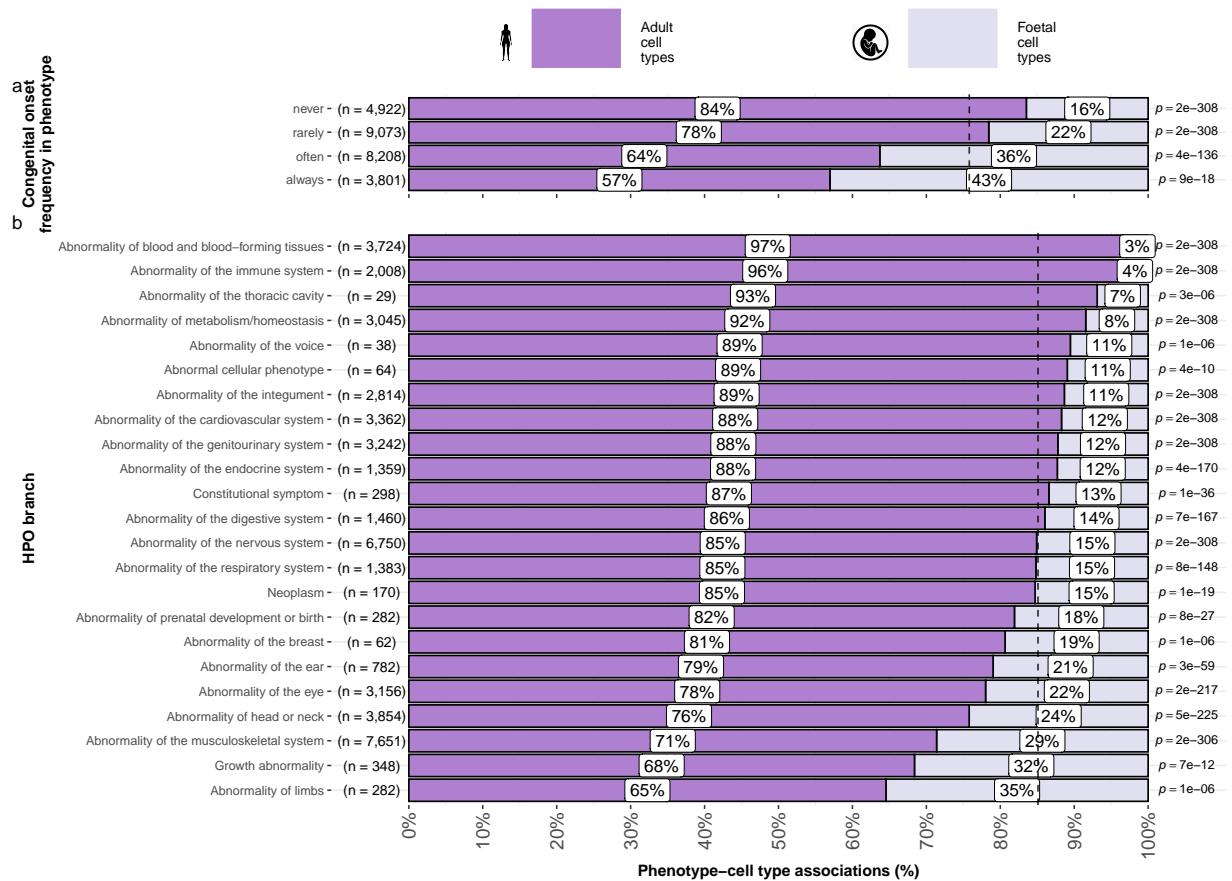
(a) Genetic disruptions to some cell types cause more clinically severe phenotypes than others. **a**, Different cell types are associated with different aspects of phenotypic severity. The dot plot shows the mean encoded frequency value for a given severity annotation (0="never", 1="rarely", 2="often",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

Figure 5

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

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

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



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

Figure 6

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

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

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

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

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

365 Specific examples of such phenotypes include ‘Excessive wrinkled skin’ and ‘Paroxysmal supraventricular
366 tachycardia’ Table 8. It is well known that ageing naturally causes a loss of skin elasticity (due to decreasing
367 collagen production) and vascular degeneration⁵⁶. Next, we were interested whether some cell types tend to
368 show strong differences in their phenotype associations between their foetal and adult forms. To test this, we
369 performed an analogous enrichment procedure as with the phenotypes, except using Cell Ontology terms and
370 the Cell Ontology graph. This analysis identified the cell type category connective tissue cell ($p=1.8 \times 10^{-3}$,
371 $\log_2(\text{fold-change})=3.2$) as the most foetal-biased cell type. No cell type categories were significantly enriched
372 for the most adult-biased cell types. This is likely due to the fact that cell types can be disrupted at different
373 stages of life, resulting in different phenotypes. Thus there the same cell types may be involved in both
374 the most foetal-biased and adult-biased phenotypes. Together, these findings serve to further validate our
375 methodology as a tool for identifying the causal cell types underlying a wide range of phenotypes.

376 Therapeutic target identification

377 In the above sections, we demonstrated how gene association databases can be used to investigate the cell
378 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^{57–59}. 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 clinicians, biomedical

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

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

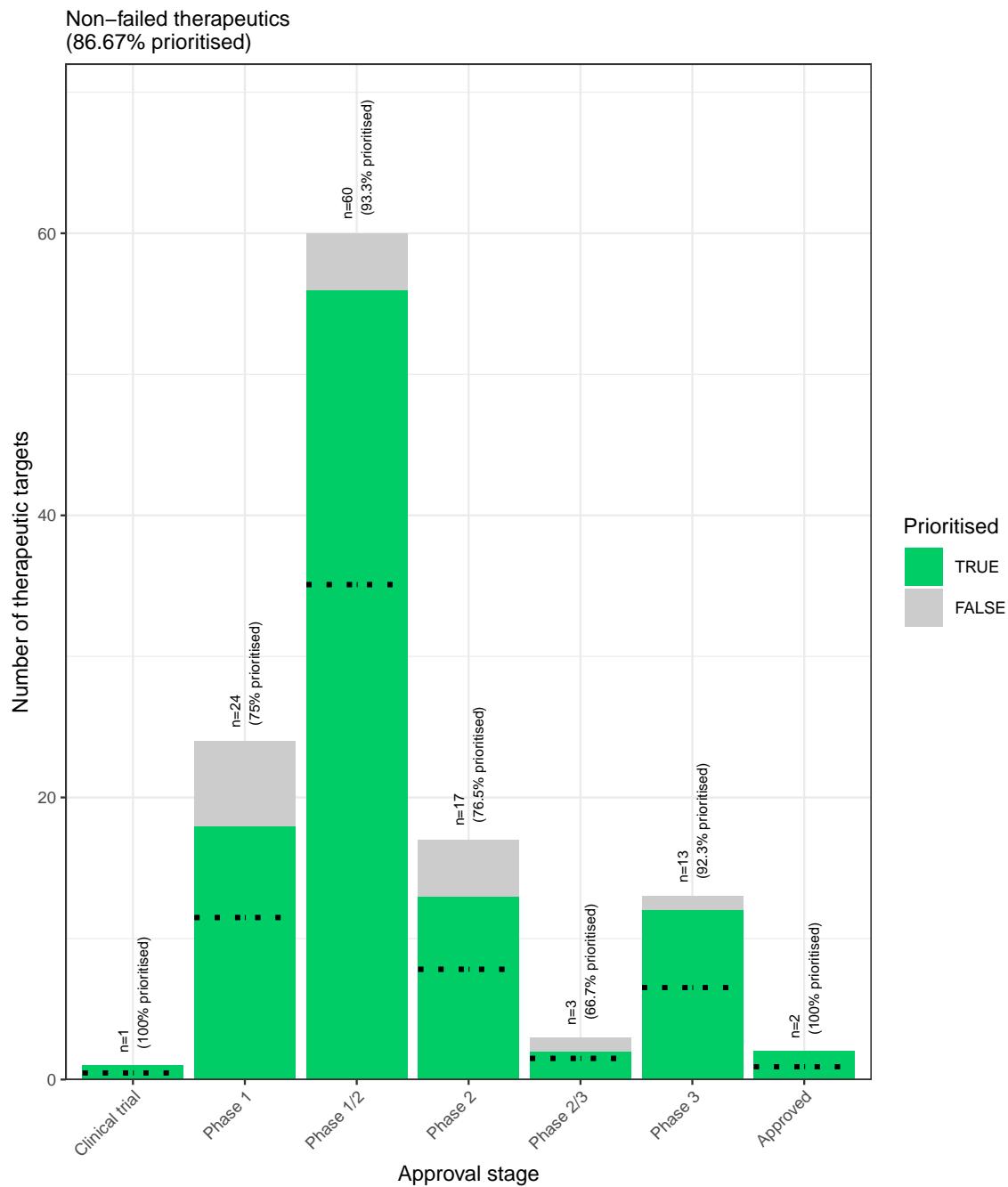
423 Therapeutic target validation

424 To determine whether the genes prioritised by our therapeutic targets pipeline were plausible, we checked
425 what percentage of gene therapy targets we recapitulated. Data on therapeutic approval status was gathered
426 from the Therapeutic Target Database (TTD; release 2025-08-07)⁶⁰. Overall, we prioritised 87% (120 total)
427 of all non-failed existing gene therapy targets (ie. those which are currently approved, investigative, or
428 undergoing clinical trials). A hypergeometric test confirmed that our prioritised targets were significantly
429 enriched for non-failed gene therapy targets ($p = 1.8 \times 10^{-5}$). For these hypergeometric tests, the background
430 gene set was composed of the union of all phenotype-associated genes in the HPO and all gene therapy
431 targets listed in TTD.

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

437 In addition to aggregate enrichment results, we also provide specific examples of successful gene therapies
438 whose cell type-specific mechanism were recapitulated by our phenotype-cell associations. In particular, our
439 pipeline nominated the gene *RPE65* within ‘retinal pigment epithelial cells’ as the top target for ‘Fundus
440 atrophy’ vision-related phenotypes that are hallmarks of ‘Leber congenital amaurosis, type II’ and ‘Se-
441 vere early-childhood-onset retinal dystrophy’. Indeed, gene therapies targeting *RPE65* within the retina of
442 patients with these rare genetic conditions are some of the most successful clinical applications of this tech-
443 nology to date, able to restore vision in many cases⁶¹. In other cases, a tissue (e.g. liver) may be known to
444 be causally involved in disease genesis, but the precise causal cell types within that tissue remain unknown
445 (e.g. hepatocytes, Kupffer cells, Cholangiocytes, Hepatic stellate cells, Natural killer cells, etc.). Tissue-level
446 investigations (e.g. using bulk transcriptomics or epigenomics) would be dominated by hepatocytes, which
447 comprise 75% of the liver. Our prioritized gene therapy targets can aid in such scenarios by providing the

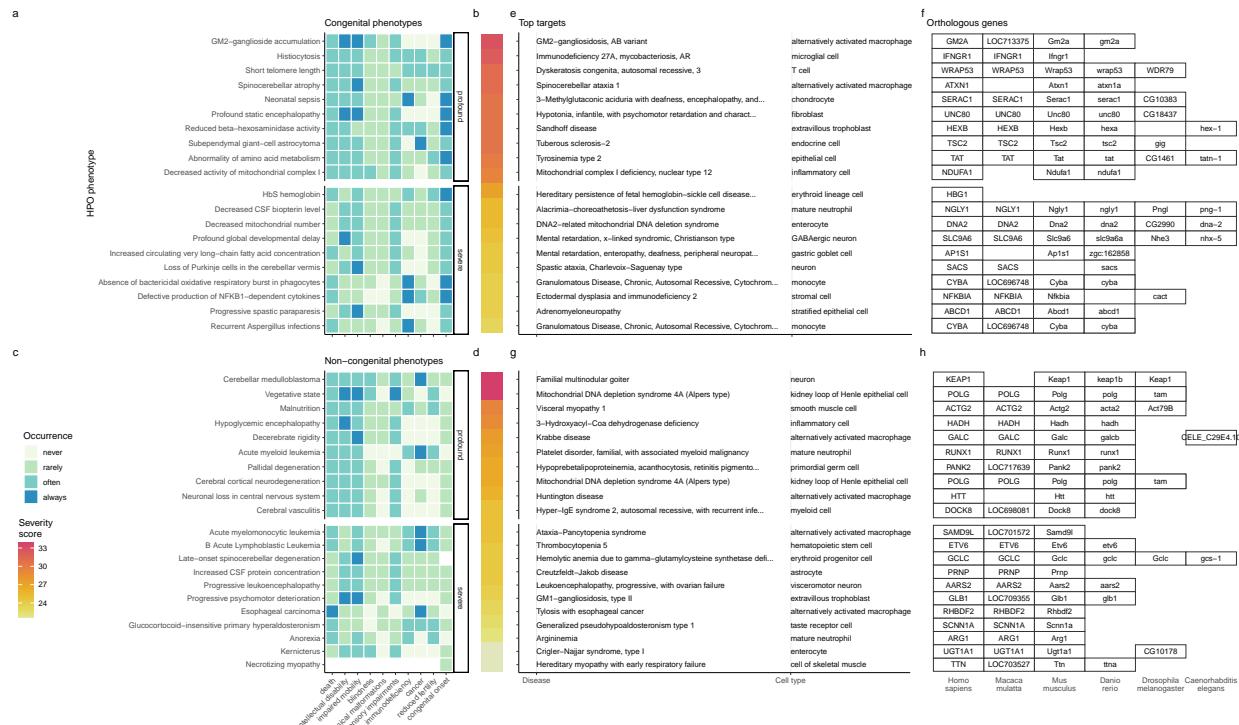
⁴⁴⁸ cell type-resolution context most likely to be causal for a given phenotype or set of phenotypes.



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

Figure 7

449 **Selected example targets**



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

Figure 8

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

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

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

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

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

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

482 Experimental model translatability

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

493 **Mappings**

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

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

502 **Discussion**

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

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

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

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

Unfortunately, there are currently only treatments available for less than 5% of RDs⁶. Novel technologies including CRISPR, prime editing, antisense oligonucleotides, viral vectors, and/or lipid nanoparticles, have been undergone significant advances in the last several years^{69–73} and proven remarkable clinical success in an increasing number of clinical applications^{74–77}. The U.S. Food and Drug Administration (FDA) recently announced an landmark program aimed towards improving the international regulatory framework to take advantage of the evolving gene/cell therapy technologies⁷⁸ with the aim of bringing dozens more therapies to patients in a substantially shorter timeframe than traditional pharmaceutical product development (typically 5–20 years with a median of 8.3 years)⁷⁹. While these technologies have the potential to revolutionise RD medicine, their successful application is dependent on first understanding the mechanisms causing each disease.

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

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

565 From our target prioritisation pipeline results, we highlight cell type-specific mechanisms for ‘GM2-
566 ganglioside accumulation’ in Tay-Sachs disease, spinocerebellar atrophy in spinocerebellar ataxia, and
567 ‘Neuronal loss in central nervous system’ in a variety of diseases (Fig. 8). Of interest, all three of these
568 neurodegenerative phenotypes involved alternatively activated (M2) macrophages. The role of macrophages
569 in neurodegeneration is complex, with both neuroprotective and neurotoxic functions, including the
570 clearance of misfolded proteins, the regulation of the blood-brain barrier, and the modulation of the immune
571 response⁸². We also recapitulated prior evidence that microglia, the resident macrophages of the nervous
572 system, are causally implicated in Alzheimer’s disease (AD) (Fig. 19)⁸³. An important contribution of our
573 current study is that we were able to pinpoint the specific phenotypes of AD caused by macrophages to
574 neurofibrillary tangles and long-tract signs (reflexes that indicate the functioning of spinal long fiber tracts).
575 Other AD-associated phenotypes were caused by other cell types (e.g. gastric goblet cells, enterocytes).

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

596 to better distinguish between signatures of highly similar cell types/subtypes⁸⁸.

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

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

614 Conclusions

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

623 Methods

624 Human Phenotype Ontology

625 The latest version of the HPO (release 2024-02-08) was downloaded from the EMBL-EBI Ontology Lookup
626 Service⁹² and imported into R using the `HPOExplorer` package. This R object was used to extract ontolog-
627 ical relationships between phenotypes as well as to assign absolute and relative ontological levels to each
628 phenotype. The latest version of the HPO phenotype-to-gene mappings and phenotype annotations were

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

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

642 Let us denote:

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

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

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

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

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

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

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

657

658

659

Tensor of Phenotype-by-Gene
evidence scores

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

660

661

662

Tensor of Disease-by-Gene
evidence scores

Disease-by-Gene-by-Phenotype
relationships

The diagram shows two input tensors: 'Tensor of Disease-by-Gene evidence scores' and 'Tensor of Phenotype-by-Gene evidence scores'. Arrows point from these tensors to the formula for calculating L_{ij} . The formula uses the sum of products of Disease (D_i), Gene (G_j), and Phenotype (P_k) scores, divided by the number of phenotypes (f), if the gene is associated with the disease (i.e., $D_i G_j \in A$). Otherwise, it is set to 1. A bracket labeled 'Disease-by-Gene-by-Phenotype relationships' spans the condition $D_i G_j \in A$.

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

664

665

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

667 Single-cell transcriptomic atlases

668 In this study, the gene by cell type specificity matrix was constructed using the Descartes Human transcriptome
 669 atlas of foetal gene expression, which contains a mixture of single-nucleus and single-cell RNA-seq
 670 data (collected with sci-RNA-seq3)³². This dataset contains 377,456 cells representing 77 distinct cell types
 671 across 15 tissues. All 121 human foetal samples ranged from 72 to 129 days in estimated postconceptual age.
 672 To independently replicate our findings, we also used the Human Cell Landscape which contains single-cell
 673 transcriptomic data (collected with microwell-seq) from embryonic, foetal, and adult human samples across
 674 49 tissues³³.

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

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

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

688

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

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

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

691

692

693 **Phenotype-cell type associations**

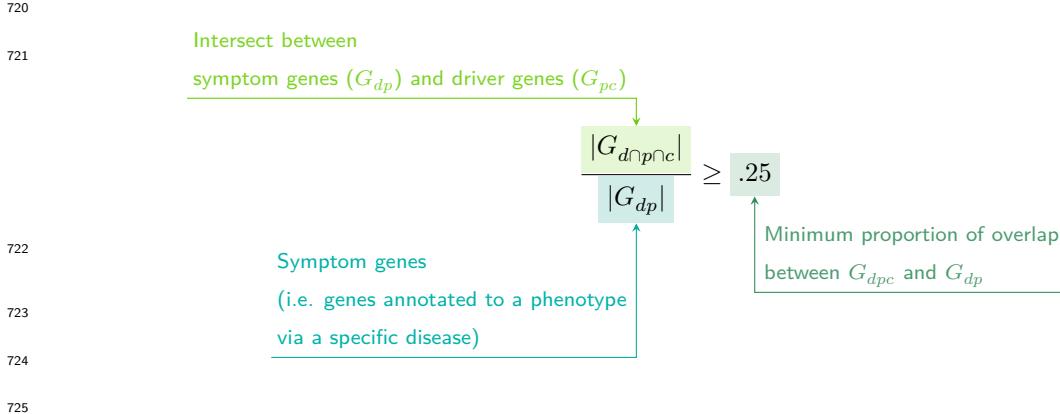
694 To test for relationships between each pairwise combination of phenotype (n=11,047) and cell type (n=201)
 695 we ran a series of univariate generalised linear models implemented via the `stats::glm` function in R. First,
 696 we filtered the gene-by-phenotype evidence score matrix (L_{ij}) and the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

705 We repeated the aforementioned procedure separately for each of the single-cell references. Once all results
 706 were generated using both cell type references (2,206,994 association tests total), we applied Benjamini-
 707 Hochberg false discovery rate⁹³ (denoted as FDR_{pc}) to account for multiple testing. Of note, we applied
 708 this correction across all results at once (as opposed to each single-cell reference separately) to ensure the
 709 FDR_{pc} was stringently controlled for across all tests performed in this study.

710 **Symptom-cell type associations**

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

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



726 Validation of expected phenotype-cell type relationships

727 We first sought to confirm that our tests (across both single-cell references) were able to recover expected
 728 phenotype-cell type relationships across seven high-level branches within the HPO (Fig. 2), including ab-
 729 normalities of the cardiovascular system, endocrine system, eye, immune system, musculoskeletal system,
 730 nervous system, and respiratory system. Within each branch the number of significant tests in a given
 731 cell type were plotted (Fig. 2b). Mappings between freeform annotations (the level at which we performed
 732 our phenotype-cell type association tests) provided by the original atlas authors and their closest CL term
 733 equivalents were provided by CellxGene³⁰. CL terms along the *x-axis* of Fig. 2b were assigned colours corre-
 734 sponding to which HPO branch showed the greatest number of enrichments (after normalising within each
 735 branch to account for differences in scale). The normalised colouring allows readers to quickly assess which
 736 HPO branch was most often associated with each cell type, while accounting for differences in the number
 737 of phenotypes across branches. We then ran a series of Analysis of Variance (ANOVA) tests to determine
 738 whether (within a given branch) a given cell type was more often enriched ($FDR < 0.05$) within that branch
 739 relative to all of the other HPO branches of an equivalent level in the ontology (including all branches not
 740 shown in Fig. 2b). After applying Benjamini-Hochberg multiple testing correction⁹³ (denoted as $FDR_{b,c}$),
 741 we annotated each respective branch-by-cell type bar according to the significance (**** : $FDR_{b,c} < 1e-04$,
 742 *** : $FDR_{b,c} < 0.001$, ** : $FDR_{b,c} < 0.01$, * : $FDR_{b,c} < 0.05$). Cell types in Fig. 2a-b were ordered along
 743 the *x-axis* according to a dendrogram derived from the CL ontology (Fig. 2c), which provides ground-truth

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

756 Validation of inter- and intra-dataset consistency

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

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

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

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

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

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

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

783 Monarch Knowledge Graph recall

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

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

797 Prioritising phenotypes based on severity

798 Only a small fraction of the the phenotypes in HPO (<1%) have metadata annotations containing informa-
799 tion on their time course, consequences, and severity. This is due to the time-consuming nature of manually
800 annotating thousands of phenotypes. To generate such annotations at scale, we previously used Generative
801 Pre-trained Transformer 4 (GPT-4), a large language model (LLM) as implemented within OpenAI’s Appli-
802 cation Programming Interface (API)³⁷. After extensive prompt engineering and ground-truth benchmarking,
803 we were able to acquire annotations on how often each phenotype directly causes intellectual disability, death,
804 impaired mobility, physical malformations, blindness, sensory impairments, immunodeficiency, cancer, re-
805 duced fertility, or is associated with a congenital onset. These criteria were previously defined in surveys
806 of medical experts as a means of systematically assessing phenotype severity⁹⁴. Responses for each metric

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

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

818 Let us denote:

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

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

Normalised Severity Score
for each phenotype

Sum of weighted annotation values
across all metrics

Numerically encoded annotation value
of metric j for phenotype p

Weight for metric j

Theoretical maximum severity score

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

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

845 The GPT-4 annotations also enabled us to assess whether foetal cell types were more often significantly
846 associated with congenital phenotypes in our Human Cell Landscape results as this single-cell reference
847 contained both adult and foetal versions of cell types (Fig. 6). To do this, we performed a chi-squared (χ^2)
848 test on the proportion of significantly associated cell types containing any of the substrings ‘fetal’, ‘fetus’,
849 ‘primordial’, ‘hESC’ or ‘embryonic’ (within cell types annotations from the original Human Cell Landscape
850 authors³³) vs. those associated without, stratified by how often the corresponding phenotype had a congenital
851 onset according to the GPT phenotype annotations (including ‘never’, ‘rarely’, ‘often’, ‘always’). In addition,
852 a series of χ^2 tests were performed within each congenital onset frequency strata, to determine whether the
853 observed proportion of foetal cell types vs. non-foetal cell types significantly deviated from the proportions
854 expected by chance.

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

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

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

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

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

878 Therapeutic target identification

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

888 Therapeutic target validation

889 To assess whether our prioritised therapeutic targets were likely to be viable, we computed the overlap
890 between our gene targets and those of existing gene therapies at various stages of clinical development
891 (Fig. 7). Gene targets were obtained for each therapy from the Therapeutic Target Database (TTD; release
892 2025-08-07) and mapped onto standardised HUGO Gene Nomenclature Committee (HGNC) gene symbols
893 using the `orthogene` R package. We stratified our overlap metrics according to whether the therapies had
894 failed (unsuccessful clinical trials or withdrawn), or were non-failed (successful or ongoing clinical trials).
895 We then conducted hypergeometric tests to determine whether the observed overlap between our prioritised
896 targets and the non-failed therapy targets was significantly greater than expected by chance (i.e. enrichment).
897 We also conducted a second hypergeometric test to determine whether the observed overlap between our
898 prioritised targets and the failed therapy targets was significantly less than expected by chance (i.e. depletion).
899 Finally, we repeated the analysis against all therapeutic targets, not just those of gene therapies, to determine
900 whether our prioritised targets had relevance to other therapeutic modalities.

901 **Experimental model translatability**

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

911 **Novel R packages**

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

920 **Rare Disease Celltyping Portal**

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

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

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

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

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

Mappings

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

Data Availability

All data is publicly available through the following resources:

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

Code Availability

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

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

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

971 Acknowledgements

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

975 Funding

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

979 References

- 980 1. Ferreira, C. R. The burden of rare diseases. *Am. J. Med. Genet. A* **179**, 885–892 (2019).
- 981 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).
- 982 3. Rare diseases BioResource.
- 983 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).
- 984 5. Molster, C. *et al.* Survey of healthcare experiences of australian adults living with rare diseases.
Orphanet J. Rare Dis. **11**, 30 (2016).
- 985 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).
- 986 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).
- 987 8. Yates, N. & Hinkel, J. The economics of moonshots: Value in rare disease drug development. *Clin.
Transl. Sci.* **15**, 809–812 (2022).

- 988 9. Nuijten, M. Pricing zolgensma - the world's most expensive drug. *J Mark Access Health Policy* **10**,
2022353 (2022).
- 989 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).
- 990 11. Gargano, M. A. *et al.* The human phenotype ontology in 2024: Phenotypes around the world. *Nucleic
Acids Res.* **52**, D1333–D1346 (2024).
- 991 12. Köhler, S. *et al.* Expansion of the human phenotype ontology (HPO) knowledge base and resources.
Nucleic Acids Res. **47**, D1018–D1027 (2019).
- 992 13. Köhler, S. *et al.* The human phenotype ontology in 2021. *Nucleic Acids Res.* **49**, D1207–D1217
(2021).
- 993 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).
- 994 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).
- 995 16. Rare diseases, common challenges. *Nat. Genet.* **54**, 215 (2022).
- 996 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).
- 997 18. Amberger, J. S. & Hamosh, A. Searching online mendelian inheritance in man (OMIM): A knowl-
edgebase of human genes and genetic phenotypes. *Curr. Protoc. Bioinformatics* **58**, 1.2.1–1.2.12
(2017).
- 998 19. McKusick, V. A. Mendelian inheritance in man and its online version, OMIM. *Am. J. Hum. Genet.*
80, 588–604 (2007).
- 999 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).
- 1000 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. Geneesk.* **152**, 518–519 (2008).
- 1001 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).
- 1002 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).
- 1003 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).

- 1004 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).
- 1005 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).
- 1006 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).
- 1007 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).
- 1008 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).
- 1009 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).
- 1010 31. Svensson, V., Veiga Beltrame, E. da & Pachter, L. A curated database reveals trends in single-cell transcriptomics. *Database* **2020**, (2020).
- 1011 32. Cao, J. *et al.* A human cell atlas of fetal gene expression. *Science* **370**, (2020).
- 1012 33. Han, X. *et al.* Construction of a human cell landscape at single-cell level. *Nature* **581**, 303–309 (2020).
- 1013 34. Cao, J. *et al.* [A human cell atlas of fetal gene expression](#). *Science* **370**, eaba7721 (2020).
- 1014 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).
- 1015 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).
- 1016 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).
- 1017 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).
- 1018 39. Diehl, A. D. *et al.* The cell ontology 2016: Enhanced content, modularization, and ontology interoperability. *J. Biomed. Semantics* **7**, 44 (2016).
- 1019 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).
- 1020 41. Heim, C. E. *et al.* Myeloid-derived suppressor cells contribute to staphylococcus aureus orthopedic biofilm infection. *J. Immunol.* **192**, 3778–3792 (2014).

- 1021 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).
- 1022 43. Stoll, H. *et al.* Staphylococcal enterotoxins Dose-Dependently modulate the generation of Myeloid-Derived suppressor cells. *Front. Cell. Infect. Microbiol.* **8**, 321 (2018).
- 1023 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).
- 1024 45. Zhou, Z., Xu, M.-J. & Gao, B. Hepatocytes: A key cell type for innate immunity. *Cell. Mol. Immunol.* **13**, 301–315 (2016).
- 1025 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).
- 1026 47. Ladhami, S. N. *et al.* Invasive meningococcal disease in patients with complement deficiencies: A case series (2008-2017). *BMC Infect. Dis.* **19**, 522 (2019).
- 1027 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).
- 1028 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).
- 1029 50. Lung, T. *et al.* The complement system in liver diseases: Evidence-based approach and therapeutic options. *J Transl Autoimmun* **2**, 100017 (2019).
- 1030 51. Reis, E. S. *et al.* Applying complement therapeutics to rare diseases. *Clin. Immunol.* **161**, 225–240 (2015).
- 1031 52. Seal, R. L. *et al.* Genenames.org: The HGNC resources in 2023. *Nucleic Acids Res.* **51**, D1003–D1009 (2023).
- 1032 53. Al-Hamoudi, W. K. Severe autoimmune hepatitis triggered by varicella zoster infection. *World J. Gastroenterol.* **15**, 1004–1006 (2009).
- 1033 54. Brewer, E. C. & Hunter, L. Acute liver failure due to disseminated varicella zoster infection. *Case Reports Hepatol* **2018**, 1269340 (2018).
- 1034 55. Eshchar, J., Reif, L., Waron, M. & Alkan, W. J. Hepatic lesion in chickenpox. A case report. *Gastroenterology* **64**, 462–466 (1973).
- 1035 56. Li, Z. *et al.* [Aging and age-related diseases: From mechanisms to therapeutic strategies](#). *Biogerontology* **22**, 165–187 (2021).
- 1036 57. Nelson, M. R. *et al.* [The support of human genetic evidence for approved drug indications](#). *Nature Genetics* **47**, 856–860 (2015).

- 1037 58. Ochoa, D. *et al.* Human genetics evidence supports two-thirds of the 2021 FDA-approved drugs. *Nature Reviews Drug Discovery* **21**, 551–551 (2022).
- 1038 59. Minikel, E. V., Painter, J. L., Dong, C. C. & Nelson, M. R. Refining the impact of genetic evidence on clinical success. *Nature* 1–6 (2024) doi:[10.1038/s41586-024-07316-0](https://doi.org/10.1038/s41586-024-07316-0).
- 1039 60. Liu, X. *et al.* The therapeutic target database: An internet resource for the primary targets of approved, clinical trial and experimental drugs. *Expert Opin. Ther. Targets* **15**, 903–912 (2011).
- 1040 61. Chiu, W. *et al.* An update on gene therapy for inherited retinal dystrophy: Experience in leber congenital amaurosis clinical trials. *International Journal of Molecular Sciences* **22**, 4534 (2021).
- 1041 62. Fenderson, B. A. Chapter 6 - developmental and genetic diseases. in *Pathology secrets (third edition)* (ed. Damjanov, I.) 98–119 (Mosby, 2009). doi:[10.1016/B978-0-323-05594-9.00006-4](https://doi.org/10.1016/B978-0-323-05594-9.00006-4).
- 1042 63. Vilcaes, A. A., Garbarino-Pico, E., Torres Demichelis, V. & Daniotti, J. L. Ganglioside synthesis by plasma membrane-associated sialyltransferase in macrophages. *International Journal of Molecular Sciences* **21**, 1063 (2020).
- 1043 64. Yohe, H. C., Coleman, D. L. & Ryan, J. L. Ganglioside alterations in stimulated murine macrophages. *Biochimica et Biophysica Acta (BBA) - Biomembranes* **818**, 81–86 (1985).
- 1044 65. Demir, S. A., Timur, Z. K., Ateş, N., Martínez, L. A. & Seyrantepe, V. GM2 ganglioside accumulation causes neuroinflammation and behavioral alterations in a mouse model of early onset tay-sachs disease. *Journal of Neuroinflammation* **17**, 277 (2020).
- 1045 66. Ferro, A., Sheeler, C., Rosa, J.-G. & Cvetanovic, M. Role of microglia in ataxias. *Journal of molecular biology* **431**, 1792–1804 (2019).
- 1046 67. Hol, E. M. & Pasterkamp, R. J. Microglial transcriptomics meets genetics: New disease leads. *Nature Reviews Neurology* 1–2 (2022) doi:[10.1038/s41582-022-00633-w](https://doi.org/10.1038/s41582-022-00633-w).
- 1047 68. Lopes, K. de P. *et al.* Atlas of genetic effects in human microglia transcriptome across brain regions, aging and disease pathologies. *bioRxiv* 2020.10.27.356113 (2020) doi:[10.1101/2020.10.27.356113](https://doi.org/10.1101/2020.10.27.356113).
- 1048 69. Bueren, J. A. & Auricchio, A. Advances and challenges in the development of gene therapy medicinal products for rare diseases. *Hum. Gene Ther.* **34**, 763–775 (2023).
- 1049 70. Bulaklak, K. & Gersbach, C. A. The once and future gene therapy. *Nat. Commun.* **11**, 5820 (2020).
- 1050 71. Godbout, K. & Tremblay, J. P. Prime editing for human gene therapy: Where are we now? *Cells* **12**, (2023).
- 1051 72. Kohn, D. B., Chen, Y. Y. & Spencer, M. J. Successes and challenges in clinical gene therapy. *Gene Ther.* **30**, 738–746 (2023).
- 1052 73. Zhao, Z., Shang, P., Mohanraju, P. & Geijzen, N. Prime editing: Advances and therapeutic applications. *Trends Biotechnol.* **41**, 1000–1012 (2023).
- 1053 74. Darrow, J. J. Luxturna: FDA documents reveal the value of a costly gene therapy. *Drug Discov. Today* **24**, 949–954 (2019).

- 1054 75. Mendell, J. R. *et al.* Single-Dose Gene-Replacement therapy for spinal muscular atrophy. *N. Engl. J. Med.* **377**, 1713–1722 (2017).
- 1055 76. Mueller, C. *et al.* 5 year expression and neutrophil defect repair after gene therapy in alpha-1 antitrypsin deficiency. *Mol. Ther.* **25**, 1387–1394 (2017).
- 1056 77. Russell, S. *et al.* Efficacy and safety of voretigene neparvovec (AAV2-hRPE65v2) in patients with RPE65-mediated inherited retinal dystrophy: A randomised, controlled, open-label, phase 3 trial. *Lancet* **390**, 849–860 (2017).
- 1057 78. Lu, C.-F. FDA takes first step toward international regulation of gene therapies to treat rare diseases. (2024).
- 1058 79. Brown, D. G., Wobst, H. J., Kapoor, A., Kenna, L. A. & Southall, N. Clinical development times for innovative drugs. *Nat. Rev. Drug Discov.* **21**, 793–794 (2022).
- 1059 80. Moffat, J. G., Vincent, F., Lee, J. A., Eder, J. & Prunotto, M. Opportunities and challenges in phenotypic drug discovery: An industry perspective. *Nat. Rev. Drug Discov.* **16**, 531–543 (2017).
- 1060 81. Zhou, Q. & Buchholz, C. J. Cell type specific gene delivery by lentiviral vectors: New options in immunotherapy. *Oncoimmunology* **2**, e22566 (2013).
- 1061 82. Gao, C., Jiang, J., Tan, Y. & Chen, S. [Microglia in neurodegenerative diseases: mechanism and potential therapeutic targets](#). *Signal Transduction and Targeted Therapy* **8**, 1–37 (2023).
- 1062 83. Mcquade, A. & Blurton-jones, M. Microglia in alzheimer’s disease : Exploring how genetics and phenotype influence risk. *Journal of Molecular Biology* 1–13 (2019) doi:[10.1016/j.jmb.2019.01.045](https://doi.org/10.1016/j.jmb.2019.01.045).
- 1063 84. Keren-shaul, H. *et al.* [A unique microglia type associated with restricting development of alzheimer ’s disease](#). *Cell* **169**, 1276–1290.e17 (2017).
- 1064 85. Deczkowska, A. *et al.* [Disease-associated microglia: A universal immune sensor of neurodegeneration](#). *Cell* **173**, 1073–1081 (2018).
- 1065 86. DiStefano, M. T. *et al.* The gene curation coalition: A global effort to harmonize gene-disease evidence resources. *Genet. Med.* **24**, 1732–1742 (2022).
- 1066 87. DiStefano, M. *et al.* P451: The gene curation coalition works to resolve discrepancies in gene-disease validity assertions. *Genetics in Medicine Open* **1**, 100498 (2023).
- 1067 88. Skene, N. G. & Grant, S. G. N. Identification of vulnerable cell types in major brain disorders using single cell transcriptomes and expression weighted cell type enrichment. *Front. Neurosci.* **10**, 16 (2016).
- 1068 89. Osmond, M. *et al.* Outcome of over 1500 matches through the matchmaker exchange for rare disease gene discovery: The 2-year experience of Care4Rare canada. *Genet. Med.* **24**, 100–108 (2022).
- 1069 90. Philippakis, A. A. *et al.* The matchmaker exchange: A platform for rare disease gene discovery. *Hum. Mutat.* **36**, 915–921 (2015).

- 1070 91. Zanello, G. *et al.* Targeting shared molecular etiologies to accelerate drug development for rare
diseases. *EMBO Mol. Med.* **15**, e17159 (2023).
- 1071 92. Côté, R. *et al.* The ontology lookup service: Bigger and better. *Nucleic Acids Res.* **38**, W155–60
(2010).
- 1072 93. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach
to multiple testing. *J. R. Stat. Soc.* (1995).
- 1073 94. Lazarin, G. A. *et al.* Systematic classification of disease severity for evaluation of expanded carrier
screening panels. *PLoS One* **9**, e114391 (2014).
- 1074 95. Solovyeva, V. V. *et al.* New approaches to tay-sachs disease therapy. *Frontiers in Physiology* **9**,
(2018).
- 1075 96. Hoffman, J. D. *et al.* Next-generation DNA sequencing of HEXA: A step in the right direction for
carrier screening. *Molecular Genetics & Genomic Medicine* **1**, 260–268 (2013).
- 1076 97. Sugiyama, K., Tagawa, S. & Toda, M. Methods for visual understanding of hierarchical system struc-
tures. *IEEE Trans. Syst. Man Cybern.* **11**, 109–125 (1981).

1077

1078

1079 **Supplementary Materials**

1080 **Supplementary Results**

1081 **Selected example targets**

1082 From our prioritised targets, we selected the following four sets of phenotypes or diseases as examples:
1083 ‘GM2-ganglioside accumulation’, ‘Spinocerebellar atrophy’, ‘Neuronal loss in central nervous system’. Only
1084 phenotypes with a GPT severity score greater than 15 were considered to avoid overplotting and to focus on
1085 the more clinically relevant phenotypes Fig. 8a-h. These examples were then selected partly on the basis of
1086 severity rankings, and partly for their relatively smaller, simpler networks than lent themselves to compact
1087 visualisations.

1088 Tay-Sachs disease (TSD) is a devastating hereditary condition in which children are born appearing healthy,
1089 which gradually degrades leading to death after 3-5 years. The underlying cause is the toxic accumulation
1090 of gangliosides in the nervous system due to a loss of the enzyme produced by *HEXA*. While this could
1091 in theory be corrected with gene editing technologies, there remain some outstanding challenges. One of
1092 which is identifying which cell types should be targeted to ensure the most effective treatments. Here
1093 we identified alternatively activated macrophages as the cell type most strongly associated with ‘GM2-
1094 ganglioside accumulation’ Fig. 8i. The role of aberrant macrophage activity in the regulation of ganglioside
1095 levels is supported by observation that gangliosides accumulate within macrophages in TSD⁶², as well as
1096 experimental evidence in rodent models^{63,64,65}. Our results not only corroborate these findings, but propose
1097 macrophages as the primary causal cell type in TSD, making it the most promising cell type to target in
1098 therapies.

1099 Another challenge in TSD is early detection and diagnosis, before irreversible damage has occurred. Our
1100 pipeline implicated extravillous trophoblasts of the placenta in ‘GM2-ganglioside accumulation’. While not
1101 necessarily a target for gene therapy (as the child is detached from the placenta after birth), checking these
1102 cells *in utero* for an absence of *HEXA* may serve as a viable biomarker as these cells normally express
1103 the gene at high levels. Early detection of TSD may lengthen the window of opportunity for therapeutic
1104 intervention⁹⁵, especially when genetic sequencing is not available or variants of unknown significance are
1105 found within *HEXA*⁹⁶.

1106 Spinocerebellar atrophy is a debilitating and lethal phenotype that occurs in diseases such as Spinocerebellar
1107 ataxia and Boucher-Nenhauser syndrome. These diseases are characterised by progressive degeneration of
1108 the cerebellum and spinal cord, leading to severe motor and cognitive impairments. Our pipeline identified
1109 M2 macrophages (labeled as the closest CL term ‘Alternatively activated macrophages’ in Fig. 8j) as the
1110 only causal cell type associated with ‘Spinocerebellar atrophy’. This strongly suggests that degeneration of
1111 cerebellar Purkinje cells are in fact downstream consequences of macrophage dysfunction, rather than being
1112 the primary cause themselves. This is consistent with the known role of macrophages, especially microglia, in

1113 neuroinflammation and other neurodegenerative conditions such as Alzheimer's and Parkinsons' disease^{66–68}.
1114 While experimental and postmortem observational studies have implicated microglia in spinocerebellar atro-
1115 phy previously⁶⁶, our results provide a statistically-supported and unbiased genetic link between known risk
1116 genes and this cell type. Therefore, targeting M2 microglia in the treatment of spinocerebellar atrophy may
1117 therefore represent a promising therapeutic strategy. This is aided by the fact that there are mouse models
1118 that perturb the ortholog of human spinocerebellar atrophy risk genes (e.g. *Atxn1*, *Pnpla6*) and reliably
1119 recapitulate the effects of this diseases at the cellular (e.g. loss of Purkinje cells), morphological (e.g. atrophy
1120 of the cerebellum, spinal cord, and muscles), and functional (e.g. ataxia) levels.

1121 Next, we investigated the phenotype 'Neuronal loss in the central nervous system'. Despite the fact that this
1122 is a fairly broad phenotype, we found that it was only significantly associated with 3 cell types (alternatively
1123 activated macrophage, macrophage, epithelial cell), specifically M2 macrophages and sinusoidal endothelial
1124 cells Fig. 8k.

1125 Skeletal dysplasia is a heterogeneous group of over 450 disorders that affect the growth and development of
1126 bone and cartilage. This phenotype can be lethal when deficient bone growth leads to the constriction of
1127 vital organs such as the lungs. Even after surgical interventions, these complications continue to arise as the
1128 child develops. Pharmacological interventions to treat this condition have largely been ineffective. While
1129 there are various cell types involved in skeletal system development, our pipeline nominated chondrocytes
1130 as the causal cell type underlying the lethal form of this condition (Fig. 19). Assuringly, we found that
1131 the disease 'Achondrogenesis Type 1B' is caused by the genes *SLC26A2* and *COL2A1* via chondrocytes.
1132 We also found that 'Platyspondylic lethal skeletal dysplasia, Torrance type'. Thus, in cases where surgical
1133 intervention is insufficient, targeting these genes within chondrocytes may prove a viable long-term solution
1134 for children suffering from lethal skeletal dysplasia.

1135 Alzheimer's disease (AD) is the most common neurodegenerative condition. It is characterised by a set of
1136 variably penetrant phenotypes including memory loss, cognitive decline, and cerebral proteinopathy. Inter-
1137 estingly, we found that different forms of early onset AD (which are defined by the presence of a specific
1138 disease gene) are each associated with different cell types via different phenotypes (Fig. 19). For example,
1139 AD 3 and AD 4 are primarily associated with cells of the digestive system ('enterocyte', 'gastric goblet
1140 cell') and are implied to be responsible for the phenotypes 'Senile plaques', 'Alzheimer disease', 'Parietal
1141 hypometabolism in FDG PET'. Meanwhile, AD 2 is primarily associated with immune cells ('alternatively
1142 activated macrophage') and is implied to be responsible for the phenotypes 'Neurofibrillary tangles', 'Long-
1143 tract signs'. This suggests that different forms of AD may be driven by different cell types and phenotypes,
1144 which may help to explain its variability in onset and clinical presentation.

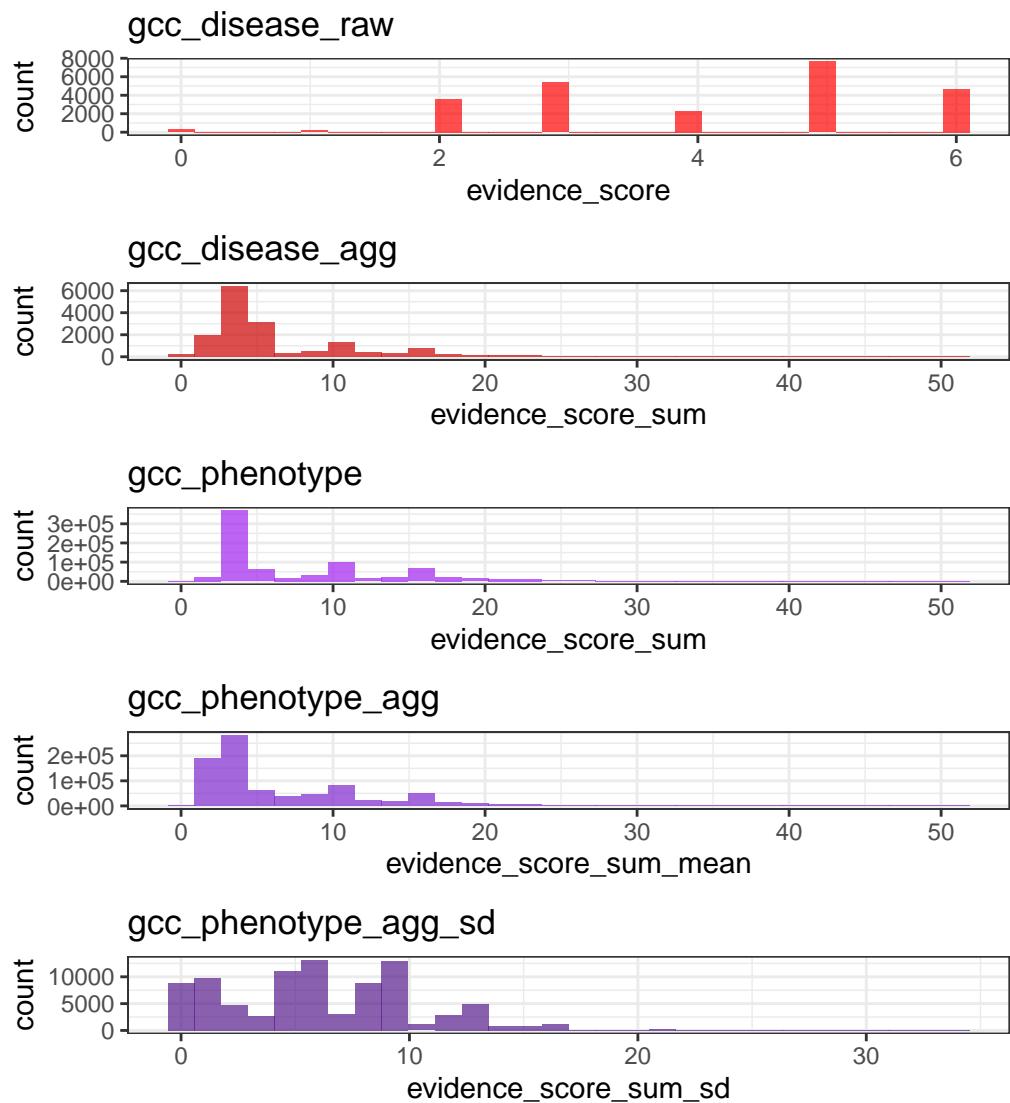
1145 Finally, Parkinson's disease (PD) is characterised by motor symptoms such as tremor, rigidity, and bradyki-
1146 nesia. However there are a number of additional phenotypes associated with the disease that span multiple
1147 physiological systems. PD 19a and PD 8 seemed to align most closely with the canonical understanding of

1148 PD as a disease of the central nervous system in that they implicated oligodendrocytes and neurons (Fig. 19).
1149 Though the reference datasets being used in this study were not annotated at sufficient resolution to distin-
1150 guish between different subtypes of neurons, in particular dopaminergic neurons. PD 19a/8 also suggested
1151 that risk variants in *LRRK2* mediate their effects on PD through both myeloid cells and oligodendrocytes
1152 by causing gliosis of the substantia nigra. The remaining clusters of PD mechanisms revolved around chon-
1153 drocytes (PD 20), amacrine cells of the eye (hereditary late-onset PD), and the respiratory/immune system
1154 (PD 14). While the diversity in cell type-specific mechanisms is somewhat surprising, it may help to explain
1155 the wide variety of cross-system phenotypes frequently observed in PD.

1156 It should be noted that the HPO only includes gene annotations for the monogenic forms of AD and PD.
1157 However it has previously been shown that there is at least partial overlap in their phenotypic and genetic
1158 aetiology with respect to their common forms. Thus understanding the monogenic forms of these diseases
1159 may shed light onto their more common counterparts.

1160 **Phenome-wide analyses discover novel phenotype-cell type associations**

1161 We visualised the putative causal relationships between genes, cell types and diseases associated with RNI as
1162 a network (Fig. 14). The phenotype ‘Recurrent Neisserial infections’ was connected to cell types through the
1163 aforementioned association test results ($FDR < 0.05$). Genes that were primarily driving these associations
1164 (i.e. genes that were both strongly linked with ‘Recurrent Neisserial infections’ and were highly specifically
1165 expressed in the given cell type) were designated as “driver genes” and retained for plotting. Across all
1166 phenotypes in the HPO, more specific phenotypes (terms in the HPO with greater IC) are not only more
1167 specific to certain cell types (Fig. 3b), but are also associated with genes that have greater cell type-specific
1168 expression within those cell types. Even so, we should note that the choice of which specificity quantiles to
1169 include is arbitrary. It should also be noted that simply because a gene is not specific to a cell type does not
1170 mean it is not important for the function of the cell type. Indeed, there are many genes that are ubiquitously
1171 expressed throughout many tissues in the body and are essential for cell function. Gene expression specificity
1172 is nevertheless a useful metric to help distinguish many hundreds of cell (sub)types with overlapping gene
1173 signatures.



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

Figure 9

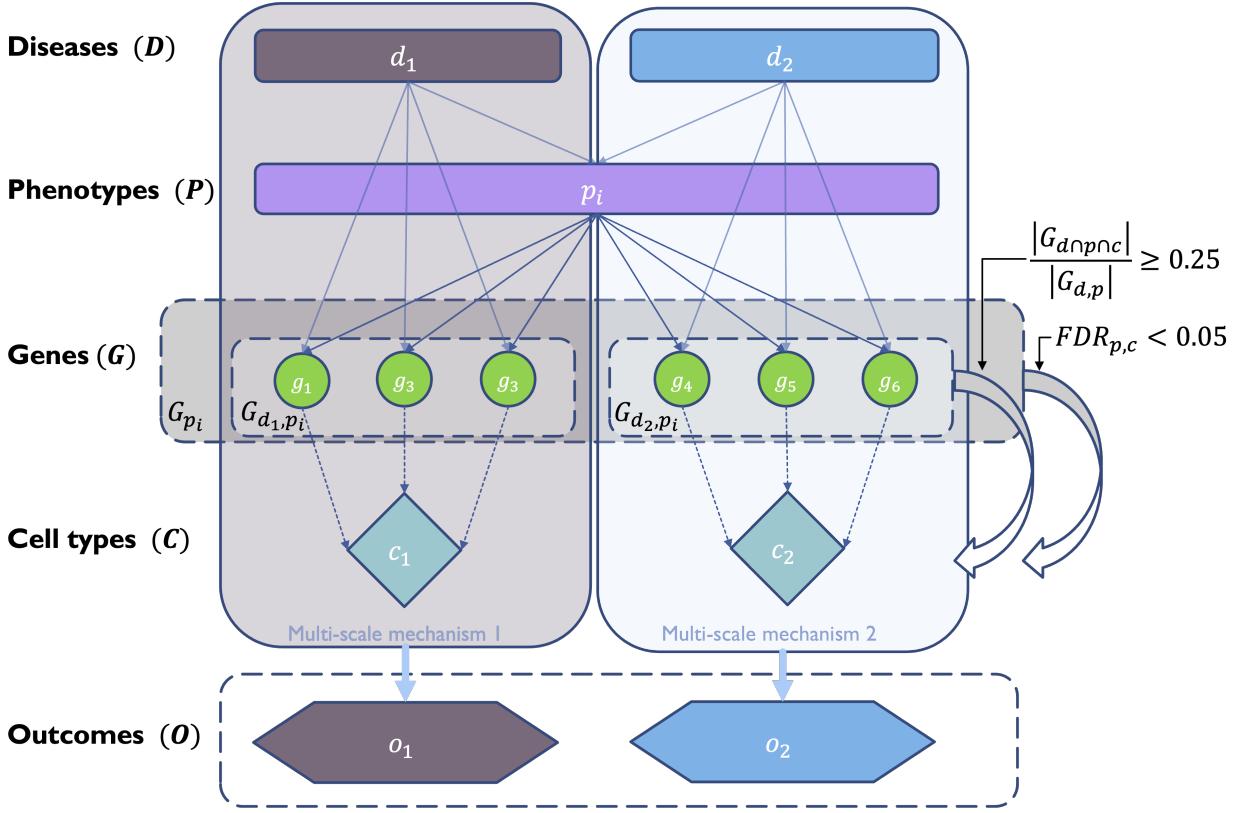
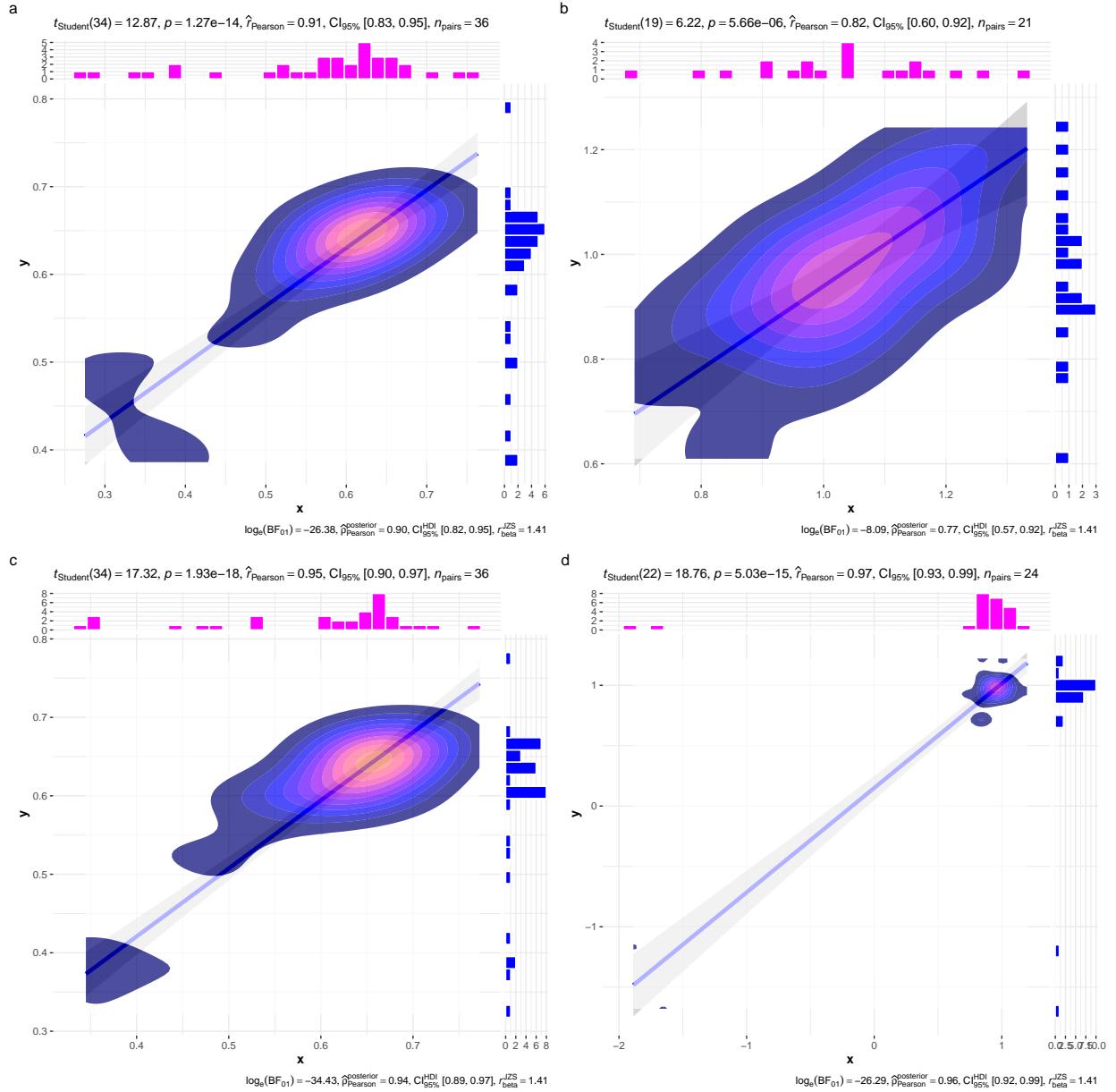
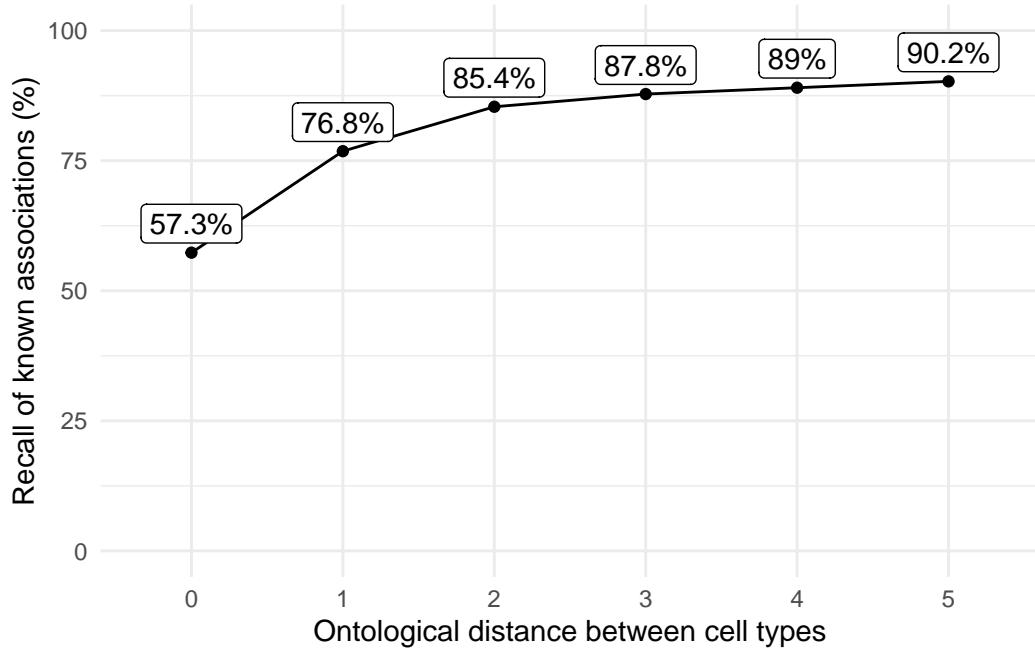


Figure 10: **Diagrammatic overview of multi-scale disease investigation strategy.** Here we provide an abstract example of differential disease aetiology across multiple scales: diseases (D), phenotypes (P), cell types (C), genes (G), and clinical outcomes (O). In the HPO, genes are assigned to phenotypes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



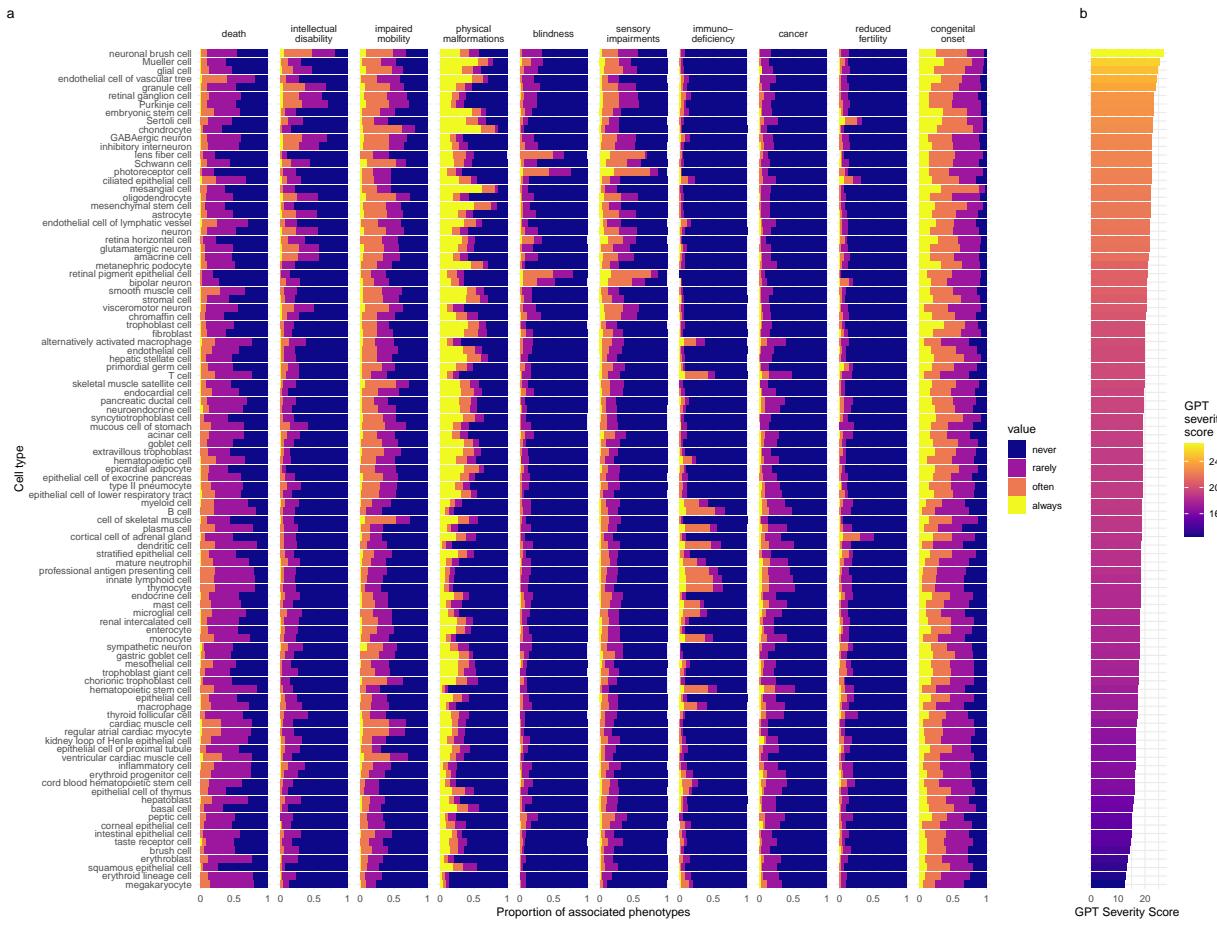
(a) Inter- and intra-dataset validation across the different CellTypeDataset (CTD) and developmental stages. Correlations are computed using Pearson correlation coefficient. Point density is plotted using a 2D kernel density estimate. **a** Correlation between the uncorrected p-values from all phenotype-cell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

Figure 11



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

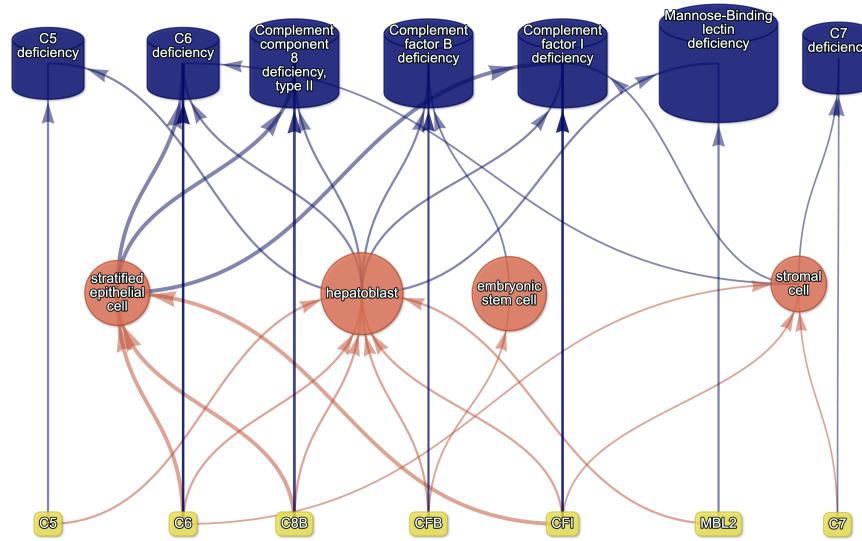
Figure 12



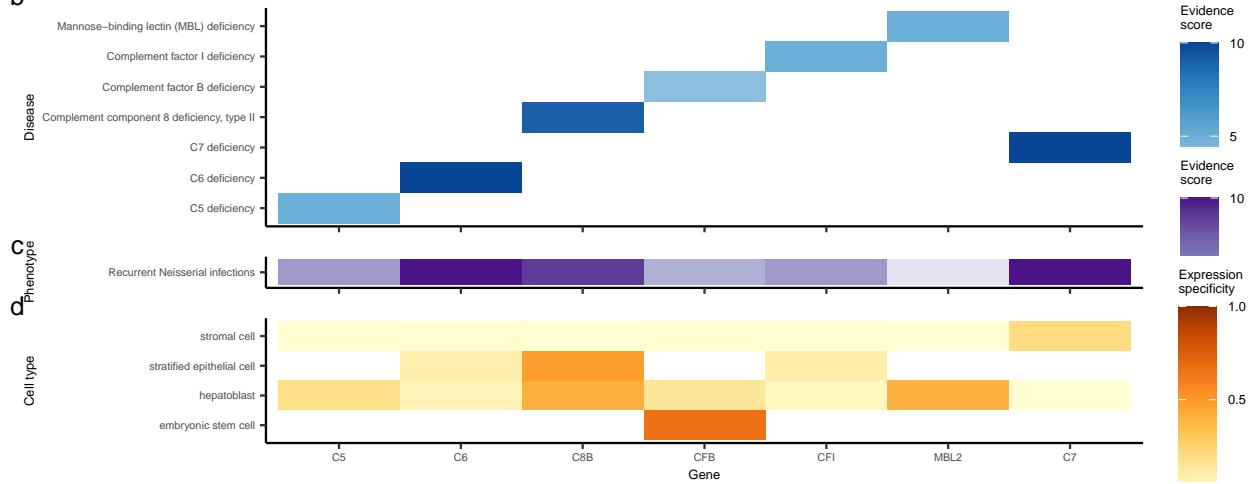
(a) **Cell types ordered by the mean severity of the phenotypes they're associated with.** **a**, The distribution of phenotype severity annotation frequencies aggregated by cell type. **b**, The composite severity score, averaged across all phenotypes associated with each cell type.

Figure 13

a

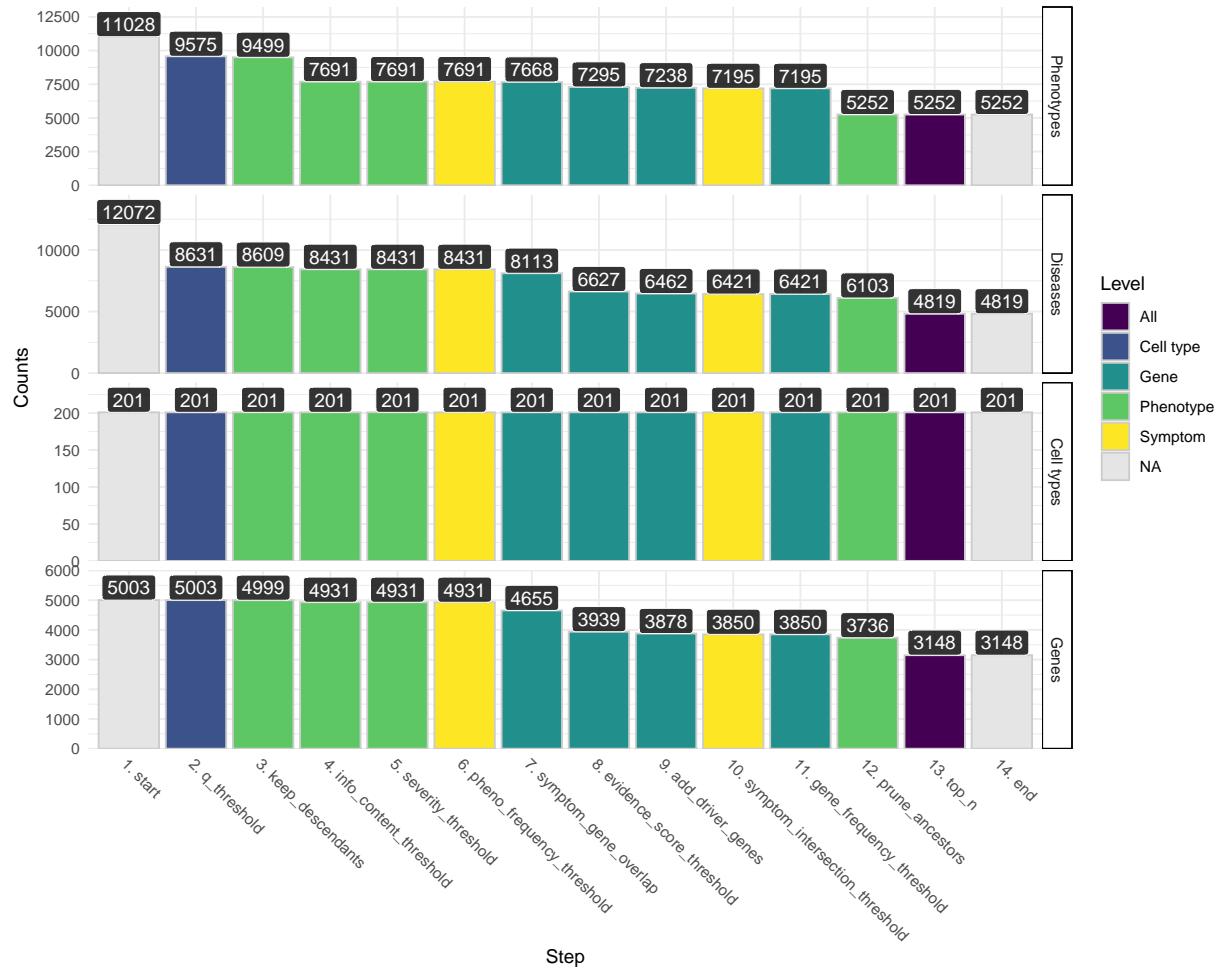


b



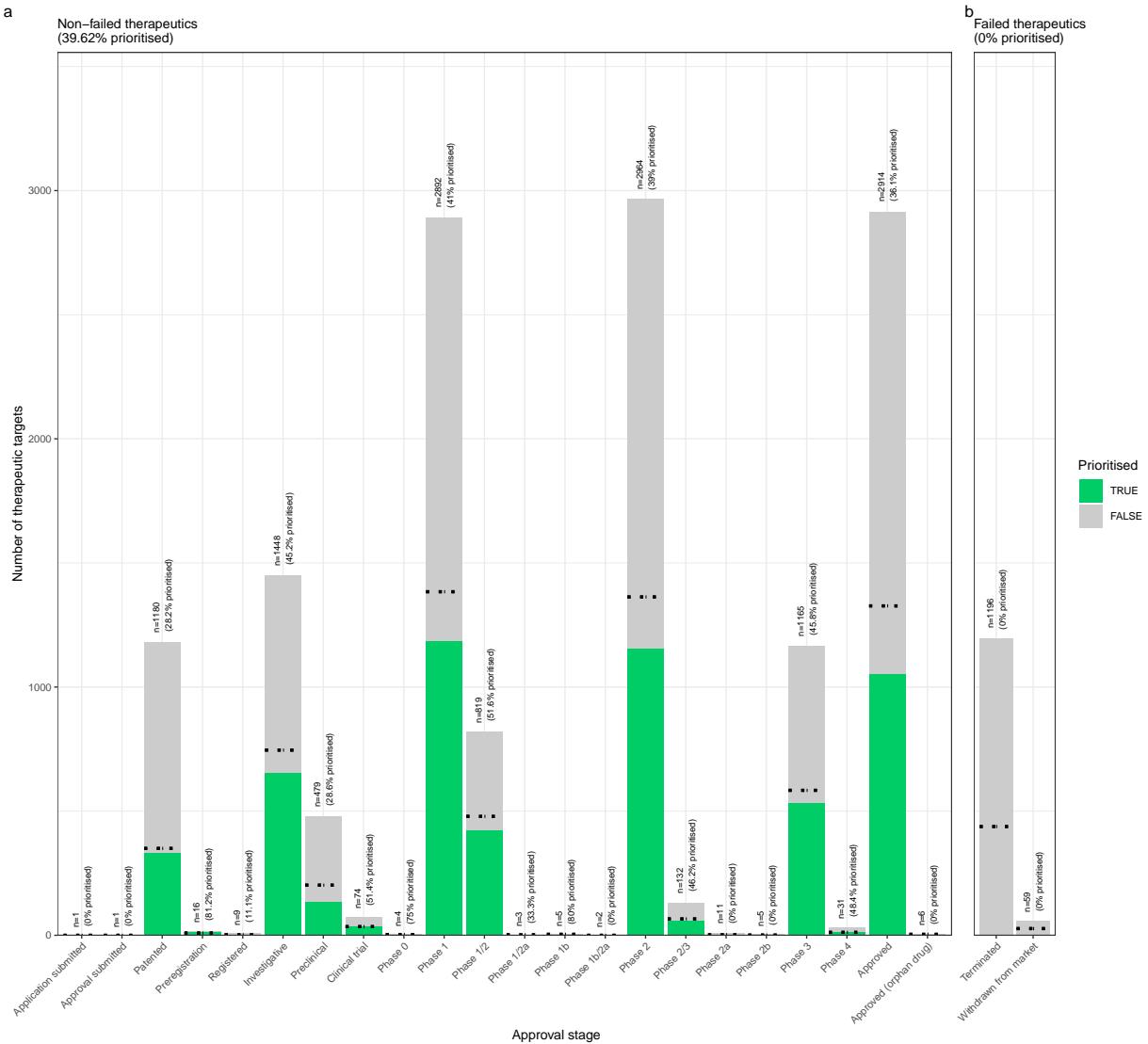
(a) **Causal network of recurrent Neisserial infections (RNI) reveals multi-scale disease aetiology.** RNI is a phenotype in seven different monogenic diseases caused by disruptions to specific complement system genes. Four cell types were significantly associated with RNI. **a**, One can trace how genes causal for RNI (yellow boxes, bottom) mediate their effects through cell types (orange circles, middle) and diseases (blue cylinders, top). Cell types are connected to RNI via association testing ($FDR < 0.05$). Genes shown here have both strong evidence for a causal role in RNI and high expression specificity in the associated cell type. Cell types can be linked to monogenic diseases via the genes specifically expressed in those cell types (i.e. are in the top 25% of cell type specificity expression quantiles). Nodes are arranged using the Sugiyama algorithm⁹⁷. **b** Expression specificity quantiles (1-40 scale) of each driver gene in each cell type (darker = greater specificity). **c** GenCC-derived evidence scores between the RNI phenotype and each gene. **d** Expression specificity (0 = least specific, 1 = most specific) of each gene in each cell type.

Figure 14



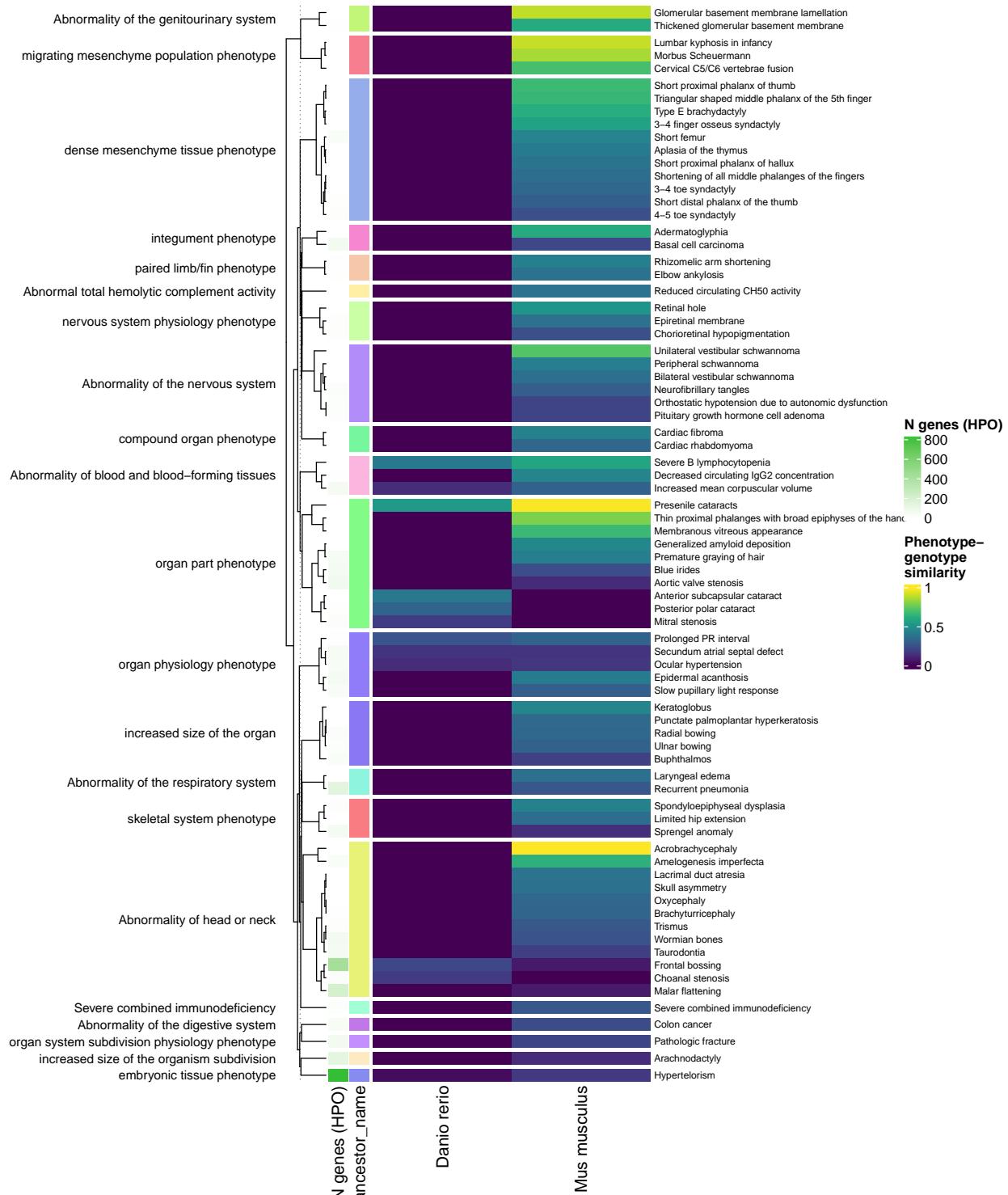
(a) **Prioritised target filtering steps.** This plot visualises the number of unique phenotype-cell type associations, cell types, genes, and phenotypes (*y-axis*) at each filtering step (*x-axis*) within the multi-scale therapeutic target prioritisation pipeline. Each step in the pipeline can be easily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

Figure 15



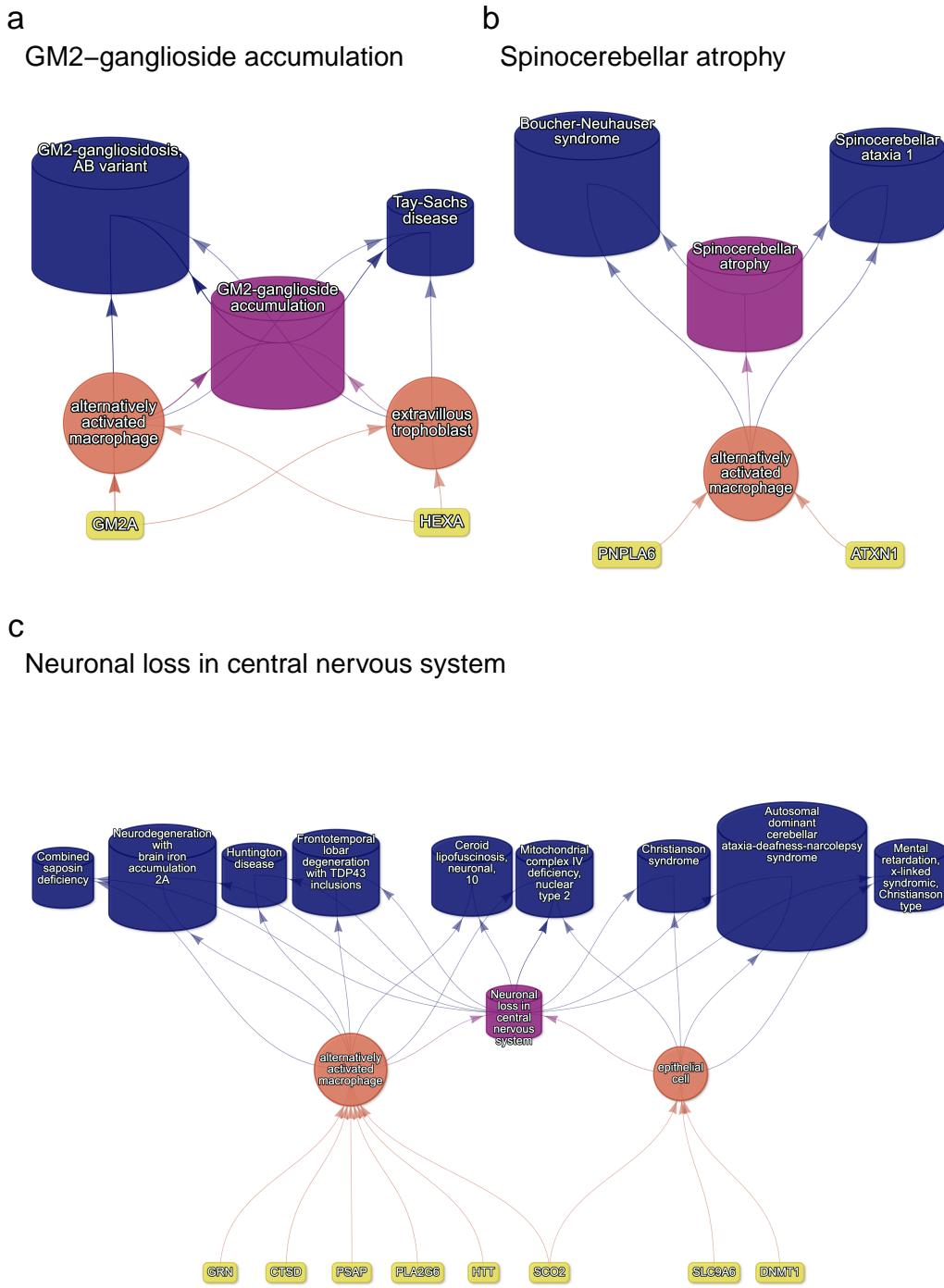
(a) **Validation of prioritised therapeutic targets.** Proportion of existing all therapy targets (documented in the Therapeutic Target Database) recapitulated by our prioritisation pipeline.

Figure 16



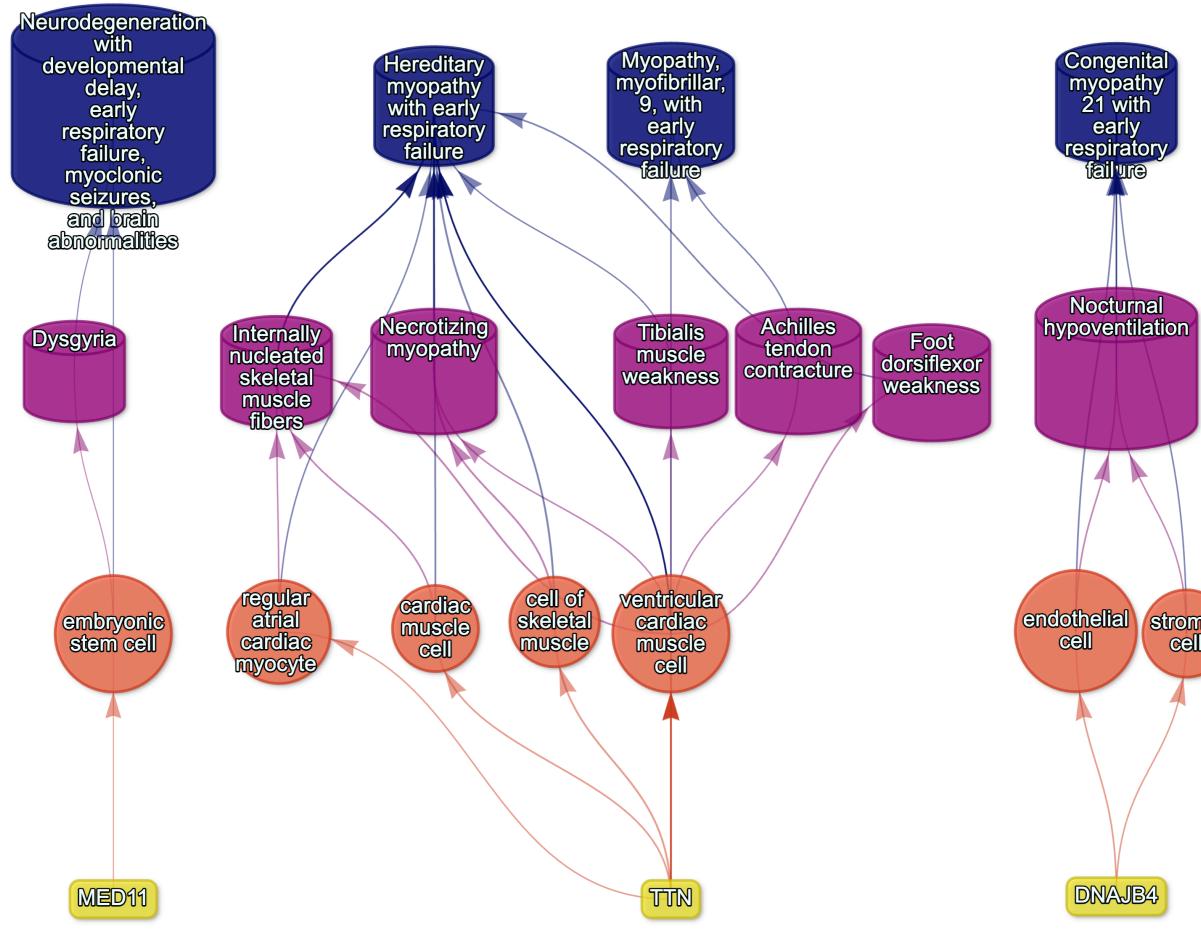
(a) **Identification of translatable experimental models.** Interspecies translatability of the top 200 human phenotypes nominated by the gene therapy prioritised pipeline. Above, the combined ontological-genotypic similarity score (SIM_{og}) is displayed as the heatmap fill colour stratified by the model organism (*x-axis*). An additional column (“n_genes_db1” on the far left) displays the total number of unique genes annotated to the phenotypic within the HPO. Phenotypes are clustered according to their ontological similarity in the HPO (*y-axis*).

Figure 17



(a) **Causal multi-scale networks reveal cell type-specific therapeutic targets.** Each disease (blue cylinders) is connected to its phenotype (purple cylinders) based on well-established clinical observations recorded within the 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 algorithm⁹⁷.

Figure 18



(a) Respiratory failure

Figure 19: **Example cell type-specific gene therapy targets for phenotypes associated with respiratory failure-related diseases.** Each disease (blue cylinders) is connected to its phenotype (purple cylinders) based on well-established clinical observations recorded within the 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 algorithm⁹⁷.

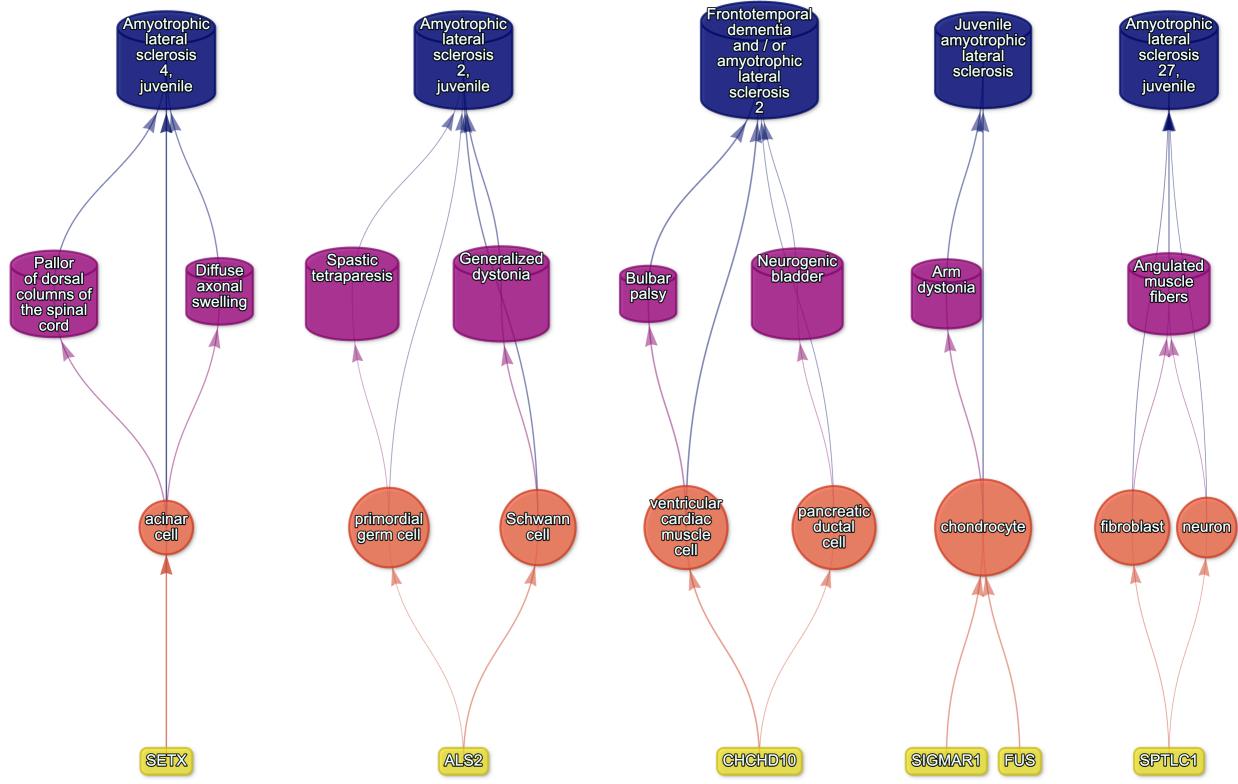


Figure 20: Causal multi-scale network for phenotypes associated with Amyotrophic Lateral Sclerosis (ALS).

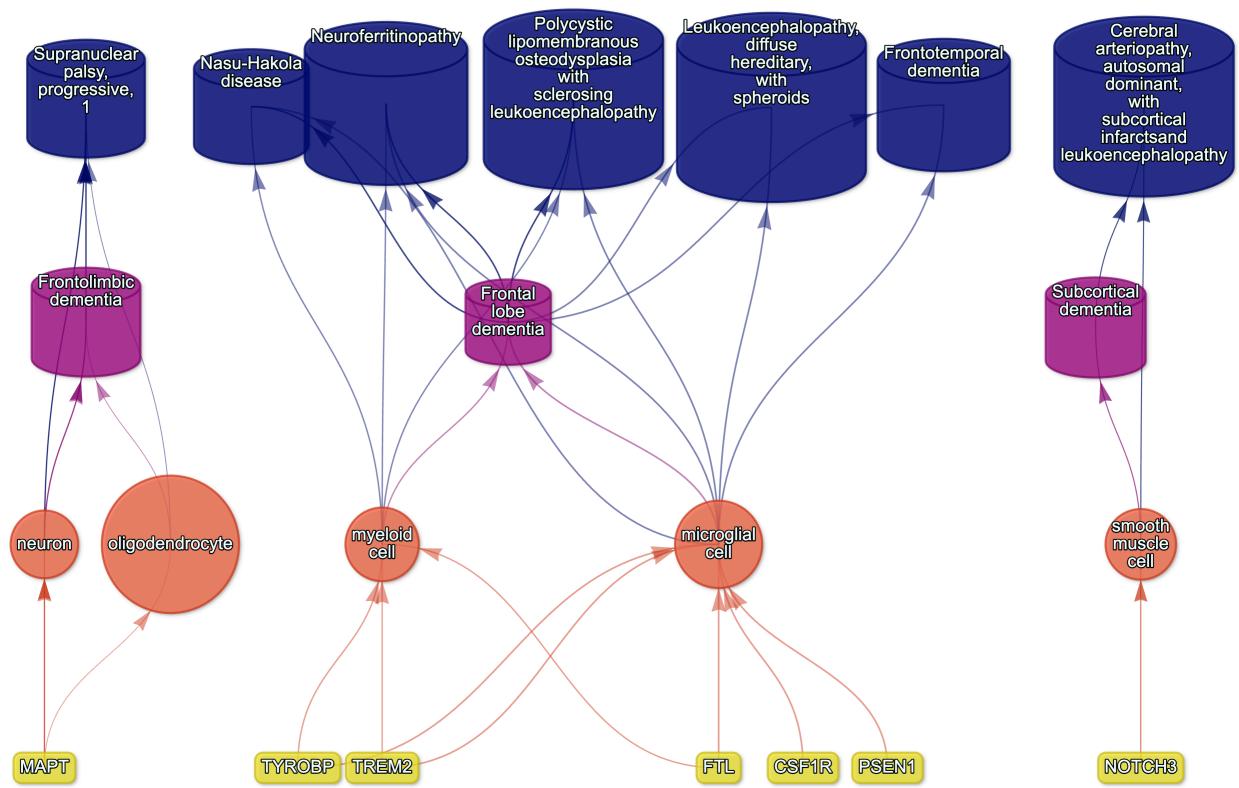


Figure 21: Causal multi-scale network for dementia phenotypes.

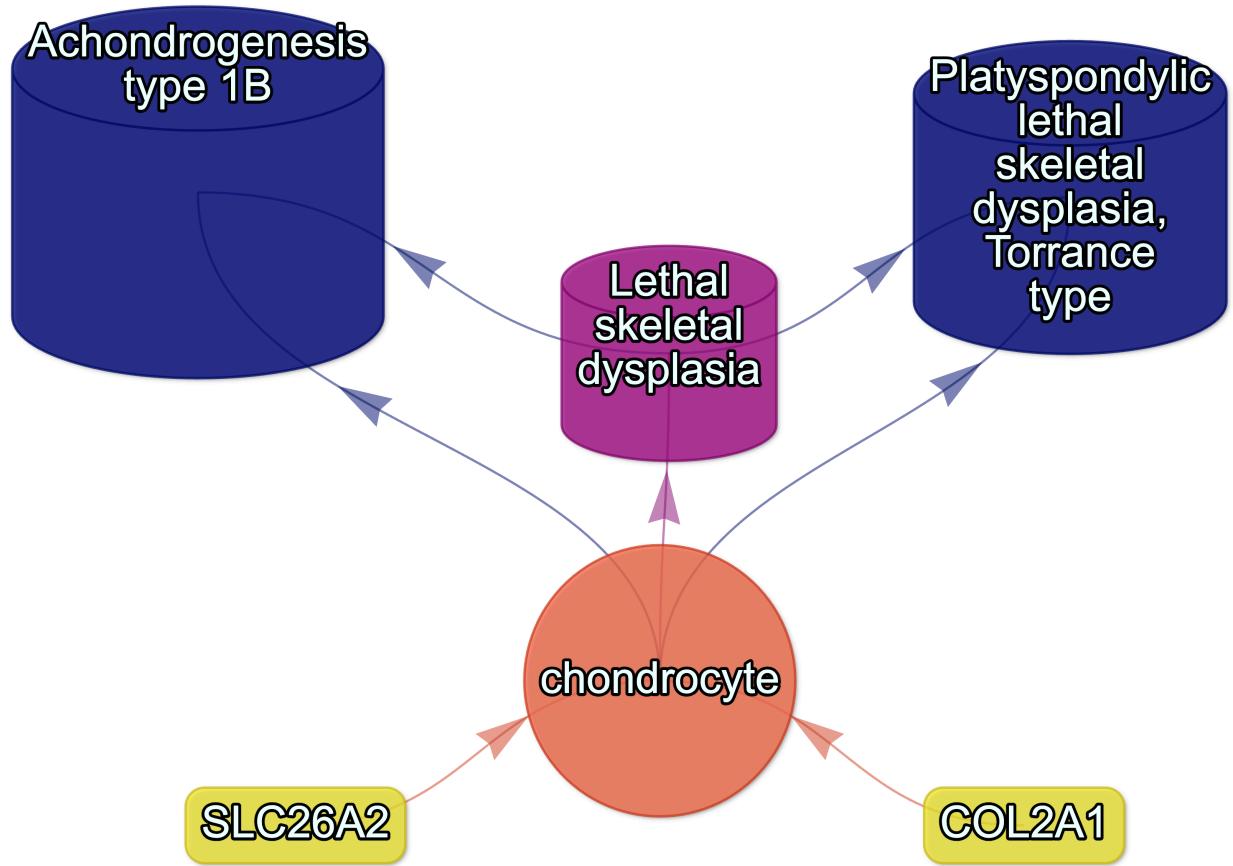


Figure 22: Causal multi-scale network for the phenotype lethal skeletal dysplasia.

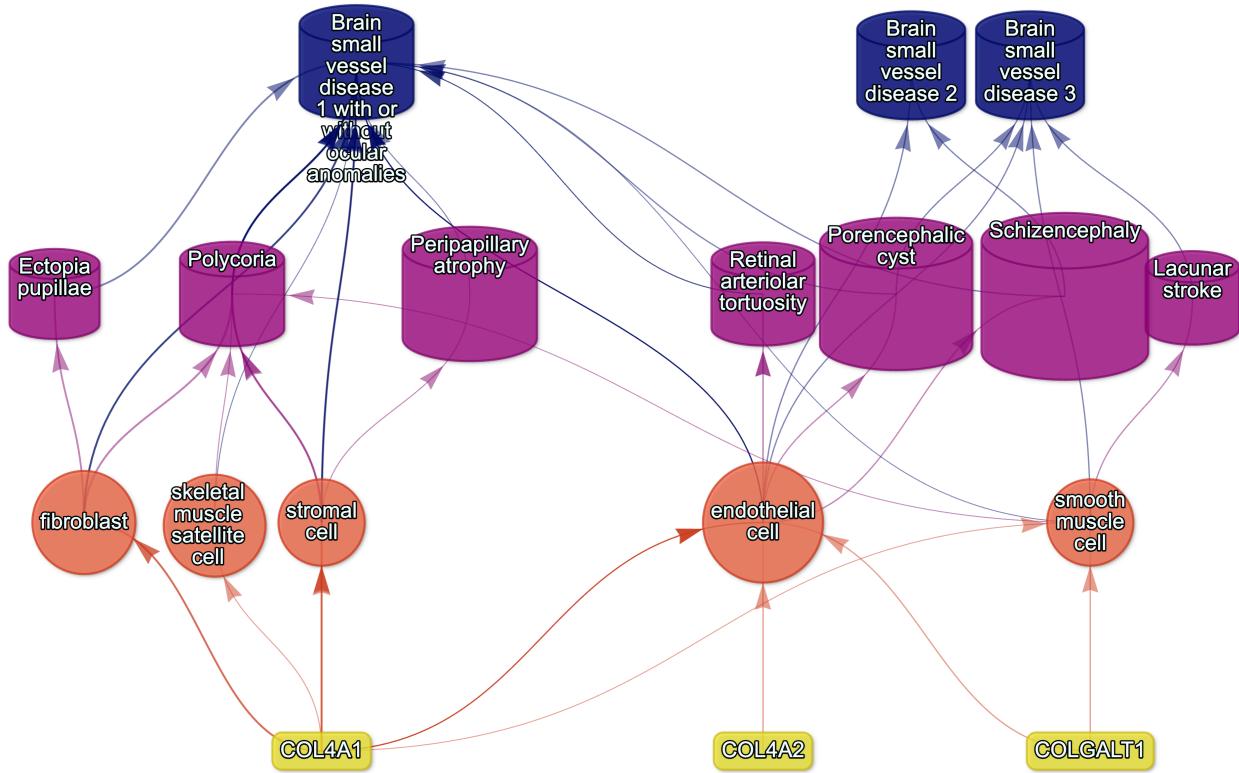


Figure 23: Causal multi-scale network for phenotypes associated with small vessel disease.

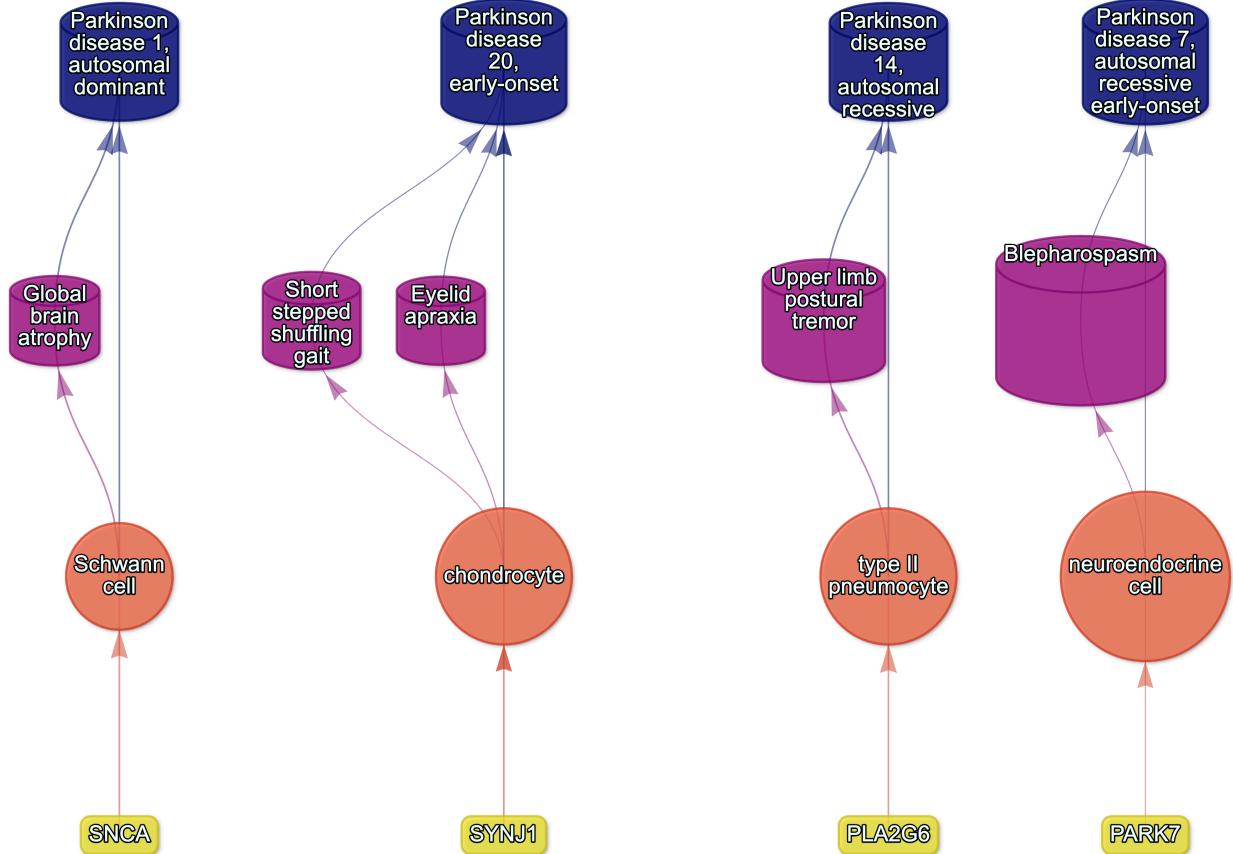


Figure 24: Causal multi-scale network for phenotypes associated with various subtypes of Parkinson's disease.

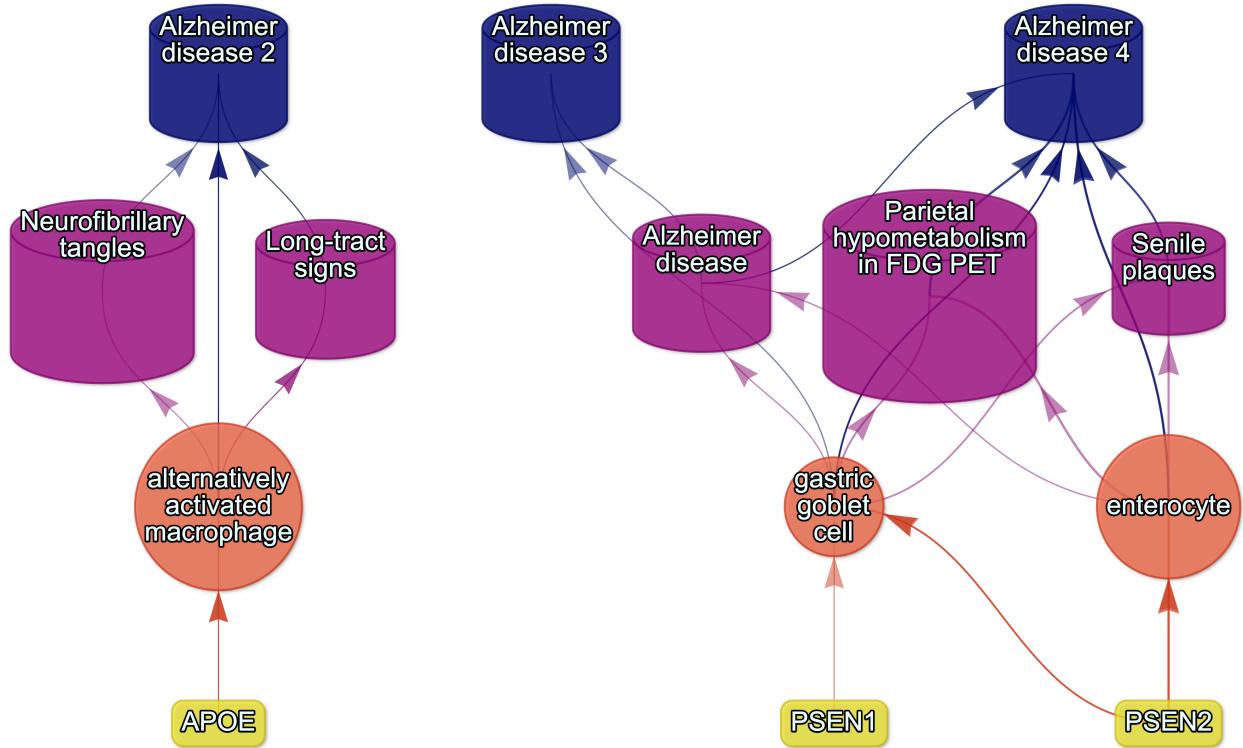


Figure 25: Causal multi-scale network for phenotypes associated with various subtypes of Alzheimer's disease.

1175 Supplementary Tables

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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