

1 Discovering Root Causal Genes 2 with High Throughput Perturbations

3 Eric V Strobl^{1*} and Eric R Gamazon²

*For correspondence:
eric.strobl@pitt.edu (ES)

4 ¹University of Pittsburgh; ²Vanderbilt University Medical Center

5

6 **Abstract** Root causal gene expression levels – or *root causal genes* for short – correspond to the
7 initial changes to gene expression that generate patient symptoms as a downstream effect.
8 Identifying root causal genes is critical towards developing treatments that modify disease near
9 its onset, but no existing algorithms attempt to identify root causal genes from data.
10 RNA-sequencing (RNA-seq) data introduces challenges such as measurement error, high
11 dimensionality and non-linearity that compromise accurate estimation of root causal effects even
12 with state-of-the-art approaches. We therefore instead leverage Perturb-seq, or high throughput
13 perturbations with single cell RNA-seq readout, to learn the causal order between the genes. We
14 then transfer the causal order to bulk RNA-seq and identify root causal genes specific to a given
15 patient for the first time using a novel statistic. Experiments demonstrate large improvements in
16 performance. Applications to macular degeneration and multiple sclerosis also reveal root causal
17 genes that lie on known pathogenic pathways, delineate patient subgroups and implicate a newly
18 defined omnigenic root causal model.

19

20 Introduction

21 *Root causes of disease* correspond to the most upstream causes of a diagnosis with strong causal
22 effects on the diagnosis. *Pathogenesis* refers to the causal cascade from root causes to the diag-
23 nosis. Genetic and non-genetic factors may act as root causes and affect gene expression as an
24 intermediate step during pathogenesis. We introduce root causal gene expression levels – or *root*
25 *causal genes* for short – that correspond to the initial changes to *gene expression* induced by genetic
26 and non-genetic root causes that have large causal effects on a downstream diagnosis (Figure 1
27 (a)). Root causal genes differ from core genes that directly cause the diagnosis and thus lie at the
28 end, rather than at the beginning, of pathogenesis (*Boyle et al., 2017*). Root causal genes also gen-
29 eralize (the expression levels of) driver genes that only account for the effects of somatic mutations
30 primarily in cancer (*Martínez-Jiménez et al., 2020*).

31 Treating root causal genes can modify disease pathogenesis in its entirety, whereas targeting
32 other causes may only provide symptomatic relief. For example, mutations in Gaucher disease
33 cause decreased expression of wild type beta-glucocerebrosidase, or the root causal gene (*Nagral,*
34 *2014*). We can give a patient blood transfusions to alleviate the fatigue and anemia associated
35 with the disease, but we seek more definitive treatments like recombinant glucocerebrosidase that
36 replaces the deficient enzyme. Enzyme replacement therapy alleviates the associated liver, bone
37 and neurological abnormalities of Gaucher disease as a downstream effect. Identifying root causal
38 genes is therefore critical for developing treatments that eliminate disease near its pathogenic
39 onset.

40 The problem is further complicated by the existence of complex disease, where a patient may
41 have multiple root causal genes that differ from other patients even within the same diagnostic
42 category (*Cano-Gamez and Trynka, 2020*). Complex diseases often have an overwhelming number

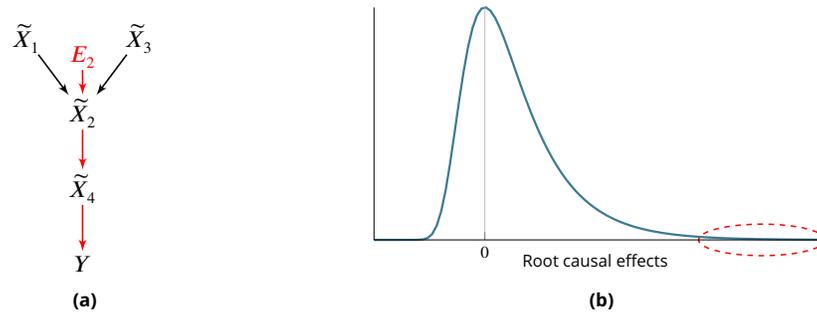


Figure 1. (a) Toy example where a variable E_2 simultaneously models genetic and non-genetic root causes that jointly have a large causal effect on a diagnose Y through gene expression \tilde{X} . E_2 first affects the gene expression level \tilde{X}_2 , or the root causal gene. The root causal gene then affects other downstream levels during pathogenesis, including the core (or direct causal) gene \tilde{X}_4 , to ultimately induce a diagnosis Y . (b) We hypothesize that the causal effects of most root causes are small, but a few are large (red ellipse), in each patient with disease. As a result, the distribution of these *root causal effects* tends to be right skewed in disease.

43 of causes but, just like a machine usually breaks down due to one or a few root causal problems, the
44 root causal genes may only represent a small subset of the genes because the causal effects of only
45 a few root causes are large (Figure 1 (b)). We thus also seek to identify *patient-specific* root causal
46 genes in order to classify patients into meaningful biological subgroups each hopefully dictated by
47 only a small group of genes.

48 No existing method identifies root causal genes from data. Many algorithms focus on discover-
49 ing associational or predictive relations, sometimes visually represented as gene regulatory net-
50 works (Costa-Silva et al., 2017; Ellington et al., 2023). Other methods even identify causal rela-
51 tions (Friedman et al., 2000; Wang et al., 2023; Wen et al., 2023; Buschur et al., 2020), but none
52 pinpoint the *first* gene expression levels that ultimately generate the vast majority of pathogen-
53 esis. Simply learning a causal graph does not resolve the issue because causal graphs do not
54 summarize the effects of *unobserved* root causes, such as unmeasured environmental changes or
55 variants, that are needed to identify all root causal genes. We therefore define the Root Causal
56 Strength (RCS) score to identify all root causal genes unique to each patient. We then design the
57 Root Causal Strength using Perturbations (RCSP) algorithm that estimates RCS from bulk RNA-seq
58 under minimal assumptions by integrating Perturb-seq, or high throughput perturbation experi-
59 ments using CRISPR-based technologies coupled with single cell RNA-sequencing (Dixit et al., 2016;
60 Adamson et al., 2016; Datlinger et al., 2017). Experiments demonstrate marked improvements
61 in performance, when investigators have access to a large bulk RNA-seq dataset and a genome-
62 wide Perturb-seq dataset from a cell line of a disease-relevant tissue. Finally, application of the
63 algorithm to two complex diseases with disparate pathogeneses recovers an *omnigenic root causal*
64 *model*, where a small set of root causal genes drive pathogenesis but impact many downstream
65 genes within each patient. As a result, nearly all gene expression levels are correlated with the
66 diagnosis at the population level.

67 Results

68 We briefly summarize the Methods in the first two subsections.

69 Definitions

70 *Differential expression analysis* identifies differences in gene expression levels between groups Y
71 (Costa-Silva et al., 2017). A gene X_i may be differentially expressed due to multiple reasons. For
72 example, X_i may cause Y , or a confounder C may explain the relation between X_i and Y such
73 that $X_i \leftarrow C \rightarrow Y$. In this paper, we take expression analysis a step further by pinpointing *causal*
74 relations from expression levels regardless of the variable type of Y (discrete or continuous). We

75 in particular seek to discover *patient-specific root causal genes* from bulk RNA-seq data, which we
76 carefully define below.

77 We represent a biological system in bulk RNA-seq as a causal graph \mathbb{G} – such as in Figure 2 (a)
78 – where p vertices \tilde{X} represent true gene expression levels in a bulk sample and Y denotes the
79 patient symptoms or diagnosis. The set \tilde{X} contains thousands of genes in practice. Directed edges
80 between the vertices in \mathbb{G} refer to direct causal relations. We assume that gene expression causes
81 patient symptoms but not vice versa so that no edge from Y is directed towards \tilde{X} . The set $\text{Pa}(\tilde{X}_i)$
82 refers to the *parents* of $\tilde{X}_i \in \tilde{X}$, or those variables with an edge directed into \tilde{X}_i . For example,
83 $\text{Pa}(\tilde{X}_2) = \{\tilde{X}_1, \tilde{X}_3\}$ in Figure 2 (a). A *root vertex* corresponds to a vertex with no parents.

84 We can associate \mathbb{G} with the structural equation $\tilde{X}_i = f_i(\text{Pa}(\tilde{X}_i), E_i)$ for each $\tilde{X}_i \in \tilde{X}$ that links
85 each vertex to its parents and error term E_i (Pearl, 2009). The error term E_i is not simply a re-
86 gression residual but instead represents the conglomeration of unobserved explanatory variables
87 that only influence \tilde{X}_i , such as unobserved transcriptional regulators, certain genetic variants and
88 specific environmental conditions. We thus also include the error terms E in the directed graph
89 of Figure 2 (b). All root vertices are error terms and vice versa. The *root causes* of Y are the error
90 terms that cause Y , or have a directed path into Y . We define the *root causal strength* (RCS) of \tilde{X}_i
91 on Y as the following absolute difference (Figure 2 (c)):

$$\begin{aligned}\Phi_i &= \left| \mathbb{E}(Y|\text{Pa}(\tilde{X}_i), E_i) - \mathbb{E}(Y|\text{Pa}(\tilde{X}_i)) \right| \\ &= \left| \mathbb{E}(Y|\text{Pa}(\tilde{X}_i), \tilde{X}_i) - \mathbb{E}(Y|\text{Pa}(\tilde{X}_i)) \right|.\end{aligned}\tag{1}$$

92 We prove the last equality in the Methods. As a result, RCS Φ_i directly measures the contribution
93 of the gene \tilde{X}_i on Y according to its error term E_i , without recovering the error term values. The
94 algorithm does not impose distributional assumptions or functional restrictions such as additive
95 noise to estimate the error term values as an intermediate step. Moreover Φ_i is patient-specific
96 because the values of $\text{Pa}(\tilde{X}_i)$ and \tilde{X}_i may differ between patients. We have $\Phi_i = 0$ when E_i is not
97 a cause of Y , and we say that the gene \tilde{X}_i is a *patient-specific root causal gene* if $\Phi_i \gg 0$, or its
98 (conditional) root causal effect is large as depicted by the red ellipse in Figure 1 (b).

99 Algorithm

100 We propose an algorithm called Root Causal Strength using Perturbations (RCSP) that estimates
101 $\Phi = \{\Phi_1, \dots, \Phi_p\}$ from genes measured in both bulk RNA-seq and Perturb-seq datasets derived
102 from possibly independent studies but from the same tissue type. We rely on bulk RNA-seq instead
103 of single cell RNA-seq in order to obtain many samples of the label Y . We focus on statistical
104 estimation rather than statistical inference because $\Phi_i > 0$ when E_i causes Y under mild conditions,
105 so we reject the null hypothesis that $\Phi_i = 0$ for many genes if many gene expression levels cause Y .
106 However, just like a machine typically breaks down due to only one or a few root causal problems,
107 we hypothesize that only a few genes have large RCS scores $\Phi_i \gg 0$ even in complex disease.

108 Estimating Φ requires access to the true gene expression levels \tilde{X} and the removal of the effects
109 of confounding. We first control for batch effects representing unwanted sources of technical vari-
110 ation such as different sequencing platforms or protocols. We however can only obtain imperfect
111 counts X from RNA sequencing even within each batch (Figure 2 (d)). Measurement error intro-
112 duces confounding as well because it prevents us from exactly controlling for the causal effects
113 of the gene expression levels. Investigators usually mitigate measurement error by normalizing
114 the gene expression levels by sequencing depth. We show in the Methods that the Poisson distri-
115 bution approximates the measurement error distribution induced by the sequencing process to
116 high accuracy (Choudhary and Satija, 2022; Sarkar and Stephens, 2021). We leverage this fact to
117 eliminate the need for normalization by sequencing depth using an asymptotic argument where
118 the library size N approaches infinity. N takes on a value of at least ten million in bulk RNA-seq,
119 but we also empirically verify that the theoretical results hold well in the Supplementary Materials.
120 We thus eliminate the Poisson measurement error and batch effects by controlling for the batches
121 B but not N in non-linear regression models.

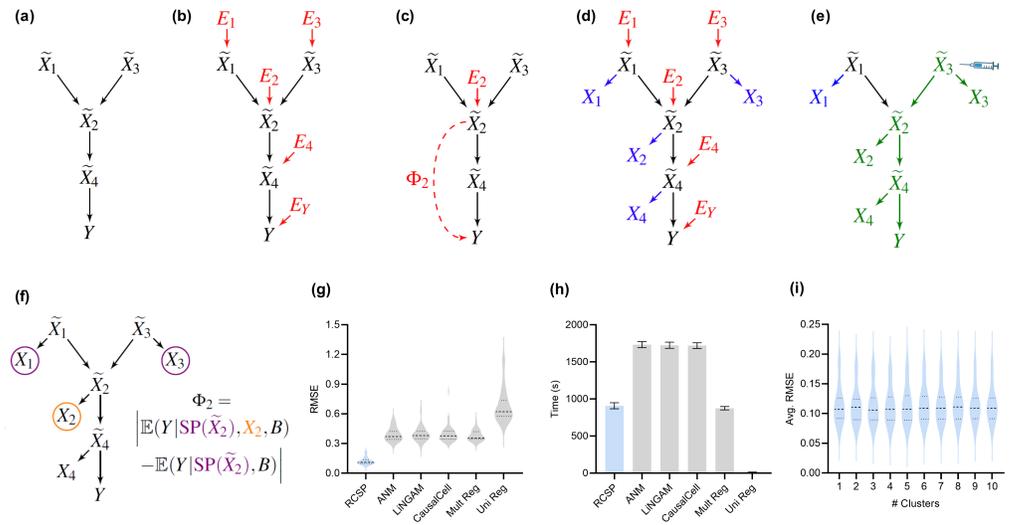


Figure 2. Method overview and synthetic data results. (a) We consider a latent causal graph over the true counts \tilde{X} . (b) We augment the graph with error terms E such that each $E_i \in E$ in red has an edge directed towards $\tilde{X}_i \in \tilde{X}$. (c) The RCS of \tilde{X}_2 , denoted by Φ_2 , quantifies the magnitude of the *conditional root causal effect*, or the strength of the causal effect from E_2 to Y conditional on $\text{Pa}(\tilde{X}_2)$. (d) We cannot observe \tilde{X} in practice but instead observe the noisy surrogates X in blue corrupted by Poisson measurement error. (e) Perturbing a variable such as \tilde{X}_3 changes the marginal distributions of downstream variables shown in green under mild conditions. (f) RCSP thus uses the perturbation data to identify (an appropriate superset of) the surrogate parents for each variable in order to compute Φ . (g) Violin plots show that RCSP achieved the smallest RMSE to the ground truth RCS values in the synthetic data. (h) RCSP also took about the same amount of time to complete as multivariate regression. Univariate regression only took 11 seconds on average, so its bar is not visible. Error bars denote 95% confidence intervals of the mean. (i) Finally, RCSP maintained low RMSE values regardless of the number of clusters considered.

122 We in particular show that Φ_i in Equation (1) is also equivalent to:

$$\Phi_i = \left| \mathbb{E}(Y | \text{SP}(\tilde{X}_i), X_i, B) - \mathbb{E}(Y | \text{SP}(\tilde{X}_i), B) \right|, \quad (2)$$

123 where $\text{SP}(\tilde{X}_i)$ refers to the *surrogate parents* of \tilde{X}_i , or the variables in X associated with $\text{Pa}(\tilde{X}_i) \subseteq \tilde{X}$.
 124 RCSP can identify (an appropriate superset of) the surrogate parents of each variable using per-
 125 turbation data because perturbing a gene changes the marginal distributions of its downstream
 126 effects – which the algorithm detects from data under mild assumptions (Figures 2 (e) and (f)).
 127 The algorithm thus only transfers the binary presence or absence of causal relations from the sin-
 128 gle cell to bulk data – rather than the exact functional relationships – in order to remain robust
 129 against discrepancies between the two data types; we empirically verify the robustness in the Sup-
 130 plementary Materials. RCSP finally performs the two non-linear regressions needed to estimate
 131 $\mathbb{E}(Y | \text{SP}(\tilde{X}_i), X_i, B)$ and $\mathbb{E}(Y | \text{SP}(\tilde{X}_i), B)$ for each Φ_i . We will compare Φ_i against Statistical Dependence
 132 (SD), a measure of correlational strength defined as $\Omega_i = |\mathbb{E}(Y | X_i, B) - \mathbb{E}(Y | B)|$ where we have re-
 133 moved the conditioning on $\text{SP}(\tilde{X}_i)$.

134 **In silico identification of root causal genes**

135 We simulated 30 bulk RNA-seq and Perturb-seq datasets from random directed graphs summa-
 136 rizing causal relations between gene expression levels. We performed single gene knock-down
 137 perturbations over 2500 genes and 100 batches. We obtained 200 cell samples from each per-
 138 turbation, and another 200 controls without perturbations. We therefore generated a total of
 139 $2501 \times 200 = 500,200$ single cell samples for each Perturb-seq dataset. We simulated 200 bulk RNA-
 140 seq samples. We compared RCSP against the Additive Noise Model (ANM) (Peters et al., 2014;

141 *Strobl and Lasko, 2023a*), the Linear Non-Gaussian Acyclic Model (LiNGAM) (*Peters et al., 2014*;
142 *Strobl and Lasko, 2022*), CausalCell (*Wen et al., 2023*), univariate regression residuals (Uni Reg),
143 and multivariate regression residuals (Multi Reg). The first two algorithms are state-of-the-art ap-
144 proaches used for error term extraction and, in theory, root causal discovery. See Methods for
145 comprehensive descriptions of the simulation setup and comparator algorithms.

146 We summarize accuracy results in Figure 2 (g) using the Root Mean Squared Error (RMSE) to
147 the ground truth Φ values. All statements about pairwise differences hold true at a Bonferonni cor-
148 rected threshold of 0.05/5 according to paired two-sided t-tests, since we compared RCSP against a
149 total of five algorithms. RCSP estimated Φ most accurately by a large margin. ANM and LiNGAM are
150 theoretically correct under their respective assumptions, but they struggle to outperform standard
151 multivariate regression due to the presence of measurement error in RNA-seq (Supplementary
152 Materials). Feature selection and causal discovery with CausalCell did not improve performance.
153 Univariate regression performed the worst, since it does not consider the interactions between
154 variables. RCSP achieved the lowest RMSE while completing in about the same amount of time as
155 multivariate regression on average (Figure 2 (h)). RCSP maintained the lowest RMSE even in the
156 cyclic case, and the performance of the algorithm remained robust to differences between the di-
157 rected graphs underlying the bulk RNA-seq and Perturb-seq data (Supplementary Materials). We
158 conclude that RCSP both scalably and accurately estimates Φ .

159 We will cluster the RCS values in real data to find patient subgroups. We therefore also per-
160 formed hierarchical clustering using Ward's method (*Ward Jr, 1963*) on the values of Φ estimated
161 by RCSP with the synthetic data. We then computed the RMSEs and averaged them within each
162 cluster. We found that RCSP maintained low average RMSE values regardless of the number of
163 clusters considered (Figure 2 (i)). We conclude that RCSP maintains accurate estimation of Φ across
164 different numbers of clusters.

165 **Oxidative stress in age-related macular degeneration**

166 We ran RCSP on a bulk RNA-seq dataset of 513 individuals with age-related macular degeneration
167 (AMD; GSE115828) and a Perturb-seq dataset of 247,914 cells generated from an immortalized
168 retinal pigment epithelial (RPE) cell line (*Ratnapriya et al., 2019*; *Replogle et al., 2022*). The Perturb-
169 seq dataset contains knockdown experiments of 2,077 genes overlapping with the genes of the
170 bulk dataset. We set the target Y to the Minnesota Grading System score, a measure of the severity
171 of AMD based on stereoscopic color fundus photographs. We always included age and sex as a
172 biological variable as covariates. We do not have access to the ground truth values of Φ in real data,
173 so we evaluated RCSP using seven alternative techniques. See Methods for a detailed rationale of
174 the evaluation of real data. RCSP outperformed all other algorithms in this dataset (Supplementary
175 Materials). We therefore only analyze the output of RCSP in detail here.

176 AMD is a neurodegenerative disease of the aging retina (*Hadziahmetovic and Malek, 2021*), so
177 age is a known root cause of the disease. We therefore determined if RCSP identified age as a
178 root cause. Note that RCSP does not need perturbation data of age to compute the RCS values
179 of age, since age has no parents in the directed graph. The algorithm estimated a heavy tailed
180 distribution of the RCS values indicating that most of the RCS values deviated away from zero
181 (Figure 3 (a)). The Deviation of the RCS (D-RCS), or the standard deviation from an RCS value of
182 zero, measures the tailedness of the distribution while preserving the unit of measurement. The
183 D-RCS of age corresponded to 0.46 – more than double that of the nearest gene (Figure 3 (d)). We
184 conclude that RCSP correctly detected age as a root cause of AMD.

185 Root causal genes typically affect many downstream genes before affecting Y . We therefore
186 expect to identify few root causal genes but many genes that correlate with Y . To evaluate this
187 hypothesis, we examined the distribution of D-RCS relative to the distribution of the Deviation of
188 Statistical Dependence (D-SD), or the standard deviation from an SD value of zero, in Figure 3 (b).
189 Notice that the histogram of D-RCS scores in Figure 3 (b) mimics a folded distribution of Figure 1 (b).
190 Thus, few D-RCS scores had large values implying the existence of only a few root causal genes. In

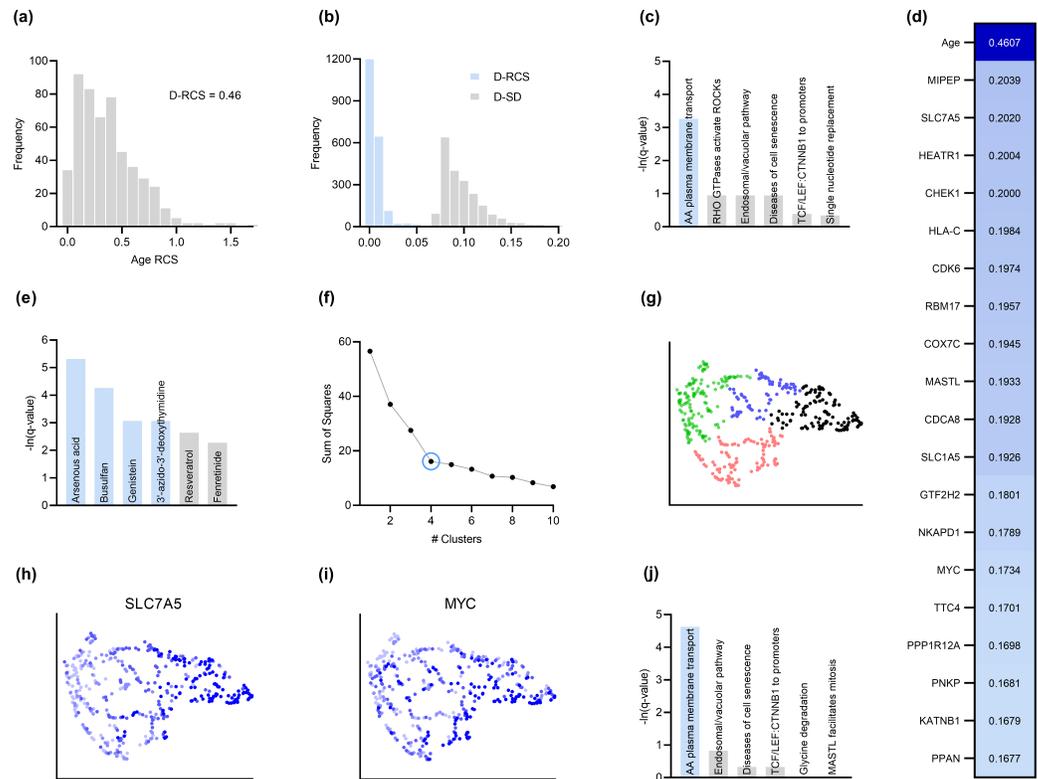


Figure 3. Analysis of AMD. (a) The distribution of the RCS scores of age deviated away from zero and had a composite D-RCS of 0.46. (b) However, the majority of gene D-RCS scores concentrated around zero, whereas the majority of gene D-SD scores concentrated around the relatively larger value of 0.10. Furthermore, the D-RCS scores of the genes in (d) mapped onto the “amino acid transport across the plasma membrane” pathway known to be involved in the pathogenesis of AMD in (c). Blue bars survived 5% FDR correction. (e) Drug enrichment analysis revealed four significant drugs, the later three of which have therapeutic potential. (f) Hierarchical clustering revealed four clear clusters according to the elbow method, which we plot by UMAP dimensionality reduction in (g). The RCS scores of the top genes in (d) increased only from the left to right on the first UMAP dimension (x-axis); we provide an example of SLC7A5 in (h) and one of three detected exceptions in (i). We therefore performed pathway enrichment analysis on the black cluster in (g) containing the largest RCS scores. (j) The amino acid transport pathway had a larger degree of enrichment in the black cluster as compared to the global analysis in (c).

191 contrast, most of the D-SD scores had relatively larger values concentrated around 0.10 implying
 192 the existence of many genes correlated with Y . We conclude that RCSP identified few root causal
 193 genes rather than many correlated genes for AMD.

194 The pathogenesis of AMD involves the loss of RPE cells. The RPE absorbs light in the back of
 195 the retina, but the combination of light and oxygen induces oxidative stress, and then a cascade
 196 of events such as immune cell activation, cellular senescence, drusen accumulation, neovascular-
 197 ization and ultimately fibrosis (*Barouch and Miller, 2007*). We therefore expect the root causal
 198 genes of AMD to include genes involved in oxidative stress during early pathogenesis. The gene
 199 MIPEP with the highest D-RCS score in Figure 3 (d) indeed promotes the maturation of oxidative
 200 phosphorylation-related proteins (*Shi et al., 2011*). The second gene SLC7A5 is a solute carrier that
 201 activates mTORC1 whose hyperactivation increases oxidative stress via lipid peroxidation (*Nachef*
 202 *et al., 2021; Go et al., 2020*). The gene HEATR1 is involved in ribosome biogenesis that is down-
 203 regulated by oxidative stress (*Turi et al., 2018*). The top genes discovered by RCSP thus identify
 204 pathways known to be involved in oxidative stress. We further verified that measurement error

205 did not explain their large D-RCS scores in Supplementary Materials.

206 We subsequently jointly analyzed the D-RCS values of all 2077 genes. We performed pathway
207 enrichment analysis that yielded one pathway “amino acid transport across the plasma membrane”
208 that passed an FDR threshold of 5% (Figure 3 (c)). The leading edge genes of the pathway included
209 the solute carriers SLC7A5 and SLC1A5. These two genes function in conjunction to increase the
210 efflux of essential amino acids out of the lysosome (*Nicklin et al., 2009; Beaumatin et al., 2019*).
211 Some of these essential amino acids like L-leucine and L-arginine activate mTORC1 that in turn
212 increases lipid peroxidation induced oxidative stress and the subsequent degeneration of the RPE
213 (*Nachef et al., 2021; Go et al., 2020*). We conclude that pathway enrichment analysis correctly
214 identified solute carrier genes involved in a known pathway promoting oxidative stress in AMD.

215 We next ran drug enrichment analysis with the D-RCS scores. The top compound arsenous
216 acid inhibits RPE proliferation (*Su et al., 2020*), but the other three significant drugs have therapeutic
217 potential (Figure 3 (e)). Busulfan decreases the requirement for intravitreal anti-VEGF injections
218 (*Dalvin et al., 2022*). Genistein is a protein kinase inhibitor that similarly attenuates neovasculariza-
219 tion (*Kinoshita et al., 2014*) and blunts the effect of ischemia on the retina (*Kamalden et al., 2011*).
220 Finally, a metabolite of the antiviral agent 3'-azido-3'-deoxythymidine inhibits neovascularization
221 and mitigates RPE degeneration (*Narendran et al., 2020*). We conclude that the D-RCS scores identified
222 promising drugs for the treatment of AMD.

223 Hierarchical clustering and UMAP dimensionality reduction on the patient-specific RCS values
224 revealed four clear clusters of patients by the elbow method on the sum of squares plot (Figures
225 3 (f) and (g), respectively). The RCS scores of most of the top genes exhibited a clear gradation
226 increasing only from the left to the right hand side of the UMAP embedding; we plot an example
227 in Figure 3 (h). We found three exceptions to this rule among the top 30 genes (example in Figure
228 3 (i) and see Supplementary Materials). RCSP thus detected genes with large RCS scores primarily
229 in the black cluster of Figure 3 (g). Pathway enrichment analysis within this cluster alone yielded
230 supra-significant results on the same pathway detected in the global analysis (Figure 3 (j) versus
231 Figure 3 (c)). Furthermore, drug enrichment analysis results by cluster confirmed that patients
232 in the black cluster with many root causal genes are likely the hardest to treat (Supplementary
233 Materials). We conclude that RCSP detected a subgroup of patients whose root causal genes have
234 large RCS scores and involve known pathogenic pathways related to oxidative stress.

235 **T cell infiltration in multiple sclerosis**

236 We next ran RCSP on 137 samples collected from CD4+ T cells of multiple sclerosis (MS; GSE137143)
237 as well as Perturb-seq data of 1,989,578 undifferentiated blast cells that can be induced to differ-
238 entiate into lymphoblasts, or the precursors of T cells and other lymphocytes (*Kim et al., 2021;*
239 *Replegle et al., 2022*). We set the target Y to the Expanded Disability Status Scale score, a measure
240 of MS severity. RCSP outperformed all other algorithms in this dataset as well (Supplementary
241 Materials).

242 MS progresses over time, and RCSP correctly detected age as a root cause of MS severity with
243 RCS values deviating away from zero (Figure 4 (a)). The distribution of gene D-RCS scores concen-
244 trated around zero with a long tail, whereas the distribution of gene D-SD scores concentrated
245 around a relatively larger value of 0.3 (Figure 4 (b)). RCSP thus detected an omnigenic root causal
246 model with a few root causal genes but many correlated genes.

247 MS is an inflammatory neurodegenerative disease that damages the myelin sheaths of nerve
248 cells in the brain and spinal cord. T cells may mediate the inflammatory process by crossing a
249 disrupted blood brain barrier and repeatedly attacking the myelin sheaths (*Fletcher et al., 2010*).
250 Damage induced by the T cells also perturbs cellular homeostasis and leads to the accumulation
251 of misfolded proteins (*Andhavarapu et al., 2019*). The root causal genes of MS thus likely include
252 genes involved in T cell infiltration across the blood brain barrier.

253 Genes with the highest D-RCS scores included MNT, CERCAM and HERPUD2 (Figure 4 (d)). MNT
254 is a MYC antagonist that modulates the proliferative and pro-survival signals of T cells after en-

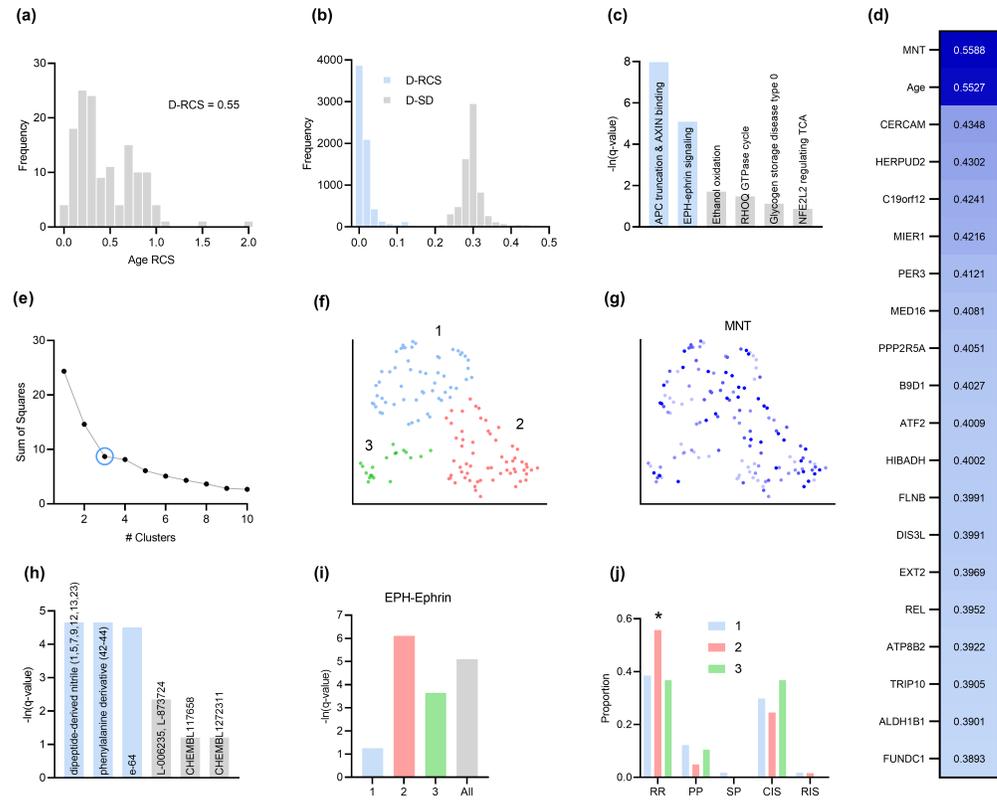


Figure 4. Analysis of MS. (a) The distribution of the RCS scores of age deviated away from zero with a composite D-RCS of 0.55. (b) The distribution of D-RCS concentrated around zero, whereas the distribution of D-SD concentrated around 0.3. (d) RCSP identified many genes with large D-RCS scores that in turn mapped onto known pathogenic pathways in MS in (c). Hierarchical clustering revealed three clusters in (e), which we plot in two dimensions with UMAP in (f). Top genes did not correlate with either dimension of the UMAP embedding; we provide an example of the MNT gene in (g). (h) Drug enrichment analysis in the green cluster implicated multiple cathepsin inhibitors. Finally, EPH-ephrin signaling survived FDR correction in (c) and was enriched in the pink cluster in (i) which contained more MS patients with the relapsing-remitting subtype in (j); subtypes include relapse-remitting (RR), primary progressive (PP), secondary progressive (SP), clinically isolated syndrome (CIS), and radiologically isolated syndrome (RIS).

255 gagement of the T cell receptor (*Gnanaprakasam and Wang, 2017*). Similarly, CERCAM is an ad-
 256 adhesion molecule expressed at high levels in microvessels of the brain that increases leukocyte
 257 transmigration across the blood brain barrier (*Starzyk et al., 2000*). HERPUD2 is involved in the
 258 endoplasmic-reticulum associated degradation of unfolded proteins (*Kokame et al., 2000*). Genes
 259 with the highest D-RCS scores thus serve key roles in known pathogenic pathways of MS.

260 We found multiple genes with high D-RCS scores in MS, in contrast to AMD where age domi-
 261 nated (Figure 4 (d) versus Figure 3 (d)). Measurement error did not account for the high scores
 262 (Supplementary Materials). We performed pathway enrichment analysis using the D-RCS scores
 263 of all genes and discovered two significant pathways at an FDR corrected threshold of 5%: “ade-
 264 nomatous polyposis coli (APC) truncation mutants have impaired AXIN binding” and “EPH-ephrin
 265 signaling” (Figure 4 (c)). APC and AXIN are both members of the Wnt signaling pathway and regu-
 266 late levels of beta-catenin (*Spink et al., 2000*). Furthermore, inhibition of Wnt/beta-catenin causes
 267 CD4+ T cell infiltration into the central nervous system via the blood brain barrier in MS (*Lengfeld
 268 et al., 2017*). Ephrins similarly regulate T cell migration into the central nervous system (*Luo et al.,
 269 2016*) and are overexpressed in MS lesions (*Sobel, 2005*). The APC-AXIN and EPH-ephrin pathways
 270 are thus consistent with the known pathophysiology of central nervous system T cell infiltration in

271 MS.

272 We subsequently performed hierarchical clustering of the RCS scores. The within cluster sum
273 of squares plot in Figure 4 (e) revealed the presence of three clusters by the elbow method. We
274 plot the three clusters in a UMAP embedding in Figure 4 (f). The clusters did not show a clear
275 relationship with MS symptom severity (Supplementary Materials) or the levels of the top most
276 genes of Figure 4 (d); we plot the MNT gene as an example in Figure 4 (g). However, further analyses
277 with additional genes revealed that the distribution of many lower ranked genes governed the
278 structure of the UMAP embedding (Supplementary Materials). The D-RCS scores of each cluster
279 also implicated different mechanisms of T cell pathology including APC-AXIN in the green cluster,
280 disturbed T cell homeostasis in the pink cluster and platelet enhanced T cell autoreactivity in the
281 blue cluster (Supplementary Materials).

282 Global drug enrichment analysis did not yield any significant drugs even at a liberal FDR thresh-
283 old of 10%. We thus ran drug enrichment analysis in each cluster of Figure 4 (f). The blue and pink
284 clusters again did not yield significant drugs. However, the third green cluster identified the cys-
285 teine cathepsin inhibitors dipeptide-derived nitriles, phenylalanine derivatives, e-64, L-006235 and
286 L-873724 (Figure 4 (h)); statistical significance of the first three held even after correcting for multi-
287 ple comparisons with the Bonferroni adjustment of 0.05/4 on the q-values. The leading edge genes
288 of the significant drugs included the cathepsins CTSL, CTSS and CTSB exclusively. These drug en-
289 richment results corroborate multiple experimental findings highlighting the therapeutic efficacy
290 of cathepsin inhibitors in a subgroup of MS patients responsive to interferon therapy (*Haves-Zburof*
291 *et al., 2011; Burster et al., 2007*).

292 Prior research has also shown that EPH-ephrin signaling is more prevalent in relapsing-remitting
293 multiple sclerosis than in other subtypes of the disease (*Golan et al., 2021*). EPH-ephrin signaling
294 survived FDR correction in our analysis (Figure 4 (c)). Furthermore, the pathway was more enriched
295 in the pink cluster than in the other two (Figure 4 (i)). The pink cluster indeed contained a higher
296 proportion of patients with the relapsing-remitting subtype (Figure 4 (j)). RCSP thus precisely iden-
297 tified the enrichment of EPH-ephrin signaling in the correct subtype of MS.

298 Discussion

299 We presented a framework for identifying root causal genes, or the gene expression levels directly
300 regulated by root causes with large causal effects on Y , by modeling the root causes using the
301 error terms of structural equation models. Each error term represents the conglomeration of un-
302 observed root causes, such as genetic variants or environmental conditions, that directly cause a
303 specific gene. We however do not have access to many of the error terms in practice, so we in-
304 troduced the root causal strength (RCS) score, or the magnitude of the conditional causal effect
305 of each error term, which we can compute using gene expression levels alone. The RCSP algo-
306 rithm computes RCS given knowledge of the causal ancestors of each variable, which we obtained
307 by Perturb-seq. RCSP only transfers the causal structure (binary cause-effect relations) from the
308 single cell to bulk data rather than the exact functional relationships in order to remain robust
309 against discrepancies between the two data types. Results with the synthetic data demonstrated
310 marked improvements over existing alternatives. The algorithm also recovered only a few root
311 causal genes that play key roles in known pathogenic pathways and implicate therapeutic drugs in
312 both AMD and MS.

313 We detected a modest number of root causal genes in both AMD and MS, but virtually all genes
314 were correlated with Y . This omnigenic model, where “omni-” refers to the nearly all genes corre-
315 lated with Y , differs from the omnigenic model involving *core genes* (*Boyle et al., 2017*). Boyle et al.
316 define core genes as genes that directly affect disease risk. The authors further elaborate that many
317 *peripheral genes* affect the functions of a modest number of core genes, so the peripheral genes
318 often explain most of disease heritability. In contrast, root causal genes may not directly cause Y
319 but lie substantially upstream of Y in the causal graph. The error terms of upstream root causal
320 genes affect many downstream genes that include both ancestors and non-ancestors of Y (Figure

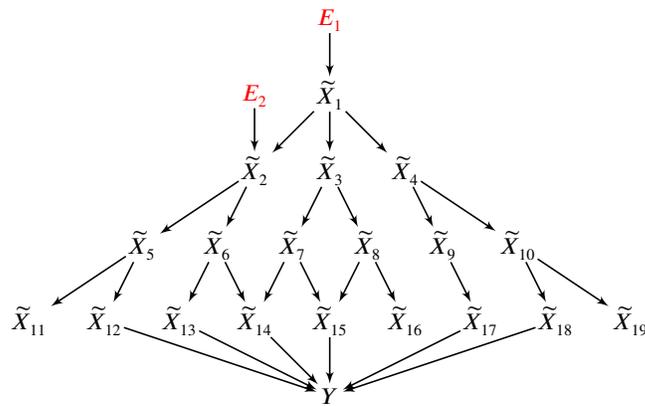


Figure 5. In this example, two root causal genes \tilde{X}_1 and \tilde{X}_2 affect many downstream genes and ultimately cause Y . Thus all genes $\tilde{X}_1, \dots, \tilde{X}_{19}$ correlate with Y , but only \tilde{X}_1 and \tilde{X}_2 have large root causal effects on Y . The omnigenic root causal model posits that only a few root causal genes affect many downstream genes, so that nearly all genes are correlated with Y . Causal genetic variants can directly cause Y or cause any gene expression level that causes Y – including those with small root causal effects – but only \tilde{X}_1 and \tilde{X}_2 have large root causal effects on Y due to genetic *and non-genetic* root causes modeled by E_1 and E_2 . In contrast, the core gene model assumes only a few direct causal genes $\tilde{X}_{12}, \tilde{X}_{13}, \tilde{X}_{14}, \tilde{X}_{15}, \tilde{X}_{17}, \tilde{X}_{18}$. These core genes do not account for the deleterious causal effects of E_1 and E_2 on $\tilde{X}_{11}, \tilde{X}_{16}$ and \tilde{X}_{19} .

321 5). These downstream genes contain traces of the root causal gene error terms that induce the
 322 many correlations with Y . The root causal model thus assumes sparsity in upstream root causal
 323 genes, whereas the core gene model assumes sparsity in the downstream direct causal genes¹.
 324 Further, each causal genetic variant tends to have only a small effect on disease risk in complex
 325 disease because the variant can directly cause Y or directly cause any causal gene including those
 326 with small root causal effects on Y ; thus, all error terms that cause Y can model genetic effects
 327 on Y . However, the root causal model further elaborates that genetic *and non-genetic factors* often
 328 combine to produce a few root causal genes with large root causal effects, where non-genetic fac-
 329 tors typically account for the majority of the large effects in complex disease. Many variants may
 330 therefore cause many genes in diseases with only a few root causal genes. We finally emphasize
 331 that the root causal model accounts for all deleterious effects of the root causal genes, whereas the
 332 core gene model only captures the deleterious effects captured by the diagnosis Y . For example,
 333 the *disease* of diabetes causes retinopathy, but retinopathy is not a part of the diagnostic criteria
 334 of diabetes. As a result, the gene expression levels that cause retinopathy but not the *diagnosis* of
 335 diabetes are not core genes, even though they are affected by the root causal genes. The sparsity
 336 of the root causal genes, the focus on the combined effects of genetic and non-genetic root causes,
 337 and the ability to account for root causal effects not represented by the target Y motivate us to
 338 use the phrase *omnigenic root causal model* in order to distinguish it from the omnigenic core gene
 339 model.

340 We identified root causal genes without imposing parametric assumptions using the RCS met-
 341 ric. Prior measures of root causal effect require restrictive functional relations, such as linear re-
 342 lations or additive noise, and continuous random variables (*Strobl and Lasko, 2022; Strobl et al.,*
 343 *2024; Strobl and Lasko, 2023a*). These restrictions ensure exact identifiability of the underlying
 344 causal graph and error terms. However, real RNA-seq is obtained from a noisy sequencing process
 345 and contains count data arguably corrupted by Poisson measurement error (*Sarkar and Stephens,*
 346 *2021*). The Poisson measurement error introduces confounding that precludes exact recovery of
 347 the underlying error terms. The one existing root causal discovery method that can handle Poisson
 348 measurement error uses single cell RNA-seq, estimates negative binomial distribution parameters
 349 and cannot scale to the thousands of genes required for meaningful root causal detection (*Strobl*

¹The omnigenic root causal model makes no statement about the number of direct causal genes, so direct causal genes may be sparse or dense.

350 **and Lasko, 2023b**). RCSP rectifies the deficiencies of these past approaches by ensuring accurate
351 root causal detection even in the presence of the counts, measurement error and high dimension-
352 ality of RNA-seq.

353 This study carries other limitations worthy of addressing in future work. The RCS score impor-
354 tantly quantifies root causal strength rather than root causal effect. As a result, the method cannot
355 be used to identify the direction of root causal effect unconditional on the parents. The root causal
356 effect and signed RCS (or expected conditional root causal effect) do not differ by much in practice
357 (Supplementary Materials), but future work may focus on exactly identifying both the strength and
358 direction of the unconditional causal effects of the error terms. Furthermore, RCS achieves patient
359 but not cell-specificity because the algorithm relies on phenotypic labels obtained from bulk RNA-
360 seq. RCSP thus cannot identify the potentially different root causal genes present within distinct
361 cell populations. Modern genome-wide Perturb-seq datasets also adequately perturb and mea-
362 sure only a few thousand, rather than all, gene expression levels. RCSP can only identify root causal
363 genes within this perturbed and measured subset. Fourth, RCSP accounts for known batch effects
364 and measurement error but cannot adjust for unknown confounding. Finally, RCSP assumes a
365 directed acyclic graph. We can transform a directed graph with cycles into an acyclic one under
366 equilibrium, but real biological distributions vary across time (*Spirtes, 1995; Bongers et al., 2021*).
367 Future work should thus aim to estimate cell-specific root causal effects under latent confounding
368 and time-varying distributions.

369 In conclusion, RCSP integrates bulk RNA-seq and Perturb-seq to identify patient-specific root
370 causal genes under a principled causal inference framework using the RCS score. RCS quantifies
371 root causal strength implicitly without requiring normalization by sequencing depth or direct ac-
372 cess to the error terms of a structural equation model. The algorithm identifies the necessary
373 causal relations to compute RCS using reliable high throughput perturbation data rather than ob-
374 servational data alone. The RCS scores often suggest an omnigenic root causal model of disease.
375 Enrichment analyses with the RCS scores frequently reveal pathogenic pathways and drug candi-
376 dates. We conclude that RCSP is a novel, accurate, scalable and disease-agnostic procedure for
377 performing patient-specific root causal gene discovery.

378 **Methods and Materials**

379 **Background on Causal Discovery**

380 We denote a singleton variable like \tilde{X}_i with italics and sets of variables like $\tilde{\mathbf{X}}$ with bold italics. We
381 can represent a causal process using a *structural equation model* (SEM) linking the $p + 1$ variables in
382 $\mathbf{Z} = \tilde{\mathbf{X}} \cup Y$ using a series of deterministic functions:

$$383 \quad Z_i = f_i(\text{Pa}(Z_i), E_i), \quad \forall Z_i \in \mathbf{Z} \quad (3)$$

384 where f_i is a function of the *parents*, or direct causes, of Z_i and an error term $E_i \in \mathbf{E}$. The error
385 terms \mathbf{E} are mutually independent. We will use the terms *vertex* and *variable* interchangeably. A
386 *root vertex* corresponds to a vertex without any parents. On the other hand, a *terminal* or *sink vertex*
387 is not a parent of any other vertex.

388 We can associate a directed graph to \mathbf{Z} by drawing a directed edge from each member of $\text{Pa}(Z_i)$
389 to Z_i for all $Z_i \in \mathbf{Z}$. A *directed path* from Z_i to Z_j corresponds to a sequence of adjacent directed
390 edges from Z_i to Z_j . If such a path exists (or $Z_i = Z_j$), then Z_i is an *ancestor* of Z_j and Z_j is a
391 *descendant* of Z_i . We collate all ancestors of Z_i into the set $\text{Anc}(Z_i)$. A *cycle* occurs when there
392 exists a directed path from Z_i to Z_j and the directed edge $Z_j \rightarrow Z_i$. A *directed acyclic graph* (DAG)
393 contains no cycles. We *augment* a directed graph by including additional vertices \mathbf{E} and drawing a
394 directed edge from each $E_i \in \mathbf{E}$ to X_i except when $X_i = E_i$ is already a root vertex. We consider an
395 augmented DAG \mathbb{G} throughout the remainder of this manuscript.

396 The vertices Z_i and Z_j are *d-connected* given $\mathbf{W} \subseteq \mathbf{Z} \setminus \{Z_i, Z_j\}$ in \mathbb{G} if there exists a path between
397 Z_i and Z_j such that every collider on the path is an ancestor of \mathbf{W} and no non-collider is in \mathbf{W} . The

397 vertices are *d-separated* if they are not d-connected. Any DAG associated with the SEM in Equation
 398 (3) also obeys the *global Markov property* where Z_i and Z_j are conditionally independent given
 399 \mathbf{W} if they are d-separated given \mathbf{W} . The term *d-separation faithfulness* refers to the converse of the
 400 global Markov property where conditional independence implies d-separation. A distribution
 401 obeys *unconditional d-separation faithfulness* when we can only guarantee d-separation faithfulness
 402 when $\mathbf{W} = \emptyset$.

403 Causal Modeling of RNA Sequencing

404 Performing causal discovery requires careful consideration of the underlying generative process.
 405 We therefore propose a causal model for RNA-seq. We differentiate between the biology and the
 406 RNA sequencing technology.

407 We represent a snapshot of a biological causal process using an SEM over $\tilde{\mathbf{X}} \cup Y$ obeying Equa-
 408 tion (3). We assume that the phenotypic target Y is a terminal vertex so that gene expression causes
 409 phenotype but not vice versa. Each $\tilde{X}_i \in \tilde{\mathbf{X}}$ corresponds to the total number of RNA molecules of
 410 a unique gene in a single cell or bulk tissue sample. The error terms model root causes that are
 411 outside of gene expression, such as genetic variation or environmental factors. Moreover, the re-
 412 lation from gene expression to Y is stochastic because $Y = f_Y(\text{Pa}(Y), E_Y)$, where E_Y introduces the
 413 stochasticity. Two individuals may therefore have the exact same error term values over $\tilde{\mathbf{X}}$ but
 414 different instantiations of Y .

415 We unfortunately cannot observe $\tilde{\mathbf{X}}$ in practice but instead measure a corrupted count \mathbf{X} using
 416 single cell or bulk RNA-seq technology. We derive the measurement error distribution from first
 417 principles. We map an exceedingly small fraction of each $\tilde{X}_i \in \tilde{\mathbf{X}}$ within a sample at unequal
 418 coverage. Let π_{ij} denote the probability of mapping one molecule of \tilde{X}_i in batch j so that $\sum_{i=1}^p \pi_{ij}$
 419 is near zero. The law of rare events ([Papoulis, 1984](#)) implies that the Poisson distribution well-
 420 approximates the library size N so that $N \sim \text{Pois}(\sum_{i=1}^p \tilde{X}_i \pi_{ij})$.

421 We write the probability of mapping \tilde{X}_i in a given sample as:

$$P_{ij} = \frac{\tilde{X}_i \pi_{ij}}{\sum_{i=1}^p \tilde{X}_i \pi_{ij}}.$$

422 This proportion remains virtually unchanged when sampling without replacement because $N \ll$
 423 $\sum_{i=1}^p \tilde{X}_i$ with small $\sum_{i=1}^p \pi_{ij}$. We can therefore approximate sampling *without* replacement by sam-
 424 pling *with* replacement using a multinomial: $\mathbf{X} \sim \text{MN}(N; P_{1j}, \dots, P_{pj})$. This multinomial and the Pois-
 425 son distribution over N together imply that the marginal distribution of each $X_i \in \mathbf{X}$ follows an
 426 independent Poisson distribution centered at $(\sum_{i=1}^p \tilde{X}_i \pi_{ij}) P_{ij} = \tilde{X}_i \pi_{ij}$, or:

$$X_i \sim \text{Pois}(\tilde{X}_i \pi_{ij}). \quad (4)$$

427 We conclude that the measurement error distribution follows a Poisson distribution to high ac-
 428 curacy. Multiple experimental results already corroborate this theoretical conclusion ([Grün et al.,
 429 2014; Sarkar and Stephens, 2021; Choudhary and Satija, 2022](#)).

430 We can represent the biology and the RNA sequencing in a single DAG over $\mathbf{X} \cup \tilde{\mathbf{X}} \cup B \cup Y$, where
 431 B denotes the batch, and Y the target variable representing patient symptoms or diagnosis. We
 432 provide a toy example in Figure 6. We draw \mathbb{G} over \mathbf{Z} in black and make each $\tilde{X}_i \in \tilde{\mathbf{X}}$ a parent of
 433 $X_i \in \mathbf{X}$ in blue. We then include the root vertex B as a parent of all members of \mathbf{X} in green. We
 434 augment this graph with the error terms of $\tilde{\mathbf{X}}$ in red and henceforth refer to the augmented DAG
 435 as \mathbb{G} . Repeated draws from the represented causal process generates a dataset.

436 No Need for Normalization by Sequencing Depth

437 We provide an asymptotic argument that eliminates the need for normalization by sequencing
 438 depth when estimating conditional expectations using bulk RNA-seq. The argument applies to the
 439 conditional expectations as a whole rather than their individual parameters.

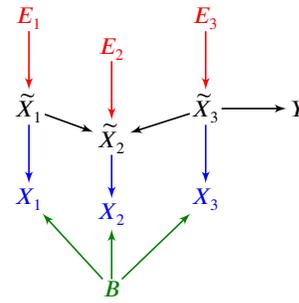


Figure 6. An example of a DAG over $X \cup \tilde{X} \cup B \cup Y$ augmented with the error terms E . The observed vertices X denote counts corrupted by batch B effects and Poisson measurement error.

440 We want to recover the causal relations between \tilde{X} by removing batch B and depth N effects
 441 from the dataset because they correspond to the sequencing process rather than the underlying
 442 biology. We first consider removing sequencing depth by finding stably expressed housekeeping
 443 genes. Let \tilde{A} denote the set of housekeeping genes where $\tilde{X}_i = \tilde{x}_i$ is a constant for each $\tilde{X}_i \in$
 444 \tilde{A} ; similarly A refers to the corresponding set with Poisson measurement error. Let $N = n$ be
 445 large enough such that $\sum_{X_i \in A} x_i > 0$ for each sample. Then dividing by $L \triangleq \sum_{X_i \in A} X_i$ controls for
 446 sequencing depth in the following sense:

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{X_i}{\sum_{X_i \in A} X_i} &= \lim_{N \rightarrow \infty} \frac{X_i/N}{\sum_{X_i \in A} X_i/N} = \frac{P_{ij}}{\sum_{X_i \in A} P_{ij}} \\ &= \frac{\tilde{X}_i \pi_{ij} / \sum_{i=1}^p \tilde{X}_i \pi_{ij}}{\sum_{\tilde{X}_i \in \tilde{A}} \tilde{x}_i \pi_{ij} / \sum_{i=1}^p \tilde{X}_i \pi_{ij}} = \frac{\tilde{X}_i \pi_{ij}}{\sum_{\tilde{X}_i \in \tilde{A}} \tilde{x}_i \pi_{ij}}, \end{aligned}$$

447 where we have divided $\tilde{X}_i \pi_{ij}$ by a constant in the last term. Thus, dividing by L removes measure-
 448 ment error within each batch as $N \rightarrow \infty$. We assume that N is so large that the approximation
 449 error is negligible. We only invoke the assumption in bulk RNA-seq, where the library size N is on
 450 the order of at least tens of millions.

451 We do not divide by L in practice because we may have $L = 0$ with finite N . We instead always
 452 include $L \cup B$ in the predictor set of downstream regressions. Conditioning on $L \cup B$ ensures that
 453 all downstream regressions mitigate depth and batch effects with adequate sequencing depth, or
 454 that $\mathbb{E}(Y|\tilde{U}, B) = \mathbb{E}(Y|U, L, B)$ for any $\tilde{U} \subseteq \tilde{X}$ as $N \rightarrow \infty$. The equality holds almost surely under a
 455 mild smoothness condition:

456 **Lemma 1.** Assume Lipschitz continuity of the conditional expectation for all $N \geq n_0$:

$$\mathbb{E} \left| \mathbb{E}(Y|\tilde{U}) - \mathbb{E}(Y|U, L, B) \right| \leq \mathbb{E} C_N \left| \tilde{U} - \frac{U}{dL} \right|,$$

457 where $d = \frac{\pi_{UB}}{\sum_{\tilde{X}_i \in \tilde{A}} \tilde{x}_i \pi_{iB}}$, $C_N \in O(1)$ is a positive constant, and we have taken an outer expectation on both
 458 sides. Then $\mathbb{E}(Y|\tilde{U}) = \lim_{N \rightarrow \infty} \mathbb{E}(Y|U, L, B)$ almost surely.

459 We delegate proofs to the Supplementary Materials unless proven here in the Methods. Note that
 460 $\lim_{N \rightarrow \infty} \frac{U}{dL} = \tilde{U}$, so the Lipschitz assumption intuitively means that accurate estimation of \tilde{U} implies
 461 accurate estimation of $\mathbb{E}(Y|\tilde{U})$. Furthermore, conditioning on the library size N instead of L can
 462 introduce spurious dependencies because N depends on all of the genes rather than just the stably
 463 expressed ones.

464 We now eliminate the need to condition on L . Note that L is a sum of independent Poisson
 465 distributions given B per Expression (4). This implies $Y \perp\!\!\!\perp L|(U, B)$ for any N , so that $\mathbb{E}(Y|\tilde{U}) =$
 466 $\lim_{N \rightarrow \infty} \mathbb{E}(Y|U, L, B) = \lim_{N \rightarrow \infty} \mathbb{E}(Y|U, B)$ almost surely. We have proved:

467 **Theorem 1.** Consider the same assumption as Lemma 1. Then $\mathbb{E}(Y|\tilde{U}) = \lim_{N \rightarrow \infty} \mathbb{E}(Y|U, B)$ almost
 468 surely, where we have eliminated the conditioning on L .

469 We emphasize again that these equalities hold for the conditional expectation but *not* for the re-
 470 gression parameters; the regression parameters do not converge in general unless we divide by L .
 471 We will only need to estimate conditional expectations in order to identify root causal genes.

472 Identifying Root Causal Genes

473 We showed how to overcome Poisson measurement error without sequencing depth normaliza-
 474 tion in the previous section. We leverage this technique to define a measure for identifying the
 475 root causal genes of Y .

476 Definitions

477 A *root cause* of Y corresponds to a root vertex that is an ancestor of Y in \mathbb{G} . All root vertices are
 478 error terms in an augmented graph. We define the *root causal effect* of any $E_i \in E$ on Y as $\Upsilon_i \triangleq$
 479 $\mathbb{P}(Y|E_i) - \mathbb{P}(Y)$ (Strobl, 2024; Strobl and Lasko, 2023c).

480 We can identify root causes using the following result:

481 **Proposition 1.** *If $E_i \perp\!\!\!\perp Y$ or $E_i \perp\!\!\!\perp Y|\text{Pa}(\tilde{X}_i)$ (or both), then E_i is a root cause of Y .*

482 We can also claim the backward direction under d-separation faithfulness. We however avoid mak-
 483 ing this additional assumption because real biological data may not arise from distributions obey-
 484 ing d-separation faithfulness in practice (Strobl, 2022).

485 Proposition 1 implies that E_i is a root cause of Y when:

$$\Delta_i \triangleq \mathbb{P}(Y|\text{Pa}(\tilde{X}_i), E_i) - \mathbb{P}(Y|\text{Pa}(\tilde{X}_i)) \neq 0.$$

486 The above quantity corresponds to the *conditional root causal effect* but not the root causal effect
 487 Υ_i due to the extra conditioning on $\text{Pa}(\tilde{X}_i)$. The two terms may also differ in direction; if $\Delta_i > 0$, then
 488 this does not imply that $\Upsilon_i > 0$, and similarly for negative values. The two variables thus represent
 489 different quantities but – in terms of priority – we would estimate Υ_i when we have nonzero Δ_i .
 490 Experimental results indicate that Υ_i and Δ_i take on similar values and agree in direction about
 491 95% of the time in practice (Supplementary Materials).

492 We now encounter two challenges. First, the quantities Υ_i and Δ_i depend on the unknown error
 493 term E_i . We can however substitute E_i with \tilde{X}_i in Δ_i due to the following result:

494 **Proposition 2.** *We have $\mathbb{P}(Y|E_i, \text{Pa}(\tilde{X}_i)) = \mathbb{P}(Y|\tilde{X}_i, \text{Pa}(\tilde{X}_i))$ under Equation (3).*

495 We can thus compute the conditional root causal effect Δ_i without access to the error terms:

$$\begin{aligned} \Delta_i &= \mathbb{P}(Y|\text{Pa}(\tilde{X}_i), E_i) - \mathbb{P}(Y|\text{Pa}(\tilde{X}_i)) \\ &= \mathbb{P}(Y|\text{Pa}(\tilde{X}_i), \tilde{X}_i) - \mathbb{P}(Y|\text{Pa}(\tilde{X}_i)). \end{aligned}$$

496 We can determine the root causal status of E_i on Y when $\Delta_i \neq 0$ per Proposition 1. Nevertheless,
 497 the term “root cause” in colloquial language refers to two concepts simultaneously: a root vertex
 498 that causes Y *and* has a large causal effect on Y . We thus say that \tilde{X}_i is a *root causal gene* of Y if
 499 $\Delta_i \gg 0$.

500 The second challenge involves computing the non-parametric probability distributions of Δ_i ,
 501 which come at a high cost. We thus define the analogous expected version by:

$$\begin{aligned} \Gamma_i &\triangleq \int y \left[p(y|\text{Pa}(\tilde{X}_i), \tilde{X}_i) - p(y|\text{Pa}(\tilde{X}_i)) \right] dy \\ &= \mathbb{E}(Y|\text{Pa}(\tilde{X}_i), \tilde{X}_i) - \mathbb{E}(Y|\text{Pa}(\tilde{X}_i)) \\ &= \mathbb{E}(Y|\text{SP}(\tilde{X}_i), X_i, B) - \mathbb{E}(Y|\text{SP}(\tilde{X}_i), B), \end{aligned}$$

502 where $p(Y)$ denotes the density of Y . Observe that if $\Delta_i = 0$, then $\Gamma_i = 0$. The converse is not
 503 true but likely to hold in real data when a change in the probability distribution also changes its
 504 expectation. The set $\text{SP}(\tilde{X}_i) \subseteq X$ denotes the *surrogate parents* of \tilde{X}_i corresponding to the variables
 505 in X associated with $\text{Pa}(\tilde{X}_i) \subseteq \tilde{X}$. The last equality holds almost surely as $N \rightarrow \infty$ by Theorem 1.

506 We call $\Phi_i \triangleq |\Gamma_i|$ the *Root Causal Strength* (RCS) of \tilde{X}_i on Y . The RCS obtains a unique value
507 $\Phi_i = \phi_{ij}$ for each patient j . We say that \tilde{X}_i is a root causal gene of Y for patient j if $\phi_{ij} \gg 0$, since
508 we posit a right skewed distribution of conditional root causal effects for each patient as in Figure
509 1 (b). We combine the RCS scores across a set of n samples using the Deviation of the RCS (D-RCS)
510 $\sqrt{\frac{1}{n} \sum_{j=1}^n \phi_{ij}^2}$, or the standard deviation of RCS from zero. We may compute D-RCS for each cluster or
511 globally across all patients depending on the context. We thus likewise say that \tilde{X}_i is a root causal
512 gene for a cluster of patients or all patients in a sample if its corresponding D-RCS score for the
513 cluster or the sample is much larger than zero, respectively. Note that we do not specify a particular
514 cutoff value for large (conditional) root causal effects, since the root causal effects likely lie on a
515 continuous graduated scale as opposed to approximately two binary values. Nevertheless, visual
516 inspection of the RCS or D-RCS histograms in disease should approximate a power law, where a
517 large mass is concentrated around zero and a long tail extends to the right similar to folding Figure
518 1 (b).

519 Algorithm

520 We now design an algorithm called Root Causal Strength using Perturbations (RCSP) that recovers
521 the RCS scores using Perturb-seq and bulk RNA-seq data.

522 Finding Surrogate Ancestors

523 Computing Φ_i for each $\tilde{X}_i \in \tilde{\mathbf{X}}$ requires access to the surrogate parents of each variable or, equiv-
524 alently, the causal graph \mathbb{G} . However, inferring \mathbb{G} using causal discovery algorithms may lead to
525 large statistical errors in the high dimensional setting (Colombo et al., 2014) and require restrictive
526 assumptions such as d-separation faithfulness (Spirtes et al., 2000) or specific functional relations
527 (Peters et al., 2014).

528 We instead directly utilize the interventional Perturb-seq data to recover a superset of the sur-
529rogate parents. We first leverage the global Markov property and equivalently write:

$$\Phi_i = \left| \mathbb{E}(Y | SA(\tilde{X}_i), X_i, B) - \mathbb{E}(Y | SA(\tilde{X}_i), B) \right|, \quad (5)$$

530 where $SA(\tilde{X}_i)$ denotes the *surrogate ancestors* of \tilde{X}_i , or the variables in \mathbf{X} associated with the ances-
531 tors of \tilde{X}_i .

532 We discover the surrogate ancestors using unconditional independence tests. For any $X_k \in \mathbf{X}$,
533 we test $X_k \perp\!\!\!\perp P_i$ by unpaired two-sided t-test, where P_i is an indicator function equal to one when
534 we perturb X_i and zero in the control samples of Perturb-seq. P_i is thus a parent of X_i alone but
535 not a child of B , so we do not need to condition on B . We use the two-sided t-test to assess for
536 independence because the t-statistic averages over cells to mimic bulk RNA-seq. If we reject the
537 null and conclude that $X_k \not\perp\!\!\!\perp P_i$, then X_k must be a descendant of P_i by the global Markov property,
538 so we include X_k into the set of surrogate descendants $SD(\tilde{X}_i)$. Curating every $X_j \in \mathbf{X}$ such that
539 $X_j \in SD(\tilde{X}_i)$ into $SA(\tilde{X}_i)$ yields the surrogate ancestors of \tilde{X}_i as desired.

540 Procedure

541 We now introduce an algorithm called Root Causal Strength using Perturbations (RCSP) that dis-
542 covers the surrogate ancestors of each variable \tilde{X} using Perturb-seq and then computes the RCS
543 of each variable using bulk RNA-seq. We summarize RCSP in Algorithm 1.

544 RCSP takes Perturb-seq and bulk RNA-seq datasets as input. The algorithm first finds the surro-
545 gate descendants of each variable in $\tilde{\mathbf{X}}$ in Line 2 in order to identify the surrogate ancestors of each
546 variable in Line 5. Access to the surrogate ancestors and the batches B allows RCSP to compute Φ_i
547 for each $X_i \in \mathbf{X}$ from the bulk RNA-seq in Line 6. The algorithm thus outputs the RCS scores Φ as
548 desired.

549 We certify RCSP as follows:

550 **Theorem 2.** (Fisher consistency) Consider the same assumption as Lemma 1. If unconditional d-separation
551 faithfulness holds, then RCSP recovers Φ almost surely as $N \rightarrow \infty$.

Algorithm 1 Root Causal Strength using Perturbations (RCSP)

Input: bulk RNA-seq data with batches B , Perturb-seq data

Output: RCS scores Φ

```
1: for each  $X_i \in X$  do
2:    $SD(\tilde{X}_i) \leftarrow$  all  $X_k \in X$  s.t.  $X_k \not\perp\!\!\!\perp P_i$  in Perturb-seq
3: end for
4: for each  $X_i \in X$  do
5:    $SA(\tilde{X}_i) \leftarrow$  all  $X_k \in X$  s.t.  $X_i \in SD(\tilde{X}_k)$ 
6:   Compute  $\Phi_i$  using Eq. (5) in bulk RNA-seq
7: end for
```

552 We engineered RCSP to only require *unconditional* d-separation faithfulness because real distribu-
553 tions may not obey full d-separation faithfulness (*Strobl, 2022*).

554 Synthetic Data

555 Simulations

556 We generated a linear SEM obeying Equation (3) specifically as $\tilde{X}_i = \tilde{X}\beta_i + E_i$ for every $\tilde{X}_i \in \tilde{X}$
557 and similarly $Y = \tilde{X}\beta_Y + E_Y$. We included $p + 1 = 2500$ variables in $\tilde{X} \cup Y$. We instantiated the
558 coefficient matrix β by sampling from a Bernoulli($2/(p - 1)$) distribution in the upper triangular
559 portion of the matrix. The resultant causal graph thus has an expected neighborhood size of 2. We
560 then randomly permuted the ordering of the variables. We introduced weights into the coefficient
561 matrix by multiplying each entry in β by a weight sampled uniformly from $[-1, -0.25] \cup [0.25, 1]$.
562 The error terms each follow a standard Gaussian distribution multiplied by 0.5. We introduced
563 batch effects by drawing each entry of the mapping efficiencies π from the uniform distribution
564 between 10 and 1000 for the bulk RNA-seq, and between 0.1 and 1 for the Perturb-seq. We set
565 $\tilde{X}_i \leftarrow \text{softplus}(\tilde{X}_i)$ and then obtained the corrupted surrogate X_i distributed $\text{Pois}(\tilde{X}_i\pi_{ij})$ for each
566 $\tilde{X}_i \in \tilde{X}$ and batch j . We chose Y uniformly at random from the set of vertices with at least one
567 parent and no children. We drew 200 samples for the bulk RNA-seq data to mimic a large but
568 common dataset size. We introduced knockdown perturbations in Perturb-seq by subtracting an
569 offset of two in the softplus function: $\tilde{X}_i \leftarrow \text{softplus}(\tilde{X}_i - 2)$. We finally drew 200 samples for the
570 control and each perturbation condition to generate the Perturb-seq data. We repeated the above
571 procedure 30 times.

572 Comparators

573 We compared RCSP against the following four algorithms:

- 574 1. Additive noise model (ANM) (*Peters et al., 2014; Strobl and Lasko, 2023a*): performs non-
575 linear regression of X_i on $\text{Pa}(X_i) \cup B$ and then regresses Y on the residuals $E \setminus E_i$ to estimate
576 $|\mathbb{E}(Y|E \setminus E_i) - \mathbb{E}(Y|X, B)|$ for each $X_i \in X$. The non-linear regression residuals are equivalent
577 to the error terms assuming an additive noise model.
- 578 2. Linear Non-Gaussian Acyclic Model (LiNGAM) (*Peters et al., 2014; Strobl and Lasko, 2022*):
579 same as above but performs linear instead of non-linear regression.
- 580 3. CausalCell (*Wen et al., 2023*): selects the top 50 genes with maximal statistical dependence to
581 Y , and then runs the Peter-Clark (PC) algorithm (*Spirtes et al., 2000*) using a non-parametric
582 conditional independence test to identify a causal graph among the top 50 genes. The algo-
583 rithm does not perform root causal inference, so we use ANM as above but condition on the
584 estimated parent sets for the top 50 genes and the ancestors inferred from the Perturb-seq
585 data otherwise.
- 586 4. Univariate regression residuals (Uni Reg): regresses Y on $X_i \cup B$ and estimates the absolute
587 residuals $|Y - \mathbb{E}(Y|X_i, B)|$ for each $X_i \in X$.

- 588 5. Multivariate regression residuals (Multi Reg): similar to above but instead computes the ab-
589 solute residuals after regressing Y on $(X \setminus X_i) \cup B$.

590 The first two methods are state-of-the-art approaches used for root causal discovery. Univariate
591 and multivariate regressions do not distinguish between predictivity and causality, but we included
592 them as sanity checks. We performed all non-linear regressions using multivariate adaptive re-
593 gression splines to control for the underlying regressor (*Friedman, 1991*). We also standardized
594 all variables before running the regressions to prevent gaming of the marginal variances in causal
595 discovery (*Reisach et al., 2021; Ng et al., 2024*). We compared the algorithms on their accuracy in
596 estimating Φ .

597 Real Data

598 Quality Control

599 We downloaded Perturb-seq datasets of retinal pigment epithelial cells from the RPE-1 cell line,
600 and undifferentiated blast cells from the K562 cell line (*Replogle et al., 2022*). We used the genome-
601 wide dataset version for the latter. We downloaded the datasets from the scPerturb database on
602 Zenodo (*Green et al., 2022*) with the same quality controls as the original paper. Replogle et al.
603 computed adjusted library sizes by equalizing the mean library size of control cells within each
604 batch. Cells with greater than a 2000 or 3000 library size, and less than 25% or 11% mitochondrial
605 RNA were kept, respectively. The parameters were chosen by plotting the adjusted library sizes
606 against the mitochondrial RNA counts and then manually setting thresholds that removed low
607 quality cells likely consisting of ambient mRNA transcripts arising from premature cell lysis or cell
608 death.

609 We next downloaded bulk RNA-seq datasets derived from patients with age-related macular
610 degeneration (AMD; GSE115828) and multiple sclerosis (MS; GSE137143) (*Ratnapriya et al., 2019;*
611 *Kim et al., 2021*). We excluded 10 individuals from the AMD dataset including one with an RNA
612 integrity number of 21.92, five missing an integrity number (all others had an integrity number of
613 less than 10), and four without a Minnesota Grading System score. We kept all samples from the
614 MS dataset derived from CD4+ T cells but filtered out genes with a mean of less than 5 counts as
615 done in the original paper.

616 We finally kept genes that were present in both the AMD bulk dataset and the RPE-1 Perturb-
617 seq dataset, yielding a final count of 513 bulk RNA-seq samples and 247,914 Perturb-seq samples
618 across 2,077 genes. We also kept genes that were present in both the MS bulk dataset and the K562
619 Perturb-seq dataset, yielding a final count of 137 bulk RNA-seq samples and 1,989,578 Perturb-seq
620 samples across 6,882 genes. We included age and sex as a biological variable as covariates for every
621 patient in both datasets in subsequent analyses.

622 Evaluation Rationale

623 We do not have access to the ground truth values of Φ in real data. We instead evaluate the RCSP
624 estimates of Φ using alternative sources of ground truth knowledge. We first assess the accuracy
625 of RCS using the control variable age as follows:

- 626 1. Determine if the RCS values of age identify age as a root cause with large causal effect in
627 diseases that progress over time.

628 Second, few root causal genes should drive pathogenesis because the effects of a few error terms
629 distribute over many downstream genes. We verify the sparsity of root causal genes as follows:

- 630 2. Determine if the distribution of D-RCS concentrates around zero more than the distribution
631 of the Deviation of Statistical Dependence (D-SD) defined as $\sqrt{\frac{1}{n} \sum_{j=1}^n \omega_{ij}^2}$ for each gene $\tilde{X}_i \in \tilde{X}$
632 where $\Omega_i = |\mathbb{E}(Y|X_i, B) - \mathbb{E}(Y|B)|$ and ω_{ij} its value for patient j .

633 Despite the sparsity of root causal genes, we still expect the root causal genes to correspond to at
634 least some known causes of disease:

- 635 3. Determine if genes with the top D-RCS scores correspond to genes known to cause the dis-
636 ease.

637 Next, the root causal genes initiate the vast majority of pathogenesis, and we often have knowledge
638 of pathogenic pathways even though we may not know the exact gene expression cascade leading
639 to disease. Intervening on root causal genes should also modulate patient symptoms. We thus
640 further evaluate the accuracy of RCSP using pathway and drug enrichment analyses as follows:

- 641 4. Determine if the D-RCS scores identify known pathogenic pathways of disease in pathway
642 enrichment analysis.
643 5. Determine if the D-RCS scores identify drugs that treat the disease.

644 Finally, complex diseases frequently involve multiple pathogenic pathways that differ between pa-
645 tients. Patients with the same complex disease also respond differently to treatment. We hence
646 evaluate the precision of RCS as follows:

- 647 6. Determine if the patient-specific RCS scores identify subgroups of patients involving different
648 but still known pathogenic pathways.
649 7. Determine if the patient-specific RCS scores identify subgroups of patients that respond dif-
650 ferently to drug treatment.

651 In summary, we evaluate RCSP in real data based on its ability to (1) identify age as a known root
652 cause, (2) suggest an omnigenic root causal model, (3) recover known causal genes, (4) find known
653 pathogenic pathways, (5) find drugs that treat the disease, and (6,7) delineate patient subgroups.

654 Enrichment Analyses

655 Multivariate adaptive regression splines introduce sparsity, but enrichment analysis performs bet-
656 ter with a dense input. We can estimate the conditional expectations of Φ using any general non-
657 linear regression method, so we instead estimated the expectations using kernel ridge regression
658 equipped with a radial basis function kernel (*Shawe-Taylor and Cristianini, 2004*). We then com-
659 puted the D-RCS across all patients for each variable in \mathbf{X} . We ran pathway enrichment analysis
660 using the fast gene set enrichment analysis (FGSEA) algorithm (*Sergushichev, 2016*) with one hun-
661 dred thousand simple permutations using the D-RCS scores and pathway information from the
662 Reactome database (version 1.86.0) (*Fabregat et al., 2017*). We likewise performed drug set en-
663 richment analysis with the Drug Signature database (version 1.0) (*Yoo et al., 2015*). We repeated
664 the above procedures for the D-RCS of all clusters identified by hierarchical clustering via Ward's
665 method (*Ward Jr, 1963*).

666 Data Availability

667 All datasets analyzed in this study have been previously published and are publicly accessible as
668 follows:

- 669 1. Bulk RNA-seq for AMD: [GSE115828](#)
670 2. Bulk RNA-seq for MS: [GSE137143](#)
671 3. Perturb-seq for the RPE-1 and K562 cell lines: DOI [10044268](#)

672 Code Availability

673 R code needed to replicate all experimental results is available on [GitHub](#).

674 Acknowledgements

675 Research reported in this report was supported by the National Human Genome Research Institute
676 of the National Institutes of Health under award numbers R01HG011138 and R35HG010718.

677 References

- 678 **Adamson B**, Norman TM, Jost M, Cho MY, Nuñez JK, Chen Y, Villalta JE, Gilbert LA, Horlbeck MA, Hein MY, et al.
679 A multiplexed single-cell CRISPR screening platform enables systematic dissection of the unfolded protein
680 response. *Cell*. 2016; 167(7):1867–1882.
- 681 **Andhavarapu S**, Mubariz F, Arvas M, Bever Jr C, Makar TK. Interplay between ER stress and autophagy: a
682 possible mechanism in multiple sclerosis pathology. *Experimental and Molecular Pathology*. 2019; 108:183–
683 190.
- 684 **Barouch FC**, Miller JW. The role of inflammation and infection in age-related macular degeneration. *Internation-*
685 *ational ophthalmology clinics*. 2007; 47(2):185–197.
- 686 **Basile MS**, Bramanti P, Mazzon E. The role of cytotoxic T-lymphocyte antigen 4 in the pathogenesis of multiple
687 sclerosis. *Genes*. 2022; 13(8):1319.
- 688 **Beaumat F**, O'prey J, Barthet VJ, Zunino B, Parvy JP, Bachmann AM, O'prey M, Kania E, Gonzalez PS, Macintosh
689 R, et al. mTORC1 activation requires DRAM-1 by facilitating lysosomal amino acid efflux. *Molecular Cell*. 2019;
690 76(1):163–176.
- 691 **Bongers S**, Forré P, Peters J, Mooij JM. Foundations of structural causal models with cycles and latent variables.
692 *The Annals of Statistics*. 2021; 49(5):2885–2915.
- 693 **Boyle EA**, Li YI, Pritchard JK. An expanded view of complex traits: from polygenic to omnigenic. *Cell*. 2017;
694 169(7):1177–1186.
- 695 **Burster T**, Beck A, Poeschel S, Øren A, Baechle D, Reich M, Roetzschke O, Falk K, Boehm BO, Youssef S, et al.
696 Interferon- γ regulates cathepsin G activity in microglia-derived lysosomes and controls the proteolytic pro-
697 cessing of myelin basic protein in vitro. *Immunology*. 2007; 121(1):82–93.
- 698 **Buschur KL**, Chikina M, Benos PV. Causal network perturbations for instance-specific analysis of single cell and
699 disease samples. *Bioinformatics*. 2020; 36(8):2515–2521.
- 700 **Cano-Gamez E**, Trynka G. From GWAS to function: using functional genomics to identify the mechanisms
701 underlying complex diseases. *Frontiers in Genetics*. 2020; 11:424.
- 702 **Choudhary S**, Satija R. Comparison and evaluation of statistical error models for scRNA-seq. *Genome Biology*.
703 2022; 23(1):27.
- 704 **Colombo D**, Maathuis MH, et al. Order-independent constraint-based causal structure learning. *Journal of*
705 *Machine Learning Research*. 2014; 15(1):3741–3782.
- 706 **Costa-Silva J**, Domingues D, Lopes FM. RNA-Seq differential expression analysis: An extended review and a
707 software tool. *PloS One*. 2017; 12(12):e0190152.
- 708 **Dalvin LA**, Olsen TW, Bakri SJ, McCullough K, Tefferi A, Al-Kali A. Busulfan treatment for myeloproliferative
709 disease may reduce injection burden in vascular endothelial growth factor-driven retinopathy. *American*
710 *Journal of Ophthalmology Case Reports*. 2022; 26:101554.
- 711 **Datlinger P**, Rendeiro AF, Schmidl C, Krausgruber T, Traxler P, Klughammer J, Schuster LC, Kuchler A, Alpar D,
712 Bock C. Pooled CRISPR screening with single-cell transcriptome readout. *Nature methods*. 2017; 14(3):297–
713 301.
- 714 **Dixit A**, Parnas O, Li B, Chen J, Fulco CP, Jerby-Arnon L, Marjanovic ND, Dionne D, Burks T, Raychowdhury R, et al.
715 Perturb-Seq: dissecting molecular circuits with scalable single-cell RNA profiling of pooled genetic screens.
716 *Cell*. 2016; 167(7):1853–1866.
- 717 **Ellington CN**, Lengerich BJ, Watkins TB, Yang J, Xiao H, Kellis M, Xing EP. Contextualized Networks Reveal
718 Heterogeneous Transcriptomic Regulation in Tumors at Sample-Specific Resolution. In: *Neural Information*
719 *and Processing Systems Workshop on Generative AI and Biology*; 2023. .
- 720 **Fabregat A**, Sidiropoulos K, Viteri G, Forner O, Marin-Garcia P, Arnau V, D'Eustachio P, Stein L, Hermjakob H.
721 Reactome pathway analysis: a high-performance in-memory approach. *BMC Bioinformatics*. 2017; 18(1):1–9.
- 722 **Fletcher JM**, Lalor S, Sweeney C, Tubridy N, Mills K. T cells in multiple sclerosis and experimental autoimmune
723 encephalomyelitis. *Clinical & Experimental Immunology*. 2010; 162(1):1–11.
- 724 **Friedman JH**. Multivariate adaptive regression splines. *The Annals of Statistics*. 1991; 19(1):1–67.

- 725 **Friedman N**, Linial M, Nachman I, Pe'er D. Using Bayesian networks to analyze expression data. In: *Proceedings*
726 *of the Fourth Annual International Conference on Computational Molecular Biology*; 2000. p. 127–135.
- 727 **Gnanaprakasam JR**, Wang R. MYC in regulating immunity: metabolism and beyond. *Genes*. 2017; 8(3):88.
- 728 **Go YM**, Zhang J, Fernandes J, Litwin C, Chen R, Wensel TG, Jones DP, Cai J, Chen Y. mTOR-initiated metabolic
729 switch and degeneration in the retinal pigment epithelium. *FASEB Journal*. 2020; 34(9):12502.
- 730 **Golan M**, Krivitsky A, Mausner-Fainberg K, Benhamou M, Vigiser I, Regev K, Kolb H, Karni A. Increased expres-
731 sion of ephrins on immune cells of patients with relapsing remitting multiple sclerosis affects oligodendro-
732 cyte differentiation. *International Journal of Molecular Sciences*. 2021; 22(4):2182.
- 733 **Green TD**, Peidli S, Shen C, Gross T, Min J, Garda S, Taylor-King JP, Marks DS, Luna A, Blüthgen N, et al. scPerturb:
734 Information Resource for Harmonized Single-Cell Perturbation Data. In: *NeurIPS 2022 Workshop on Learning*
735 *Meaningful Representations of Life*; 2022. .
- 736 **Grün D**, Kester L, Van Oudenaarden A. Validation of noise models for single-cell transcriptomics. *Nature*
737 *Methods*. 2014; 11(6):637–640.
- 738 **Hadziahmetovic M**, Malek G. Age-related macular degeneration revisited: From pathology and cellular stress
739 to potential therapies. *Frontiers in Cell and Developmental Biology*. 2021; 8:612812.
- 740 **Haves-Zbuorof D**, Paperna T, Gour-Lavie A, Mandel I, Glass-Marmor L, Miller A. Cathepsins and their endoge-
741 nous inhibitors cystatins: expression and modulation in multiple sclerosis. *Journal of Cellular and Molecular*
742 *Medicine*. 2011; 15(11):2421–2429.
- 743 **Kamalden T**, Ji D, Fawcett R, Osborne N. Genistein blunts the negative effect of ischaemia to the retina caused
744 by an elevation of intraocular pressure. *Ophthalmic Research*. 2011; 45(2):65–72.
- 745 **Kim K**, Pröbstel AK, Baumann R, Dyckow J, Landefeld J, Kogl E, Madireddy L, Loudermilk R, Eggers EL, Singh
746 S, et al. Cell type-specific transcriptomics identifies neddylation as a novel therapeutic target in multiple
747 sclerosis. *Brain*. 2021; 144(2):450–461.
- 748 **Kinoshita S**, Noda K, Tagawa Y, Inafuku S, Dong Y, Fukuhara J, Dong Z, Ando R, Kanda A, Ishida S. Genistein
749 attenuates choroidal neovascularization. *The Journal of Nutritional Biochemistry*. 2014; 25(11):1177–1182.
- 750 **Kokame K**, Agarwala KL, Kato H, Miyata T. Herp, a new ubiquitin-like membrane protein induced by endoplas-
751 mic reticulum stress. *Journal of Biological Chemistry*. 2000; 275(42):32846–32853.
- 752 **Lengfeld JE**, Lutz SE, Smith JR, Diaconu C, Scott C, Kofman SB, Choi C, Walsh CM, Raine CS, Agalliu I, et al.
753 Endothelial Wnt/ β -catenin signaling reduces immune cell infiltration in multiple sclerosis. *Proceedings of*
754 *the National Academy of Sciences*. 2017; 114(7):E1168–E1177.
- 755 **Luo H**, Broux B, Wang X, Hu Y, Ghannam S, Jin W, Larochelle C, Prat A, Wu J. EphrinB1 and EphrinB2 regu-
756 late T cell chemotaxis and migration in experimental autoimmune encephalomyelitis and multiple sclerosis.
757 *Neurobiology of Disease*. 2016; 91:292–306.
- 758 **Martínez-Jiménez F**, Muiños F, Sentís I, Deu-Pons J, Reyes-Salazar I, Arnedo-Pac C, Mularoni L, Pich O, Bonet J,
759 Kranas H, et al. A compendium of mutational cancer driver genes. *Nature Reviews Cancer*. 2020; 20(10):555–
760 572.
- 761 **Nachef M**, Ali AK, Almutairi SM, Lee SH. Targeting SLC1A5 and SLC3A2/SLC7A5 as a potential strategy to
762 strengthen anti-tumor immunity in the tumor microenvironment. *Frontiers in immunology*. 2021; 12:624324.
- 763 **Nagraal A**. Gaucher disease. *Journal of Clinical and Experimental Hepatology*. 2014; 4(1):37–50.
- 764 **Narendran S**, Pereira F, Yerramothu P, Apicella I, Wang Sb, Varshney A, Baker KL, Marion KM, Ambati M, Ambati
765 VL, et al. A clinical metabolite of azidothymidine inhibits experimental choroidal neovascularization and
766 retinal pigmented epithelium degeneration. *Investigative ophthalmology & visual science*. 2020; 61(10):4–4.
- 767 **Ng I**, Huang B, Zhang K. Structure learning with continuous optimization: A sober look and beyond. In: *Causal*
768 *Learning and Reasoning PMLR*; 2024. p. 71–105.
- 769 **Nicklin P**, Bergman P, Zhang B, Triantafellow E, Wang H, Nyfeler B, Yang H, Hild M, Kung C, Wilson C, et al.
770 Bidirectional transport of amino acids regulates mTOR and autophagy. *Cell*. 2009; 136(3):521–534.
- 771 **Olsen TW**, Feng X. The Minnesota Grading System of eye bank eyes for age-related macular degeneration.
772 *Investigative Ophthalmology and Visual Science*. 2004; 45(12):4484–4490.

- 773 **Orian JM**, D'Souza CS, Kocovski P, Krippner G, Hale MW, Wang X, Peter K. Platelets in multiple sclerosis: early
774 and central mediators of inflammation and neurodegeneration and attractive targets for molecular imaging
775 and site-directed therapy. *Frontiers in Immunology*. 2021; 12:620963.
- 776 **Papoulis A**. Probability, Random Variables and Stochastic Processes. McGraw-Hill; 1984.
- 777 **Pearl J**. Causality. Cambridge University Press; 2009.
- 778 **Peters J**, Mooij JM, Janzing D, Schölkopf B. Causal discovery with continuous additive noise models. *Journal of*
779 *Machine Learning Research*. 2014; .
- 780 **Ratnapriya R**, Sosina OA, Starostik MR, Kwicklis M, Kapphahn RJ, Fritsche LG, Walton A, Arvanitis M, Gieser L,
781 Pietraszkiewicz A, et al. Retinal transcriptome and eQTL analyses identify genes associated with age-related
782 macular degeneration. *Nature Genetics*. 2019; 51(4):606–610.
- 783 **Reisach A**, Seiler C, Weichwald S. Beware of the simulated DAG! causal discovery benchmarks may be easy to
784 game. *Advances in Neural Information Processing Systems*. 2021; 34:27772–27784.
- 785 **Replogle JM**, Saunders RA, Pogson AN, Hussmann JA, Lenail A, Guna A, Mascibroda L, Wagner EJ, Adelman
786 K, Lithwick-Yanai G, et al. Mapping information-rich genotype-phenotype landscapes with genome-scale
787 Perturb-seq. *Cell*. 2022; 185(14):2559–2575.
- 788 **Sarkar A**, Stephens M. Separating measurement and expression models clarifies confusion in single-cell RNA
789 sequencing analysis. *Nature Genetics*. 2021; 53(6):770–777.
- 790 **Sergushichev A**. An algorithm for fast preranked gene set enrichment analysis using cumulative statistic cal-
791 culation. *BioRxiv*. 2016; 60012:1–9.
- 792 **Shawe-Taylor J**, Cristianini N. Kernel Methods for Pattern Analysis. Cambridge University Press; 2004.
- 793 **Shi Y**, Qu J, Zhang D, Zhao P, Zhang Q, Tam POS, Sun L, Zuo X, Zhou X, Xiao X, et al. Genetic variants at 13q12. 12
794 are associated with high myopia in the Han Chinese population. *The American Journal of Human Genetics*.
795 2011; 88(6):805–813.
- 796 **Sobel RA**. Ephrin A receptors and ligands in lesions and normal-appearing white matter in multiple sclerosis.
797 *Brain Pathology*. 2005; 15(1):35–45.
- 798 **Spink KE**, Polakis P, Weis WI. Structural basis of the Axin–adenomatous polyposis coli interaction. *The EMBO*
799 *journal*. 2000; 19(10):2270–2279.
- 800 **Spirtes P**, Glymour C, Scheines R. Causation, Prediction, and Search. 2nd ed. MIT press; 2000.
- 801 **Spirtes P**. Directed cyclic graphical representations of feedback models. In: *Proceedings of the Eleventh Confer-*
802 *ence on Uncertainty in Artificial Intelligence*; 1995. p. 491–498.
- 803 **Starzyk RM**, Rosenow C, Frye J, Leismann M, Rodzinski E, Putney S, Tuomanen EI. Cerebral cell adhesion
804 molecule: a novel leukocyte adhesion determinant on blood-brain barrier capillary endothelium. *The Journal*
805 *of Infectious Diseases*. 2000; 181(1):181–187.
- 806 **Strobl EV**. Causal discovery with a mixture of DAGs. *Machine Learning*. 2022; p. 1–25.
- 807 **Strobl EV**. Counterfactual Formulation of Patient-Specific Root Causes of Disease. *Journal of Biomedical Infor-*
808 *matics*. 2024; .
- 809 **Strobl EV**, Lasko TA. Identifying patient-specific root causes of disease. In: *Proceedings of the 13th ACM Interna-*
810 *tional Conference on Bioinformatics, Computational Biology and Health Informatics*; 2022. p. 1–10.
- 811 **Strobl EV**, Lasko TA. Identifying patient-specific root causes with the heteroscedastic noise model. *Journal of*
812 *Computational Science*. 2023; 72:102099.
- 813 **Strobl EV**, Lasko TA. Root Causal Inference from Single Cell RNA Sequencing with the Negative Binomial. In:
814 *Proceedings of the 14th ACM International Conference on Bioinformatics, Computational Biology and Health Infor-*
815 *matics BCB '23*, New York, NY, USA: Association for Computing Machinery; 2023. .
- 816 **Strobl EV**, Lasko TA. Sample-specific root causal inference with latent variables. In: *Conference on Causal*
817 *Learning and Reasoning* PMLR; 2023. p. 895–915.

- 818 **Strobl EV**, Lasko TA, Gamazon ER. Mitigating pathogenesis for target discovery and disease subtyping. *Com-*
819 *puters in Biology and Medicine*. 2024; 171:108122.
- 820 **Su Y**, Wang F, Hu Q, Qu Y, Han Y. Arsenic trioxide inhibits proliferation of retinal pigment epithelium by down-
821 *regulating expression of extracellular matrix and p27*. *International Journal of Clinical and Experimental*
822 *Pathology*. 2020; 13(2):172.
- 823 **Turi Z**, Senkyrikova M, Mistrik M, Bartek J, Moudry P. Perturbation of RNA Polymerase I transcription machinery
824 *by ablation of HEATR1 triggers the RPL5/RPL11-MDM2-p53 ribosome biogenesis stress checkpoint pathway*
825 *in human cells*. *Cell Cycle*. 2018; 17(1):92–101.
- 826 **Wang L**, Trasanidis N, Wu T, Dong G, Hu M, Bauer DE, Pinello L. Dictys: dynamic gene regulatory network
827 *dissects developmental continuum with single-cell multiomics*. *Nature Methods*. 2023; 20(9):1368–1378.
- 828 **Ward Jr JH**. Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Associ-*
829 *ation*. 1963; p. 236–244.
- 830 **Wen Y**, Huang J, Guo S, Elyahu Y, Monsonego A, Zhang H, Ding Y, Zhu H. Applying causal discovery to single-cell
831 *analyses using CausalCell*. *Elife*. 2023; 12:e81464.
- 832 **Yoo M**, Shin J, Kim J, Ryall KA, Lee K, Lee S, Jeon M, Kang J, Tan AC. DSigDB: drug signatures database for gene
833 *set analysis*. *Bioinformatics*. 2015; 31(18):3069–3071.

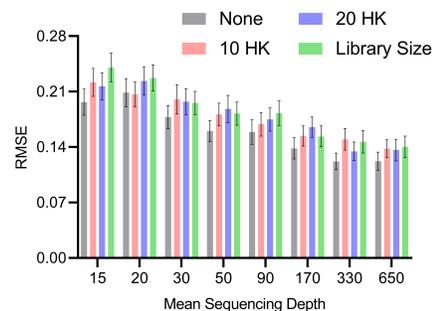
834 Supplementary Materials

835 Additional Synthetic Data Results

836 Normalization by Sequencing Depth

837 We theoretically showed that RCS does not require normalization by sequencing depth in the Meth-
838 ods using an asymptotic argument. We tested this claim empirically by drawing 200 bulk RNA-seq
839 samples from random DAGs as in the Methods but over $p + 1 = 250$ variables. We varied the mean
840 sequencing depth N/p of each gene from 15, 20, 30, 50, 90, 170, 330 to 650 counts; multiplying
841 N/p by p recovers the library size N . We only included one batch in the bulk RNA-seq in order
842 to isolate the effect of sequencing depth. We compared no normalization, normalization by 10
843 housekeeping genes, normalization by 20 housekeeping genes, and normalization by library size.
844 We repeated each experiment 100 times and thus generated a total of $100 \times 4 \times 8 = 3200$ datasets.

845 We plot the results in Supplementary Figure 1. All methods improved with increasing mean
846 sequencing depth as expected. The no normalization strategy performed the best at low mean
847 sequencing depths, followed by the housekeeping genes and then total library size. The result
848 even held with a small library size of $N = 15 \times 249 = 3735$ at the smallest mean sequencing depth of
849 15, suggesting that the asymptotic argument holds well in bulk RNA-seq where N/p is often greater
850 than 500 and N greater than the tens of millions. However, the average RMSEs of all normalization
851 methods became more similar as sequencing depth increased. We conclude that no normalization
852 exceeds or matches the accuracy of other strategies. We therefore do not normalize by sequencing
853 depth in subsequent analyses.

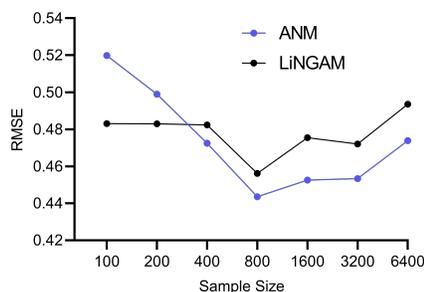


Supplementary Figure 1. Mean RMSE to the ground truth RCS values across different mean sequencing depths and normalization strategies. The no normalization strategy achieved low RMSEs at lower mean sequencing depths, but the performances of all methods converged as the mean sequencing depths increased. Error bars denote 95% confidence intervals of the mean RMSE.

854 Functional Causal Models and Measurement Error

855 The experiments in the Results section quantify the accuracies of the algorithms in estimating Φ .
856 However, the functional causal models ANM and LiNGAM also estimate the error terms as an in-
857 termediate step, whereas RCSP does not. We therefore also investigated the accuracies of ANM
858 and LiNGAM in estimating the error term values.

859 Theoretical results suggest that ANM and LiNGAM cannot consistently estimate the error terms
860 in RNA-seq due to the Poisson measurement error. We empirically tested this hypothesis by sam-
861 pling from bulk RNA-seq data as in the Methods but with $p + 1 = 100$ and a batch size of one in
862 order to isolate the effect of measurement error. We repeated the experiment 100 times for bulk
863 RNA-seq sample sizes of 100, 200, 400, 800, 1600 and 3200. We plot the results in Supplementary
864 Figure 2. The accuracies of ANM and LiNGAM did not improve beyond an RMSE of 0.44 to the
865 ground truth error term values even with a large sample size of 6400. We conclude that ANM and
866 LiNGAM cannot estimate the error terms accurately in the presence of measurement error even
867 with large sample sizes.

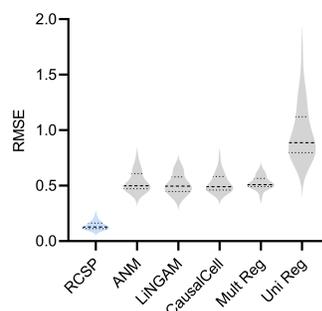


Supplementary Figure 2. Mean RMSE values to the ground truth error term values across different sample sizes. The accuracies of ANM and LiNGAM do not improve with increasing sample sizes.

868 Cyclic Causal Graphs

869 We also evaluated the algorithms on directed graphs with cycles. We generated a linear SEM over
870 $p + 1 = 1000$ variables in $\tilde{X} \cup Y$. We sampled the coefficient matrix β from a Bernoulli($1/(p - 1)$) distri-
871 bution but did not restrict the non-zero coefficients to the upper triangular portion of the matrix.
872 We then proceeded to permute the variable ordering and weight each entry as in the Methods for
873 the DAG. We repeated this procedure 30 times and report the results in Supplementary Figure 3.

874 RCSP again outperformed all other algorithms even in the cyclic case. The results suggest that
875 conditioning on the surrogate ancestors also estimates the RCS well even in the cyclic case. How-
876 ever, we caution that an error term E_i can affect the ancestors of \tilde{X}_i when cycles exist. As a result,
877 the RCS may not isolate the causal effect of the error term and thus not truly coincide with the
878 notion of a root causal effect in cyclic causal graphs.



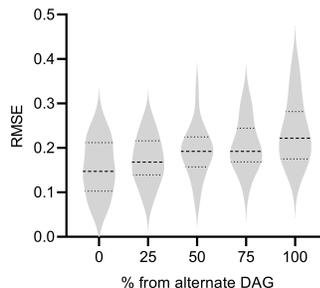
Supplementary Figure 3. RCSP achieved the lowest RMSE in cyclic graphs as well. However, error terms can influence ancestors in the cyclic case, so the interpretation of the RCS remains unclear when cycles exist.

879 DAG Incongruence

880 We next assessed the performance of RCSP when the DAG underlying the Perturb-seq data differs
881 from the DAG underlying the bulk RNA-seq data. We considered a mixture of two random DAGs in
882 bulk RNA-seq, where one of the DAGs coincided with the Perturb-seq DAG and second alternate
883 DAG did not. We instantiated and simulated samples from each DAG as per the previous subsec-
884 tion. We generated 0%, 25%, 50%, 75%, and 100% of the bulk RNA-seq samples from the alternate
885 DAG, and the rest from the Perturb-seq DAG. We ideally would like to see the performance of RCSP
886 degrade gracefully, as opposed to abruptly, as the percent of samples derived from the alternate
887 DAG increases.

888 We summarize results in Supplementary Figure 4. As expected, RCSP performed the best when
889 we drew all samples from the same underlying DAG for Perturb-seq and bulk RNA-seq. However,
890 the performance of RCSP also degraded slowly as the percent of samples increased from the al-

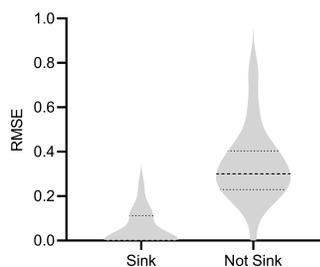
891 ternate DAG. We conclude that RCSP can accommodate some differences between the underlying
892 DAGs in Perturb-seq and bulk RNA-seq with only a mild degradation in performance.



Supplementary Figure 4. The performance of RCSP degrades gracefully as the percent of samples from the alternate DAG increases.

893 Non-Sink Target

894 We finally considered the scenario where Y is a non-sink (or non-terminal) vertex. If Y is a parent
895 of a gene expression level, then we cannot properly condition on the parents because modern
896 Perturb-seq datasets usually do not intervene on Y or measure Y . We therefore empirically inves-
897 tigated the degradation in performance resulting from a non-sink target Y , in particular for gene
898 expression levels where Y is a parent. We again simulated 200 samples from bulk RNA-seq and
899 each condition of Perturb-seq with a DAG over 1000 vertices, an expected neighborhood size of 2
900 and a non-sink target Y . We then removed the outgoing edges from Y and resampled the DAG with
901 a sink target. We compare the results of RCSP for both DAGs in gene expression levels where Y is
902 a parent. We plot the results in Supplementary Figure 5. As expected, we observe a degradation
903 in performance when Y is not terminal, where the mean RMSE increased from 0.045 to 0.342. We
904 conclude that RCSP is sensitive to violations of the sink target assumption.



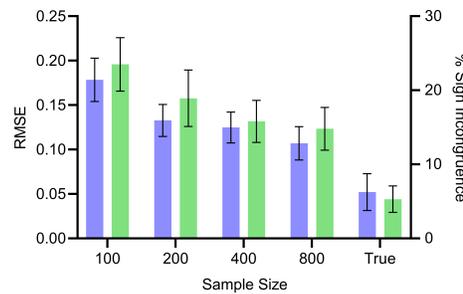
Supplementary Figure 5. Results with a sink or non-sink target Y . RCSP estimated the RCS scores less accurately with a non-sink target indicating that the algorithm is sensitive to violations of the sink target assumption.

905 Root Causal Effect versus Conditional Root Causal Effect

906 We compared the expected and unconditional root causal effects $\Omega_i \triangleq \mathbb{E}(Y|E_i) - \mathbb{E}(Y)$ to the ex-
907 pected and conditional root causal effects or, equivalently, the signed RCS scores Γ . These root
908 causal effects and conditional root causal effects are not the same, but they are similar. We empiri-
909 cally investigated the differences between the estimated values of Γ and the true values of Ω using
910 the RMSE and also the percent of samples with incongruent signs; Γ and Ω have incongruent signs
911 if one is positive and the other is negative. We again drew 200 bulk RNA-seq samples from random
912 DAGs as in the Methods over $p + 1 = 250$ variables with one batch. We varied the bulk RNA-seq
913 sample size from 100, 200, 400 to 800. We also compared true Γ against true Ω by estimating the

914 two to negligible error using 20,000 samples of \tilde{X} . We repeated each experiment 100 times and
915 thus generated a total of $100 \times 5 = 500$ datasets.

916 We summarize the results in Supplementary Figure 6. The estimated Γ values approached the
917 true Ω values with increasing sample sizes. The true Γ values did not converge exactly to the true Ω
918 values, but the RMSE remained low at 0.05 and the two values differed in sign only around 5.3% of
919 the time. Increasing the number of samples of \tilde{X} to 50,000 did not change performance, confirming
920 that we reached the floor. We conclude that the empirical results replicate the theoretical results
921 because Γ and Ω do not match exactly. However, the two quantities take on similar values and
922 their signs matched around 95% of the time in practice.



Supplementary Figure 6. Mean RMSE (blue, left) and percent sign incongruence (green, right) of the expected root causal effects and signed RCS values, respectively. The RMSE continues to decrease with increasing sample size but reaches a floor of around 0.05. Similarly, the percent sign incongruence decreases but reaches a floor of around 5%.

923 **Additional Results for Age-related Macular Degeneration**

924 **Algorithm Comparisons**

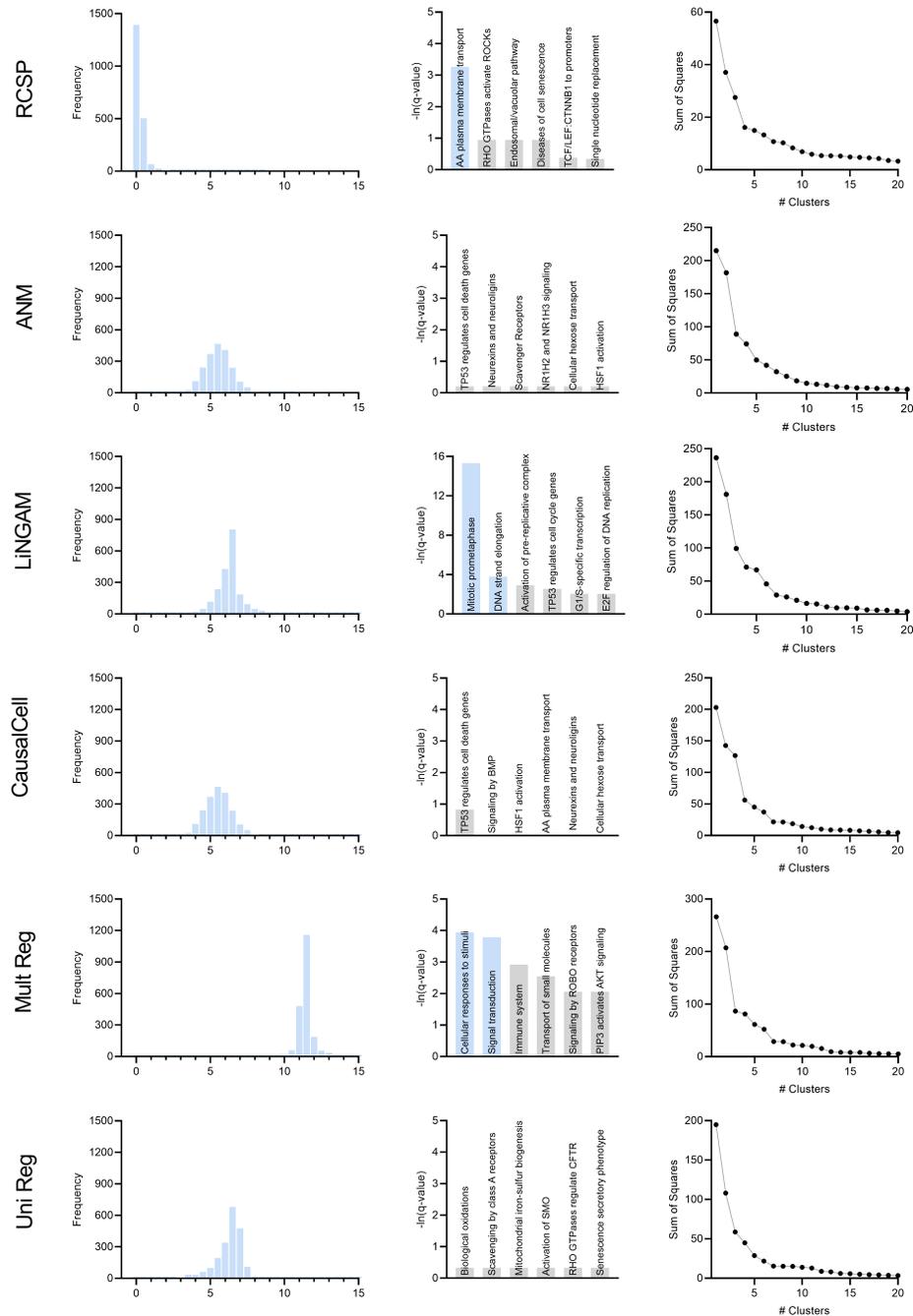
925 We say that an algorithm performs well in real data if it simultaneously (1) identifies a sparse set
926 of root causal genes, (2) recovers known pathogenic pathways with high specificity measured by
927 the sparsity of leading edge genes, and (3) clusters patients into clear subgroups.

928 We compared the algorithms with the AMD data. We summarize the results in Supplementary
929 Figure 7 plotted on the next page. The figure contains 6 rows and 3 columns. Similar to the D-RCS,
930 we can compute the standard deviation of the output of each algorithm from zero for each gene.
931 The first column in Supplementary Figure 7 denotes the histograms of these standard deviations
932 across the genes. We standardized the outputs to have mean zero and unit variance. We then
933 added the minimum value so that all histograms begin at zero; note that the bars at zero are not
934 visible for many algorithms, since only a few genes attained standard deviations near the minimum.
935 If an algorithm accurately identifies root causal genes, then it should only identify a few genes with
936 large conditional root causal effects under the omnigenic root causal model. The RCSP algorithm
937 had a histogram with large probability mass centered around zero with a long tail to the right. The
938 standard deviations of the outputs of the other algorithms attained large values for nearly all genes.
939 Incorporating feature selection and causal discovery with CausalCell introduced more outliers in
940 the histogram of ANM. We conclude that only RCSP detected an omnigenic root causal model.

941 We plot the results of pathway enrichment analysis in the second column of Supplementary
942 Figure 7. RCSP, LiNGAM and univariate regression detected pathways related to oxidative stress
943 in AMD. However, the “mitotic prometaphase” and “DNA strand elongation” pathways in blue for
944 LiNGAM involved 94 and 27 leading edge genes, respectively. The “cellular responses to stimuli”
945 and “signal transduction” pathways for multivariate regression also involved 253 and 282 leading
946 edge genes. In contrast, the “amino acid plasma membrane transport” pathway for RCSP involved
947 two leading edge genes. We conclude that RCSP identified a known pathogenic pathway of AMD
948 with the fewest number of leading edge genes.

949 We finally plot the clustering results in the third column of Supplementary Figure 7. The RCSP
 950 sum of squares plot revealed a sharp elbow at four groups of patients, whereas the other plots
 951 did not reveal a clear number of categories using the elbow method. We conclude that only RCSP
 952 identified clear subgroups of patients in AMD.

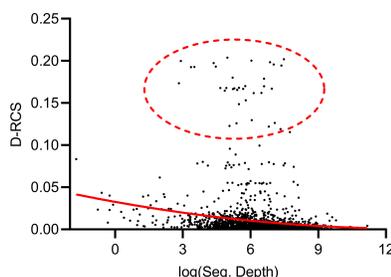
953 In summary, RCSP detected a small set of root causal genes, identified pathogenic pathways
 954 with maximal specificity and discovered distinguishable patient subgroups. We therefore conclude
 955 that RCSP outperformed all other algorithms in the AMD dataset.



Supplementary Figure 7. Comparison of the algorithms in age-related macular degeneration.

956 Effect of Sequencing Depth

957 Theorem 1 states that RCS scores may exhibit bias with insufficient sequencing depth. The genes
 958 with large D-RCS scores may therefore simply have low sequencing depths. To test this hypothesis,
 959 we plotted sequencing depth against D-RCS scores. Consistent with Theorem 1, we observed a
 960 small negative correlation between D-RCS and sequencing depth ($\rho = -0.16$, $p=2.04E-13$), and D-
 961 RCS scores exhibited greater variability at the lowest sequencing depths (Supplementary Figure 8).
 962 However, genes with the largest D-RCS scores had mean sequencing depths interspersed between
 963 20 and 3000. We conclude that genes with the largest D-RCS scores had a variety of sequencing
 964 depths ranging from low to high.



Supplementary Figure 8. Mean sequencing depth of each gene plotted against their D-RCS scores in AMD. Genes with the largest D-RCS scores (red ellipse) had a variety of sequencing depths.

965 Biological Results

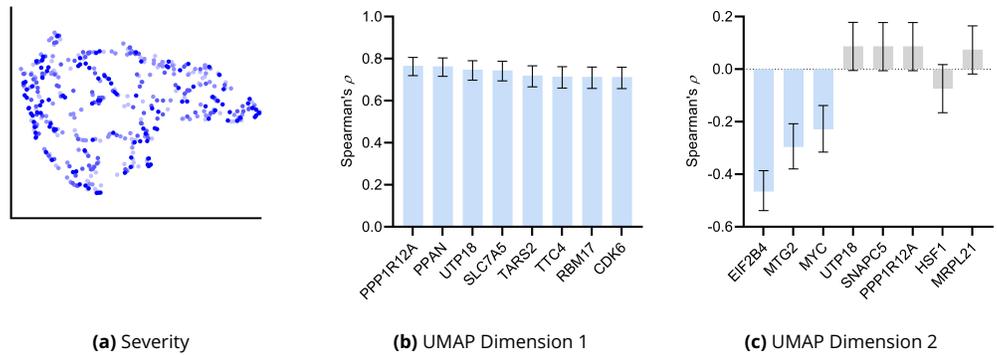
966 We provide the full pathway enrichment analysis results in Supplementary Table 1 corresponding
 967 to Figure 3 (c). We summarize pathway enrichment analysis of the black cluster of Figure 3 (g) in
 968 Figure 3 (j). However, analyses of the blue, green and pink clusters did not yield significant pathways
 969 even at a liberal FDR threshold of 10%.

970 We examined whether the clusters of Figure 3 (g) differentiate dry and wet macular degenera-
 971 tion. Wet macular degeneration is associated with the highest Minnesota Grading System (MGS)
 972 score of 4 (*Olsen and Feng, 2004*). We plotted the UMAP embedding against MGS (Supplementary
 973 Figure 9 (a)). None of the two UMAP dimensions correlated significantly with the MGS score (5%
 974 uncorrected threshold by Spearman's correlation test). These results and the large RCS scores of
 975 age in Figure 3 (a) seem to support the hypothesis that wet macular degeneration is a more severe
 976 type of dry macular degeneration. However, MGS does not differentiate between wet macular de-
 977 generation and late stage dry macular degeneration involving geographical atrophy. We therefore
 978 cannot separate late stage dry and wet macular degeneration using the RCS scores alone.

Pathway	p-value	q-value	Effect Size	Leading Edge
Amino acid transport across the plasma membrane	2.44e-05	0.038	0.995	8140,6510
RHO GTPases Activate ROCKs	2.09e-03	0.388	0.976	4659,5500
Endosomal/Vacuolar pathway	2.32e-03	0.388	0.998	3107
Diseases of Cellular Senescence	2.97e-03	0.388	0.997	1021
Binding of TCF/LEF:CTNNB1 to target gene promoters	6.52e-03	0.680	0.993	4609
APEX1-Indep. Resolution of AP Sites via Nucleotide Replacement	7.28e-03	0.712	0.980	11284,7515
MASTL Facilitates Mitotic Progression	1.59e-02	0.978	0.911	84930,983
PI5P Regulates TP53 Acetylation	1.94e-02	0.978	0.980	79837
Formation of Incision Complex in GG-NER	2.24e-02	0.978	0.791	2966,9978,2967
Glycine degradation	2.24e-02	0.978	0.977	1738
Prefoldin mediated transfer of substrate to CCT/TriC	3.96e-02	0.978	0.787	5203,5201,10576

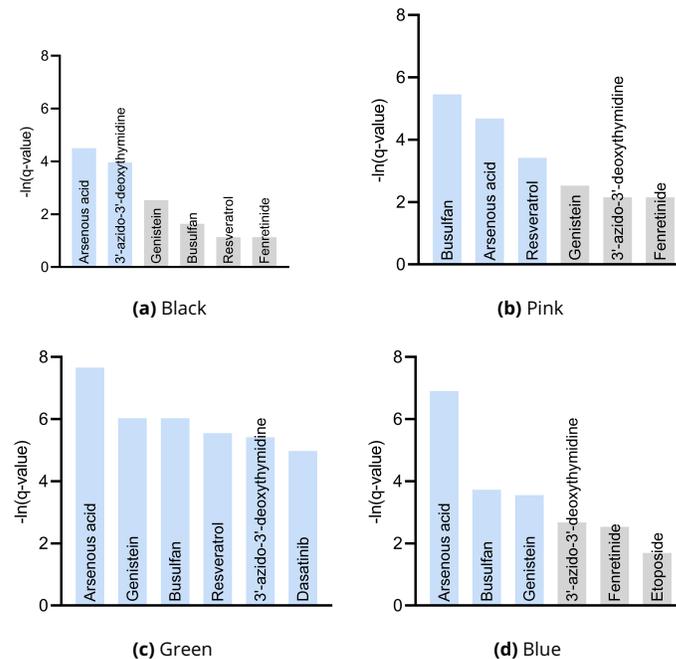
Supplementary Table 1. Full pathway enrichment analysis results for all patients in the AMD dataset. We list the Entrez gene IDs of up to the top three leading edge genes in the right-most column.

979 We correlated the two UMAP dimensions with the top 30 genes ranked by their RCS scores. We
 980 plot genes with the highest correlation to the first and second UMAP dimensions in Supplementary
 981 Figures 9 (b) and 9 (c), respectively. Many genes correlated with the first dimension, but only three
 982 genes correlated with the second at an FDR threshold of 5%.



Supplementary Figure 9. Additional UMAP embedding results for AMD. (a) The UMAP dimensions did not correlate with AMD severity as assessed by the MGS score. Many genes correlated with the first UMAP dimension in (b), but only three genes correlated with the second UMAP dimension in (c). Blue bars passed an FDR threshold of 5%, and error bars denote 95% confidence intervals.

983 We finally performed drug enrichment analysis in each of the four clusters in Figure 3 (g). We
 984 summarize the results in Supplementary Figure 10. Only two drugs – and one potentially thera-
 985 peutic option – passed FDR correction in patients in the black cluster with the most identified root
 986 causal genes according to the RCS scores. In contrast, enrichment analysis identified many drugs
 987 in patients in the green cluster with the lowest RCS scores and thus relatively few root causal genes.
 988 The pink and blue clusters yielded moderate results. We conclude that drug enrichment analysis
 989 expectedly identified more drugs for patients on the left hand side of the UMAP embedding with
 990 fewer root causal genes than on the right hand side with many simultaneous root causal genes.



Supplementary Figure 10. Drug enrichment analysis results by cluster in Figure 3 (g). The analyses recovered similar drugs across clusters, but the results for the green cluster in (c) were supra-significant.

991 **Additional Results for Multiple Sclerosis**

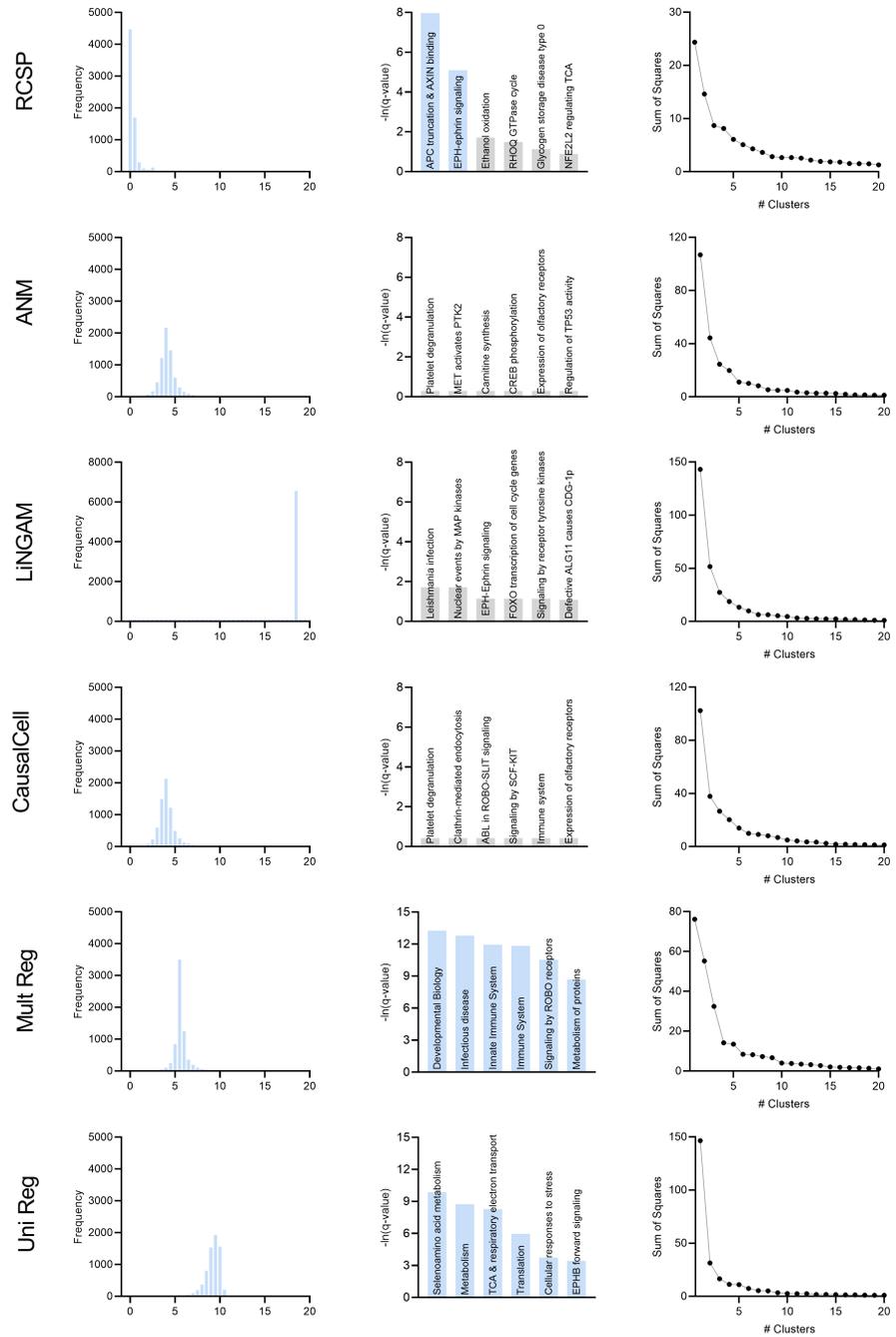
992 Algorithm Comparisons

993 We compared the algorithms using the MS data with the same criteria used for the AMD dataset.
994 We summarize the results in Supplementary Figure 11 plotted on the next page. Only the his-
995 togram of RCSP had large probability mass centered around zero as shown in the first column.
996 The histogram of LiNGAM contained many outliers, so it appears to spike around a value of 18.
997 The histograms of ANM and CausalCell were again near identical. We conclude that only the his-
998 togram of RCSP supported an omnigenic root causal model in MS.

999 We performed pathway enrichment analysis on the algorithm outputs and summarize the
1000 results in the second column of Supplementary Figure 11. The functional causal models ANM,
1001 LiNGAM and CausalCell did not identify significant pathways at an FDR corrected threshold of 0.05.
1002 In contrast, multivariate and univariate regression both identified many significant pathways in
1003 blue with no specific link to the blood brain barrier. The top six significant pathways for multivari-
1004 ate and univariate regression involved 112 to 831 and 18 to 545 leading edge genes, respectively.
1005 In contrast, the two significant pathways of RCSP involved only 2 and 9 leading genes. We conclude
1006 that RCSP detected pathogenic pathways of MS with the sparsest set of leading edge genes.

1007 We finally clustered the algorithm outputs into patient subgroups. We list the sum of squares
1008 plots in the third column of Supplementary Figure 11. Univariate regression did not differentiate
1009 between the patients because it detected one dominating cluster. RCSP and multivariate regres-
1010 sion identified clear subgroups according to the elbow method, whereas the sum of squares plots
1011 for ANM, LiNGAM and CausalCell showed no clear cutoffs. We conclude that only RCSP and multi-
1012 variate regression identified clear patient subgroups in MS.

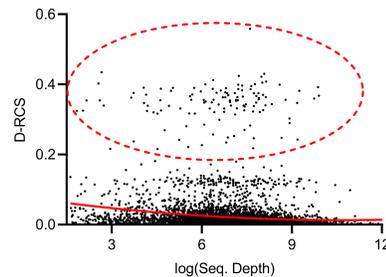
1013 In summary, only RCSP simultaneously detected an omnigenic root causal model, identified
1014 pathogenic pathways with high specificity and discovered clear patient subgroups. We therefore
1015 conclude that RCSP also outperformed all other algorithms in the MS dataset.



Supplementary Figure 11. Comparison of the algorithms in multiple sclerosis.

1016 Effect of Sequencing Depth

1017 We plot sequencing depth against the D-RCS scores of each gene similar to the AMD dataset. We
 1018 again observed a small negative correlation ($\rho = -0.136$, $p < 2.2E-16$), indicating that genes with low
 1019 sequencing depths had slightly higher D-RCS scores on average (Supplementary Figure 12). How-
 1020 ever, genes with the largest D-RCS scores again had a variety of sequencing depths. We conclude
 1021 that sequencing depth has minimal correlation with the largest D-RCS scores.



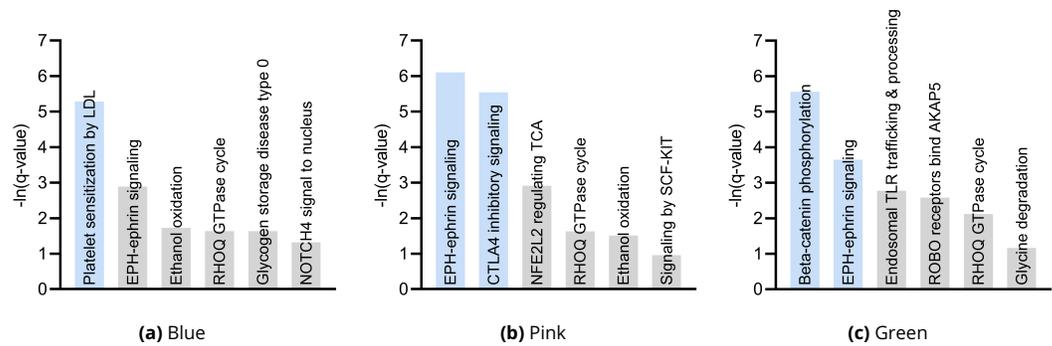
Supplementary Figure 12. Mean sequencing depth of each gene plotted against their D-RCS scores in MS. Genes with the largest D-RCS scores (red ellipse) again had a variety of sequencing depths.

1022 Biological Results

1023 We provide the full global pathway enrichment analysis results for MS in Supplementary Table 2.
 1024 Pathway enrichment analysis of the individual clusters in Figure 4 (f) consistently implicated EPH-
 1025 ephrin signaling among the top two pathways. However, each cluster also involved one separate
 1026 additional pathway (Supplementary Figure 13). The green cluster involved the same APC-AXIN path-
 1027 way as the global analysis via beta-catenin. On the other hand, the blue cluster involved “platelet
 1028 sensitization by LDL.” Low density lipoprotein enhances platelet aggregation. Platelet degranula-
 1029 tion in turn drives the generation of autoreactive T cells in the peripheral circulation during distur-
 1030 bance of the blood brain barrier (*Orian et al., 2021*). Finally, CTLA4 regulates T-cell homeostasis
 1031 and inhibits autoimmunity for the pink cluster (*Basile et al., 2022*). The D-RCS scores of each cluster
 1032 thus implicate different mechanisms of T cell pathology.

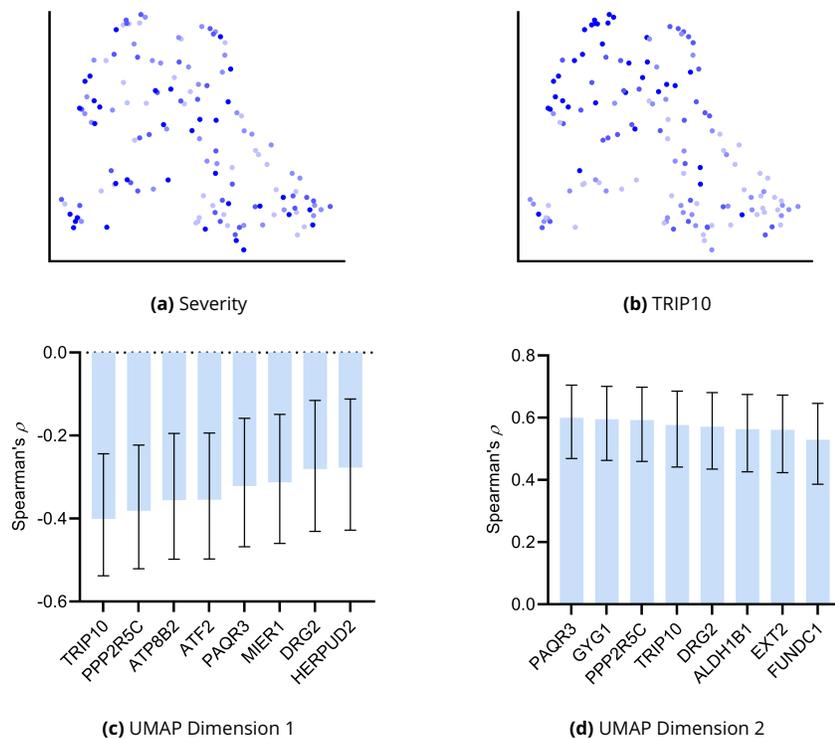
Pathway	p-value	q-value	Effect Size	Leading Edge
APC truncation mutants have impaired AXIN binding	1.91e-06	3.45e-4	0.960	5525,5527
EPH-ephrin signaling	4.23e-05	6.12e-3	0.826	8874,102,8976
Ethanol oxidation	2.02e-03	0.182	0.967	219,128
RHOQ GTPase cycle	2.72e-03	0.226	0.793	9322,8874,10395
Glycogen storage disease type 0 (muscle GYS1)	4.32e-03	0.322	0.996	2992
NFE2L2 regulating TCA cycle genes	6.31e-03	0.414	0.970	4199,3417
C6 deamination of adenosine	7.42e-03	0.414	0.981	103,104
Ion channel transport	7.63e-03	0.414	0.728	57198,540,55515
Synthesis of IP3 and IP4 in the cytosol	7.65e-03	0.414	0.904	3633,805,23236
Diseases associated with glycosaminoglycan metabolism	8.21e-03	0.414	0.894	2132,11285,3339
Signaling by SCF-KIT	8.67e-03	0.414	0.794	7006,5578,3815

Supplementary Table 2. Full pathway enrichment analysis results for all patients in the MS dataset. We again list up to the top three leading edge genes in the right-most column.



Supplementary Figure 13. Pathway enrichment analysis results by cluster consistently revealed EPH-ephrin signaling as well as an additional pathway implicating T cell pathology.

1033 The severity of MS, as assessed by the Expanded Disability Status Scale (EDSS) score, did not
 1034 correlate with either dimension of the UMAP embedding (Supplementary Figure 14 (a)). The top
 1035 genes in Figure 4 (d) such as MNT and CERCAM also did not correlate. However, lower ranked genes
 1036 such as TRIP10 did (Supplementary Figure 14 (b)). An expanded correlation analysis with the top
 1037 30 genes revealed significant correlations across a variety of lower ranked genes (Supplementary
 1038 Figures 14 (c) and 14 (d)). We conclude that the distribution of lower ranked genes govern the
 1039 structure of the UMAP embedding in Figure 4 (f).



Supplementary Figure 14. Additional analyses of the UMAP embedding for MS. (a) The UMAP dimensions did not correlate with MS severity as assessed by EDSS. However, lower ranked genes such as TRIP10 correlated with both dimensions in (b). We expanded the analysis to the top 30 genes and plot the genes with the highest correlations to UMAP dimension one and two in (c) and (d), respectively.

1040 **Proofs**

1041 **Lemma 1.** Assume Lipschitz continuity of the conditional expectation for all $N \geq n_0$:

$$\mathbb{E} \left| \mathbb{E}(Y|\tilde{U}) - \mathbb{E}(Y|U, L, B) \right| \leq \mathbb{E} C_N \left| \tilde{U} - \frac{U}{dL} \right|, \quad (6)$$

1042 where $d = \frac{\pi_{UB}}{\sum_{\tilde{X}_i \in \tilde{A}} \tilde{x}_i \pi_i B}$, $C_N \in O(1)$ is a positive constant, and we have taken an outer expectation on both
 1043 sides. Then $\mathbb{E}(Y|\tilde{U}) = \lim_{N \rightarrow \infty} \mathbb{E}(Y|U, L, B)$ almost surely.

1044 *Proof.* We can write the following sequence:

$$\begin{aligned} \mathbb{E} \left| \mathbb{E}(Y|\tilde{U}) - \lim_{N \rightarrow \infty} \mathbb{E}(Y|U, L, B) \right| &= \mathbb{E} \lim_{N \rightarrow \infty} \left| \mathbb{E}(Y|\tilde{U}) - \mathbb{E}(Y|U, L, B) \right| \\ &\leq \mathbb{E} \lim_{N \rightarrow \infty} C_N \left| \tilde{U} - \frac{U}{dL} \right| \leq C \mathbb{E} \left| \tilde{U} - \frac{1}{d} \lim_{N \rightarrow \infty} \frac{U}{L} \right| = C \mathbb{E} \left| \tilde{U} - \frac{1}{d} \tilde{U} d \right| = 0, \end{aligned}$$

1045 where we have applied Expression (6) at the first inequality. We have $C_N \leq C$ for all $N \geq n_0$ in the
 1046 second inequality because $C_N \in O(1)$. With the above bound, choose $a > 0$ and invoke the Markov
 1047 inequality:

$$\mathbb{P} \left(\left| \mathbb{E}(Y|\tilde{U}) - \lim_{N \rightarrow \infty} \mathbb{E}(Y|U, L, B) \right| \geq a \right) \leq \frac{1}{a} \mathbb{E} \left| \mathbb{E}(Y|\tilde{U}) - \lim_{N \rightarrow \infty} \mathbb{E}(Y|U, L, B) \right| = 0.$$

1048 The conclusion follows because we chose a arbitrarily. □

1049 **Proposition 1.** If $E_i \perp\!\!\!\perp Y$ or $E_i \perp\!\!\!\perp Y | \text{Pa}(\tilde{X}_i)$ (or both), then E_i is a root cause of Y .

1050 *Proof.* If $E_i \perp\!\!\!\perp Y$ or $E_i \perp\!\!\!\perp Y | \text{Pa}(\tilde{X}_i)$ (or both), then E_i and Y are d-connected by the global Markov
 1051 property. Since E_i is a root vertex, the d-connection implies that there exists a directed path from
 1052 E_i to Y . □

1053 **Proposition 2.** We have $\mathbb{P}(Y|E_i, \text{Pa}(\tilde{X}_i)) = \mathbb{P}(Y|\tilde{X}_i, \text{Pa}(\tilde{X}_i))$ under Equation (3).

1054 *Proof.* We can write:

$$\mathbb{P}(Y|E_i, \text{Pa}(\tilde{X}_i)) = \mathbb{E}_{\tilde{X}_i | E_i, \text{Pa}(\tilde{X}_i)} \mathbb{P}(Y|E_i, \tilde{X}_i, \text{Pa}(\tilde{X}_i)) = \mathbb{P}(Y|E_i, \tilde{X}_i, \text{Pa}(\tilde{X}_i)) = \mathbb{P}(Y|\tilde{X}_i, \text{Pa}(\tilde{X}_i)).$$

1055 The second equality follows because \tilde{X}_i is a constant given E_i and $\text{Pa}(\tilde{X}_i)$. The third equality follows
 1056 by the global Markov property because Y is a terminal vertex. □

1057 **Theorem 2.** (Fisher consistency) Consider the same assumption as Lemma 1. If unconditional d-separation
 1058 faithfulness holds, then RCSP recovers Φ almost surely as $N \rightarrow \infty$.

1059 *Proof.* If $X_k \perp\!\!\!\perp P_i$ in Line 2 of Algorithm 1, then X_k is a descendant of the root vertex P_i under
 1060 the global Markov property. Similarly, if X_k is a descendant of P_i , then X_k is d-connected to P_i
 1061 so $X_k \perp\!\!\!\perp P_i$ by unconditional d-separation faithfulness. Hence, $\text{SD}(\tilde{X}_i)$ contains only and all the
 1062 surrogate descendants of \tilde{X}_i for each $\tilde{X}_i \in \tilde{X}$. This in turn implies that $\text{SA}(\tilde{X}_i)$ in Line 5 of Algorithm
 1063 1 contains only and all the surrogate ancestors of \tilde{X}_i . Hence, RCSP now has access to the correct
 1064 set $\text{SA}(\tilde{X}_i)$ as well as B for each $\tilde{X}_i \in \tilde{X}$. We finally invoke Theorem 1 to conclude that RCSP recovers
 1065 Φ almost surely as $N \rightarrow \infty$. □