

Transcriptome-Wide Root Causal Inference

Eric V. Strobl^{1,*} and Eric R. Gamazon²

¹University of Pittsburgh

²Vanderbilt University Medical Center

ABSTRACT

Root causal genes correspond to the first gene expression levels perturbed during pathogenesis by genetic or non-genetic factors. Targeting root causal genes has the potential to alleviate disease entirely by eliminating pathology near its onset. No existing algorithm discovers root causal genes from observational data alone. We therefore propose the Transcriptome-Wide Root Causal Inference (TWRCI) algorithm that identifies root causal genes and their causal graph using a combination of genetic variant and unperturbed bulk RNA sequencing data. TWRCI uses a novel competitive regression procedure to annotate cis and trans-genetic variants to the gene expression levels they directly cause. The algorithm simultaneously recovers a causal ordering of the expression levels to pinpoint the underlying causal graph and estimate root causal effects. TWRCI outperforms alternative approaches across a diverse group of metrics by directly targeting root causal genes while accounting for distal relations, linkage disequilibrium, patient heterogeneity and widespread pleiotropy. We demonstrate the algorithm by uncovering the root causal mechanisms of two complex diseases, which we confirm by replication using independent genome-wide summary statistics.

1 Introduction

Genetic and non-genetic factors can influence gene expression levels to ultimately cause disease. Root causal gene expression levels – or *root causal genes* for short – correspond to the *initial* changes to *gene expression* that ultimately generate disease as a downstream effect¹. Root causal genes differ from core genes that directly cause the phenotype and thus lie at the end, rather than at the beginning, of pathogenesis². Root causal genes also generalize driver genes that only account for the effects of somatic mutations primarily in protein coding sequences in cancer³.

Discovering root causal genes is critical towards identifying drug targets that modify disease near its pathogenic onset and thus mitigate downstream pathogenesis in its entirety⁴. The problem is complicated by the existence of complex disease, where the causal effects of the root causal genes may differ between patients even within the same diagnostic category. However, the recently defined *omnigenic root causal model* posits that only a few root causal genes affect nearly all downstream genes to initiate the vast majority of pathology in each patient¹. We thus more specifically seek to identify *personalized* root causal genes specific to any given individual.

Only one existing algorithm accurately identifies personalized root causal genes¹, but the algorithm requires access to genome-wide Perturb-seq data, or high throughput perturbations with single cell RNA sequencing readout^{5,6}. Perturb-seq is currently expensive and difficult to obtain in many cell types. We instead seek a method that can uncover personalized root causal genes directly from widespread observational (or non-experimental) datasets.

We make the following contributions in this paper:

1. We introduce the conditional root causal effect (CRCE) that measures the causal effect of the genetic and non-genetic factors, which directly affect a gene expression level, on the phenotype.
2. We propose a novel strategy called Competitive Regression that provably annotates both cis and trans-genetic variants to the gene expression level or phenotype they directly cause without conservative significance testing.
3. We create an algorithm called Transcriptome-Wide Root Causal Inference (TWRCI) that uses the annotations to reconstruct a

personalized causal graph summarizing the CRCEs of gene expression levels from a combination of genetic variant and bulk RNA sequencing observational data.

4. We show with confirmatory replication that TWRCI identifies only a few root causal genes that accurately distinguish subgroups of patients even in complex diseases – consistent with the omnigenic root causal model.

We provide an example of the output of TWRCI in Figure 1. TWRCI annotates both cis and trans genetic variants to the expression level or phenotype they *directly* cause. We prove that the direct causal annotations allow the algorithm to uniquely reconstruct the causal graph between the gene expression levels that cause the phenotype as well as estimate their CRCEs. The algorithm summarizes the CRCEs in the graph by weighing and color-coding each vertex, where vertex size correlates with magnitude, green induces disease and red prevents disease. TWRCI thus provides a succinct summary of root causal genes and their root causal effect sizes specific to a given patient using observational data alone. TWRCI outperforms combinations of existing algorithms across all subtasks: annotation, graph reconstruction and CRCE estimation. No existing algorithm performs all subtasks simultaneously.

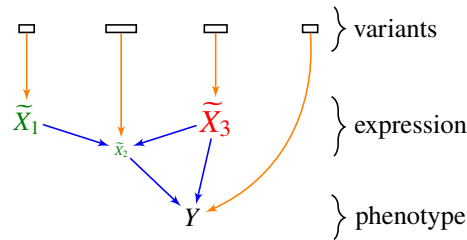


Figure 1. Toy example of a root causal mechanism inferred by TWRCI for a specific patient. Rectangles denote sets of genetic variants, potentially in linkage disequilibrium even between sets. Each set of variants directly causes a gene expression level in \tilde{X} or the phenotype Y . Larger lettered vertices denote larger CRCE magnitudes and colors refer to their direction – green is a positive CRCE and red is negative. TWRCI simultaneously **annotates**, **reconstructs** and estimates the **CRCEs**.

2 Results

2.1 Overview of TWRCI

2.1.1 Setup

We seek to identify not just causal but *root causal* genes. We must therefore carefully define the generative process. We consider a set of variants S , the transcriptome \tilde{X} and the phenotype Y . We represent the generative causal process using a directed graph like in Figure 2 (a), where the variants cause the transcriptome, and the transcriptome causes the phenotype. Directed edges denote direct causal relations between variables. In practice, the sets S and \tilde{X} contain millions and thousands of variables, respectively. As described in Methods 4.2, we cannot measure the values of \tilde{X} exactly using RNA sequencing but instead measure values X corrupted by Poisson measurement error and batch effects.

2.1.2 Variable Selection

Simultaneously handling millions of variants and thousands of gene expression levels currently requires expensive computational resources. Moreover, most variants and gene expression levels do not inform the discovery of root causal genes for a particular phenotype Y . The Transcriptome-Wide Root Causal Inference (TWRCI) thus first performs variable selection by eliminating variants and gene expression levels unnecessary for root causal inference.

TWRCI first identifies variants $T \subseteq S$ associated with Y using widely available summary statistics at a liberal α threshold, such as $5e-5$, in order to capture many causal variants. The algorithm then uses individual-level data – where each individual has variant data, bulk gene expression data from the relevant tissue and phenotype data (variant-expression-phenotype) – to learn a regression model predicting X from T . TWRCI identifies the subset of expression levels $\tilde{R} \subseteq \tilde{X}$ that it can predict better

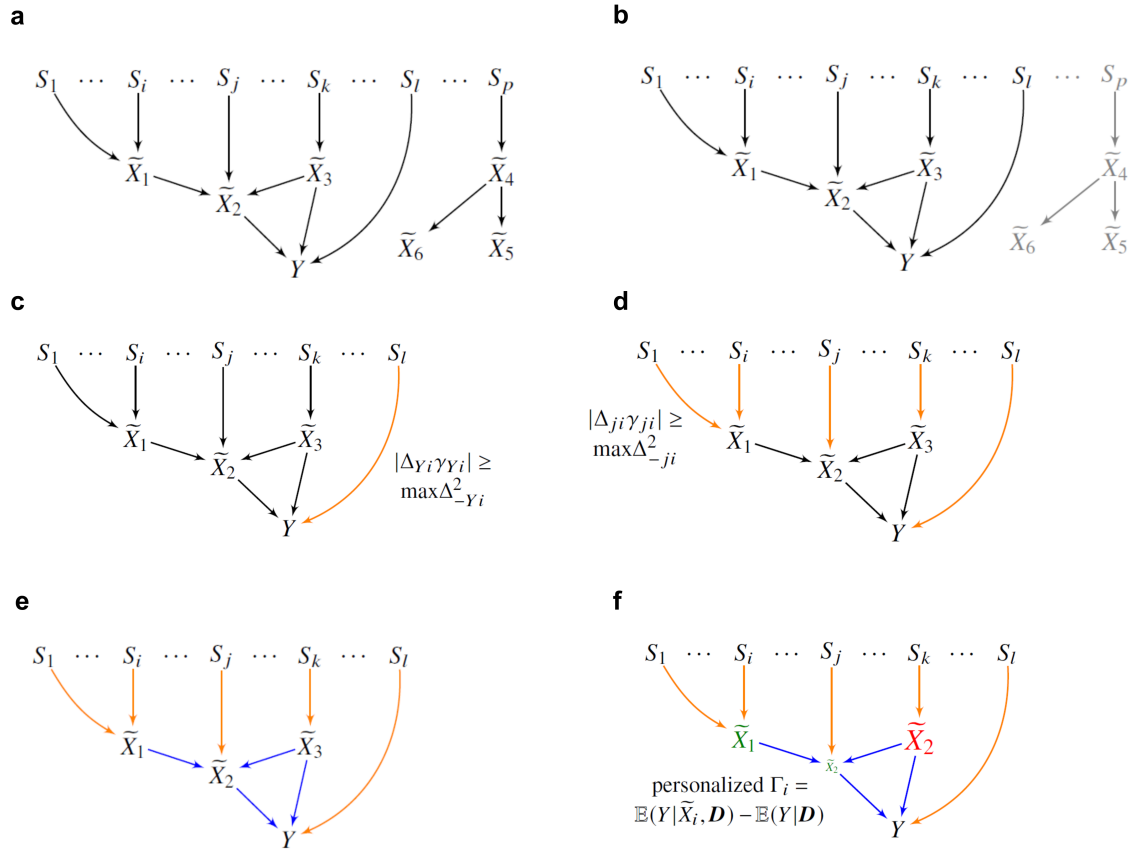


Figure 2. Overview of the TWRCI algorithm. (a) We redraw Figure 1 in more detail. We do not have access to the underlying causal graph in practice. (b) TWRCI first performs variable selection by only keeping variants and gene expression levels correlated with Y (and their common cause confounders) as shown in black. (c) The algorithm then uses Competitive Regression to find the variants that directly cause Y in orange. (d) TWRCI iteratively repeats Competitive Regression for each gene expression level as well – again shown in orange. (e) The algorithm next performs causal discovery to identify the causal relations between the gene expression levels and the phenotype in blue. (f) Finally, TWRCI weighs each vertex $\tilde{X}_i \in \tilde{\mathbf{X}}$ by the magnitude of Γ_i and color codes the vertex by its direction (green is positive, red is negative). TWRCI thus ultimately recovers a causal graph like the one shown in Figure 1.

than chance. We refer the reader to Methods 4.4.2 for details on the discovery of additional nuisance variables required to mitigate confounding. We prove that $T \cup \tilde{\mathbf{R}}$ retains all of the causes of Y in $S \cup \tilde{\mathbf{X}}$.

2.1.3 Annotation by Competitive Regression

We want to annotate both cis and trans-variants to the gene expression level that they directly cause in $\tilde{\mathbf{R}}$. We also want to annotate variants to the phenotype Y in order to account for horizontal pleiotropy, where variants bypass $\tilde{\mathbf{R}}$ and directly cause Y . TWRCI achieves both of these feats through a novel process called *Competitive Regression*.

TWRCI accounts for horizontal pleiotropy by applying Competitive Regression to the phenotype. We do not restrict the theoretical results detailed in Methods to linear models, but linear models trained on genotype data currently exhibit competitive performance^{7,8}. TWRCI therefore trains debiased linear ridge regression models⁹ predicting Y from $T \cup \tilde{\mathbf{R}}$ without requiring Gaussian distributions; let γ_{Yi} refer to the coefficient for T_i in the regression model. Similarly, let Δ_{-Yi} correspond to the matrix of coefficients for T_i in the regression models predicting $\tilde{\mathbf{R}}$ (but not Y) from T ; notice that we have not conditioned on the gene expression levels in this case. If T_i directly causes Y , then it will predict Y given $T \setminus T_i$ and given $T \cup \tilde{\mathbf{R}} \setminus T_i$ (i.e., Δ_{Yi} and γ_{Yi} will both be non-zero), but T_i will not predict any gene expression level given $T \setminus T_i$ (i.e., $\max \Delta_{-Yi}^2$ will be zero). We also prove

the converse direction in Methods 4.4.4. TWRCI therefore annotates T_i to Y , if $|\Delta_{Y_i\gamma_{Y_i}}|$ deviates away from zero even after conditioning on gene expression so that $|\Delta_{Y_i\gamma_{Y_i}}| \geq \max \Delta_{-Y_i}^2$ – i.e., $\Delta_{Y_i\gamma_{Y_i}}$ “beats” $\Delta_{-Y_i}^2$ in a competitive process (Figure 2 (c)).

We prove in Methods 4.4.3 that Competitive Regression successfully recovers the direct causes of Y , denoted by $S_Y \subseteq T$, so long as Y is a sink vertex that does not cause any other variable. We also require analogues of two standard assumptions used in instrumental variable analysis: relevance and exchangeability¹⁰. In this paper, *relevance* means that at least one variant in T directly causes each gene expression level in \tilde{R} ; the assumption usually holds because T contains orders of magnitude more variants than entries in \tilde{R} . On the other hand, *exchangeability* assumes that T and other sets of direct causal variants not in T share no latent confounders (details in Methods 4.4.2); this assumption holds approximately due to the weak causal relations emanating from variants to gene expression and the phenotype. Exchangeability also weakens as T grows larger.

We further show that Competitive Regression can recover $S_i \subseteq T$, or the direct causes of $\tilde{R}_i \in \tilde{R}$ in T , when \tilde{R}_i causes Y and turns into a sink vertex after removing Y from consideration (Methods 4.4.4). As a result, TWRCI removes Y and appends it to the empty ordered set K to ensure that some $\tilde{R}_i \in \tilde{R}$ is now a sink vertex. We introduce a statistical criterion in Methods 4.4.4 that allows TWRCI to find the sink vertex \tilde{R}_i after removing Y . TWRCI then annotates S_i to \tilde{R}_i again using Competitive Regression (Figure 2 (d)), removes R_i from R and appends R_i to the front of K . The algorithm iterates until it has removed all variables from $R \cup Y$ and placed them into the causal order K .

2.1.4 Causal Discovery and CRCE Estimation

Annotation only elucidates the direct causal relations from variants to gene expression, but it does not recover the causal relations between gene expression or the causal relations from gene expression to the phenotype. We want TWRCI to recover the *entire* biological mechanism from variants all the way to the phenotype.

TWRCI thus subsequently runs a causal discovery algorithm with the causal order K to uniquely identify the causal graph over $\tilde{R} \cup Y$ (Figure 2 (e)). The algorithm also estimates the personalized or *conditional root causal effect* (CRCE) of gene expression levels that cause Y :

$$\begin{aligned} \Gamma_i &= \mathbb{E}(Y | \widehat{E_i} \cup \widehat{S_i}, D) - \mathbb{E}(Y | D), \\ &= \mathbb{E}(Y | \tilde{X}_i, D) - \mathbb{E}(Y | D), \end{aligned} \tag{1}$$

where we choose $D \subseteq \tilde{R} \cup T$ carefully to ensure that the second equality holds (Methods 4.3). The CRCE Γ_i of $\tilde{X}_i \in \tilde{R}$ thus measures the causal effect of the genetic factors S_i and the non-genetic factors E_i on Y that perturb \tilde{X}_i first. The CRCE values differ between patients, so TWRCI can recover different causal graphs by weighing each vertex according to the patient-specific CRCE values $\Gamma = \gamma$ (Figure 2 (f)). The gene \tilde{X}_i is a *personalized root causal gene* if $|\gamma_i| > 0$. The *omnigenic root causal model* posits that $|\gamma| \gg 0$ for only a small subset of genes in each patient even in complex disease.

2.2 TWRCI accurately annotates, reconstructs and estimates in silico

No existing algorithm recovers personalized root causal genes from observational data alone. However, existing algorithms can annotate variants using different criteria and reconstruct causal graphs from observational data. We therefore compared TWRCI against state of the art algorithms in annotation and causal graph reconstruction using 100 semi-synthetic datasets with real variant data but simulated gene expression and phenotype data (Methods 4.6).

Many different annotation methods exist with different objectives. Most methods nevertheless annotate variants by at least considering proximity to the transcription start site (TSS), with the hope that variants near the TSS of a gene will *directly* affect that gene’s expression level; for example, a variant in the exonic region of a gene may compromise its mRNA stability, while a variant in the promoter region may affect its transcription rate. We thus compare a diverse range of methods in *direct causal* annotation, or assigning variants to the gene expression levels they directly cause. This criterion accommodates other annotation objectives from a mathematical perspective as well – solving direct causation automatically solves causation, colocalization and

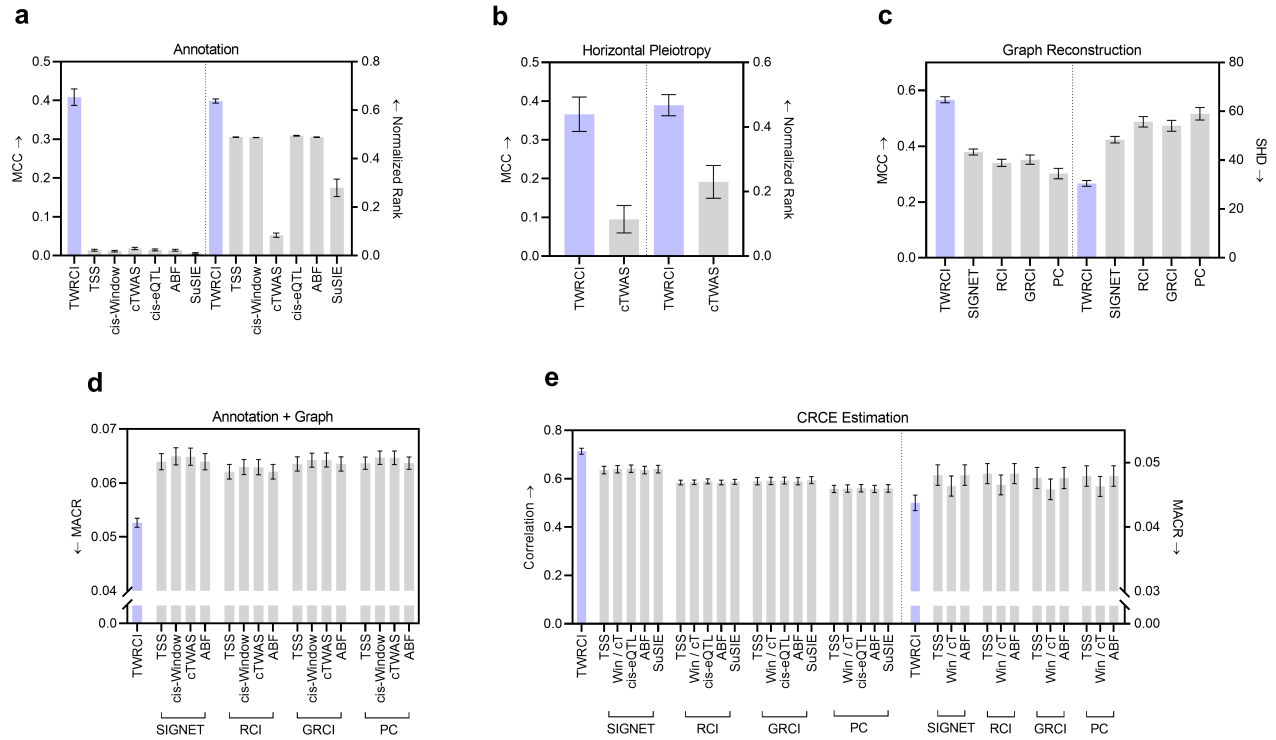


Figure 3. Semi-synthetic data results in terms of (a) direct causal annotation, (b) annotation focused on horizontal pleiotropy only, (c) graph reconstruction, (d) combined annotation and graph reconstruction, and (e) CRCE estimation accuracy. Four of the graphs summarize two evaluation metrics. Arrows near the y-axis denote whether a higher (upward arrow) or a lower (downward arrow) score is better. We do not plot the results of cis-eQTL and SuSIE in (d) and (e) when they exhibit much worse performance. The cis-window and cTWAS algorithms have the exact same CRCE estimates in (e) because accounting for horizontal pleiotropy in cTWAS does not change the conditioning set D in Equation (1); we thus denote cis-Window and cTWAS as Win/cT for short. TWRCI in purple outperformed all algorithms across all nine evaluation metrics. Error bars correspond to 95% confidence intervals.

correlation as progressively more relaxed cases. We in particular compare nearest TSS, a one mega-base cis-window¹, the causal transcriptome-wide association study (cTWAS)¹¹, the maximally correlated gene within the cis-window (cis-eQTL)¹², colocalization with approximate Bayes factors (ABF)¹³, and colocalization with Sum of Single Effects model (SuSIE)¹⁴. We then performed causal graph reconstruction using SIGNET^{15,16}, RCI¹⁷, GRCI¹⁸ and the PC algorithm^{19,20}. We evaluated TWRCI against all combinations of annotation and graph reconstruction methods. See Methods 4.5 and 4.8 for a detailed description of comparator algorithms and evaluation metrics, respectively. All statements about empirical results mentioned below hold at a Bonferroni corrected threshold of 0.05 divided by the number of comparator algorithms according to two-sided paired t-tests.

We first summarize the accuracy results for annotation of direct causes only. All existing annotation algorithms utilize heuristics such as location, correlation or colocalization to infer causality. Only TWRCI provably identifies the direct causes of each gene expression level (Theorem 1 in Methods 4.4.6). Empirical results corroborate this theoretical conclusion. TWRCI achieved the highest accuracy as assessed by Matthew’s correlation coefficient (MCC) to the true direct causal variants of each gene expression level and phenotype (Figure 3 (a) left). The algorithm also ranked the ground truth direct causal variants the highest by assigning the ground truth causal variants larger regression coefficient magnitudes than non-causal variants (Figure 3 (a) right). Both TWRCI and cTWAS account for horizontal pleiotropy, but TWRCI again outperformed cTWAS even when we

¹If multiple genes were present in the window, then we assigned the variant to the gene with the nearest TSS.

only compared the true and inferred variants that directly cause the phenotype using MCC and the normalized rank (Figure 3 (b)). We conclude that TWRCI annotated the genetic variants to their direct effects most accurately.

We obtained similar results with causal graph reconstruction. TWRCI obtained the best performance according to the highest MCC and the lowest structural hamming distance (SHD) to the ground truth causal graphs (Figure 3 (c)). We then assessed the performance of combined annotation and graph reconstruction using the mean absolute correlation of the residuals (MACR), or the mean absolute correlation between the *indirect* causes of a gene expression level and the residual gene expression level obtained after partialing out the inferred *direct* causes; if an algorithm annotates and reconstructs accurately, then each gene expression level should not correlate with its indirect causes after partialing out its direct causes, so the MACR should attain a small value. TWRCI accordingly achieved the lowest MACR as compared to all possible combinations of existing algorithms (Figure 3 (d)). The cis-eQTL and SuSIE algorithms obtained MACR values greater than 0.3 because many cis-variants did not correlate or colocalize with the expression level of the gene with the nearest TSS; we thus do not plot the results of these algorithms. We conclude that TWRCI used annotations to reconstruct the causal graph most accurately by provably accounting for both cis and trans-variants.

We finally analyzed CRCE estimation accuracy. Computing the CRCE requires access to the inferred annotations and causal graph. We therefore again evaluated TWRCI against all possible combinations of existing algorithms. The CRCE estimates of TWRCI attained the largest correlation to the ground truth CRCE values (Figure 3 (e) left). Further, if an algorithm accurately estimates the components $\mathbb{E}(Y|\tilde{X}_i, \mathbf{D})$ and $\mathbb{E}(Y|\mathbf{D})$ of the CRCE in Equation (1), then the residual $Y - \mathbb{E}(Y|\tilde{X}_i, \mathbf{D})$ should not correlate with $S_i \cap T$. TWRCI accordingly obtained the lowest mean absolute correlation of these residuals (MACR) against all combinations of algorithms (Figure 3 (e) right). The cis-eQTL and SuSIE algorithms again attained much worse MACR values above 0.4 because they failed to annotate many causal variants to their gene expression levels. We conclude that TWRCI outperformed existing methods in CRCE estimation. TWRCI therefore annotated, reconstructed and estimated the most accurately according to all nine evaluation criteria. The algorithm also completed within about 3 minutes for each dataset (Supplementary Figure 1).

2.3 Chronic and exaggerated immunity in COPD

We next ran the algorithms using summary statistics of a large GWAS of COPD²¹ consisting of 13,530 cases and 454,945 controls of European ancestry. We downloaded individual variant-expression-phenotype data of lung tissue from GTEx²² with 96 cases and 415 controls. We also replicated results using an independent GWAS consisting of 4,017 cases and 162,653 controls of East Asian ancestry²¹. COPD is a chronic inflammatory condition of the airways or the alveoli that leads to persistent airflow obstruction²³. Exposure to respiratory infections or environmental pollutants can also trigger acute on chronic inflammation called COPD exacerbations that worsen the obstruction.

2.3.1 Accuracy

We first compared the accuracy of the algorithms in variant annotation, graph reconstruction and CRCE estimation. We can compute the MACR metrics – representing two of the nine evaluation criteria used in the previous section – with real data. We summarize the MACR for simultaneous variant annotation and graph reconstruction averaged over ten nested cross-validation folds in Figure 4 (a). TWRCI achieved the lowest MACR out of all combinations of algorithms within about 3 minutes (Supplementary Figure 2 (b) and (c)). Performance differed primarily by the annotation method rather than the causal discovery algorithm. Conservative annotation algorithms, such as colocalization by SuSIE, again failed to achieve a low MACR because they frequently failed to annotate at least one variant to every gene expression level. MACR values for CRCE estimation followed a similar pattern (Figure 4 (b)) because accurate annotation and reconstruction enabled accurate downstream CRCE estimation.

We next followed^{11,24} and downloaded a set of silver standard genes enriched in genes that cause COPD. The KEGG database does not contain a pathway for COPD, so we downloaded the gene set from the DisGeNet database instead (UMLS C0024117, curated)^{25,26}. Many silver standard genes are causal but not *root* causal for COPD. If an algorithm truly identifies

root causal genes, then partialing out the root causal genes from all of the downstream non-root causal genes and the phenotype should explain away the vast majority of the causal effect between the non-root causal genes and the phenotype according to the omnigenic root causal model. We therefore computed another MACR metric, the mean absolute correlation between the residuals of the silver standard genes and the residuals of the phenotype after partialing out the inferred root causal genes. TWRCI again obtained the lowest MACR value (Figure 4 (c)). We conclude that TWRCI identified the root causal genes most accurately according to known causal genes in COPD.

2.3.2 Horizontal pleiotropy and trans-variants

We studied the output of TWRCI in detail to gain insight into important issues in computational genomics. Previous studies have implicated the existence of widespread horizontal pleiotropy in many diseases²⁸. TWRCI can annotate variants directly to the phenotype, so we can use TWRCI to assess the existence of widespread pleiotropy. The variable selection step of TWRCI identified fourteen gene expression levels surviving false discovery rate (FDR) correction at a liberal 10% threshold; eight of these levels ultimately caused the phenotype, including two psoriasis susceptibility genes, a complement protein and five MHC class II genes. TWRCI annotated 13.7% of the variants that cause COPD directly to the phenotype, despite competition for variants between the phenotype and the eight gene expression levels (Figure 4 (d)). Many variants thus directly cause COPD by bypassing expression. We conclude that TWRCI successfully identified widespread horizontal pleiotropy in COPD. In contrast, cTWAS failed to identify any variants that bypass gene expression because all variants had very small effects on the phenotype, especially after accounting for gene expression; as a result, no variants ultimately had a posterior inclusion probability greater than 0.8 according to cTWAS.

TWRCI annotates both cis and trans-variants, so we examined the locations of the annotated variants relative to the TSS for each of the eight causal genes. Most of the variants lying on the same chromosome as the TSS fell within a one megabase distance from the TSS (Figure 4 (e) blue). However, 78% of the variants were located on different chromosomes. We thus compared the variants annotated to causal genes by TWRCI against a previously published list of trans-eQTLs associated with any phenotype in a large-scale search²⁹ (Methods 4.7.3). Variants annotated by TWRCI were located 1.94 times closer to trans-eQTLs than expected by chance (10,000 permutations, $p < 0.001$, 95% CI [1.93,1.95]). We next examined the effect sizes of the variants that cause the phenotype. We regressed the phenotype on variants inferred to directly or indirectly cause the phenotype using linear ridge regression. We then computed the moving average of the magnitudes of the regression coefficients over different distances from the TSS. The magnitudes remained approximately constant with increasing distance from the TSS (Figure 4 (e) red). Moreover, the magnitudes for variants located on different chromosomes did not converge to zero (dotted line). We thus conclude that trans-variants play a significant role in modulating gene expression to cause COPD.

2.3.3 Root Causal Mechanism

We next analyzed the output of TWRCI to elucidate the root causal mechanism of COPD. The pathogenesis of COPD starts with inhaled irritants that trigger an exaggerated and persistent activation of inflammatory cells such as macrophages, T cells and B cells²³. These cells in turn regulate a variety of inflammatory mediators that promote alveolar wall destruction, abnormal tissue repair and mucous hypersecretion obstructing airflow. The root causal genes of COPD therefore likely involve genes mediating chronic and exaggerated inflammation in the lung.

Eight of the fourteen gene expression levels ultimately caused the COPD phenotype in the causal graph reconstructed by TWRCI (Figure 4 (f)). The graph contained five MHC class II genes that present extracellular peptide antigens to CD4+ T cells in the adaptive immune response³⁰. Subsequent activation of T cell receptors regulates a variety of inflammatory mediators and cytokines³¹. Moreover, the complement fragment C4a³² as well as the psoriasis susceptibility genes PSORS1C1 and PSORS1C2³³ help initiate and maintain the exaggerated inflammatory response seen in COPD. The recovered causal graph thus implicates chronic exaggerated inflammation as the root causal mechanism of COPD. TWRCI replicated these results by again discovering C4A and the MHC class II genes in an independent GWAS dataset composed of individuals of East Asian ancestry (Supplementary Figure 3 (a)).

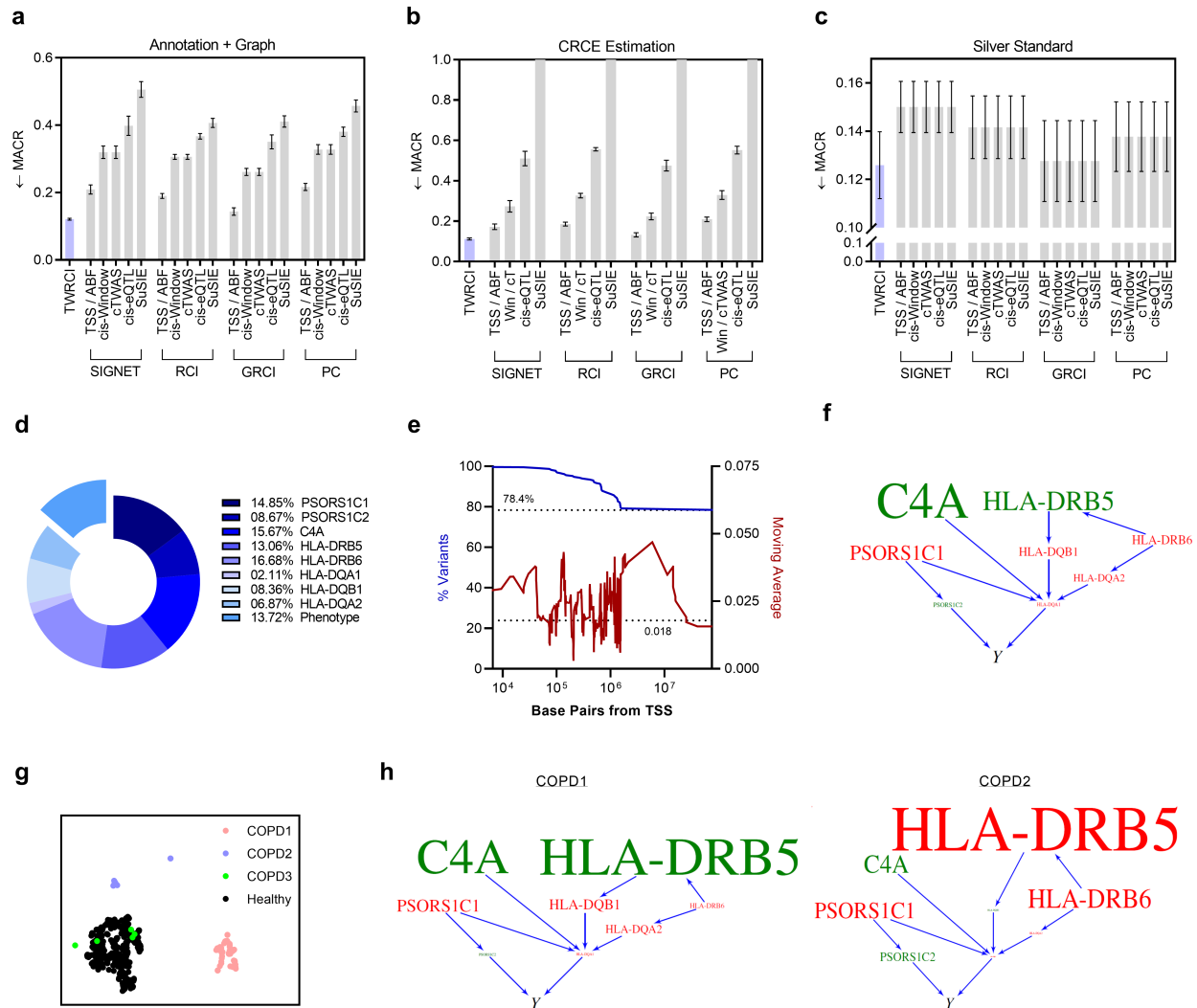


Figure 4. Results for COPD. (a) TWRCI outperformed all other combinations of algorithms in direct causal annotation and graph reconstruction by achieving the lowest MACR; error bars correspond to one standard error of the mean in accordance with the one standard error rule of cross-validation²⁷. (b) TWRCI similarly achieved the lowest MACR for CRCE estimation. (c) Silver standard genes exhibited the smallest correlation with the phenotype after partialing out the root causal genes inferred by TWRCI. (d) More than 13% of the causal variants exhibited horizontal pleiotropy. TWRCI annotated the remaining causal variants to eight gene expression levels. (e) TWRCI assigned approximately 78% of the causal variants to genes located on different chromosomes. Most causal variants annotated to a gene on the same chromosome fell within a one megabase distance from the TSS (blue, left). The average magnitude of the regression coefficients remained approximately constant with increasing distance from the TSS (red, right); the dotted line again corresponds to variants on different chromosomes. (f) The COPD-wide causal graph revealed multiple MHC class II genes as root causal. (g) UMAP dimensionality reduction revealed two clusters of COPD patients well-separated from the healthy controls. (h) The directed graphs highlighted different root causal genes within each of the two clusters.

We finally analyzed the personalized CRCE estimates in more detail. We can decompose the CRCE estimate of each gene into genetic and non-genetic components according to Equation (1). The genetic variants explained only 6.4% of the estimated variance of the CRCE for HLA-DRB5, 1.4% for C4A and <1% for the other six causal genes. We conclude that non-genetic factors account for nearly all of the explained variance in the CRCE estimates. We then performed UMAP dimensionality reduction³⁴ on the causal gene expression levels. Hierarchical clustering with Ward's method³⁵ yielded three clear clusters of patients with COPD (Figure 4 (g)) according to the elbow method on the sum of squares plot (Supplementary Figure 2 (a)). UMAP differentiated two of the COPD clusters from healthy controls, each with different mean CRCE estimates (Figure 4 (h) directed graphs). For example, HLA-DRB5 had a large positive CRCE in cluster one but a large negative CRCE in cluster two. The CRCE estimates thus differentiated multiple subgroups of patients consistent with the known pathobiology of COPD; we likewise obtained similar results in the second GWAS dataset (Supplementary Figure 3 (b) and (c)).

2.4 Oxidative stress in ischemic heart disease

We also ran the algorithms on summary statistics of ischemic heart disease (IHD) consisting of 31,640 cases and 187,152 controls from Finland³⁶. We used variant-expression-phenotype data of whole blood from GTEx²² with 113 cases and 547 controls. We used whole blood because IHD arises from narrowing or obstruction of the coronary arteries most commonly secondary to atherosclerosis with transcription products released into the bloodstream³⁷. We replicated the results using an independent set of GWAS summary statistics from 20,857 cases and 340,337 controls from the UK Biobank³⁸.

2.4.1 Accuracy

We compared the algorithms in variant annotation, graph reconstruction and CRCE estimation accuracy. TWRCI achieved the lowest MACR in both cases (Figure 5 (a) and (b)) within about one hour (Supplementary Figure 4 (b) and (c)). Cis-eQTLs and colocalization with SuSIE failed to annotate many variants because many trans-variants again predicted gene expression. We obtained similar results with a set of silver standard genes downloaded from the KEGG database (hsa05417)³⁹, where TWRCI outperformed all other algorithms (Figure 5 (c)).

2.4.2 Horizontal pleiotropy and trans-variants

The genetic variants predicted 27 gene expression levels at an FDR threshold of 10% with six genes inferred to cause the phenotype. We plot the six genes in the directed graph recovered by TWRCI in Figure 5 (f). TWRCI sorted approximately 8-23% of the causal variants to each of the six genes (Figure 5 (d)). Moreover, TWRCI annotated approximately 17% of the causal variants directly to the phenotype supporting widespread horizontal pleiotropy in IHD. In contrast, cTWAS again did not detect any variants that directly cause the phenotype with a posterior inclusion probability greater than 0.8.

We analyzed the inferred causal effects of cis and trans-variants. Only 7.4% of the annotated variants were located on the same chromosome, and those on the same chromosome were often located over 10 megabases from the TSS (Figure 5 (e) blue). Moreover, variants annotated by TWRCI were located 4.46 times closer to a published list of trans-eQTLs than expected by chance (10,000 permutations, $p = 0.0014$, 95% CI [4.39,4.52]). The magnitudes of the regression coefficients remained approximately constant with increasing distance from the TSS and converged to 0.002 – rather than to zero – on different chromosomes (Figure 5 (e) red). We conclude that trans-variants also play a prominent role in IHD.

2.4.3 Root Causal Mechanism

We next examined the root causal genes of IHD. IHD is usually caused by atherosclerosis, where sites of disturbed laminar flow and altered shear stress trap low-density lipoprotein (LDL)⁴⁰. Reactive oxygen species then oxidize LDL and stimulate an inflammatory response. T cells in turn stimulate macrophages that ingest the oxidized LDL. The macrophages then develop into lipid-laden foam cells that form the initial fatty streak of an eventual atherosclerotic plaque. We therefore expect the root causal genes of IHD to involve oxidative stress and the inflammatory response.

TWRCI identified MRPL1, TRBV6-2 and FAM241B as the top three root causal genes (Figure 5 (f)). MRPL1 encodes a mitochondrial ribosomal protein that helps synthesize complex proteins involved in the respiratory chain⁴¹. Deficiency of

MRPL1 can lead to increased oxidative stress. TRBV6-2 encodes a T-cell receptor beta variable involved in the inflammatory response and accumulation of T-cells in the atherosclerotic plaque⁴². Moreover, knocking out FAM241B induces the cytoplasmic buildup of large lysosome-derived vacuoles that generate foam cells⁴³. We conclude that the root causal genes identified by TWRCI correspond to known genes involved in the pathogenesis of IHD. Finally, TWRCI rediscovered MRPL1 in a second independent GWAS dataset (Supplementary Figure 5 (a)).

We next dissected the CRCE estimates in detail. The annotated variants explained less than 1.5% of the CRCE variance for MRPL1, TRBV6-2 and FAM241B (Figure 5 (g)). Non-genetic factors therefore account for the vast majority of the CRCE variance. UMAP dimensionality reduction and then hierarchical cluster on the causal genes discovered by TWRCI revealed two clusters of IHD patients (Supplementary Figure 4 (a)). The largest of the two clusters lied distal to the cluster of healthy controls (Figure 5 (h)). Furthermore, the FAM241B, TRBV6-2 and MRPL1 genes retained the largest mean CRCEs in this cluster (Figure 5 (i)). TWRCI likewise replicated the large mean CRCE estimate for MRPL1 in the independent GWAS dataset (Supplementary Figure 5 (a) and (b)). We conclude that the CRCE estimates also identify genes that differentiate patient subgroups in IHD.

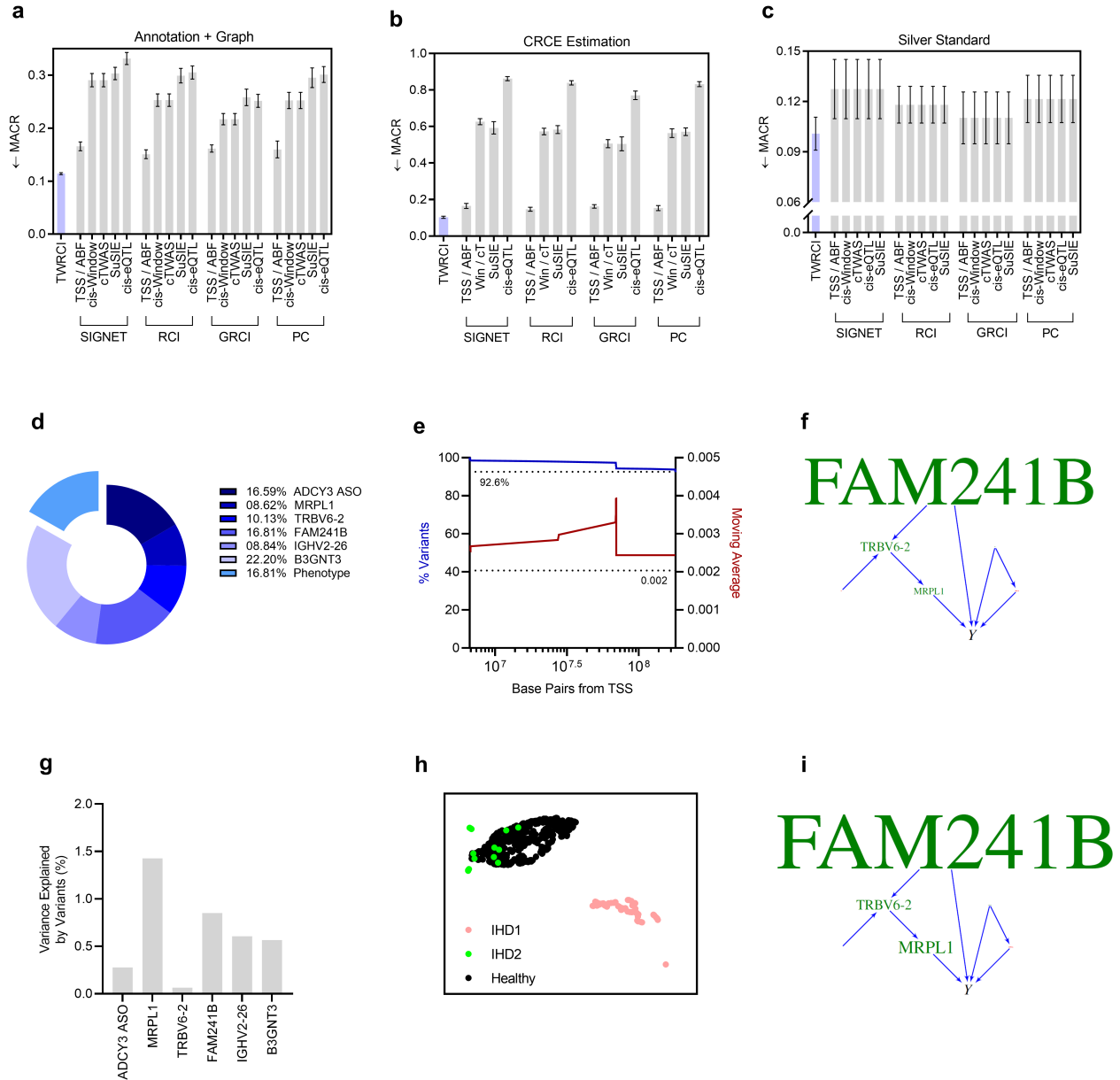


Figure 5. Results for IHD. (a) TWRCI again outperformed all other algorithms in combined annotation and graph reconstruction by achieving the lowest MACR. (b) TWRCI also estimated the CRCEs most accurately relative to all possible combinations of the other algorithms. (c) TWRCI outperformed all other algorithms with a silver standard set of genes causally involved in atherosclerosis. (d) TWRCI annotated varied numbers of variants to six causal expression levels as well as the phenotype. (e) Nearly all of the annotated variants were located distal to the TSS (blue), and the magnitudes of their causal effects did not consistently increase or decrease on average with greater distance from the TSS (red). (f) TWRCI estimated the largest mean CRCEs for MRPL1, TRBV6-2 and FAM241B. (g) The annotated variants only explained a small proportion (<1.5%) of the variance for all CRCE estimates. (h) UMAP dimensionality reduction identified one cluster of patients clearly separated from healthy controls. (i) The mean CRCEs of MRPL1, TRBV6-2 and FAM241B remained the largest in this cluster.

3 Discussion

We introduced the CRCE of a gene, a measure of the causal effect of the genetic and non-genetic factors that directly cause a gene expression level on a phenotype. We then created the TWRCI algorithm that estimates the CRCE of each gene after simultaneously annotating variants and reconstructing the causal graph for improved statistical power. TWRCI annotates, reconstructs and estimates more accurately than alternative algorithms across multiple semi-synthetic and real datasets. Applications of TWRCI to COPD and IHD revealed succinct sets of root causal genes consistent with the known pathogenesis of each disease, which we verified by replication. Furthermore, clustering delineated patient subgroups whose pathogeneses were dictated by different root causal genes.

Our experimental results highlight the importance of incorporating trans-variants in statistical analysis. TWRCI annotated many variants distal to the TSS of each gene. These trans-variants improved the ability of the algorithm to learn models of gene regulation consistent with the correlations in the data according to the MACR criteria. Moreover, variants annotated by TWRCI were located closer to the positions of a previously published list of trans-eQTLs than expected by chance²⁹. In contrast, nearest TSS, cis-windows, cTWAS, cis-eQTLs and the colocalization methods all rely on cis-variants that did not overlap with many GWAS hits both in the COPD and IHD datasets. Most GWAS hits likely lie distal to the TSSs in disease due to natural selection against cis-variants with large causal effects on gene expression⁴⁴. As a result, algorithms that depend solely on cis-variants can fail to detect a large proportion of variants that cause disease in practice.

TWRCI detected widespread horizontal pleiotropy accounting for 13-17% of the causal variants in both the COPD and IHD datasets. Previous studies have detected horizontal pleiotropy in around 20% of causal variants even after considering thousands of gene expression levels as well²⁸. Moreover, many of the variants annotated to the phenotype by TWRCI correlated with gene expression (Supplementary Figures 2 (d) and 4 (d)). Accounting for widespread horizontal pleiotropy thus mitigates pervasive confounding between gene expression levels and the phenotype.

The cTWAS algorithm did not detect widespread pleiotropy in the real datasets. The algorithm also underperformed TWRCI in the semi-synthetic data, even when we restricted the analyses to variants that directly cause the phenotype. We obtained these results because cTWAS relies on the SuSIE algorithm to identify pleiotropic variants. However, pleiotropic variants usually exhibit weak causal relations to the phenotype, so most of these variants do not achieve a large posterior inclusion probability in practice. Algorithms that depend on *absolute* measures of certainty, such as posterior probabilities or p-values, miss many causal variants with weak causal effects in general. TWRCI therefore instead annotates variants by relying on *relative* certainty via a novel process called Competitive Regression, which we showed leads to more consistent causal models across multiple metrics.

We re-emphasize that TWRCI is the only algorithm that accurately recovers *root* causal genes *initiating* pathogenesis. Other methods such as colocalization and cTWAS identify causal genes *involved* in pathogenesis, regardless of whether the genes are root causal or not root causal. As a result, only TWRCI inferred a few genes with large CRCE magnitudes even in complex diseases. Moreover, genes with non-zero CRCE magnitudes explained away most of the causal effects of the non-root causal genes in the silver standards. Both of these results are consistent with the omnigenic root causal model, or the hypothesis that only a few root causal genes initiate the vast majority of pathology in each patient even in complex disease by affecting a very large number of downstream genes¹.

Recall that the above root causal genes differ from driver genes and core genes. Root causal genes generalize driver genes by accounting for all of the factors that directly influence gene expression levels across all diseases, rather than just somatic mutations in cancer³. Accounting for both genetic and non-genetic factors is especially important when non-genetic factors explain the majority of the variance in the root causal effects, as we saw in COPD and IHD. Finally, root causal genes differ from core genes, or the gene expression levels that directly cause a phenotype, by focusing on the beginning rather than the end of pathogenesis². Root causal genes may affect the expression levels of downstream genes so that many genes are differentially expressed between patients and healthy controls including many core genes. A few root causal genes can therefore increase the number of core genes.

TWRCI provably identifies root causal genes and attains high empirical accuracy, but the algorithm carries several limitations.

The algorithm cannot accommodate cycles or directed graphs with different directed edges, even though cycles may exist and direct causal relations may differ between patient populations in practice⁴⁵. TWRCI also estimates CRCE values at the patient-specific level, but the CRCEs may also vary between different cell types. Finally, the algorithm uses linear rather than non-linear models to quantify the causal effects of the variants on gene expression or the phenotype. Future work should therefore consider relaxing the single DAG constraint and accommodating non-linear relations. Future work will also focus on scaling the method to millions of genetic variants without feature selection.

In summary, we introduced an algorithm called TWRCI for accurate estimation and interpretation of the CRCE using personalized causal graphs. TWRCI empirically discovers only a few gene expression levels with large CRCE magnitudes even within different patient subgroups of complex disease in concordance with the omnigenic root causal model⁴⁶. We conclude that TWRCI is a novel, accurate and disease agnostic procedure that couples variant annotation with graph reconstruction to identify root causal genes using observational data alone.

Acknowledgements

Research reported in this manuscript was supported by (1) the National Human Genome Research Institute of the National Institutes of Health under award numbers R01HG011138 and R35HG010718, and (2) the National Institute of General Medical Sciences of the National Institutes of Health under award number R01GM140287.

References

1. Strobl, E. V. & Gamazon, E. R. Discovering root causal genes with high throughput perturbations. *eLife* (2024, in press).
2. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An expanded view of complex traits: from polygenic to omnigenic. *Cell* **169**, 1177–1186 (2017).
3. Martínez-Jiménez, F. *et al.* A compendium of mutational cancer driver genes. *Nat. Rev. Cancer* **20**, 555–572 (2020).
4. Strobl, E. V., Lasko, T. A. & Gamazon, E. R. Mitigating pathogenesis for target discovery and disease subtyping. *Comput. Biol. Medicine* 108122 (2024).
5. Dixit, A. *et al.* Perturb-seq: dissecting molecular circuits with scalable single-cell rna profiling of pooled genetic screens. *Cell* **167**, 1853–1866 (2016).
6. Replogle, J. M. *et al.* Mapping information-rich genotype-phenotype landscapes with genome-scale perturb-seq. *Cell* **185**, 2559–2575 (2022).
7. Huang, C. *et al.* Personal transcriptome variation is poorly explained by current genomic deep learning models. *Nat. Genet.* **55**, 2056–2059 (2023).
8. Sasse, A. *et al.* Benchmarking of deep neural networks for predicting personal gene expression from dna sequence highlights shortcomings. *Nat. Genet.* **55**, 2060–2064 (2023).
9. Zhang, Y. & Politis, D. N. Ridge regression revisited: Debiasing, thresholding and bootstrap. *The Annals Stat.* **50**, 1401–1422 (2022).
10. Lousdal, M. L. An introduction to instrumental variable assumptions, validation and estimation. *Emerg. Themes Epidemiol.* **15**, 1 (2018).
11. Zhao, S. *et al.* Adjusting for genetic confounders in transcriptome-wide association studies improves discovery of risk genes of complex traits. *Nat. Genet.* **56**, 336–347 (2024).
12. Rockman, M. V. & Kruglyak, L. Genetics of global gene expression. *Nat. Rev. Genet.* **7**, 862–872 (2006).
13. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).

14. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection in regression, with application to genetic fine mapping. *J. Royal Stat. Soc. Ser. B: Stat. Methodol.* **82**, 1273–1300 (2020).
15. Chen, C., Ren, M., Zhang, M. & Zhang, D. A two-stage penalized least squares method for constructing large systems of structural equations. *J. Mach. Learn. Res.* **19**, 1–34 (2018).
16. Jiang, Z. *et al.* Signet: transcriptome-wide causal inference for gene regulatory networks. *Sci. Reports* **13**, 19371 (2023).
17. Strobl, E. V. & Lasko, T. A. Identifying patient-specific root causes of disease. In *Proceedings of the 13th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*, 1–10 (2022).
18. Strobl, E. V. & Lasko, T. A. Identifying patient-specific root causes with the heteroscedastic noise model. *J. Comput. Sci.* **72**, 102099 (2023).
19. Spirtes, P., Glymour, C. & Scheines, R. *Causation, Prediction, and Search* (MIT press, 2000), 2nd edn.
20. Wen, Y. *et al.* Applying causal discovery to single-cell analyses using causalcell. *Elife* **12**, e81464 (2023).
21. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human phenotypes. *Nat. Genet.* **53**, 1415–1424 (2021).
22. Consortium, G. The gtex consortium atlas of genetic regulatory effects across human tissues. *Science* **369**, 1318–1330 (2020).
23. Barnes, P. J. Inflammatory mechanisms in patients with chronic obstructive pulmonary disease. *J. Allergy Clin. Immunol.* **138**, 16–27 (2016).
24. Zhou, D. *et al.* A unified framework for joint-tissue transcriptome-wide association and mendelian randomization analysis. *Nat. Genet.* **52**, 1239–1246 (2020).
25. Caramori, G. *et al.* Copd immunopathology. In *Seminars in Immunopathology*, vol. 38, 497–515 (Springer, 2016).
26. Piñero, J. *et al.* The disgenet knowledge platform for disease genomics: 2019 update. *Nucleic Acids Res.* **48**, D845–D855 (2020).
27. Breiman, L. *Classification and regression trees* (Routledge, 2017).
28. Verbanck, M., Chen, C.-Y., Neale, B. & Do, R. Detection of widespread horizontal pleiotropy in causal relationships inferred from mendelian randomization between complex traits and diseases. *Nat. Genet.* **50**, 693–698 (2018).
29. Vösa, U. *et al.* Large-scale cis-and trans-eqtl analyses identify thousands of genetic loci and polygenic scores that regulate blood gene expression. *Nat. Genet.* **53**, 1300–1310 (2021).
30. Nurwidya, F., Damayanti, T. & Yunus, F. The role of innate and adaptive immune cells in the immunopathogenesis of chronic obstructive pulmonary disease. *Tuberc. Respir. Dis.* **79**, 5 (2016).
31. West, E. E., Kolev, M. & Kemper, C. Complement and the regulation of t cell responses. *Annu. Rev. Immunol.* **36**, 309–338 (2018).
32. Detsika, M., Palamaris, K., Dimopoulou, I., Kotanidou, A. & Orfanos, S. The complement cascade in lung injury and disease. *Respir. Res.* **25**, 20 (2024).
33. Li, X. *et al.* Association between psoriasis and chronic obstructive pulmonary disease: a systematic review and meta-analysis. *PLoS One* **10**, e0145221 (2015).
34. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).
35. Ward Jr, J. H. Hierarchical grouping to optimize an objective function. *J. Am. Stat. Assoc.* 236–244 (1963).
36. Elsworth, B. *et al.* The mrc ieu opengwas data infrastructure. *BioRxiv* 2020–08 (2020).

37. Jensen, R. V., Hjortbak, M. V. & Bøtker, H. E. Ischemic heart disease: an update. In *Seminars in Nuclear Medicine*, vol. 50, 195–207 (Elsevier, 2020).
38. Bycroft, C. *et al.* The uk biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
39. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. Kegg: new perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–D361 (2017).
40. Batty, M., Bennett, M. R. & Yu, E. The role of oxidative stress in atherosclerosis. *Cells* **11**, 3843 (2022).
41. Gan, X. *et al.* Tag-mediated isolation of yeast mitochondrial ribosome and mass spectrometric identification of its new components. *Eur. J. Biochem.* **269**, 5203–5214 (2002).
42. Saigusa, R., Winkels, H. & Ley, K. T cell subsets and functions in atherosclerosis. *Nat. Rev. Cardiol.* **17**, 387–401 (2020).
43. Lenk, G. M. *et al.* Crispr knockout screen implicates three genes in lysosome function. *Sci. Reports* **9**, 9609 (2019).
44. Mostafavi, H., Spence, J. P., Naqvi, S. & Pritchard, J. K. Systematic differences in discovery of genetic effects on gene expression and complex traits. *Nat. Genet.* **55**, 1866–1875 (2023).
45. Strobl, E. V. Causal discovery with a mixture of dags. *Mach. Learn.* 1–25 (2022).
46. Yang, W. *et al.* Promoter-sharing by different genes in human genome—cpne1 and rbm12 gene pair as an example. *BMC Genomics* **9**, 1–16 (2008).
47. Lauritzen, S. L., Dawid, A. P., Larsen, B. N. & Leimer, H.-G. Independence properties of directed markov fields. *Networks* **20**, 491–505 (1990).
48. Choudhary, S. & Satija, R. Comparison and evaluation of statistical error models for scrna-seq. *Genome Biol.* **23**, 27 (2022).
49. Sarkar, A. & Stephens, M. Separating measurement and expression models clarifies confusion in single-cell rna sequencing analysis. *Nat. Genet.* **53**, 770–777 (2021).
50. Bowley, A. The standard deviation of the correlation coefficient. *J. Am. Stat. Assoc.* **23**, 31–34 (1928).
51. Stone, M. Cross-validators choice and assessment of statistical predictions. *J. Royal Stat. Soc. Ser. B (Methodological)* **36**, 111–133 (1974).
52. Storey, J. D. The positive false discovery rate: a bayesian interpretation and the q-value. *The Annals Stat.* **31**, 2013–2035 (2003).
53. Colombo, D., Maathuis, M. H. *et al.* Order-independent constraint-based causal structure learning. *J. Mach. Learn. Res.* **15**, 3741–3782 (2014).
54. Cristianini, N. & Shawe-Taylor, J. *An introduction to support vector machines and other kernel-based learning methods* (Cambridge University Press, 2000).
55. Strobl, E. V., Zhang, K. & Visweswaran, S. Approximate kernel-based conditional independence tests for fast non-parametric causal discovery. *J. Causal Inference* **7**, 20180017 (2019).
56. Shimizu, S., Hoyer, P. O., Hyvärinen, A., Kerminen, A. & Jordan, M. A linear non-gaussian acyclic model for causal discovery. *J. Mach. Learn. Res.* **7** (2006).
57. Hoyer, P., Janzing, D., Mooij, J. M., Peters, J. & Schölkopf, B. Nonlinear causal discovery with additive noise models. *Adv. Neural Inf. Process. Syst.* **21** (2008).
58. Danecek, P. *et al.* Twelve years of samtools and bcftools. *Gigascience* **10**, giab008 (2021).
59. Westra, H.-J. *et al.* Systematic identification of trans eqtls as putative drivers of known disease associations. *Nat. Genet.* **45**, 1238–1243 (2013).

60. Matthews, B. W. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophys. Acta (BBA)-Protein Struct.* **405**, 442–451 (1975).
61. Acid, S. & de Campos, L. M. Searching for bayesian network structures in the space of restricted acyclic partially directed graphs. *J. Artif. Intell. Res.* **18**, 445–490 (2003).

4 Methods

4.1 Background on Causal Discovery

Causal discovery refers to the process of discovering causal relations from data. We let italicized letters such as Z_i denote a singleton random variable and bold italicized letters such as \mathbf{Z} denote sets of random variables. Calligraphic letters such as \mathcal{Z} refer to sets of sets.

We consider a set of p endogenous variables \mathbf{Z} . We represent a causal process over \mathbf{Z} using a *structural equation model* (SEM) consisting of a series of deterministic functions:

$$Z_i = f_i(\text{Pa}(Z_i), E_i) \quad \forall Z_i \in \mathbf{Z}, \quad (2)$$

where $\text{Pa}(Z_i) \subseteq \mathbf{Z} \setminus Z_i$ denotes the *parents*, of direct causes, of Z_i and $E_i \in \mathbf{E}$ an *exogenous variable*, also called an *error* or a *noise term*. We assume that the variables in \mathbf{E} are mutually independent. The set $\text{Ch}(Z_i)$ refers to the *children*, or direct effects, of Z_i where $Z_j \in \text{Ch}(Z_i)$ if and only if $Z_i \in \text{Pa}(Z_j)$.

We can associate an SEM with a *directed graph* \mathbb{G} by a drawing a directed edge from Z_j to Z_i when $Z_j \in \text{Pa}(Z_i)$. We thus use the words *variable* and *vertex* interchangeably. A *root vertex* in \mathbb{G} refers to a vertex without any parents, whereas a *sink* or *terminal vertex* refers to a vertex without any children. A *path* between Z_0 and Z_n corresponds to an ordered sequence of distinct vertices $\langle Z_0, \dots, Z_n \rangle$ such that Z_i and Z_{i+1} are adjacent for all $0 \leq i \leq n-1$. In contrast, a *directed path* from Z_0 to Z_n corresponds to an ordered sequence of distinct vertices $\langle Z_0, \dots, Z_n \rangle$ such that $Z_i \in \text{Pa}(Z_{i+1})$ for all $0 \leq i \leq n-1$. We say that Z_j is an *ancestor* of Z_i , and likewise that Z_i is a *descendant* of Z_j , if there exists a directed path from Z_j to Z_i (or $Z_j = Z_i$). We collect all ancestors of Z_j into the set $\text{Anc}(Z_j)$, and all its non-descendants into the set $\text{Nd}(Z_j)$. We write $Z_i \in \text{Anc}(\mathbf{A})$ when Z_i is an ancestor of any variable in \mathbf{A} , and likewise $\text{Nd}(\mathbf{A})$ for the non-descendants. The variable Z_j *causes* Z_i if Z_j is an ancestor of Z_i and $Z_j \neq Z_i$. A *root cause* of Z_i corresponds to a root vertex that also causes Z_i .

A *cycle* exists in \mathbb{G} when Z_j causes Z_i and vice versa. A *directed acyclic graph* (DAG) corresponds to a directed graph without cycles. A *collider* corresponds to Z_j in the triple $Z_i \rightarrow Z_j \leftarrow Z_k$. Two vertices Z_i and Z_j are *d-connected* given $\mathbf{W} \subseteq \mathbf{Z} \setminus \{Z_i, Z_j\}$ if there exists a path between Z_i and Z_j such that no non-collider is in \mathbf{W} and all colliders are ancestors of \mathbf{W} . We denote d-connection by $Z_i \perp\!\!\!\perp_d Z_j | \mathbf{W}$ for shorthand. The two vertices are *d-separated* given \mathbf{W} , likewise denoted by $Z_i \perp\!\!\!\perp_d Z_j | \mathbf{W}$, if they are not d-connected. The *Markov boundary* of Z_i , denoted by $\text{Mb}(Z_i)$, corresponds to the not necessarily unique but smallest set of variables in $\mathbf{Z} \setminus Z_i$ such that $Z_i \perp\!\!\!\perp_d (\mathbf{Z} \setminus \text{Mb}(Z_i)) | \text{Mb}(Z_i)$. A path is *blocked* by \mathbf{W} if \mathbf{W} contains at least one non-collider on the path or does not contain an ancestor of a collider (or both).

A probability density that obeys an SEM associated with the DAG \mathbb{G} also factorizes according to the graph:

$$p(\mathbf{Z}) = \prod_{i=1}^p p(Z_i | \text{Pa}(Z_i)).$$

Any density that factorizes as above obeys the *global Markov property*, where Z_i and Z_j are conditionally independent given \mathbf{W} , or $Z_i \perp\!\!\!\perp Z_j | \mathbf{W}$, if $Z_i \perp\!\!\!\perp_d Z_j | \mathbf{W}^{47}$. A density obeys *d-separation faithfulness* when the converse holds: if $Z_i \perp\!\!\!\perp Z_j | \mathbf{W}$, then $Z_i \perp\!\!\!\perp_d Z_j | \mathbf{W}$. The Markov boundary of Z_i uniquely corresponds to the parents, children and parents of the children (or *spouses*) of Z_i under d-separation faithfulness.

4.2 Causal Modeling of Variants, Gene Expression and the Phenotype

We divide the set of random variables \mathbf{Z} into disjoint sets $Y \cup S \cup L \cup \tilde{X}$ corresponding to the phenotype Y , q genetic variants S , latent variables L modeling linkage disequilibrium (LD) and m gene expression levels \tilde{X} . We model the causal process over \mathbf{Z} using the following SEM associated with a DAG \mathbb{G} :

$$\begin{aligned}
 L_i &= f_i(\text{Pa}(L_i), E_i), & \forall L_i \in L \\
 S_j &= f_j(\text{Pa}(S_j), E_j), & \forall S_j \in S \\
 \tilde{X}_k &= f_k(\text{Pa}(\tilde{X}_k), E_k), & \forall \tilde{X}_k \in \tilde{X}, \\
 Y &= f_Y(\text{Pa}(Y), E_Y),
 \end{aligned}
 \tag{3}$$

where $\text{Pa}(L_i) \subseteq L$, $\text{Pa}(S_j) \subseteq L$, $\text{Pa}(\tilde{X}_k) \subseteq (\tilde{X} \cup S)$ and $\text{Pa}(Y) \subseteq (\tilde{X} \cup S)$ for any latent variable, any genetic variant, any gene expression level and the phenotype, respectively. In other words, linkage disequilibrium L generates variants S , and variants and gene expression generate other gene expression levels \tilde{X} and the phenotype Y (example in Figure 6 (a)). We assume that Y is a sink vertex such that gene expression and variants cause Y but not vice versa.

Let S_i denote the direct causes of \tilde{X}_i in S . We require $S_i \neq \emptyset$ for all $\tilde{X}_i \in \tilde{X}$ so that at least one variant directly causes each gene expression level. We also assume that any single variant can only *directly* cause one gene expression level or the phenotype (but not both). Investigators have reported only a few rare exceptions to this latter assumption in the literature⁴⁶. A variant may however indirectly cause many gene expression levels.

We unfortunately cannot measure the exact values of gene expression using RNA sequencing (RNA-seq) technology. Numerous theoretical and experimental investigations have revealed that RNA-seq suffers from independent Poisson measurement error^{48,49}:

$$X_i \sim \text{Pois}(\tilde{X}_i \pi_{ij}),$$

where π_{ij} denotes the *mapping efficiency* of \tilde{X}_i in batch j . We thus sample $Y \cup S \cup L \cup X \cup B$ from the DAG like the one shown in Figure 6 (b) in practice, where B denotes the batch. With slight abuse of terminology, we will still call \tilde{X}_i a *sink vertex* if it has only one child X_i .

We can perform consistent regression under Poisson measurement error. Let $N = \sum_{i=1}^m X_i$ denote the library size and let $\tilde{N}_j = \sum_{i=1}^m \tilde{X}_i \pi_{ij}$ denote the true unobserved total gene expression level weighted by the mapping efficiencies in batch j . Also let $\tilde{U} \subseteq \tilde{X}$ and $V \subseteq S$ refer to any subset of gene expression levels and variants, respectively. The following result holds:

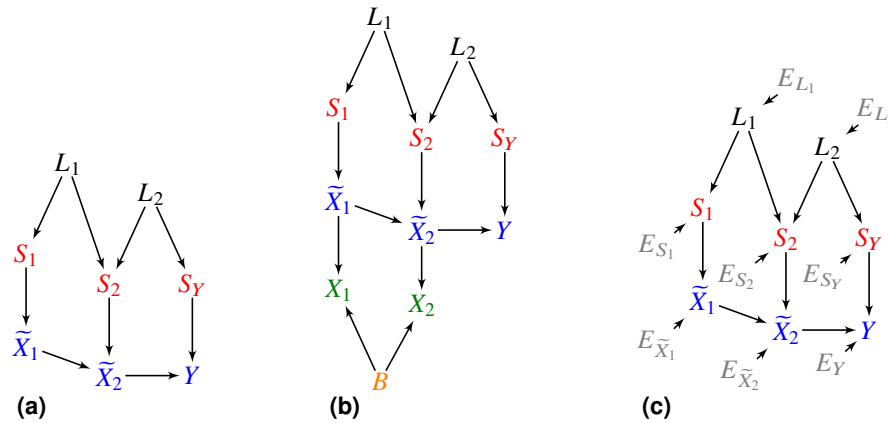


Figure 6. (a) An example of a DAG over \mathbf{Z} . In (b), the additional vertices X denote counts corrupted by batch B effects and Poisson measurement error. (c) We can also augment the DAG in (a) with root vertex error terms E .

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

$$\mathbb{E} \left| \mathbb{E}(Z_i | \tilde{U}, \mathbf{V}) - \mathbb{E}(Z_i | \mathbf{U}, \mathbf{V}, B) \right| \leq \mathbb{E} C_N \left| \tilde{U} - \frac{\mathbf{U}}{N} \frac{\tilde{N}_B}{\pi_{UB}} \right|,$$

where $C_N \in O(1)$ is a positive constant, and we have taken an outer expectation on both sides. Then $\mathbb{E}(Z_i | \tilde{U}, \mathbf{V}) = \lim_{N \rightarrow \infty} \mathbb{E}(Z_i | \mathbf{U}, \mathbf{V}, B)$ almost surely.

We delegate proofs to the Supplementary Materials. Intuitively, $\frac{\mathbf{U}}{N} \frac{\tilde{N}_B}{\pi_{UB}}$ approaches \tilde{U} as the library size increases, so the above lemma states that accurate estimation of \tilde{U} implies accurate estimation of $\mathbb{E}(Z_i | \tilde{U}, \mathbf{V})$. We can thus consistently estimate any conditional expectation $\mathbb{E}(Z_i | \tilde{U}, \mathbf{V})$ using $\mathbb{E}(Z_i | \mathbf{U}, \mathbf{V}, B)$ when the library size approaches infinity. We only apply the asymptotic argument to bulk RNA-seq, where the library size is on the order of at least tens of millions. We henceforth implicitly assume additional conditioning on B whenever regressing to or on bulk RNA-seq data in order to simplify notation.

4.3 Conditional Root Causal Effects

We define the root causal effect of a gene expression level on the phenotype Y . We focus on Equation (3) with the endogenous variables \mathbf{Z} and the exogenous variables \mathbf{E} . If the error terms \mathbf{E} are mutually independent, then we can *augment* the associated DAG \mathbb{G} with \mathbf{E} by drawing a directed edge from each E_{Z_i} to its direct effect Z_i (Figure 6 (c)). We denote the resultant graph by \mathbb{G}' , where we always have $E_{Z_i} \in \text{Pa}_{\mathbb{G}'}(Z_i)$ and the subscript emphasizes the augmented DAG; if we do not place a subscript, then we refer to the original DAG \mathbb{G} . Only the error terms are root vertices in \mathbb{G}' , so only exogenous variables that cause Y can be root causes of Y .

The *root causal effect* of Z_i on Y given the exogenous variables \mathbf{E} is the causal effect of its direct causes in \mathbf{E} on Y :

$$\mathbb{P}(Y | \text{Pa}_{\mathbb{G}'}(Z_i) \cap \mathbf{E}) - \mathbb{P}(Y) = \mathbb{P}(Y | E_{Z_i}) - \mathbb{P}(Y). \quad (4)$$

The variable Z_i is the first variable in \mathbf{Z} affected by E_{Z_i} , and Z_i may in turn causally affect Y . The exogenous variable E_{Z_i} models the effects of environmental, epigenetic and other *non-genetic* factors on Z_i because the set of endogenous variables $\mathbf{Z} = Y \cup \mathbf{S} \cup \mathbf{L} \cup \tilde{\mathbf{X}}$ includes the *genetic factors* \mathbf{S} . The root causal effect is a special case of the *conditional root causal effect* (CRCE) given the exogenous variables \mathbf{E} :

$$\mathbb{P}(Y | E_{Z_i}, \mathbf{D}) - \mathbb{P}(Y | \mathbf{D})$$

where (1) $\mathbf{D} \subseteq \text{Nd}_{\mathbb{G}'}(Z_i) \setminus (E_{Z_i} \cup Z_i)$ and (2) $Y \perp\!\!\!\perp_d E_{Z_i} | Z_i \cup \mathbf{D}$. The first condition ensures that \mathbf{D} does not block any directed path from Z_i to Y . The second ensures that \mathbf{D} eliminates any confounding between E_{Z_i} and Y . The first condition actually implies the second in this case because \mathbf{E} are root vertices. If we set $\mathbf{D} = \emptyset$, then we recover the unconditional root causal effect in Equation (4).

We are however interested in identifying the causal effects of both genetic *and* non-genetic factors on Y through gene expression $\tilde{\mathbf{X}}$ with potential confounding between members of \mathbf{S} due to LD. We therefore expand the set of exogenous variables to $\mathbf{E} \cup \mathbf{S}$ representing the non-genetic and genetic factors, respectively. We define the conditional root causal effect of $\tilde{X}_i \in \tilde{\mathbf{X}}$ given the exogenous variables $\mathbf{E} \cup \mathbf{S}$ as:

$$\begin{aligned} & \mathbb{P}(Y | \text{Pa}_{\mathbb{G}'}(\tilde{X}_i) \cap (\mathbf{E} \cup \mathbf{S}), \mathbf{D}) - \mathbb{P}(Y | \mathbf{D}) \\ &= \mathbb{P}(Y | E_i \cup S_i, \mathbf{D}) - \mathbb{P}(Y | \mathbf{D}), \end{aligned}$$

where we write $E_{\tilde{X}_i} \cup S_{\tilde{X}_i}$ as $E_i \cup S_i$ to prevent cluttering of notation. The set $E_i \cup S_i$ thus refers to the direct causes of \tilde{X}_i in $\mathbf{E} \cup \mathbf{S}$. The above conditional root causal effect measures the causal effect of the root vertices \mathbf{E} on Y as they pass through $E_i \cup S_i$ to \tilde{X}_i .

We can likewise choose any \mathbf{D} such that $\mathbf{D} \subseteq \text{Nd}_{\mathcal{G}'}(\tilde{X}_i) \setminus (E_i \cup \mathbf{S}_i \cup \tilde{X}_i)$ and $Y \perp\!\!\!\perp_d (E_i \cup \mathbf{S}_i) | \tilde{X}_i \cup \mathbf{D}$. We choose \mathbf{D} carefully to satisfy these two conditions as well as elicit favorable mathematical properties by setting $\mathbf{D} = \tilde{\mathbf{V}}_i \cup (\mathbf{T} \setminus \mathbf{S}_i)$, where $\tilde{\mathbf{V}}_i = \text{Pa}_{\mathcal{G}'}(\tilde{X}_i) \cap \tilde{\mathbf{X}}$ and $\mathbf{T} = \{\mathbf{S}_i \in \mathbf{S} : \mathbf{S}_i \perp\!\!\!\perp_d Y\}$. This particular choice of \mathbf{D} allows us to write:

$$\begin{aligned} \Psi_i &= \mathbb{P}(Y | E_i \cup \mathbf{S}_i, \mathbf{D}) - \mathbb{P}(Y | \mathbf{D}), \\ &= \mathbb{P}(Y | \tilde{X}_i, \mathbf{D}) - \mathbb{P}(Y | \mathbf{D}), \end{aligned}$$

so that we do not need to recover E_i as an intermediate step. We prove the second equality in Proposition 1 of the Supplementary Materials under *exchangeability*, or no latent confounding by L between any two entries of $\mathbf{T} \cup \{\mathbf{S}_i \setminus \mathbf{T} : \tilde{X}_i \notin \text{Anc}(Y)\}$; this union corresponds to a set of sets including \mathbf{T} and each entry of $\{\mathbf{S}_i \setminus \mathbf{T} : \tilde{X}_i \notin \text{Anc}(Y)\}$ in the set. Exchangeability holds approximately in practice due to the weak causal relations emanating from variants to gene expression and the phenotype. Moreover, the assumption weakens with more variants in \mathbf{T} . Now the first gene expression level in $\tilde{\mathbf{X}}$ affected by $E_i \cup \mathbf{S}_i$ is \tilde{X}_i . We thus call \tilde{X}_i a *root causal gene* if \tilde{X}_i also causes Y such that $\Psi_i \neq 0$.

We finally focus on the expected version of Ψ_i to enhance computational speed, improve statistical efficiency and overcome Poisson measurement error according to Lemma 1:

$$\Gamma_i = \mathbb{E}(Y | \tilde{X}_i, \mathbf{D}) - \mathbb{E}(Y | \mathbf{D}), \quad (5)$$

The *omnigenic root causal model* posits that $|\gamma| \gg 0$ for only a small subset of gene expression levels in each patient with $\Gamma = \gamma$. We thus seek to estimate the values γ for each patient. We use the acronym CRCEs to specifically refer to Γ from here on.

4.4 Algorithm

4.4.1 Strategy Overview

We seek to accurately annotate, reconstruct and estimate the CRCEs using (1) summary statistics as well as (2) linked variant-expression-phenotype data. We summarize the proposed Transcriptome-Wide Root Causal Inference (TWRCI) algorithm in Algorithm 1. TWRCI first uses summary statistics to identify variants \mathbf{T} associated with the phenotype at a liberal α threshold in Line 1. The algorithm also identifies gene expression levels $\mathbf{R} \subseteq \mathbf{X}$ predictable by \mathbf{T} in Line 1 from the variant-expression-phenotype data. TWRCI then annotates non-overlapping sets of variants to the phenotype in Line 2 and each gene expression level in Line 3 using a novel process called Competitive Regression; we prove that annotated variants include all of the direct causes in \mathbf{T} . TWRCI identifies the causal ordering \mathbf{K} among $\tilde{\mathbf{R}}$ during the annotation process. The algorithm finally recovers the directed graph uniquely given \mathbf{K} in Line 4 and estimates the CRCE of each gene inferred to cause Y using the estimated graph $\hat{\mathbf{G}}$ and the annotations \mathcal{P} in Line 5. TWRCI can thus weigh and color-code each node in $\hat{\mathbf{G}}$ that causes Y by the CRCE estimates for each patient. We will formally prove that TWRCI is sound and complete at the end of this subsection.

Algorithm 1 Transcriptome-Wide Root Causal Inference (TWRCI)

Input: summary statistics, $\mathbf{S} \cup \mathbf{X} \cup Y$

Output: $\mathcal{P}, \mathbf{K}, \hat{\mathbf{G}}, \Gamma$

- 1: $\mathbf{T}, \mathbf{R}, \mathbf{N} \leftarrow$ Variable selection with Algorithm 2
 - 2: $\mathbf{P}_Y \leftarrow$ Annotate some variants in \mathbf{T} to Y using Algorithm 3
 - 3: $\mathcal{P}, \mathbf{K} \leftarrow$ Annotate remaining variants in \mathbf{T} to gene expression levels and obtain the causal order using Algorithm 4
 - 4: $\hat{\mathbf{G}} \leftarrow$ Recover DAG using Algorithms 5 and 6
 - 5: $\Gamma \leftarrow$ Compute CRCE of each gene inferred to cause Y using $\hat{\mathbf{G}}$ and \mathcal{P}
-

4.4.2 Variable Selection

We summarize the variable selection portion of TWRCI in Algorithm 2. TWRCI first reduces the number of variants using summary statistics by only keeping variants with a significant association to the phenotype at a very liberal α threshold (Line 1); we use $5e-5$, or a three orders of magnitude increase from the usual threshold of $5e-8$. We do not employ clumping or other

pre-processing methods that may remove more variants from consideration. Let \mathbf{T} denote the variants that survive this screening step so that $\mathbf{T} = \{S_i \in \mathbf{S} : S_i \not\perp_d Y\}$.

The variable selection algorithm then identifies the gene expression levels predictable by \mathbf{T} using the variant-expression-phenotype data in Line 2. We operationalize this step by linearly regressing \mathbf{X} on \mathbf{T} using half of the samples, and then testing whether the predicted level \widehat{X}_i and the true level X_i linearly correlate in the second half for each $X_i \in \mathbf{X}^{50}$. This sample splitting procedure ensures proper control of the Type I error rate⁵¹. We keep gene expression levels $\mathbf{R} \subseteq \mathbf{X}$ that achieve a q-value below a liberal FDR threshold of 10%⁵². We say that \mathbf{T} is *relevant* if it contains at least one variant that directly causes each member of $\widetilde{\mathbf{R}}$. We finally repeat the above procedure after regressing out \mathbf{R} from $\mathbf{X} \setminus \mathbf{R}$ and \mathbf{T} in Line 3 in order to identify $\widetilde{\mathbf{N}}$, or all parents of $\widetilde{\mathbf{R}}$ in $\widetilde{\mathbf{X}} \setminus \widetilde{\mathbf{R}}$. We call \mathbf{N} the set of *nuisance variables*, since we will need to condition on them, but they do not contain the ancestors of Y . Algorithm 2 formally identifies the necessary ancestors needed for downstream inference:

Lemma 2. *Assume d -separation faithfulness and relevance. Then, (1) $\mathbf{T} \cup \widetilde{\mathbf{R}}$ contains all of the ancestors of Y in $\mathbf{S} \cup \widetilde{\mathbf{X}}$, and (2) $(\text{Mb}(\widetilde{\mathbf{R}}_i) \cap \widetilde{\mathbf{X}}) \subseteq (\widetilde{\mathbf{R}} \setminus \widetilde{\mathbf{R}}_i) \cup \widetilde{\mathbf{N}}$ for any $\widetilde{\mathbf{R}}_i \in \widetilde{\mathbf{R}}$.*

Algorithm 2 Variable Selection

Input: summary statistics, $\mathbf{S} \cup \mathbf{X} \cup Y$

Output: $\mathbf{T}, \mathbf{R}, \mathbf{N}$

- 1: $\mathbf{T} \leftarrow S_i \in \mathbf{S}$ such that $S_i \not\perp Y$ using summary statistics
 - 2: $\mathbf{R} \leftarrow X_i \in \mathbf{X}$ such that $X_i \not\perp \mathbf{T}$ using variant-expression-phenotype data
 - 3: $\mathbf{N} \leftarrow X_i \in \mathbf{X} \setminus \mathbf{R}$ such that $X_i \not\perp \mathbf{T} | \mathbf{R}$ using variant-expression-phenotype data
-

4.4.3 Annotation for Horizontal Pleiotropy

TWRCI next annotates the associated variants \mathbf{T} to their direct effects in $\widetilde{\mathbf{R}} \cup Y$. The algorithm first annotates a sink vertex and then gradually works its way up the DAG until it annotates the final root vertex.

TWRCI assumes that Y is a sink vertex, so it first annotates to Y . A variant exhibits *horizontal pleiotropy* if it directly causes Y . We propose a novel competitive regression scheme to annotate all members of $\mathbf{T} \cap \text{Pa}(Y) = \mathbf{S}_Y$ to Y .

We mildly assume equality in conditional expectation implies equality in conditional distribution and vice versa. Let $\widetilde{\mathbf{Q}} = \widetilde{\mathbf{R}} \cup Y$ and likewise $\mathbf{Q} = \mathbf{R} \cup Y$. We also mildly assume that the following *contribution scores* exist and are finite: $\Delta_{ij} = \mathbb{E}|\partial \mathbb{E}(Q_i | \mathbf{T}) / \partial T_j|$ and $\gamma_{ij} = \mathbb{E}|\partial \mathbb{E}(Q_i | \widetilde{\mathbf{Q}} \setminus \widetilde{Q}_i, \mathbf{T}) / \partial T_j|$. The scores correspond to the variable coefficients in linear regression. We use the contribution scores to annotate any $T_j \in \mathbf{T}$ such that $|\Delta_{Yj} \gamma_{Yj}| \geq \max_{\mathbf{R}_j} \Delta_{\mathbf{R}_j}^2$ to Y , since this set of variants corresponds to a superset of \mathbf{S}_Y by the following result:

Corollary 1. *Under d -separation faithfulness, relevance and exchangeability, $|\Delta_{Yj} \gamma_{Yj}| \geq \max_{\mathbf{R}_j} \Delta_{\mathbf{R}_j}^2$ if and only if $T_j \notin \text{Anc}(\mathbf{R})$ or $T_j \in \text{Pa}(Y)$ (or both).*

The proof follows directly from Lemma 3 in the Supplementary Materials.

The Competitive Regression (CR) algorithm summarized in Algorithm 3 computes the contribution scores in order to annotate variants to Y . Let Δ_{-i} denote the removal of the i^{th} row from Δ corresponding to $Q_i = Y$. We use debiased linear ridge regression⁹ to compute Δ in Line 1 and γ_i in Line 2. CR compares the two quantities and outputs the set $\mathbf{P}_i = \{T_j : |\Delta_{ij} \gamma_{ij}| \geq \max_{-ij} \Delta_{-ij}^2\}$, or a superset of $\mathbf{S}_i \cap \mathbf{T}$ not including any other variants with children in $\widetilde{\mathbf{Q}} \setminus \widetilde{Q}_i$ according to Corollary 1, in Line 3.

4.4.4 Annotation and Causal Order

The CR algorithm requires the user to specify a known sink vertex. We drop this assumption by integrating CR into the Annotation and Causal Order (ACO) algorithm that automatically finds a sink vertex at each iteration.

ACO takes $\mathbf{R}, \mathbf{N}, Y, \mathbf{T}, \mathbf{P}_Y$ as input as summarized in Algorithm 4. The algorithm constructs a causal ordering over $\mathbf{R} \cup Y$ in \mathbf{K} by iteratively eliminating a sink vertex from \mathbf{R} and appending it to the front of \mathbf{K} . ACO also instantiates a list \mathcal{P} and assigns

Algorithm 3 Competitive Regression (CR)

Input: T, Q, N, \tilde{Q}_i

Output: P_i

- 1: $\Delta \leftarrow$ Matrix of coefficients with rows obtained after regressing Q_j on T for all $Q_j \in Q$
 - 2: $\gamma_i \leftarrow$ Row vector of coefficients obtained after regressing Q_i on T and $(\tilde{Q} \setminus \tilde{Q}_i) \cup \tilde{N}$
 - 3: $P_i \leftarrow \{T_j : |\Delta_{ij}\gamma_{ij}| \geq \max \Delta_{-ij}^2\}$
-

Algorithm 4 Annotation and Causal Order (ACO)

Input: R, N, Y, T, P_Y

Output: K, \mathcal{P}

- 1: $\mathcal{P} \leftarrow$ Empty list
 - 2: $K \leftarrow Y; O \leftarrow P_Y$
 - 3: **repeat**
 - 4: $\Delta \leftarrow$ Contributions after regressing R on $T \setminus O$
 - 5: $C \leftarrow \emptyset$
 - 6: **for all** $R_i \in R$ **do**
 - 7: $\gamma_i \leftarrow$ Contributions after regressing R_i on $(\tilde{R} \setminus \tilde{R}_i) \cup \tilde{N} \cup (T \setminus O)$
 - 8: $U_i \leftarrow \{T_j : |\Delta_{ij}\gamma_{ij}| \geq \max \Delta_{-ij}^2\}$
 - 9: **if** $U_i = \emptyset$ **then**
 - 10: $C_i \leftarrow \infty$
 - 11: **else**
 - 12: $C_i \leftarrow$ Measure of dependence between R_i and $T \setminus (O \cup U_i)$ given $(\tilde{R} \setminus \tilde{R}_i) \cup \tilde{N} \cup U_i$
 - 13: **end if**
 - 14: **end for**
 - 15: $R_i \leftarrow$ Most independent variable in R according to C
 - 16: $K \leftarrow$ Append R_i to the front of K
 - 17: $R \leftarrow R \setminus R_i$
 - 18: $P_i \leftarrow U_i$
 - 19: $O \leftarrow O \cup P_i$
 - 20: **until** $R = \emptyset$
-

genetic variants $P_i = \{T_j : |\Delta_{ij}\gamma_{ij}| \geq \max \Delta_{-ij}^2\} \in \mathcal{P}$ to each gene expression level $R_i \in R$ in Lines 8 and 18 using the following generalization of Corollary 1:

Lemma 3. *Assume d -separation faithfulness, relevance and exchangeability. Further assume that \tilde{Q}_i is a sink vertex. Then, $|\Delta_{ij}\gamma_{ij}| \geq \max \Delta_{-ij}^2$ if and only if $T_j \notin \text{Anc}(Q \setminus Q_i)$ or $T_j \in \text{Pa}(\tilde{Q}_i)$ (or both).*

The set P_i is thus again a superset of $S_i \cap T$, and any additional variants in P_i do not directly cause another gene expression level or the phenotype.

ACO determines whether \tilde{R}_i is indeed a sink vertex from data using the following result:

Lemma 4. *\tilde{R}_i is a sink vertex if and only if $R_i \perp\!\!\!\perp (T \setminus U_i) | (\tilde{R} \setminus \tilde{R}_i) \cup \tilde{N} \cup U_i$ in Line 12 of ACO under d -separation faithfulness, relevance and exchangeability.*

ACO practically determines whether any \tilde{R}_i is indeed a sink vertex post variable elimination by first computing the residuals F_i after regressing R_i on $\tilde{R} \setminus \tilde{R}_i$, the nuisance variables \tilde{N} and the identified variants U_i . A sink vertex \tilde{R}_i has residuals F_i that are uncorrelated with the variants in $T \setminus (O \cup U_i)$ in Line 12 by Lemma 4, so ACO can identify the sink vertex \tilde{R}_i in Line 15 as the variable with the smallest absolute linear correlation. The algorithm then appends R_i to the front of K and eliminates R_i from R in Lines 16 and 17, respectively. ACO finally adds U_i to \mathcal{P} in Line 18, so U_i can be removed from T of the next iteration through O . We formally prove the following result:

Lemma 5. Under d -separation faithfulness, relevance and exchangeability, ACO recovers the correct causal order \mathbf{K} over $\tilde{\mathbf{R}}$ and $(S_i \cap T) \subseteq P_i$ for all $R_i \in \tilde{\mathbf{R}}$.

4.4.5 Causal Graph Discovery

TWRCI uses the causal order \mathbf{K} and the annotations \mathcal{P} to perform causal discovery. The algorithm runs the (stabilized) skeleton discovery procedure of the Peter-Clark (PC) algorithm to identify the presence or absence of edges between any two gene expression levels (Algorithm 5)^{19,53}. We modify the PC algorithm so that it tests whether R_i and R_j are conditionally independent given $P_i \cup \tilde{N}$ and subsets of the neighbors of \tilde{R}_i in $\tilde{\mathbf{R}} \setminus \tilde{R}_i$ in Line 12 to ensure that we condition on all parents of \tilde{R}_i . Finally, we orient the edges using the causal order \mathbf{K} in Line 19 to uniquely recover the DAG over $\tilde{\mathbf{R}}$:

Lemma 6. Under d -separation faithfulness, relevance and exchangeability, the graph discovery algorithm outputs the true sub-DAG over $\tilde{\mathbf{R}}$ given a conditional independence oracle, \mathbf{K} and \mathcal{P} .

We next include the phenotype Y into the causal graph. We often only have a weak causal effect from gene expression and variants to the phenotype. We therefore choose to detect any causal relation to Y rather than just direct causal relations using Algorithm 6. Algorithm 6 only conditions on $\tilde{V}_i \cup P_i$ in Line 4 to discover both direct and indirect causation in concordance with the following result:

Lemma 7. Under d -separation faithfulness, relevance and exchangeability, \tilde{R}_i causes Y – and likewise the vertices $S_i \cup E_i$ cause Y – if and only if $Y \not\perp\!\!\!\perp \tilde{R}_i | \tilde{V}_i \cup P_i$.

Algorithm 5 Graph Discovery

Input: \mathbf{R} , N , \mathcal{P} , \mathbf{K} , type I error rate α

Output: DAG $\hat{\mathbf{G}}$ over $\tilde{\mathbf{R}}$

```

1: Form a fully connected undirected graph  $\hat{\mathbf{G}}$  over  $\tilde{\mathbf{R}}$ 
2:  $l \leftarrow -1$ 
3: repeat
4:   Let  $l = l + 1$ 
5:   repeat
6:     for each  $\tilde{R}_i \in \tilde{\mathbf{R}}$  do
7:        $\text{Adj}_{\hat{\mathbf{G}}}(\tilde{R}_i) \leftarrow$  Vertices adjacent to  $\tilde{R}_i$  in  $\hat{\mathbf{G}}$ 
8:     end for
9:     Select a new ordered pair of vertices  $(\tilde{R}_i, \tilde{R}_j)$  that are adjacent in  $\hat{\mathbf{G}}$  and satisfy  $|\text{Adj}_{\hat{\mathbf{G}}}(\tilde{R}_i) \setminus \tilde{R}_j| \geq l$ 
10:    repeat
11:      Choose a new set  $\mathbf{W} \subseteq \text{Adj}_{\hat{\mathbf{G}}}(\tilde{R}_i) \setminus \tilde{R}_j$  with  $|\mathbf{W}| = l$ 
12:      Test whether  $R_i$  and  $R_j$  are independent given  $\tilde{\mathbf{W}} \cup \tilde{N} \cup P_i$  to obtain p-value  $p$ 
13:      if  $p > \alpha$  then
14:        Delete the edge  $\tilde{R}_i - \tilde{R}_j$  from  $\hat{\mathbf{G}}$ 
15:      end if
16:    until  $\tilde{R}_i$  and  $\tilde{R}_j$  are no longer adjacent in  $\hat{\mathbf{G}}$  or all such subsets with  $|\mathbf{W}| = l$  have been considered
17:  until all ordered pairs of adjacent vertices  $(\tilde{R}_i, \tilde{R}_j)$  in  $\hat{\mathbf{G}}$  with  $|\text{Adj}_{\hat{\mathbf{G}}}(\tilde{R}_i) \setminus \tilde{R}_j| \geq l$  have been considered
18: until all pairs of adjacent vertices  $(\tilde{R}_i, \tilde{R}_j)$  in  $\hat{\mathbf{G}}$  satisfy  $|\text{Adj}_{\hat{\mathbf{G}}}(\tilde{R}_i) \setminus \tilde{R}_j| \leq l$ 
19: Orient the edges of  $\hat{\mathbf{G}}$  according to the causal order  $\mathbf{K}$ 

```

4.4.6 Conditional Root Causal Effect Estimation

TWRCI finally estimates the CRCEs of the genes that cause Y given the recovered graph $\hat{\mathbf{G}}$ and the annotations \mathcal{P} . We estimate the two conditional expectations in Equation (5) using kernel ridge regression⁵⁴. We embed X_i and $\text{Pa}_{\hat{\mathbf{G}}}(\tilde{R}_i) = \tilde{V}_i$ using a radial basis function kernel but embed $T \setminus P_i$ using a normalized linear kernel. We normalize the latter to prevent the linear kernel from dominating the radial basis function kernel, since the variables in $T \setminus P_i$ typically far outnumber those in \tilde{V}_i .

We now integrate all steps of TWRCI by formally proving that TWRCI is sound and complete:

Algorithm 6 CRCE Graph Discovery

Input: $R, N, Y, \mathcal{P}, \widehat{G}$ over \widetilde{R} , type I error rate α

Output: DAG \widehat{G} over $\widetilde{R} \cup Y$

- 1: Add vertex Y in \widehat{G}
 - 2: Draw a directed edge from each vertex in \widetilde{R} to Y in \widehat{G}
 - 3: **for each** $\widetilde{R}_i \in \widetilde{R}$ **do**
 - 4: Test whether R_i and Y are independent given $\text{Pa}_{\widehat{G}}(\widetilde{R}_i) \cup \widetilde{N} \cup \mathcal{P}_i$ to obtain p-value p
 - 5: **if** $p > \alpha$ **then**
 - 6: Delete the edge $\widetilde{R}_i \rightarrow Y$ from \widehat{G}
 - 7: **end if**
 - 8: **end for**
-

Theorem 1. (Fisher consistency) Under d -separation faithfulness, relevance and exchangeability, TWRCI identifies all of the direct causal variants of $Y \cup (\text{Anc}(Y) \cap \widetilde{X})$, the unique causal graph over $Y \cup (\text{Anc}(Y) \cap \widetilde{X})$ and the CRCEs of $\text{Anc}(Y) \cap \widetilde{X}$ almost surely as $N \rightarrow \infty$ with Lipschitz continuous conditional expectations and a conditional independence oracle.

We perform conditional independence testing by correlating the regression residuals of smooth non-linear transformations of the gene expression levels and phenotype⁵⁵. As a result, Lemma 1 also enables accurate conditional independence testing over subsets of $T \cup \widetilde{R} \cup \widetilde{N}$, even though we only have access to $T \cup R \cup N$.

4.4.7 Time Complexity

We analyze the time complexity of TWRCI in detail. TWRCI can admit different regression procedures, so we will assume that each regression takes $O(c^3)$ time, where c denotes the dimensionality of the conditioning set typically much larger than the sample size n . Most regression procedures satisfy the requirement.

TWRCI first runs Algorithm 2 which requires $O(q)$ time in Line 1 with summary statistics, $O(q^3 m)$ time in Line 2 with at most m regressions on T , and $O(m^3(m+q))$ time for at most $m+q$ regressions on \widetilde{R} in Line 3. Algorithm 2 thus takes $O(m^4 + m^3 q) + O(q^3 m)$ time in total.

TWRCI next annotates to Y using Algorithm 3 which takes $O(q^3 m) + O((m+q)^3)$ time for Lines 1 and 2, respectively. Annotation to Y therefore carries a total time complexity of $O(m^3 q^3)$. TWRCI then runs Algorithm 4. Each iteration of the repeat loop in Line 3 of Algorithm 4 takes $O(q^3)$ time for the regression in Line 4 and $O(m(m+q)^3)$ time for the at most m regressions in Line 7. The repeat loop iterates at most m times, so Algorithm 4 has a total time complexity of $O(m(q^3 + m(m+q)^3)) = O(m^5 q^3)$.

Algorithm 5 dominates Algorithm 6 in time during the causal graph discovery portion of TWRCI. Algorithm 5 runs in $O(m^e(m+q)^3) = O(m^{e+3} q^3)$ time, where e denotes the maximum neighborhood size¹⁹. Finally, CRCE estimation in Line 5 requires $O(2m(m+q)^3) = O(m^4 q^3)$ time for at most $2m$ regressions on expression levels and variants. Thus TWRCI in total requires $O(m^4 + m^3 q) + O(q^3 m) + O(m^3 q^3) + O(m^5 q^3) + O(m^{e+3} q^3) + O(m^4 q^3) = O(m^5 q^3) + O(m^{e+3} q^3)$ time. We conclude that the ACO and Graph Discovery sub-algorithms dominate the time complexity of TWRCI.

4.5 Comparators

We compared TWRCI against state of the art algorithms enumerated below.

Annotation:

1. Nearest TSS: annotates each variant to its closest gene according to the TSS.
2. Cis-window: annotates a variant to a gene if the variant lies within a one megabase window of the TSS. If a variant lies in multiple windows, then we assign the variant to the closest TSS.
3. Causal transcriptome-wide association study (cTWAS)¹¹: annotates variants to genes using cis-windows and then accounts for horizontal pleiotropy using the Sum of Single Effects (SuSIE) algorithm.

4. Cis-eQTLs¹²: annotates a variant to a gene if (1) the variant lies in the cis-window of the gene per above, and (2) the variant correlates most strongly with that gene expression level relative to the other levels.
5. Colocalization with approximate Bayes factors¹³: annotates each variant to the gene expression level with the highest colocalization probability according to approximate Bayes factors. We *could not* differentiate this method from cis-windows using the MACR criteria for the real data (Methods 4.8), since the algorithm always assigns higher approximate Bayes factors to cis-variants.
6. Colocalization with SuSIE^{13,14}: same as above but with probabilities determined according to SuSIE. We *could* differentiate this method from cis-windows using the MACR criteria for the real data.

Causal Graph Reconstruction:

1. SIGNET^{15,16}: predicts gene expression levels from variants using ridge regression and then recovers the genetic ancestors of each expression level by running the adaptive LASSO on the predicted expression levels. The method thus assumes linearity.
2. RCI¹⁷: assumes a linear non-Gaussian acyclic model⁵⁶, and recovers the causal order by maximizing independence between gene expression level residuals obtained from linear regression.
3. GRCI¹⁸: same as above but assumes an additive noise model⁵⁷ and uses non-linear regression.
4. PC/CausalCell²⁰: runs the stabilized PC algorithm^{19,53} on the gene expression levels using a non-parametric conditional independence test⁵⁵.

4.6 Semi-Synthetic Data

The causal graph reconstruction algorithms all require a variable selection step with gene expression data, since they cannot scale to the tens of thousands of genes with the neighborhood sizes seen in practice^{1,20}. We therefore assessed the performance of the algorithms independent of variable selection by first instantiating a DAG directly over $\tilde{\mathcal{Q}}$ with $p = 30$ variables including 29 gene expression levels and a single phenotype. We generated a linear SEM obeying Equation (3) such that $\tilde{Q}_i = \tilde{\mathbf{Q}}\beta_i + S_i\theta_i + E_i$ for every $\tilde{Q}_i \in \tilde{\mathcal{Q}}$ with $E_i \sim N(0, 1/25)$ to enable detection of weak causal effects from variants. We drew the coefficient matrix β from a Bernoulli($2/(p-1)$) in the upper triangular portion of the matrix and then randomly permuted the ordering of the variables. The resultant DAG has an expected neighborhood size of 2. We then weighted the coefficient matrix between the gene expression levels and phenotype by sampling uniformly from $[-1, -0.25] \cup [0.25, 1]$.

We instantiated the variants \mathbf{T} and θ as follows. We downloaded summary statistics from a wide variety of IEU datasets listed in Table 1 and filtered variants at a liberal α threshold of $5e-5$. We selected a variant to be closest to the TSS of each gene uniformly at random and assigned direct causal variants to the 29 gene expression levels with probability proportional to the inverse of the absolute distance from the closest variant plus one. As a result, variants closer to the TSS are more likely to have a direct causal effect on the gene expression level. We assigned the remaining variants to the phenotype. We sampled \mathbf{T} by bootstrap from the GTEx version 8²² individual-level genotype data and the weights θ uniformly from $[-0.15, -0.05] \cup [0.05, 0.15]$ because variants usually have weak causal effects.

We converted the above linear SEM to a non-linear one by setting $\tilde{Q}_i \leftarrow \text{softplus}(\tilde{Q}_i)$ for each $\tilde{Q}_i \in \tilde{\mathcal{Q}}$. We obtained each measurement error corrupted surrogate R_i by sampling from $\text{Pois}(\tilde{R}_i\pi_{i1})$ for each $\tilde{R}_i \in \tilde{\mathbf{R}}$. We drew the mapping efficiencies $\pi_{\cdot 1}$ for a single batch from the uniform distribution between 100 and 10000 for the bulk RNA sequencing data. We repeated the entirety of the above procedure 100 times to generate 100 independent variant-expression-phenotype datasets. We ran TWRCI and all combinations of the comparator algorithms on each dataset.

4.7 Real Data

4.7.1 Data Availability

All real datasets analyzed in this study have been previously published and are publicly accessible. The COPD datasets include:

Dataset	Trait	# Variants	# Cases	# Controls
ieu-b-5067	Alzheimer's	669	954	487331
ieu-b-4967	Appendicitis	736	4604	481880
ieu-b-4972	Endocarditis	300	1080	485404
ieu-b-4971	Cholecystitis	691	4052	482432
ieu-b-4973	Lower respiratory tract infection	1116	14135	472349
ieu-a-1187	Major depression	566	135458	344901
ieu-b-4965	Colorectal cancer	1791	5657	372016
ieu-b-5063	Upper respiratory tract infection	451	2795	483689
ieu-b-4956	Lymphoid leukemia	988	760	372016
ieu-b-4953	Liver cell carcinoma	517	168	372016

Table 1. Variant data used during semi-synthetic data generation.

1. Summary statistics: [ebi-a-GCST90018807](#)
2. Individual level variant and phenotype data: [GTEx V8 Protected Access Data](#)
3. Gene expression data: [GTEx V8 Lung](#)
4. Replication summary statistics: [ebi-a-GCST90018587](#)

The IHD datasets include:

1. Summary statistics: [finn-b-I9_ISCHHEART](#)
2. Individual level variant and phenotype data: [GTEx V8 Protected Access Data](#)
3. Gene expression data: [GTEx V8 Whole Blood](#)
4. Replication summary statistics: [ukb-d-I9_IHD](#)

4.7.2 Quality Control

We selected variants T at an α threshold of $5e-5$ for both the COPD and IHD summary statistics. We harmonized the variant data of the IEU and GTEx datasets by lifting the GTEx variant data from the hg38 to hg19 build using the liftover command in BCFtools version 1.18⁵⁸. We ensured that the reference and alternative alleles matched in both datasets after lifting for every variant. We removed gene expression levels with a mean count of less than five. We subjected the gene expression data to an inverse hyperbolic sine transformation to mitigate the effects of outliers. We regressed out the first 5 principal components, sequencing platform (Illumina HiSeq 2000 or HiSeq X), sequencing protocol (PCR-based or PCR-free) and sex from all variables in the linked GTEx variant-expression-phenotype data. Then, we either included age as a covariate for algorithms that accept a nuisance covariate, or regressed out age from the expression and phenotype data for algorithms that do not accept a nuisance covariate.

4.7.3 Comparison to trans-eQTLs

TWRCI annotated many trans-variants in both of the real datasets. Other authors have proposed *trans-eQTLs* as variants that lie distal to the TSS and correlate with at least one reported phenotype in the Catalog of Published GWAS⁵⁹. TWRCI annotates variants based on direct causality rather than correlation and an overlap with another phenotype. However, we hypothesized that the variants discovered by TWRCI should still lie close to at least a subset of the trans-eQTLs. To test this hypothesis, we downloaded trans-eQTL results from the [eQTLGen database](#)²⁹. We then standardized the positions of the variants within each chromosome by their standard deviation to account for variable chromosome length and polymorphism density. Next, we computed the nearest neighbor distances between the variants annotated to causal genes by TWRCI and the trans-eQTLs. We used the median of these normalized distances M as a robust statistic of central tendency.

We used a permutation test to test the null hypothesis that the variants annotated to causal genes by TWRCI are distributed arbitrarily far from the trans-eQTLs. We recomputed the median statistic 10,000 times after permuting the positions of the trans-eQTL variants. The p-value corresponds to the proportion of permuted statistics smaller than M . We reject the null hypothesis – and thus conclude that the variants annotated to causal genes by TWRCI lie close to trans-eQTLs – when the p-value falls below 0.05.

4.8 Metrics

We evaluated the accuracy of the algorithms using the nine metrics listed below for the synthetic data. We evaluated annotation quality using the following two metrics:

1. Matthew's Correlation Coefficient (MCC)⁶⁰ between the estimated annotations and the ground truth direct causal variants. Larger is better.
2. Rank of the estimated coefficients $\hat{\theta}$ normalized by the rank of the ground truth coefficients θ . Larger is better.

We also computed the above two quantities only using the variants that directly cause the phenotype in order to evaluate the ability of the algorithms to account for horizontal pleiotropy (3. and 4.). We evaluated the causal graph reconstruction quality using the following two metrics:

5. Structural Hamming Distance (SHD)⁶¹ between the estimated and the ground truth causal graph. Smaller is better.
6. MCC between the estimated and the ground truth causal graph. Larger is better.

We evaluated combined annotation and graph reconstruction quality using Lemma 4:

7. Mean absolute correlation of the residuals (MACR) defined as the mean absolute correlation between (a) the variants T and ancestral gene expression levels, and (b) the gene expression residuals after partialing out the inferred parents. Smaller is better under the global Markov property and exchangeability. If the algorithm infers no direct causal variants in T and no parents in \hat{G} for some \tilde{R}_i , then this situation violates the relevance assumption, where at least one variant in T directly causes \tilde{R}_i . We thus set the absolute correlation of \tilde{R}_i to one in this case.

We assessed the accuracy in CRCE estimation using the following metrics:

8. Root mean squared error between the estimated CRCE and the ground truth CRCE averaged over all gene expression levels. We do not have access to the ground truth CRCE, so we estimate it to negligible error with kernel ridge regression using the ground truth parents. Smaller is better.
9. MACR between (a) the residuals $Y - \hat{E}(Y|\tilde{R}_i, D)$ and (b) the inferred set P_i , which should be zero under the global Markov property and exchangeability. Smaller is better. We again set the absolute correlation to one for \tilde{R}_i if the algorithm infers no direct causal variants and no parents in \hat{G} under relevance.

We can compute the MACR metrics 7. and 9. on real data, so we evaluate the algorithms using these two metrics in the IHD and COPD datasets. We also have access to silver standard sets of genes known to be causally involved in disease from either the DisGeNet²⁶ or KEGG database³⁹. We therefore compute a third MACR metric with the real data:

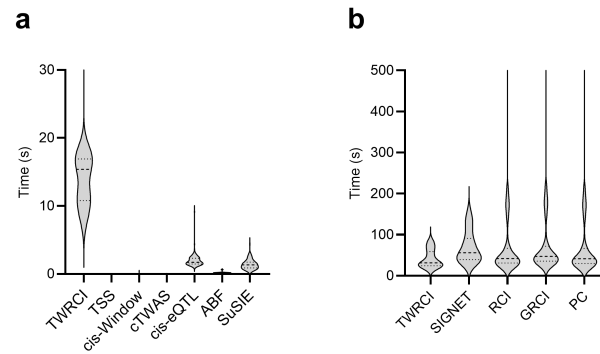
10. A causal gene should at least correlate with the phenotype, so we first correlate the silver standard genes with the phenotype and only keep silver standard genes with a significant correlation ($p < 0.05$ uncorrected). We then compute a MACR metric between (a) the kept silver standard genes after partialing out genes with non-zero CRCEs and (b) the phenotype after partialing out genes with non-zero CRCEs.

4.9 Code Availability

R code needed to replicate all experimental results is available on [GitHub](https://github.com).

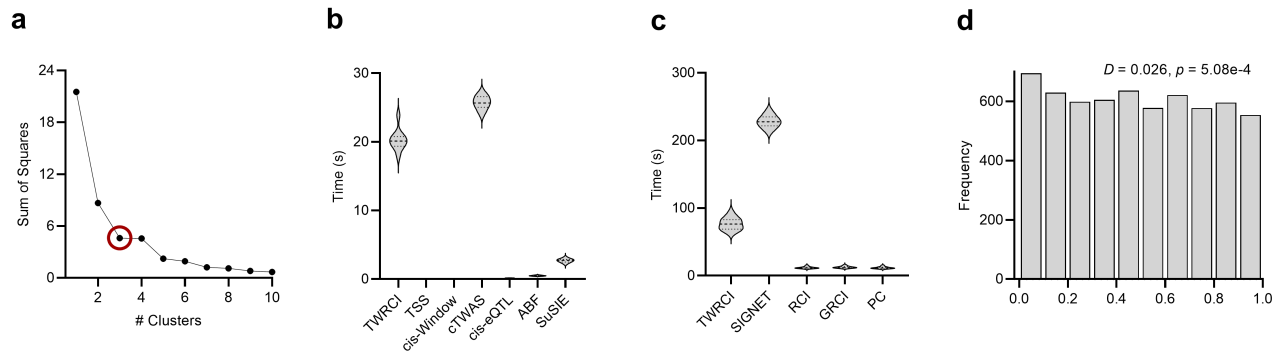
5 Supplementary Materials

5.1 Additional Semi-Synthetic Data Results

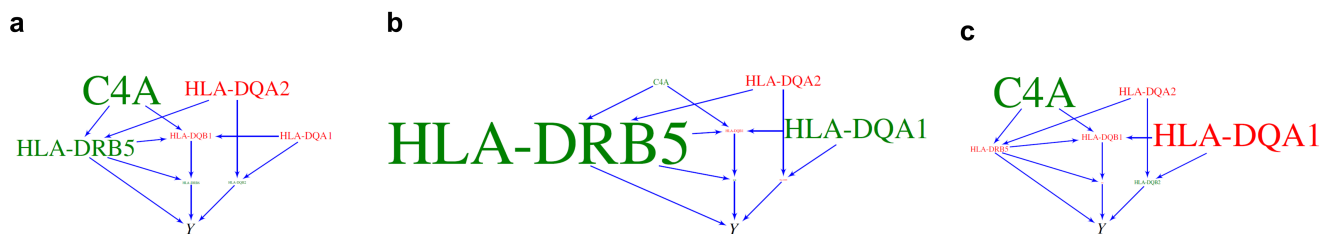


Supplementary Figure 1. Timing results for the semi-synthetic datasets split into the variant annotation and graph reconstruction portions because they took the longest by far. (a) TWRCI took the longest time during annotation, but (b) all algorithms spent the majority of the time in causal graph reconstruction over \tilde{R} in congruence with the time complexity results of Methods 4.4.7. TWRCI completed within about 3 minutes overall.

5.2 Additional COPD Data Results

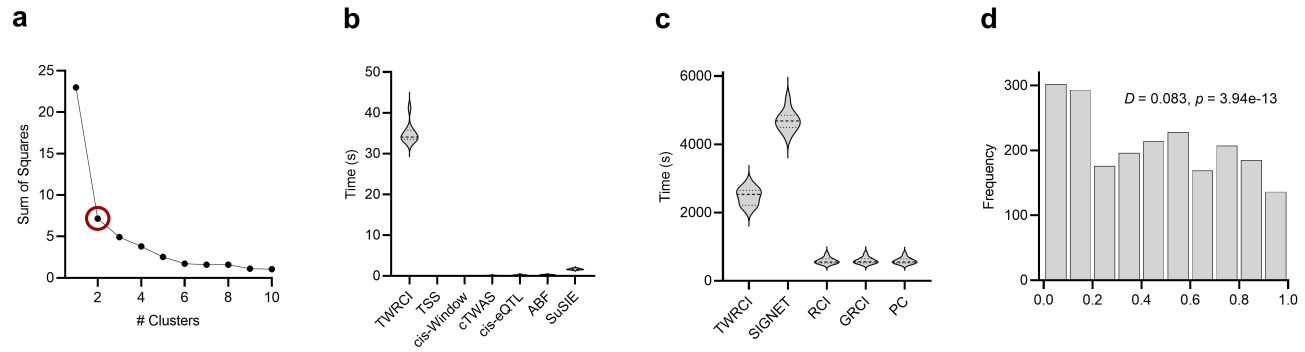


Supplementary Figure 2. Additional results for COPD. (a) Sum of squares plot for hierarchical clustering using Ward's method revealed three clusters according to the elbow method, or the cluster size with the maximum distance from the imaginary line drawn between the first and last cluster sizes. TWRCI took the second longest time to complete in annotation (b) and the second longest time to complete in graph reconstruction (c). RCI, GRCI and PC all took a much smaller amount of time to reconstruct the causal graph because they ignore the genetic variants. (d) Histogram of Pearson correlation test p-values computed between variants annotated to the phenotype and gene expression levels. The p-values did not follow a uniform distribution according to the Kolomogorov-Smirnov test with statistic D indicating the presence of confounding between the variants annotated to the phenotype and gene expression.

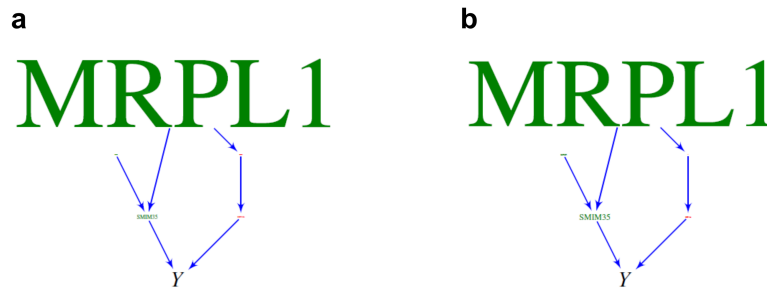


Supplementary Figure 3. Replication results in an independent set of individuals of East Asian ancestry (dataset ebi-a-GCST90018587). We summarize results for (a) all patients, (b) cluster one in Figure 4 (g) and (c) cluster two. TWRCI again identified C4A and multiple MHC class II genes involving the adaptive immune system.

5.3 Additional IHD Data Results



Supplementary Figure 4. Additional results for IHD. (a) Sum of squares plot revealed two clusters according to the elbow method. (b) TWRCI took the longest to annotate, but (c) the timing results for graph reconstruction dominated in this case. Methods using SIGNET thus took the longest overall in this dataset. (d) Histogram of Pearson correlation test p-values were again non-uniform, indicating confounding between the variants annotated to the phenotype and gene expression.



Supplementary Figure 5. Replication results in an independent set of patients from the UK Biobank (dataset ukb-d-I9_IHD). We summarize results for (a) all patients, and (b) cluster one. TWRCI again identified MRPL1 as a root causal gene with a large positive CRCE.

5.4 Proofs

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

$$\mathbb{E} \left| \mathbb{E}(Z_i | \tilde{\mathbf{U}}, \mathbf{V}) - \mathbb{E}(Z_i | \mathbf{U}, \mathbf{V}, B) \right| \leq \mathbb{E} C_N \left| \tilde{\mathbf{U}} - \frac{\mathbf{U}}{N} \frac{\tilde{N}_B}{\pi_{\mathbf{U}B}} \right|,$$

where $C_N \in O(1)$ is a positive constant, and we have taken an outer expectation on both sides. Then $\mathbb{E}(Z_i | \tilde{\mathbf{U}}, \mathbf{V}) = \lim_{N \rightarrow \infty} \mathbb{E}(Z_i | \mathbf{U}, \mathbf{V}, B)$ almost surely.

Proof. We can write the following sequence:

$$\begin{aligned} & \mathbb{E} \left| \mathbb{E}(Z_i | \tilde{\mathbf{U}}, \mathbf{V}) - \lim_{N \rightarrow \infty} \mathbb{E}(Z_i | \mathbf{U}, \mathbf{V}, B) \right| = \mathbb{E} \lim_{N \rightarrow \infty} \left| \mathbb{E}(Z_i | \tilde{\mathbf{U}}, \mathbf{V}) - \mathbb{E}(Z_i | \mathbf{U}, \mathbf{V}, B) \right| \\ & \leq \mathbb{E} \lim_{N \rightarrow \infty} C_N \left| \tilde{\mathbf{U}} - \frac{\mathbf{U}}{N} \frac{\tilde{N}_B}{\pi_{\mathbf{U}B}} \right| \leq \mathbb{E} C \left| \tilde{\mathbf{U}} - \lim_{N \rightarrow \infty} \frac{\mathbf{U}}{N} \frac{\tilde{N}_B}{\pi_{\mathbf{U}B}} \right| = \mathbb{E} C \left| \tilde{\mathbf{U}} - \frac{\tilde{\mathbf{U}} \pi_{\mathbf{U}B}}{\tilde{N}_B} \frac{\tilde{N}_B}{\pi_{\mathbf{U}B}} \right| = C \mathbb{E} |\tilde{\mathbf{U}} - \tilde{\mathbf{U}}| = 0, \end{aligned}$$

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

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

The conclusion follows because we chose a arbitrarily. □

Proposition 1. We have $\mathbb{P}(Y | E_i \cup S_i, \mathbf{D}) - \mathbb{P}(Y | \mathbf{D}) = \mathbb{P}(Y | \tilde{X}_i, \mathbf{D}) - \mathbb{P}(Y | \mathbf{D})$ under exchangeability.

Proof. We can write:

$$\begin{aligned} \mathbb{P}(Y | E_i \cup S_i, \mathbf{D}) &= \mathbb{P}(Y | E_i, \tilde{\mathbf{V}}_i, \mathbf{T} \cup S_i) = \mathbb{E}_{\tilde{X}_i | E_i, \tilde{\mathbf{V}}_i, \mathbf{T}, S_i} \mathbb{P}(Y | \tilde{X}_i, E_i, \tilde{\mathbf{V}}_i, \mathbf{T} \cup S_i) \\ &= \mathbb{P}(Y | \tilde{X}_i, E_i \cup S_i, \tilde{\mathbf{V}}_i, \mathbf{T} \setminus S_i) = \mathbb{P}(Y | \tilde{X}_i, \tilde{\mathbf{V}}_i, \mathbf{T} \setminus S_i) = \mathbb{P}(Y | \tilde{X}_i, \mathbf{D}). \end{aligned}$$

The third equality follows because \tilde{X}_i is a constant given E_i and $\text{Pa}(\tilde{X}_i) = \tilde{\mathbf{V}}_i \cup S_i$. For the fourth equality, all paths between S_i and Y are blocked by $\tilde{X}_i \cup \tilde{\mathbf{V}}_i \cup \mathbf{T} \setminus S_i$ under exchangeability. We thus have $Y \perp\!\!\!\perp_d (E_i \cup S_i) | (\tilde{X}_i, \tilde{\mathbf{V}}_i, \mathbf{T} \setminus S_i)$ and $Y \perp\!\!\!\perp (E_i \cup S_i) | (\tilde{X}_i, \tilde{\mathbf{V}}_i, \mathbf{T} \setminus S_i)$ by the global Markov property. □

Lemma 2. Assume d -separation faithfulness and relevance. Then, (1) $\mathbf{T} \cup \tilde{\mathbf{R}}$ contains all of the ancestors of Y in $S \cup \tilde{\mathbf{X}}$, and (2) $\text{Mb}(\tilde{R}_i) \cap \tilde{\mathbf{X}} \subseteq (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}}$ for any $\tilde{R}_i \in \tilde{\mathbf{R}}$.

Proof. We first prove (1). If S_i is an ancestor of Y , then $S_i \not\perp\!\!\!\perp_d Y$, so $S_i \not\perp\!\!\!\perp Y$ by d -separation faithfulness. It follows that $S_i \in \mathbf{T}$ by Line 1 of Algorithm 2. If \tilde{X}_i is an ancestor of Y , then so is $S_i \subseteq \mathbf{T}$. Hence $\tilde{X}_i \not\perp\!\!\!\perp_d \mathbf{T}$, so $\tilde{X}_i \not\perp\!\!\!\perp \mathbf{T}$ and $\tilde{X}_i \in \tilde{\mathbf{R}}$ by d -separation faithfulness and Line 2, respectively. We chose S_i and \tilde{X}_i arbitrarily, so the set $\mathbf{T} \cup \tilde{\mathbf{R}}$ contains all of the ancestors of Y in $S \cup \tilde{\mathbf{X}}$.

We now prove (2). We need to show that $(\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}}$ contains the parents, children and spouses of any \tilde{R}_i , provided that these relatives are also in $\tilde{\mathbf{X}}$. Note that $\tilde{R}_i \not\perp\!\!\!\perp_d \mathbf{T}$ by Line 2 under the global Markov property. Hence, the parents and children of \tilde{R}_i in $\tilde{\mathbf{X}}$ are also d -connected to \mathbf{T} and hence dependent on \mathbf{T} under d -separation faithfulness. It follows that $\tilde{\mathbf{R}} \setminus \tilde{R}_i$ contains all of the parents and children of \tilde{R}_i also by Line 2. Next, suppose $\tilde{\mathbf{R}} \setminus \tilde{R}_i$ does not contain a spouse of \tilde{R}_i , which we denote by \tilde{X}_j . Then we have $\tilde{X}_j \rightarrow \tilde{R}_i \leftarrow S_i$ and $S_i \in \mathbf{T}$ under relevance. Hence $X_j \not\perp\!\!\!\perp_d \mathbf{T} | \tilde{\mathbf{R}}$, so $X_j \not\perp\!\!\!\perp \mathbf{T} | \tilde{\mathbf{R}}$ by d -separation faithfulness and $X_i \in \mathbf{N}$ by Line 3. It follows that $\tilde{\mathbf{N}} \cup (\tilde{\mathbf{R}} \setminus \tilde{R}_i)$ contains all of the spouses of \tilde{R}_i . We conclude that $(\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}}$ contains all members of $\text{Mb}(\tilde{R}_i) \cap \tilde{\mathbf{X}}$ of any $\tilde{R}_i \in \tilde{\mathbf{R}}$. □

Lemma 8. Under d -separation faithfulness, relevance and exchangeability, (1) $T_j \notin \text{Anc}(Q_i)$ if and only if $Q_i \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$ and (2) $Q_i \perp\!\!\!\perp T_j | (\tilde{\mathcal{Q}} \setminus \tilde{Q}_i, \tilde{\mathbf{N}}, \mathbf{T} \setminus T_j)$ and $Q_i \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$ if and only if $T_j \in \text{Pa}(\tilde{Q}_i)$.

Proof. For the first statement and forward direction, if $T_j \notin \text{Anc}(Q_i)$, then $Q_i \perp\!\!\!\perp_d T_j | \mathbf{T} \setminus T_j$ under exchangeability, so $Q_i \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$ by the global Markov property. For the backward direction, if $Q_i \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$, then $Q_i \perp\!\!\!\perp_d T_j | \mathbf{T} \setminus T_j$ by d -separation faithfulness. No directed path can thus exist from T_j to Q_i , so $T_j \notin \text{Anc}(Q_i)$.

We next address the second statement. The backward direction follows immediately from d -separation faithfulness. For the forward direction, if $Q_i \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$, then $T_j \in \text{Anc}(Q_i)$ from statement (1). Furthermore, if $Q_i \perp\!\!\!\perp T_j | (\tilde{\mathcal{Q}} \setminus \tilde{Q}_i, \tilde{\mathbf{N}}, \mathbf{T} \setminus T_j)$ then $Q_i \perp\!\!\!\perp_d T_j | (\tilde{\mathcal{Q}} \setminus \tilde{Q}_i, \tilde{\mathbf{N}}, \mathbf{T} \setminus T_j)$ under the global Markov property. Note that $(\tilde{\mathcal{Q}} \setminus \tilde{Q}_i) \cup \tilde{\mathbf{N}} \cup \mathbf{T}$ contains $\text{Mb}(\tilde{Q}_i) \cap \tilde{\mathbf{X}}$ by Lemma 2 under d -separation faithfulness and relevance. Therefore, if T_j is not in the Markov boundary of \tilde{Q}_i , then $(\tilde{\mathcal{Q}} \setminus \tilde{Q}_i) \cup \tilde{\mathbf{N}} \cup \mathbf{T} \setminus T_j$ contains $\text{Mb}(\tilde{Q}_i) \cap (\tilde{\mathbf{X}} \cup \mathbf{T})$. As a result, all paths between T_j and \tilde{Q}_i are blocked by $(\tilde{\mathcal{Q}} \setminus \tilde{Q}_i) \cup \tilde{\mathbf{N}} \cup \mathbf{T} \setminus T_j$ under exchangeability. We thus arrive at the contradiction $Q_i \perp\!\!\!\perp_d T_j | (\tilde{\mathcal{Q}} \setminus \tilde{Q}_i, \tilde{\mathbf{N}}, \mathbf{T} \setminus T_j)$. It follows that T_j must be in the Markov boundary of \tilde{Q}_i and therefore can only be a parent or a spouse of \tilde{Q}_i (or both). If T_j is a spouse but not a parent of \tilde{Q}_i , then we arrive at another contradiction that $T_j \notin \text{Anc}(Q_i)$. Hence $T_j \in \text{Pa}(\tilde{Q}_i)$. \square

Lemma 3. Assume d -separation faithfulness, relevance and exchangeability. Further assume that \tilde{Q}_i is a sink vertex. Then, $|\Delta_{ij}\gamma_{ij}| \geq \max \Delta_{-ij}^2$ if and only if $T_j \notin \text{Anc}(\mathcal{Q} \setminus Q_i)$ or $T_j \in \text{Pa}(\tilde{Q}_i)$ (or both).

Proof. Assume $|\Delta_{ij}\gamma_{ij}| \geq \max \Delta_{-ij}^2$ for the forward direction. We have two cases. If $|\Delta_{ij}\gamma_{ij}| > 0$, then $Q_i \perp\!\!\!\perp T_j | (\tilde{\mathcal{Q}} \setminus \tilde{Q}_i, \tilde{\mathbf{N}}, \mathbf{T} \setminus T_j)$ and $Q_i \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$, so $T_j \in \text{Pa}(\tilde{Q}_i)$ by Lemma 8. If $|\Delta_{ij}\gamma_{ij}| = 0$, then $\max \Delta_{-ij}^2 = 0$, so $Q_k \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$ for all $Q_k \in \mathcal{Q} \setminus Q_i$. We conclude that $T_j \notin \text{Anc}(\mathcal{Q} \setminus Q_i)$ by again invoking Lemma 8.

For the backward direction, if $T_j \notin \text{Anc}(\mathcal{Q} \setminus Q_i)$, then $Q_k \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$ for all $Q_k \in \mathcal{Q} \setminus Q_i$ by Lemma 8. Thus $\max \Delta_{-ij}^2 = 0$ so $|\Delta_{ij}\gamma_{ij}| \geq \max \Delta_{-ij}^2$. If $T_j \in \text{Pa}(\tilde{Q}_i)$, then $T_j \notin \text{Anc}(\mathcal{Q} \setminus Q_i)$ because \tilde{Q}_i is a sink vertex. Hence $Q_k \perp\!\!\!\perp T_j | \mathbf{T} \setminus T_j$ for all $Q_k \in \mathcal{Q} \setminus Q_i$ by Lemma 8, so $\max \Delta_{-ij}^2 = 0$. We conclude that $|\Delta_{ij}\gamma_{ij}| \geq \max \Delta_{-ij}^2$. \square

Lemma 4. \tilde{R}_i is a sink vertex if and only if $R_i \perp\!\!\!\perp (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$ in Line 12 of ACO under d -separation faithfulness, relevance and exchangeability.

Proof. Assume that \tilde{R}_i is a sink vertex for the forward direction. We have two cases:

1. If $\tilde{R}_i \in \text{Anc}(Y)$, then $\text{Pa}(\tilde{R}_i) \subseteq (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup U_i$ by the first statement of Lemma 2 and Lemma 3. Note that $\tilde{\mathbf{N}}$ Hence, $R_i \perp\!\!\!\perp_d (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$ because \tilde{R}_i is a sink vertex, and $R_i \perp\!\!\!\perp (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$ in Line 12 by the global Markov property.
2. If $\tilde{R}_i \notin \text{Anc}(Y)$, then $(\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$ contains all of the parents of \tilde{R}_i in $\tilde{\mathbf{X}}$ and \mathbf{T} by Lemma 2 and Lemma 3, respectively. Moreover, the other direct causal variants of \tilde{R}_i , or $S_i \setminus \mathbf{T}$, share no latent confounders with \mathbf{T} or any other direct causal variant set excluding \mathbf{T} by exchangeability. Hence, we also have $R_i \perp\!\!\!\perp_d (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$, and $R_i \perp\!\!\!\perp (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$ in Line 12 by the global Markov property.

We have exhausted all possibilities and thus conclude that $R_i \perp\!\!\!\perp (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$.

For the backward direction, assume $R_i \perp\!\!\!\perp (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$ so that $R_i \perp\!\!\!\perp_d (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$ by d -separation faithfulness. Assume for a contradiction that \tilde{R}_i is not a sink vertex. Then there exists a path $\tilde{R}_i \rightarrow \tilde{R}_j \leftarrow T_k$ for some $T_k \in S_j \cap \mathbf{T}$ by relevance. We thus have $\Delta_{ik} = 0$ but $\Delta_{jk}^2 > 0$, so $T_k \notin U_i$ and $T_k \in \mathbf{T} \setminus U_i$. We arrive at the contradiction $R_i \perp\!\!\!\perp_d (\mathbf{T} \setminus U_i) | (\tilde{\mathbf{R}} \setminus \tilde{R}_i) \cup \tilde{\mathbf{N}} \cup U_i$. The variable \tilde{R}_i must therefore be a sink vertex. \square

Lemma 5. Under d -separation faithfulness, relevance and exchangeability, ACO recovers the correct causal order \mathbf{K} over $\tilde{\mathbf{R}}$ and $(S_i \cap \mathbf{T}) \subseteq P_i$ for all $R_i \in \tilde{\mathbf{R}}$.

Proof. We use proof by induction. Base: Suppose $\tilde{\mathbf{R}}$ contains one variable \tilde{R}_i . Then $\mathbf{K} = (R_i, Y)$ because R_i is trivially the most independent variable in \mathbf{R} according to \mathbf{C} of Line 15. The variable R_i is a sink vertex after Y is eliminated, so we have $(S_i \cap T) \subseteq P_i$ under d-separation faithfulness, relevance and exchangeability by Lemma 3. Step: Assume that the conclusion holds when $\tilde{\mathbf{R}}$ contains $p - 1$ variables. We need to prove the statement when $\tilde{\mathbf{R}}$ contains p variables. Assume for now that \tilde{R}_p is an arbitrary sink vertex in $\tilde{\mathbf{R}}$. Lemma 3 then guarantees $|\Delta_{p,j}\gamma_{p,j}| \geq \max \Delta_{-p,j}^2$ for each $S_j \in S_p \cap T$ in Line 8 under d-separation faithfulness, relevance and exchangeability. We thus have $(S_p \cap T) \subseteq P_p$ and no variant of any other parent set is in P_p . Finally, the measure of dependence C_p in Line 15 identifies R_p as a sink vertex by Lemma 4. ACO thus eliminates R_p from \mathbf{R} and appends it to the front of \mathbf{K} . The conclusion follows by the inductive hypothesis. \square

Lemma 6. *Under d-separation faithfulness, relevance and exchangeability, the graph discovery algorithm outputs the true sub-DAG over $\tilde{\mathbf{R}}$ given a conditional independence oracle, \mathbf{K} and \mathcal{P} .*

Proof. The set $\tilde{\mathbf{N}} \cup \tilde{\mathbf{R}}$ contains all of the parents of any $\tilde{R}_i \in \tilde{\mathbf{R}}$ in $\tilde{\mathbf{X}}$ by Lemma 2. Furthermore, P_i contains all of the parents of \tilde{R}_i in T for any $\tilde{R}_i \in \tilde{\mathbf{R}}$ by Lemma 4. The stabilized skeleton discovery procedure of the PC algorithm thus recovers all and only the undirected edges in the true DAG over $\tilde{\mathbf{R}}$ under d-separation faithfulness and exchangeability⁵³. The conclusion follows because ACO recovers the true causal order over $\tilde{\mathbf{R}}$ also by Lemma 5, so Algorithm 5 infers the true sub-DAG uniquely over $\tilde{\mathbf{R}}$ in Line 19. \square

Lemma 7. *Under d-separation faithfulness, relevance and exchangeability, \tilde{R}_i causes Y – and likewise the vertices $S_i \cup E_i$ cause Y – if and only if $Y \not\perp\!\!\!\perp \tilde{R}_i | \tilde{\mathbf{V}}_i \cup P_i$.*

Proof. Recall that $(S_i \cap T) \subseteq P_i$ by Lemma 5 under d-separation faithfulness, relevance and exchangeability.

Now if \tilde{R}_i causes Y , then there exists a directed path from \tilde{R}_i to Y so $Y \not\perp\!\!\!\perp_d \tilde{R}_i | \tilde{\mathbf{V}}_i \cup P_i$. We then have $Y \not\perp\!\!\!\perp \tilde{R}_i | \tilde{\mathbf{V}}_i \cup P_i$ by d-separation faithfulness.

For the backward direction, assume that \tilde{R}_i does not cause Y . All paths between \tilde{R}_i and Y are blocked by $\tilde{\mathbf{V}}_i \cup P_i$ under exchangeability. Thus \tilde{R}_i and Y are d-separated given $\tilde{\mathbf{V}}_i \cup P_i$. We invoke the global Markov property to conclude that $Y \perp\!\!\!\perp \tilde{R}_i | \tilde{\mathbf{V}}_i \cup P_i$. \square

Theorem 1. (Fisher consistency) *Under d-separation faithfulness, relevance and exchangeability, TWRCI identifies all of the direct causal variants of $Y \cup (\text{Anc}(Y) \cap \tilde{\mathbf{X}})$, the unique causal graph over $Y \cup (\text{Anc}(Y) \cap \tilde{\mathbf{X}})$ and the CRCEs of $\text{Anc}(Y) \cap \tilde{\mathbf{X}}$ almost surely as $N \rightarrow \infty$ with Lipschitz continuous conditional expectations and a conditional independence oracle.*

Proof. Lemma 2 ensures that $T \cup \tilde{\mathbf{R}}$ from Line 1 of Algorithm 1 contains all of the ancestors of Y in $S \cup \tilde{\mathbf{X}}$. Thus $(\text{Anc}(Y) \cap \tilde{\mathbf{X}}) \subseteq \tilde{\mathbf{R}}$ and $(\text{Anc}(Y) \cap S) \subseteq T$. TWRCI identifies $S_Y \subseteq P_Y$ in Line 2 by Lemma 3 under d-separation faithfulness, relevance and exchangeability. The algorithm also identifies $S_i \subseteq P_i$ for each $\tilde{R}_i \in \text{Anc}(Y)$ under d-separation faithfulness, relevance and exchangeability in Line 3 by invoking Lemma 5. Furthermore, TWRCI recovers the causal order over $\tilde{\mathbf{R}}$ via Lemma 5. TWRCI thus uniquely recovers the sub-DAG over $\tilde{\mathbf{R}}$ in Line 4 by Lemma 6 and then correctly includes Y in the graph by Lemma 7. TWRCI finally identifies the CRCEs of $\text{Anc}(Y) \cap \tilde{\mathbf{X}}$ almost surely in Line 5 given the recovered DAG over $\tilde{\mathbf{R}} \cup Y$ and \mathcal{P} by Lemma 1. \square