# A statistical physics approach for disease module detection

Supplemental Material

Xu-Wen Wang[1], Dandi Qiao[1], Michael H. Cho[1], Dawn L. DeMeo[1], Edwin K. Silverman[1], Yang-Yu Liu[1]

[1]*Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts 02115, USA*

# 1. Existing disease module detection method used for comparison

Numerous computational methods have been developed to identify disease modules (de Weerd et al. 2020; Ulitsky et al. 2010; Bersanelli et al. 2016) by projecting the molecular profiles into the network and assigning scores/weights to nodes/edges. Those methods can be briefly classified into three categories: (1) Permutation-based methods (e.g., Hierarchical HotNet (Reyna et al. 2018) and DOMINO (Levi et al. 2021)), which involve a network randomization process to correct for typical PPI network biases. (2) Seed-genes-based methods (e.g., DIAMOnD (Ghiassian et al. 2015)), which "grow" the disease module from the seed genes in a local agglomerative process. (3) Heuristic-based methods, e.g., ModuleDiscoverer (Vlaic et al. 2018), ROBUST (Bernett et al. 2022), attempt to find optimal subgraphs within the molecular interaction network. We compared our RFIM method with the following five existing methods.

## 1.1 DIAMOnD.

DIseAse Module Detection (DIAMOnD) algorithm (Ghiassian et al. 2015) starts from a set of seed genes and prioritizes the other proteins of the interactome for their putative disease relevance. Therefore, the identified disease module might significantly rely on the initial seed genes. In this work, we used the implementation of DIAMOnD incorporated in R (Team 2013) package MODifieR (de Weerd et al. 2020) which can simultaneously use the PPI network and gene-wise $p$-values. The number of genes in the disease module in DIAMOnD set to be 5% of the total genes, just the same as the RFIM. We used the default values for other parameters: seed weight (10); The $p$-value cutoff for differentially expressed genes is 0.05 and these significantly expressed genes were selected as seed genes and the seed genes were excluded from the disease module in MODifieR (de Weerd et al. 2020).

## 1.2 ModuleDiscoverer.

ModuleDiscoverer (Vlaic et al. 2018) uses a randomization heuristic-based approximation of the community structure based on maximal clique enumeration problem. ModuleDiscoverer first approximates the underlying community structure by iterative enumeration of gene cliques from random seed genes. Then, the union of all significantly enriched cliques are ensembled into a large module (Vlaic et al. 2018). Number of permutations is set to be 1,000. Number of times the algorithm is repeated is set to be 3.

## 1.3 DOMINO.

DIMINO (Levi et al. 2021) (Discovery of active Modules In Networks using Omics) finds disjoint connected subnetworks that the active genes are over-represented based on a permutation-based method that empirically evaluate GO terms. The main steps of DOMINO are: (1) Partition the network into disjoint and highly connected slices. (2) Detect relevant slices including over-represented active genes. (3) For each slice, refine it to sub-slice and repartition slice into putative modules. (4) Final modules are those with over-represented by active genes. The active genes were selected as genes with $p$-value lower than $10^{-4}$.

## 1.4 ROBUST.

ROBUST (Bernett et al. 2022) is a disease module mining method via enumeration of diverse prize-collecting Steiner trees. In naïve approach, disease module is identified by running mining method many times on shuffled input and return the subgraph induced by nodes contained in many of the return modules. Robust overcome the runtime limitation of naïve approach by enumerating pairwise diverse rather than merely pairwise non-identical disease modules (Bernett et al. 2022). We used the default parameters for ROBUST: initial fraction 0.25; reduction factor 0.9; number of Steiner trees to be computed 30 and threshold 0.1.

## 1.5 Hierarchical HotNet.

Hierarchical HotNet (Reyna et al. 2018) is an algorithm that finds a hierarchy of altered subnetworks and identify statistically significant subnetworks in the hierarchy. Hierarchical HotNet (1) combines network topology and node scores, (2) defines a similarity matrix from network using a random walk, (3) construct a hierarchy of clusters consisting of strongly connected components and (4) assesses the statistical significance of clusters in the hierarchy. As Hierarchical HotNet requires much longer computation time, we reduced the number of permutations into 20.

# 2. Existing optimal subgraphs identification methods.

This type of optimization procedure, i.e., finding optimal subgraphs within the molecular interaction network, turns out to be computationally intractable in general due to their NP-hard nature. For example, MWCS and PCST are two classical formalisms of this optimization procedure.

## 2.1 Maximum-Weight Connected Subgraph (MWCS).

In cases where we have no prior-knowledge of the edge weights that capture the confidence level of the protein-protein interactions, we focus on the node weights (i.e., the gene-wise $p$-values). The optimization problem is referred to as the MWCS problem, where given a graph $G(V, E)$ with node weights $z_i$, we try to find the connected subgraph with maximum weight:

$$\max_{A \subseteq G} S_A, \text{ with } S_A \equiv \frac{Z_A - \mu_k}{\sigma_k} \text{ and } Z_A \equiv \frac{1}{\sqrt{k}} \sum_{i \in A} z_i \qquad [1]$$

Here $z_i \equiv \Phi^{-1}(1 - p_i)$, $p_i$ is the $p$-value of node (gene) $i$ obtained from a genome-wide association study (GWAS) and $\Phi^{-1}$ denotes the inverse normal cumulated distribution function (CDF) that converts the $p$-values to $Z$-scores (node weights). $Z_A$ is the aggregate score of subgraph $A$ with size $k = |A|$, and $S_A$ is the normalized score of subgraph $A$, calculated by comparing $Z_A$ to the distribution of aggregate scores for node sets of size $k$ randomly chosen from the graph. Two heuristic algorithms have been developed to solve the MWCS problem approximately: (i) the simulated annealing algorithm (Van Laarhoven and Aarts 1987) and (ii) the greedy search algorithm (Nacu et al. 2007). The exact solution of the MWCS problem can only be obtained by solving the corresponding integer linear programming (ILP) problem, which is NP-hard in general.

Dittrich *et al.* introduced a scalable scoring function with false discovery rate (FDR) as a meaningful adjustment parameter (Dittrich et al. 2008). The additivity of this logarithmic score enables them to exactly solve the MWCS problem using an ILP approach (Dittrich et al. 2008; El-Kebir and Klau 2014). A fundamental limitation of the MWCS formalism is that edge weights are completely ignored, which might cause serious problems in the presence of many false positive interactions. Moreover, the MWCS formalism could be confounded by nodes with high degrees, such as *TP53* and *UBC*, which interact with hundreds of other proteins in the cellular molecular interactome and hence may make them more likely to be selected in the optimal subgraphs.

## 2.2 Prize-Collecting Steiner Tree (PCST).

Molecular interaction data (especially PPIs) typically contain large numbers of false positives that can diminish the predictive power of integrative approaches. To optimally exploit the large amount of human interaction data that have been generated in the past decade, the detection of optimal subgraphs of the functional network has been formalized as a PCST problem (Bailly-Bechet et al. 2011): Given a graph $G(V, E)$ with positive node weights (also called node prizes) $b_i$ and positive edge weights (also called edge costs) $c_e$, find the connected subgraph $A$ that maximizes the following objective function:

$$\max_{A \subseteq G} S_A, \text{ with } S_A \equiv \lambda \sum_{i \in V(A)} b_i - \sum_{e \in E(A)} c_e \qquad [2]$$
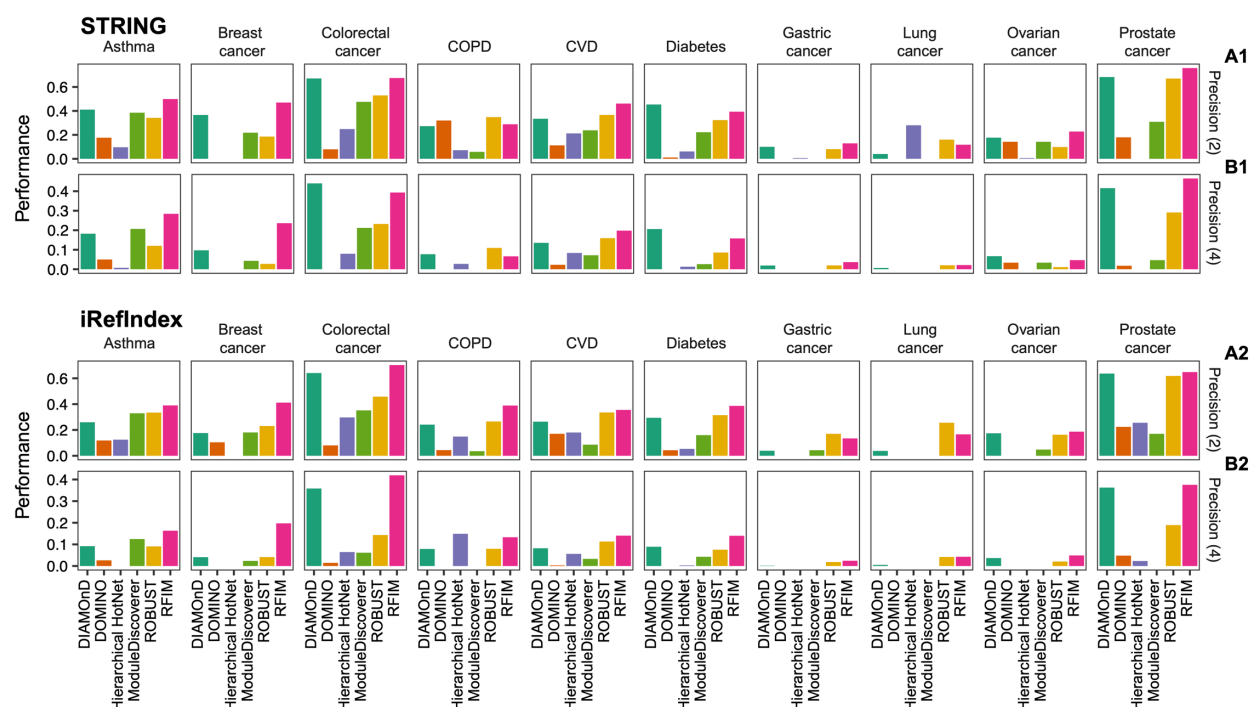
Here, the edge costs $c_e$'s are chosen such that edges with high confidence interactions (e.g., protein interactions verified in small-scale experiments or found in many large-scale datasets) have lower values with respect to low confidence ones (Bailly-Bechet et al. 2011; Biazzo et al. 2012). The node prizes are computed as $b_i = -\log p_i$ with $p_i$ the *p*-value of node (gene) *i*. (Note that $b_i$ is, by definition, positive.) The control parameter $\lambda$ regulates the trade-off between the edge costs and node prizes, and its value indirectly controls the size of the optimal subgraph $A$. In spite of its apparent simplicity, the PCST problem is computationally intractable (NP-hard) (Bailly-Bechet et al. 2011; Biazzo et al. 2012). A fundamental limitation of the PCST formalism is that there is no clear theoretical criterion to choose the control parameter $\lambda$. In practice, one has to solve the PCST problem for a wide range of $\lambda$ values and then use some heuristics to select the best $\lambda$ value (Bailly-Bechet et al. 2011; Gitter et al. 2014; Balbin et al. 2013). Moreover, by definition, the output of the PCST formalism is always a tree, i.e., an acyclic connected subgraph. Recently, the PCST problem has been extended to the prize-collecting Steiner forest (PCSF) problem by introducing another control parameter, where the optimal subgraph is a set of disjoint tree graphs (Gitter et al. 2014). However, the disease module is not necessarily a tree or a forest graph.
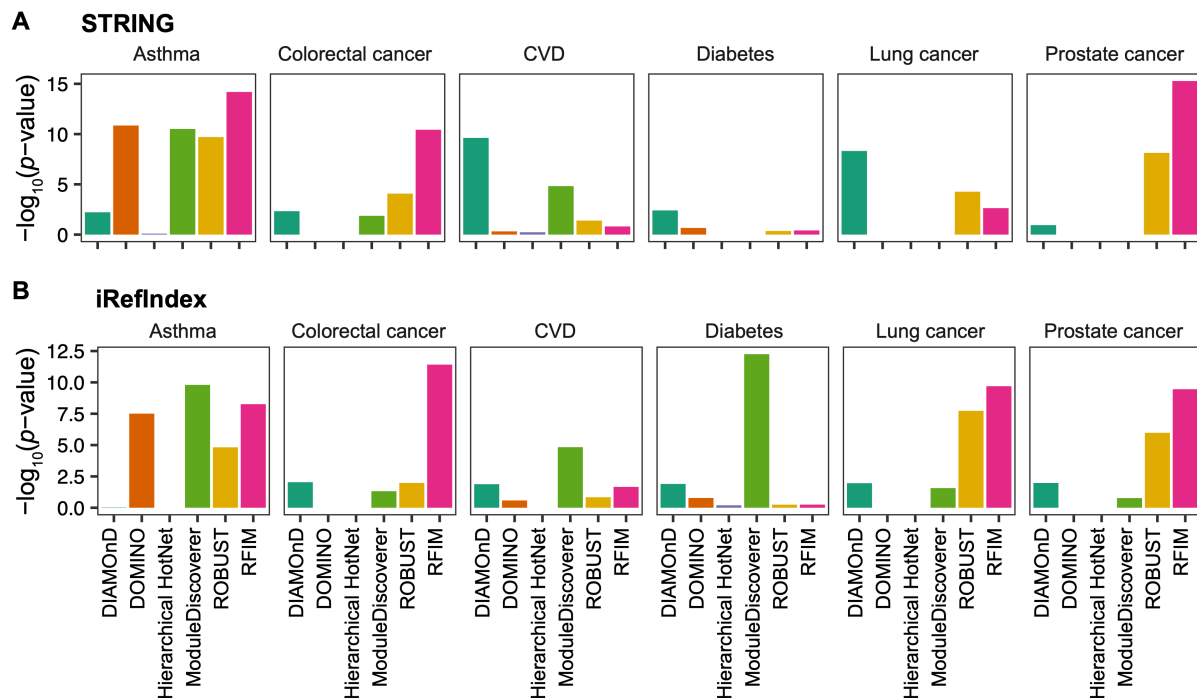
# References

Bailly-Bechet M, Borgs C, Braunstein A, Chayes J, Dagkessamanskaia A, François J-M, Zecchina R. 2011. Finding undetected protein associations in cell signaling by belief propagation. *Proc Natl Acad Sci* **108**: 882–887.

Balbin OA, Prensner JR, Sahu A, Yocum A, Shankar S, Malik R, Fermin D, Dhanasekaran SM, Chandler B, Thomas D. 2013. Reconstructing targetable pathways in lung cancer by integrating diverse omics data. *Nat Commun* **4**: 1–13.

Bernett J, Krupke D, Sadegh S, Baumbach J, Fekete SP, Kacprowski T, List M, Blumenthal DB. 2022. Robust disease module mining via enumeration of diverse prize-collecting Steiner trees ed. L. Cowen. *Bioinformatics* **38**: 1600–1606.

Bersanelli M, Mosca E, Remondini D, Castellani G, Milanesi L. 2016. Network diffusion-based analysis of high-throughput data for the detection of differentially enriched modules. *Sci Rep* **6**: 34841.

Biazzo I, Braunstein A, Zecchina R. 2012. Performance of a cavity-method-based algorithm for the prize-collecting Steiner tree problem on graphs. *Phys Rev E* **86**: 026706.

de Weerd HA, Badam TVS, Martínez-Enguita D, Åkesson J, Muthas D, Gustafsson M, Lubovac-Pilav Z. 2020. MODifieR: an Ensemble R Package for Inference of Disease Modules from Transcriptomics Networks ed. L. Cowen. *Bioinformatics* **36**: 3918–3919.

Dittrich MT, Klau GW, Rosenwald A, Dandekar T, Müller T. 2008. Identifying functional modules in protein–protein interaction networks: an integrated exact approach. *Bioinformatics* **24**: i223–i231.

El-Kebir M, Klau GW. 2014. Solving the maximum-weight connected subgraph problem to optimality. *ArXiv Prepr ArXiv14095308*.

Ghiassian SD, Menche J, Barabási A-L. 2015. A DIseAse MOdule Detection (DIAMOnD) Algorithm Derived from a Systematic Analysis of Connectivity Patterns of Disease Proteins in the Human Interactome ed. A. Rzhetsky. *PLOS Comput Biol* **11**: e1004120.

Gitter A, Braunstein A, Pagnani A, Baldassi C, Borgs C, Chayes J, Zecchina R, Fraenkel E. 2014. Sharing information to reconstruct patient-specific pathways in heterogeneous diseases. In *Biocomputing 2014*, pp. 39–50, World Scientific.

Kolberg L, Raudvere U, Kuzmin I, Vilo J, Peterson H. 2020. gprofiler2 -- an R package for gene list functional enrichment analysis and namespace conversion toolset g:Profiler. *F1000Research* **9**: ELIXIR-709.

Levi H, Elkon R, Shamir R. 2021. DOMINO: a network-based active module identification algorithm with reduced rate of false calls. *Mol Syst Biol* **17**: e9593.

Nacu Ş, Critchley-Thorne R, Lee P, Holmes S. 2007. Gene expression network analysis and applications to immunology. *Bioinformatics* **23**: 850–858.

Piñero J, Bravo À, Queralt-Rosinach N, Gutiérrez-Sacristán A, Deu-Pons J, Centeno E, García-García J, Sanz F, Furlong LI. 2016. DisGeNET: a comprehensive platform integrating information on human disease-associated genes and variants. *Nucleic Acids Res* gkw943.

Reyna MA, Leiserson MD, Raphael BJ. 2018. Hierarchical HotNet: identifying hierarchies of altered subnetworks. *Bioinformatics* **34**: i972–i980.

Team RC. 2013. R: A language and environment for statistical computing.

Ulitsky I, Krishnamurthy A, Karp RM, Shamir R. 2010. DEGAS: De Novo Discovery of Dysregulated Pathways in Human Diseases ed. T. Ravasi. *PLoS ONE* **5**: e13367.

Van Laarhoven PJ, Aarts EH. 1987. Simulated annealing. In *Simulated annealing: Theory and applications*, pp. 7–15, Springer.

Vlaic S, Conrad T, Tokarski-Schnelle C, Gustafsson M, Dahmen U, Guthke R, Schuster S. 2018. ModuleDiscoverer: Identification of regulatory modules in protein-protein interaction networks. *Sci Rep* **8**: 433.

Yu G, He Q-Y. 2016. ReactomePA: an R/Bioconductor package for reactome pathway analysis and visualization. *Mol Biosyst* **12**: 477–479.
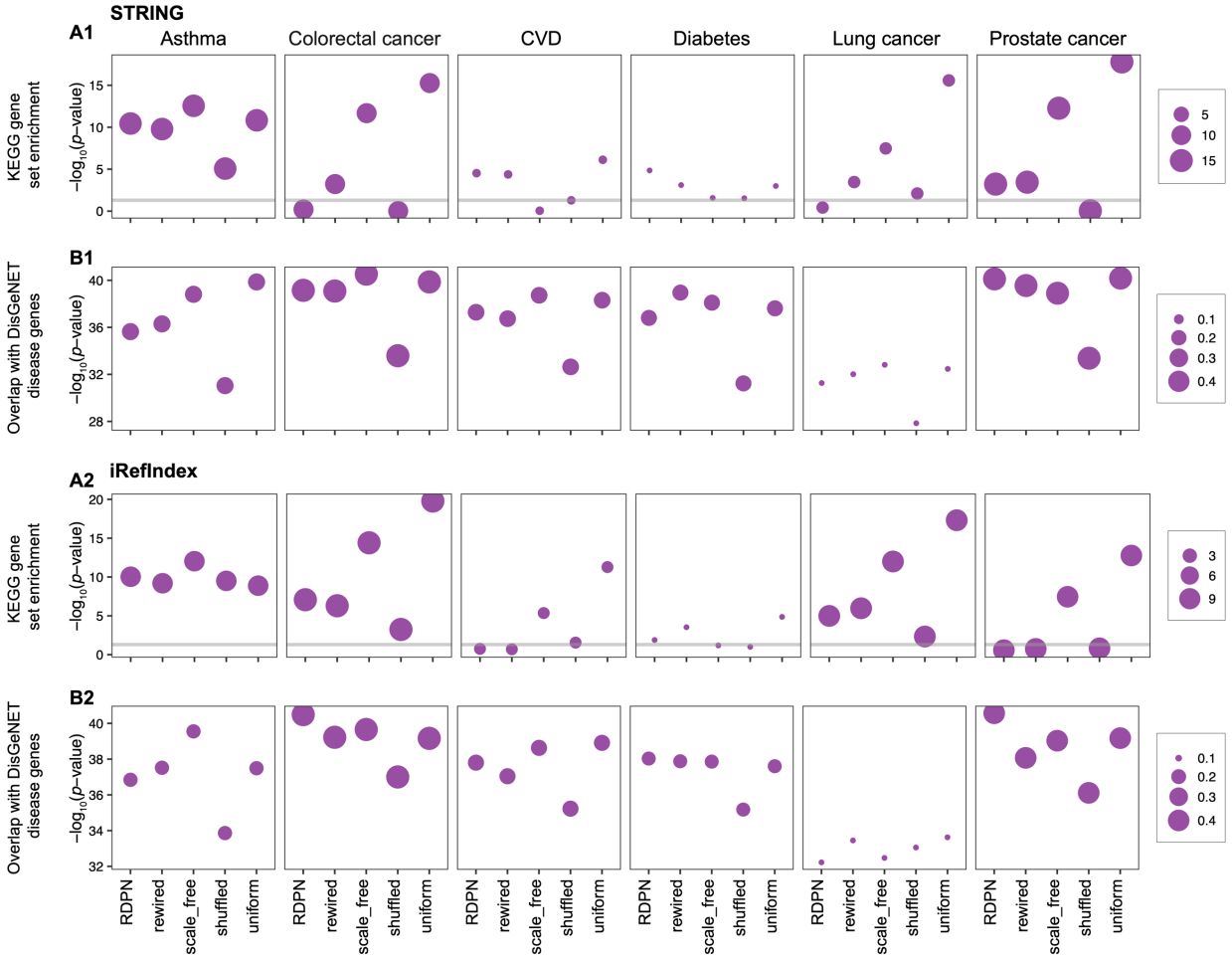
# Supplemental Figures



**Supplemental Fig S1: Enrichment analysis for genes in disease module detected from integrating the GWAS with two interactomes for each phenotype and method.** We obtained the enriched pathways for disease genes using the ReactomePA package (Yu and He 2016) whose *p*-values are lower than 0.05 cutoff adjust by False Discovery Rate (FDR). Then, we extracted the disease-associated genes of each phenotype using the DisGeNET (Piñero et al. 2016) database and calculated the precision defined as fraction of enriched pathways with at least two (A) or four (B) disease-associated genes over all enriched pathways.
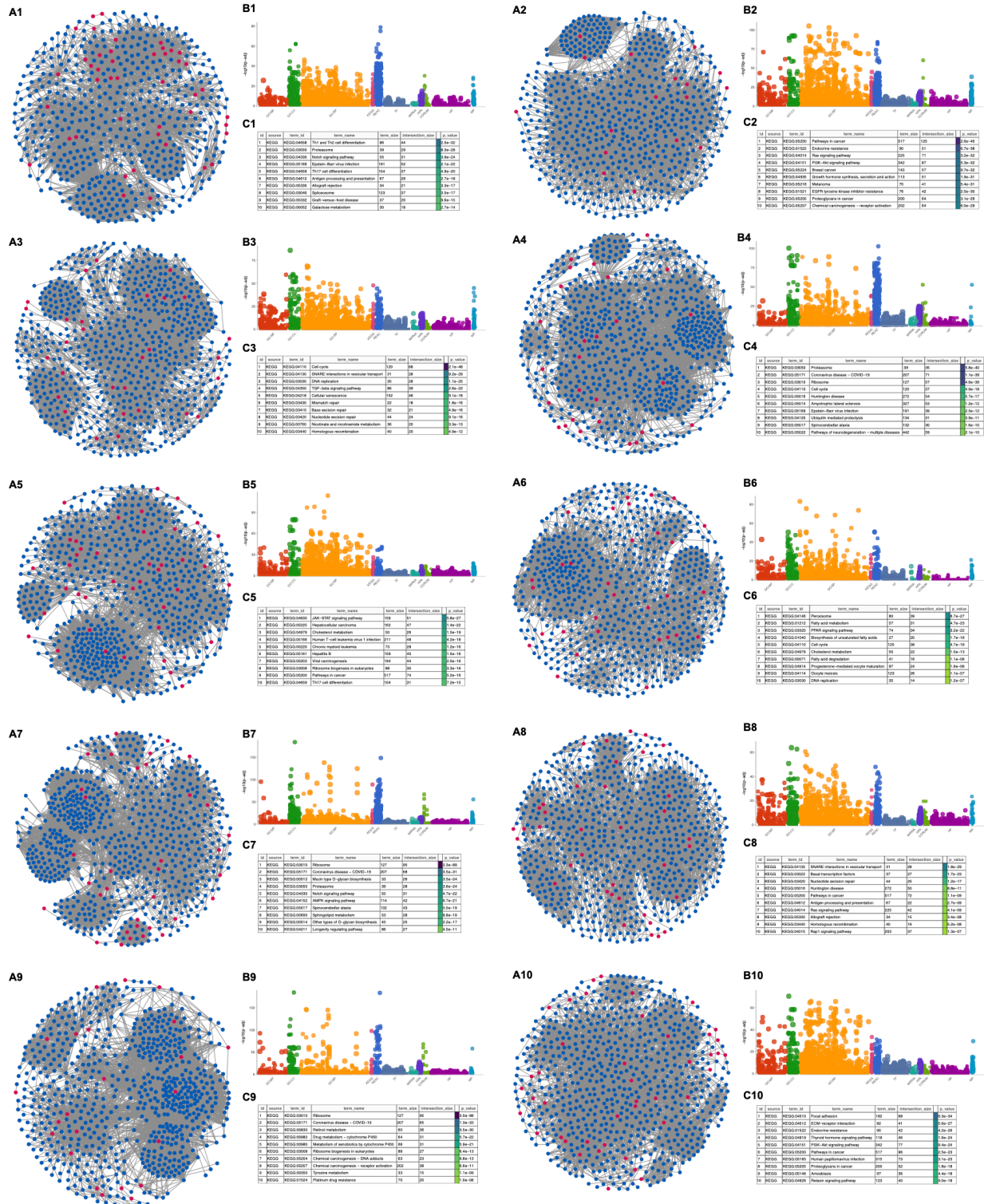
**Supplemental Fig S2: KEGG enrichment analysis for genes in disease module detected from integrating the GWAS with two interactomes for each phenotype and method.** We first extracted the enriched pathways of genes in disease module with *p*-value cutoff 0.05. Then, we computed the average *p*-values of pathways related to the phenotype in KEGG database. KEGG pathways: asthma: hsa05310; breast cancer: hsa05224; lung cancer: hsa05222, hsa05223; colorectal cancer: hsa05210; gastric cancer: hsa05226; prostate cancer: hsa05215; CVD: hsa05410, hsa05412, hsa05414, hsa05416; diabetes: hsa04940 and hsa04930. Note that there are no KEGG pathways associated with COPD and ovarian cancer; and there are no significant KEGG pathways (with *p*<0.05, based on GSEA) associated with breast cancer and gastric cancer. Hence, those four diseases were not considered here.

**Supplemental Fig S3: Meaningfulness test of disease module identified by RFIM.** Two meaningfulness tests were performed: (A1-A2) KEGG gene set enrichment analysis (GSEA). We compared the *p*-value of gene set in the disease module identified from the original PPI network with that identified from 10 randomized PPI networks generated by each of the following five random network generators: degree preserving (RDPN), expected degree preserving (rewired), topology preserving (shuffled), scale-free (scale_free) and uniform (uniform). (B1-B2) Overlap with DisGeNET disease gene. We compared the overlap genes with DisGeNET genes in the original PPI network with that in 10 randomized PPI networks generated by each of the above five random network generators. Then we showed the *p*-values using one-sided one-sample t-test. (A): STRING and (B) iRefIndex. Node size represents the *p*-value. KEGG pathways: asthma: hsa05310; breast cancer: hsa05224; lung cancer: hsa05222, hsa05223; colorectal cancer: hsa05210; gastric cancer: hsa05226; prostate cancer: hsa05215; CVD: hsa05410, hsa05412, hsa05414, hsa05416; diabetes: hsa04940 and hsa04930. Note that there are no KEGG pathways associated with COPD and ovarian cancer; and there are no significant KEGG pathways (with *p*<0.05, based on GSEA) associated with breast cancer and gastric cancer. Hence, those four diseases were not considered here.

**Supplemental Fig S4: Subnetworks of genes in the disease module and enriched functional terms using STRING interactome. (A1-A10):** Subnetwork of genes in disease modules of asthma, breast cancer, colorectal cancer, COPD, CVD, diabetes, gastric cancer, lung cancer, ovarian cancer and prostate cancer. Genes with *p*-values lower (higher) than 0.001 are colored with

red (blue). **(B1-B10):** Enriched functional terms of genes in disease modules of asthma, breast cancer, COPD, CVD, diabetes and lung cancer using gprofier2 (Kolberg et al. 2020, 2). **(C1-C10):** Top-10 enriched KEGG pathways of genes in the disease modules of asthma, breast cancer, colorectal cancer, COPD, CVD, diabetes, gastric cancer, lung cancer, ovarian cancer and prostate cancer.

**Supplemental Fig S5: Subnetworks of genes in the disease module and enriched functional terms using iRefIndex interactome. (A1-A10):** Subnetwork of genes in disease modules of asthma, breast cancer, colorectal cancer, COPD, CVD, diabetes, gastric cancer, lung cancer, ovarian cancer and prostate cancer. Genes with *p*-values lower (higher) than 0.001 are colored with red (blue). **(B1-B10):** Enriched functional terms of genes in disease modules of asthma, breast

cancer, COPD, CVD, diabetes and lung cancer using gprofier2 (Kolberg et al. 2020, 2). **(C1-C10):** Top-10 enriched KEGG pathways of genes in the disease modules of asthma, breast cancer, colorectal cancer, COPD, CVD, diabetes, gastric cancer, lung cancer, ovarian cancer and prostate cancer.



**Supplemental Fig S6: Overlap between disease modules identified by different methods. (A):** STRING. **(B):** iRefIndex.