# An enhanced disease network with robust cross-phenotype relationships via variant frequency-inverse phenotype frequency
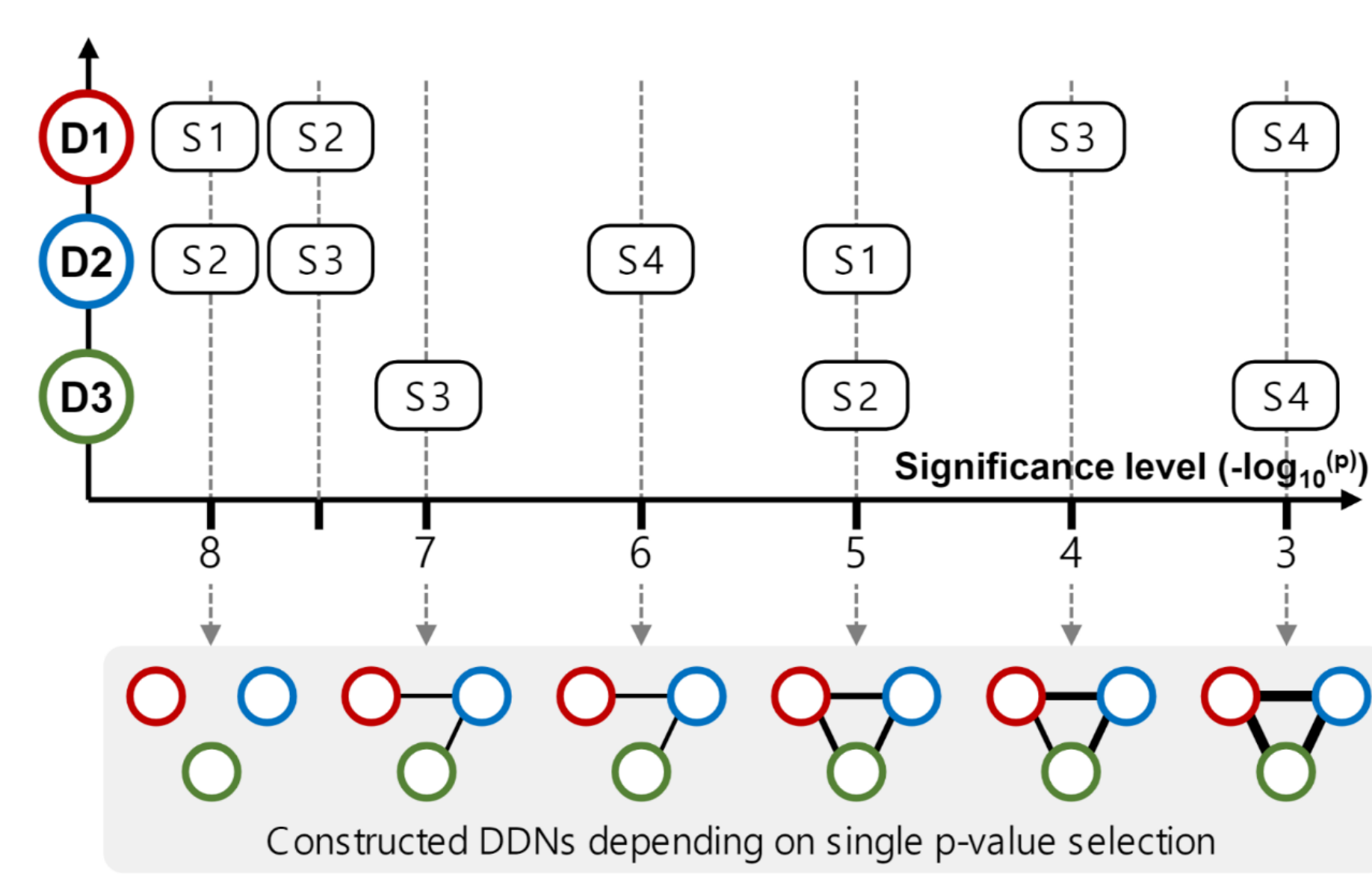
**Yonghyun Nam, Vivek Sriram, Dokyoon Kim**
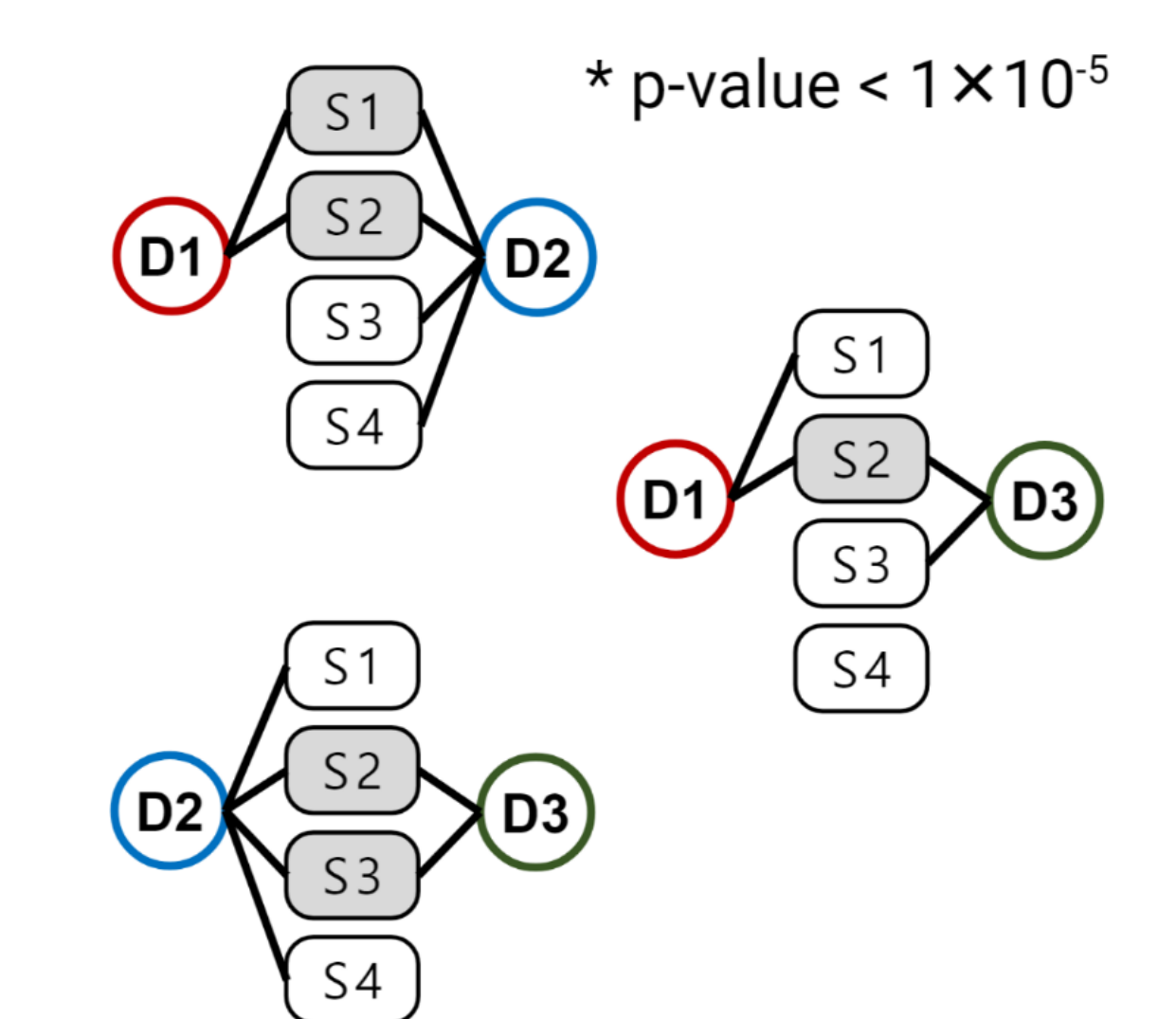Department of biostatistics, epidemiology & informatics

* Contact: Yonghyun Nam, Ph.D.,(Yonghyun.Nam@pennmedicine.upenn.edu)

## Abstract

**Motivation:** Biobank-scaled phenome-wide association studies (PheWAS) identify associations between multiple phenotypes and common genetic variants. Disease-disease networks (DDNs) can intuitively observe relationships across many phenotypes. However, DDNs constructed from PheWAS summary data will highlight different disease associations depending on the chosen significance threshold. Thus, we proposed a novel method, variant frequency-inverse phenotype frequency (VF-IPF), to increase significance and decrease uncertainty in associations when constructing DDNs. **Results:** We constructed an enhanced DDN (eDDN) among 421 phenotypes from UK Biobank (UKBB) PheWAS summary data. We then used VF-IPF to enhance connections across diseases. To validate the enhanced cross phenotype relationships in eDDN, graph-based semi-supervised learning was applied to predict scores of disease co-occurrence. Ground truth co-occurrences were generated from UKBB inpatient data. Compared to conventional DDNs built from individual significance levels, the eDDN showed improved performance in identifying cross-phenotype associations. To demonstrate its clinical significance, we further validated our predictions using myocardial infarction as an index disease of interest against co-occurrence information from the Penn Medicine Biobank.

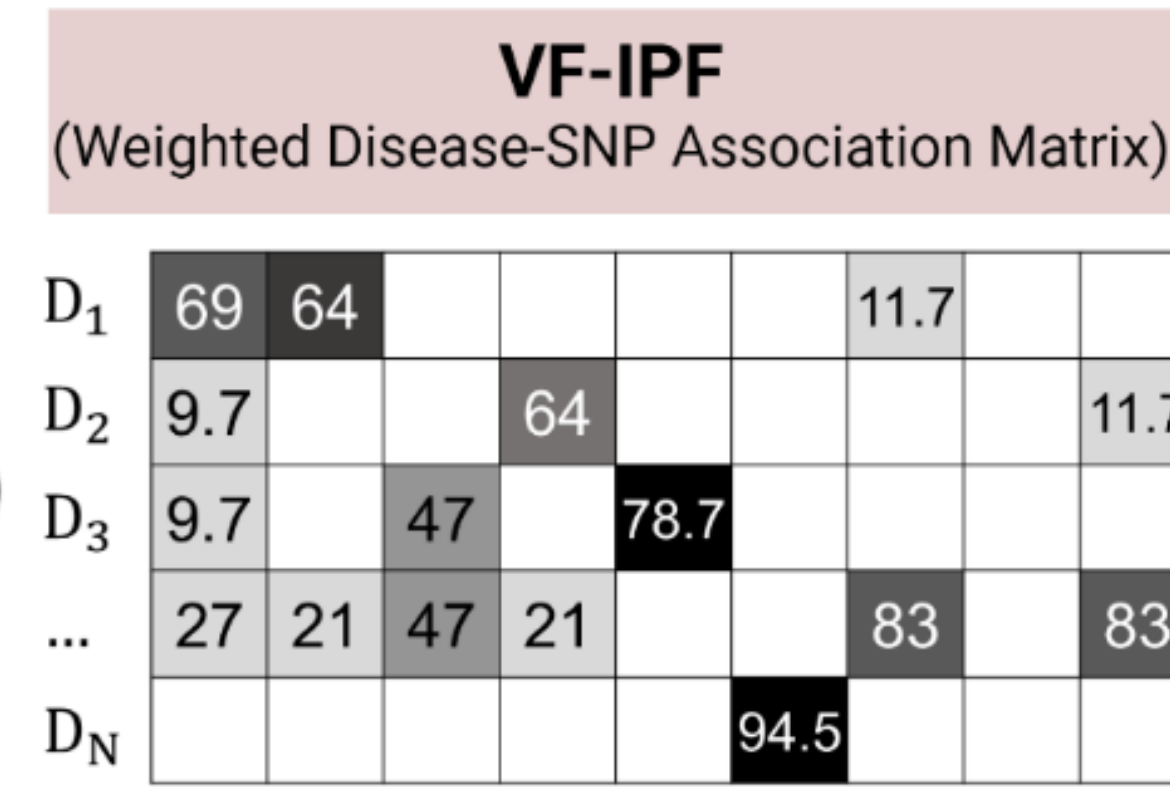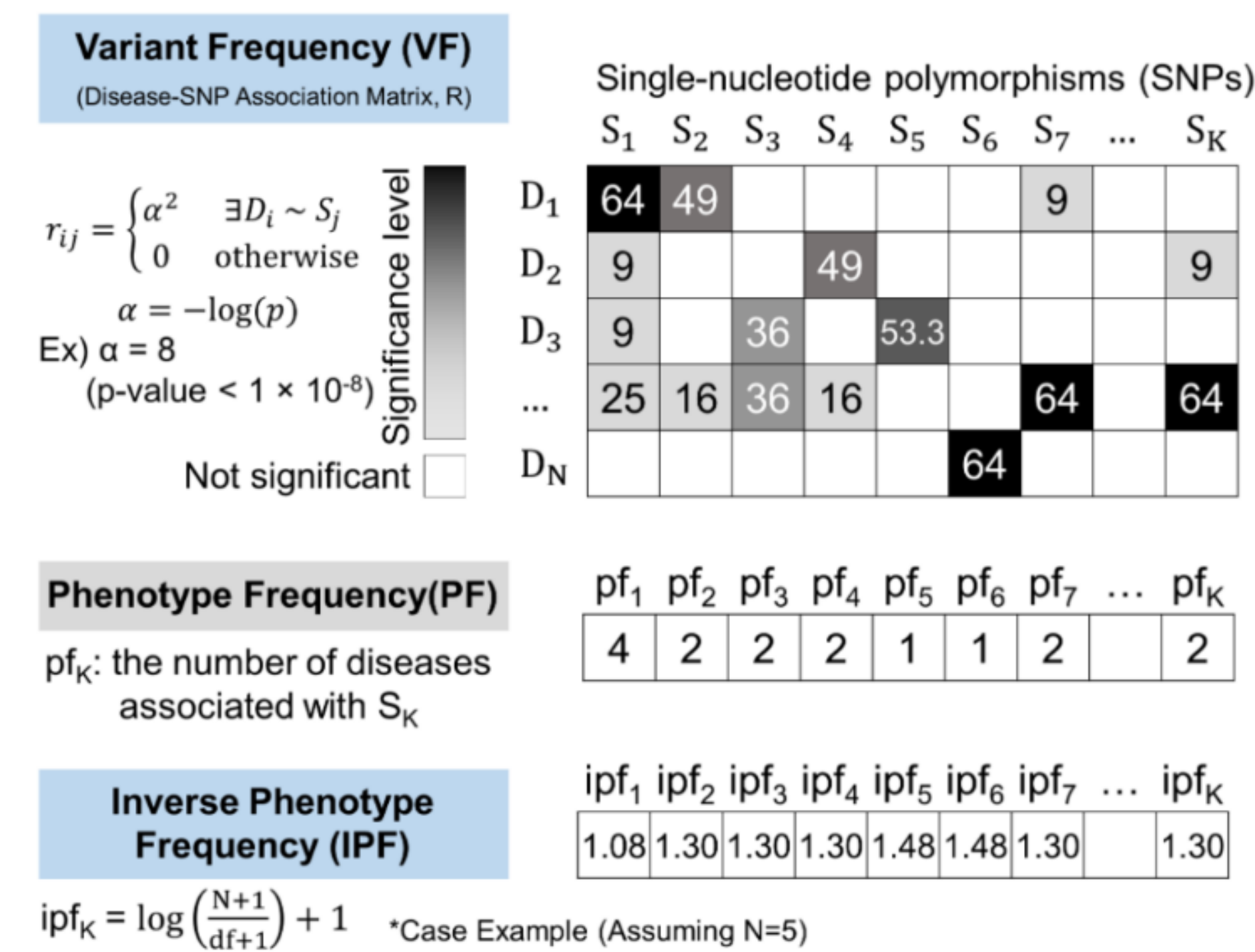**a** Disease - SNP associations in PheWAS data



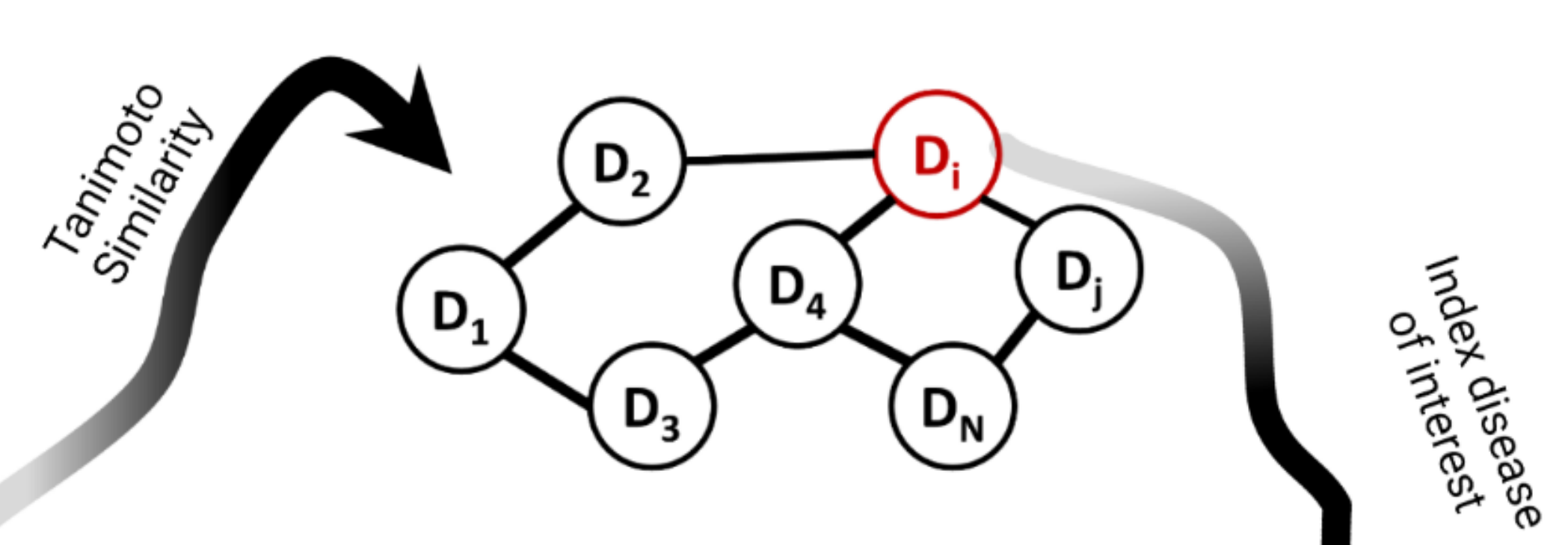**b** Disease-Disease associations by shared component hypothesis



## Proposed Method: Enhanced DDN with VF-IPF

- We develop a novel method to **increase the significance of associations** and **decrease uncertainty of associations** simultaneously in disease-SNP associations regardless of the arbitrary selection of significance levels
- We propose **the enhanced disease-disease network (eDDN)**, which can discover more scalable and robust relationships across multiple phenotypes from biobank-scaled PheWAS data

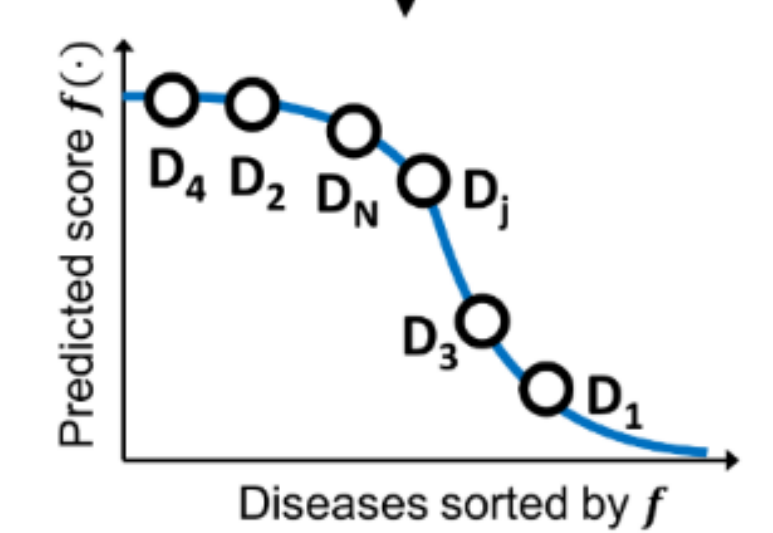**a** Variant frequency-Inverse phenotype frequency (VF-IPF)
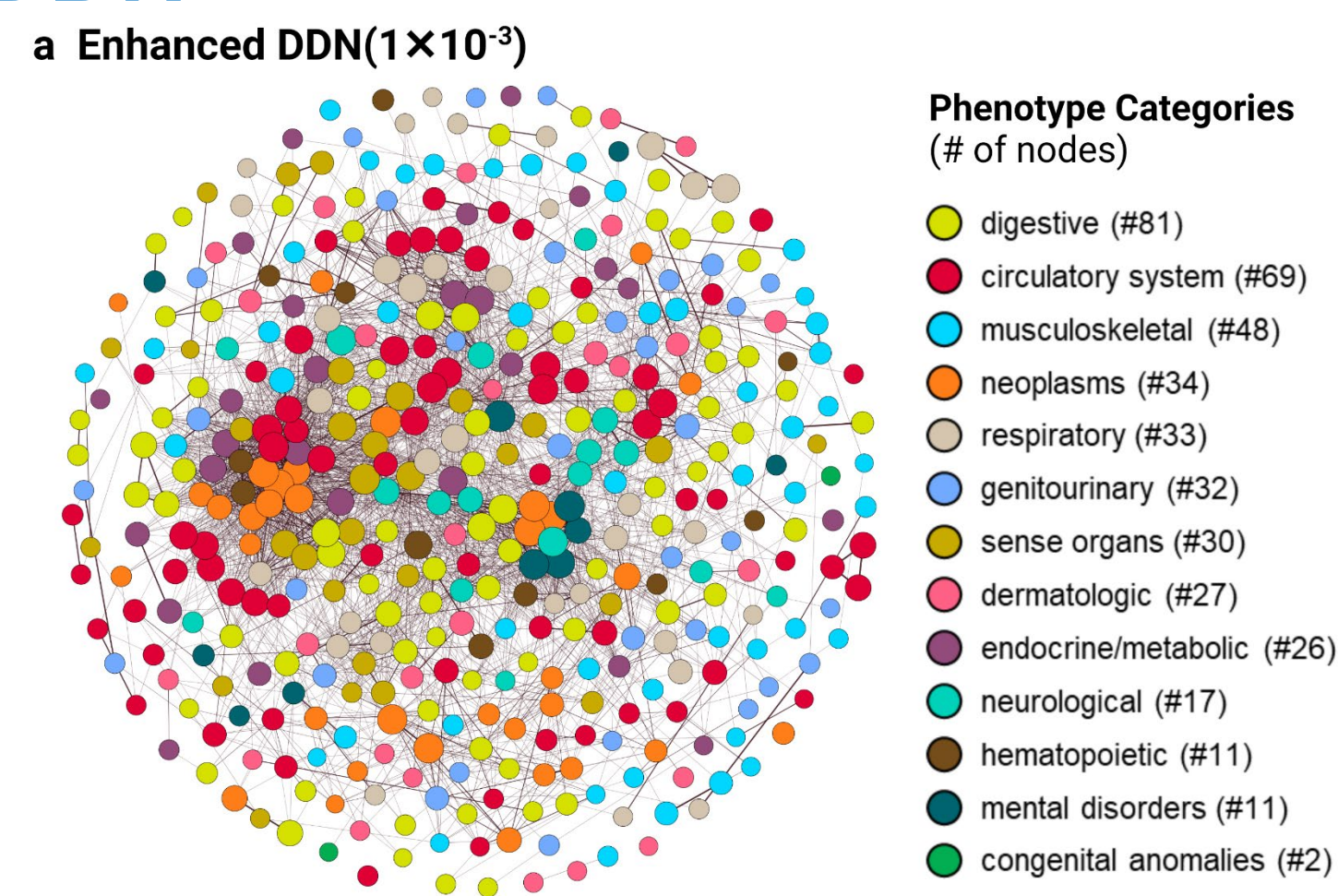


**b** Construction enhanced disease-disease network

**c** Predicting co-occurrence diseases using graph-based SSL

Initial setup for graph-based SSL
- Index disease of interest: $D_i$
- Labeled disease: $y(D_i)=1$
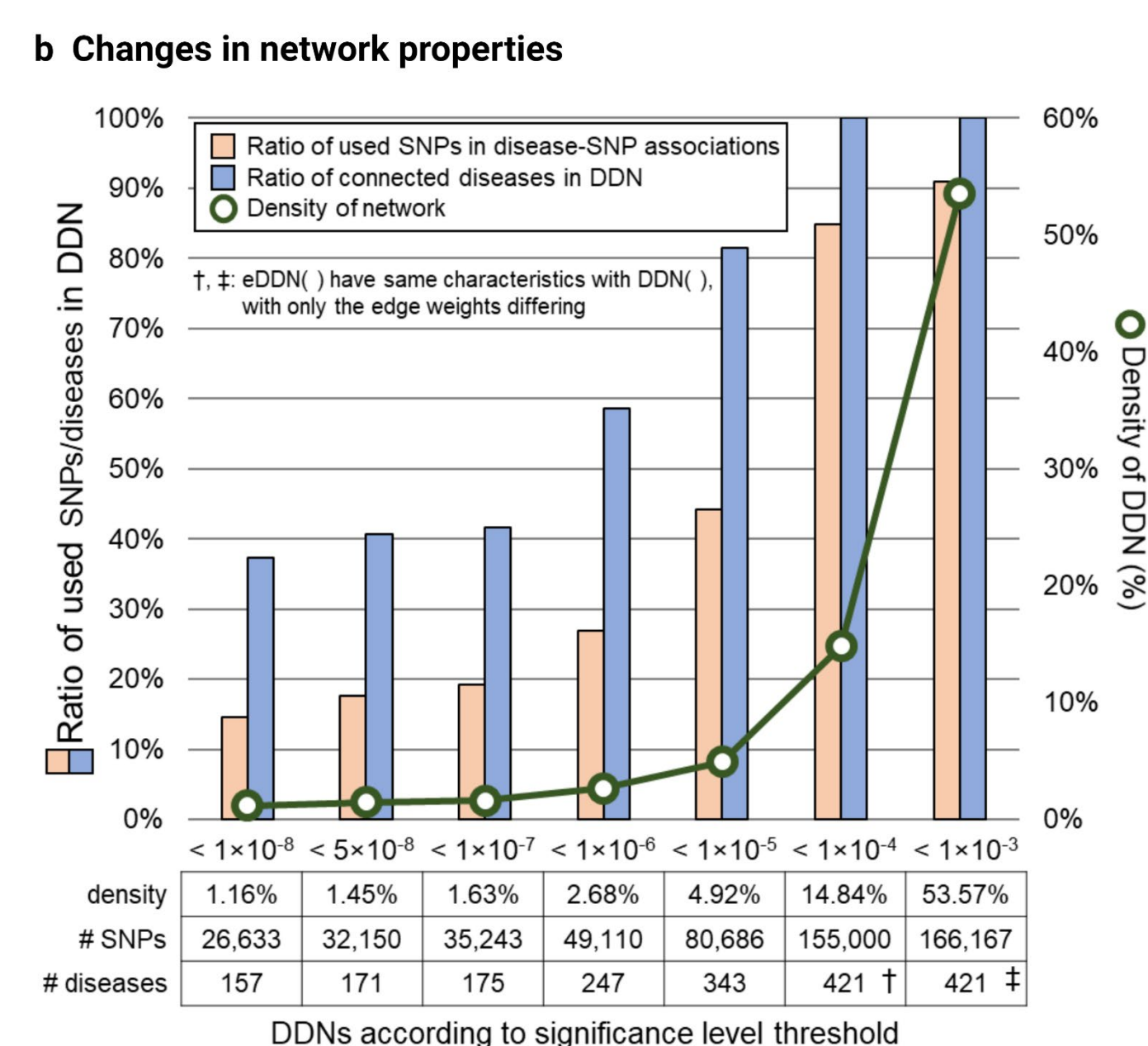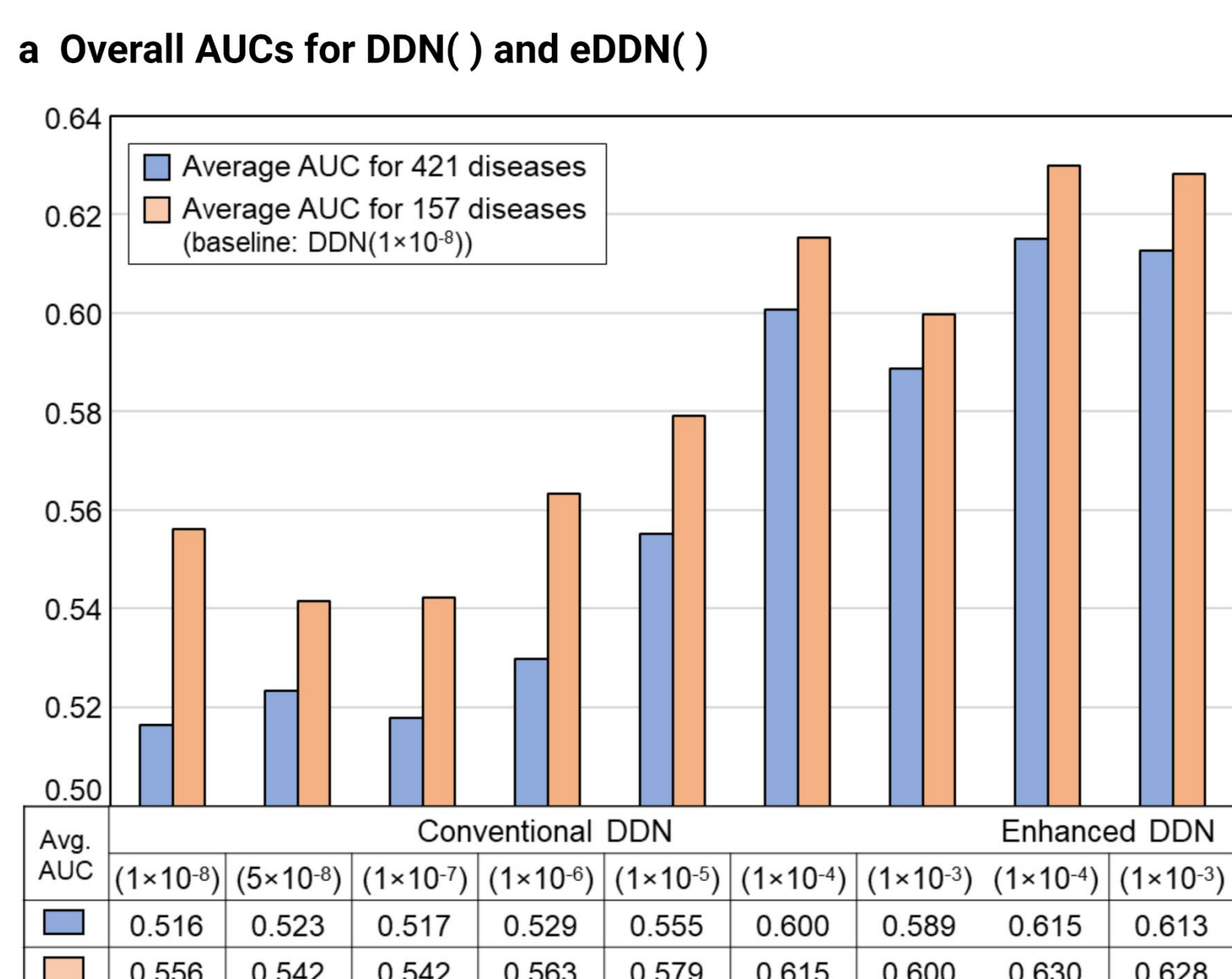- Unlabeled disease: $y(D_1, D_2, .., D_N)=0$

## Result (1): Construction of eDDN

- The eDDN were constructed using UKBB PheWAS summary statistics
  - Number of Phenotypes: 421

**a** Enhanced DDN($1\times10^{-3}$)



Phenotype Categories (# of nodes)
- digestive (#81)
- circulatory system (#69)
- musculoskeletal (#48)
- neoplasms (#34)
- respiratory (#33)
- genitourinary (#32)
- sense organs (#30)
- dermatologic (#27)
- endocrine/metabolic (#26)
- neurological (#17)
- hematopoietic (#11)
- mental disorders (#11)
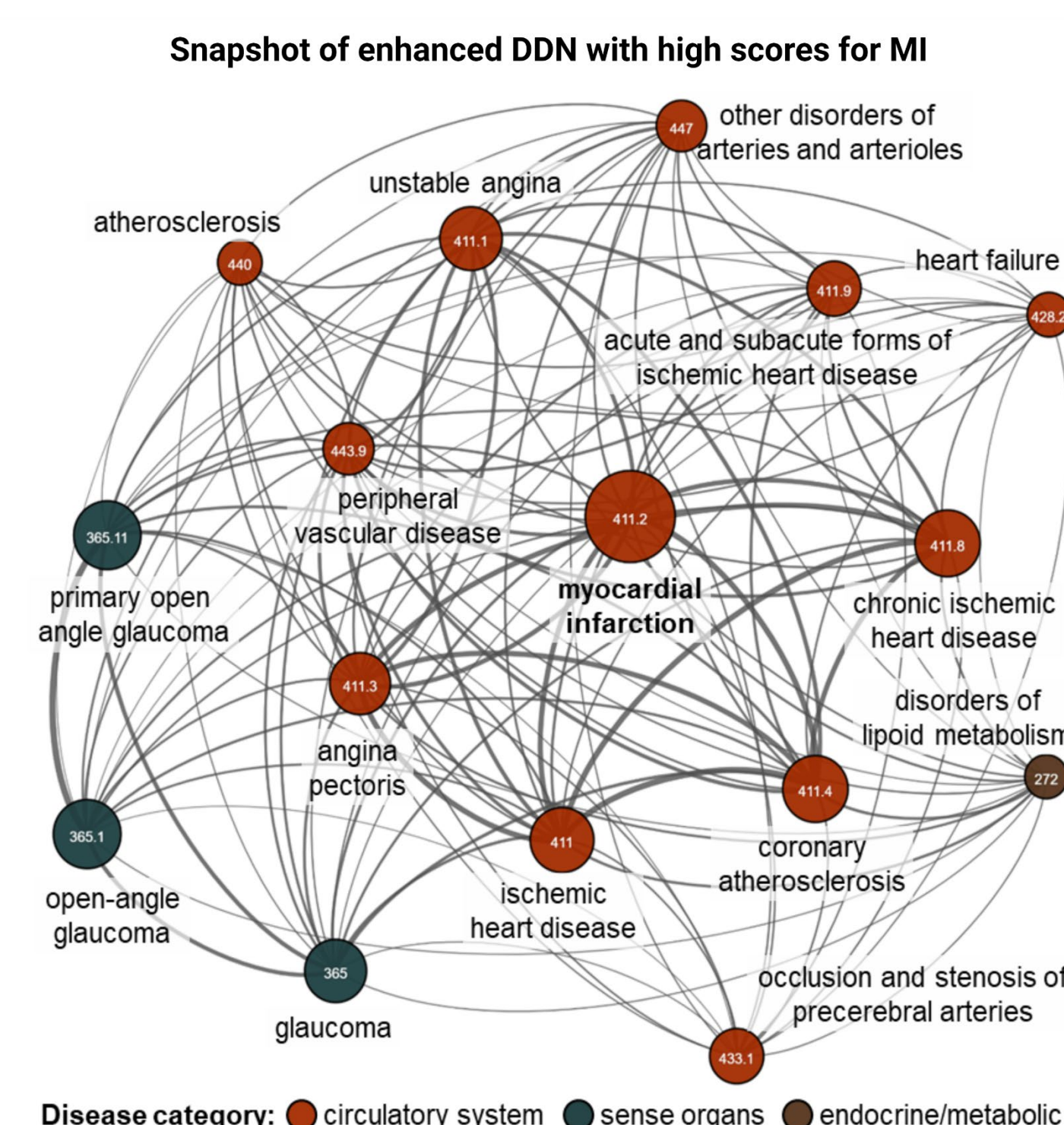- congenital anomalies (#2)

## Result (2): Performance Comparison

- **Task:** Predicting co-occurrence scoring
  - Ground truth: Comorbidity measurements calculated by UKBB inpatient EHR data
  - Performance measurement: AUC
- The eDDN showed the highest average AUC of 0.630 in co-occurrence prediction when compared with the other eight DDNs

**a** Overall AUCs for DDN( ) and eDDN( )



| Avg. AUC | Conventional DDN | | | | | | Enhanced DDN | |
|---|---|---|---|---|---|---|---|---|---|
| | ($1\times10^{-8}$) | ($5\times10^{-8}$) | ($1\times10^{-7}$) | ($1\times10^{-6}$) | ($1\times10^{-5}$) | ($1\times10^{-4}$) | ($1\times10^{-3}$) | ($1\times10^{-4}$) | ($1\times10^{-3}$) |
| | 0.516 | 0.523 | 0.517 | 0.529 | 0.555 | 0.600 | 0.589 | 0.615 | 0.613 |
| | 0.556 | 0.542 | 0.542 | 0.563 | 0.579 | 0.615 | 0.600 | 0.630 | 0.628 |

**b** Changes in network properties



†, ‡: eDDN( ) have same characteristics with DDN( ), with only the edge weights differing

| | < $1\times10^{-8}$ | < $5\times10^{-8}$ | < $1\times10^{-7}$ | < $1\times10^{-6}$ | < $1\times10^{-5}$ | < $1\times10^{-4}$ | < $1\times10^{-3}$ |
|---|---|---|---|---|---|---|---|
| density | 1.16% | 1.45% | 1.63% | 2.68% | 4.92% | 14.84% | 53.57% |
| # SNPs | 26,633 | 32,150 | 35,243 | 49,110 | 80,686 | 155,000 | 166,167 |
| # diseases | 157 | 171 | 175 | 247 | 343 | 421 † | 421 ‡ |

DDNs according to significance level threshold

## Result (3): Clinical implementation of co-occurrence scoring

- This result shows a proof-of-concept study for demonstrating the utility of eDDN.
- Transferring the discovered knowledge in genetic information to clinical significance is essential in translational bioinformatics.
- The results can give significance to the discovered information and provide evidence that can be used for clinical implementation.
- Predicted results were verified using UKBB EHR in the analysis as well as comorbidity information obtained from Penn Medicine Biobank (PMBB), supplementing the evidence that can be used to utilize the clinical implications.

Snapshot of enhanced DDN with high scores for MI



Disease category: ● circulatory system ● sense organs ● endocrine/metabolic

| Phenotype description (PheCode) | Score ($f$) | Relative Risk [95% CI] | |
|---|---|---|---|
| | | UKBB | PMBB |
| Coronary atherosclerosis (411.4) | 0.9884 | 13.014 [13.011, 13.017] | 3.439 [3.437, 3.442] |
| Chronic ischemic heart disease (411.8) | 0.9883 | 13.542 [13.539, 13.545] | 6.400 [6.390, 6.410] |
| Open-angle glaucoma (365.1) | 0.9904 | 1.848 [1.807, 1.890] | 1.133 [1.084, 1.183] |
| Primary open angle glaucoma (365.11) | 0.9901 | 1.857 [1.815, 1.899] | 1.555 [1.386, 1.745] |
| Ischemic Heart Disease (411) | 0.9867 | 11.283 [11.281, 11.284] | 3.488 [3.485, 3.490] |
| Unstable angina (411.1) | 0.9852 | 14.569 [14.558, 14.580] | 5.164 [5.146, 5.181] |
| Glaucoma (365) | 0.9844 | 1.943 [1.931, 1.956] | 1.331 [1.315, 1.348] |
| Angina pectoris (411.3) | 0.9832 | 9.448 [9.445, 9.451] | 4.255 [4.243, 4.267] |
| Occlusion and stenosis of precerebral arteries (433.1) | 0.9723 | 6.116 [6.068, 6.165] | 2.809 [2.799, 2.819] |
| acute and subacute forms of ischemic heart disease (411.9) | 0.9721 | 21.922 [21.885, 21.958] | 5.104 [5.036, 5.174] |
| Peripheral vascular disease (443.9) | 0.9656 | 8.273 [8.255, 8.292] | 3.227 [3.219, 3.234] |
| Other disorders of arteries and arterioles (447) | 0.9640 | 5.584 [5.493, 5.676] | 2.259 [2.249, 2.269] |
| Atherosclerosis (440) | 0.9422 | 10.366 [10.084, 10.656] | 3.053 [3.041, 3.065] |
| Heart failure (428.2) | 0.9376 | 12.325 [12.314, 12.335] | 3.948 [3.915, 3.981] |
| Disorders of lipid metabolism (272) | 0.9325 | 5.301 [5.300, 5.303] | 1.718 [1.717, 1.720] |

## Conclusion

- The eDDN can have more reliable explanatory power for disease-disease associations
- It can be robust even in the presence of false positive disease-SNP associations according to less stringent thresholds (p-value < $1\times10^{-3}$).