*Gene Function Prediction*

# Exploiting Ontology Graph for Predicting Sparsely Annotated Gene Function

Sheng Wang[1,†], Hyunghoon Cho[2,†], ChengXiang Zhai[1], Bonnie Berger[2,3] and Jian Peng[1,*]

[1] Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL, USA.
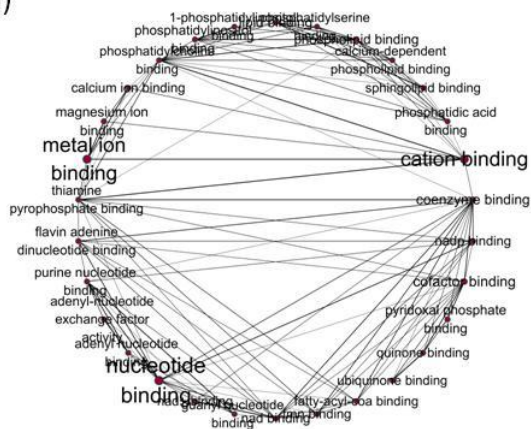
[2] Computer Science and Artificial Intelligence Laboratory, MIT, Cambridge, MA, USA

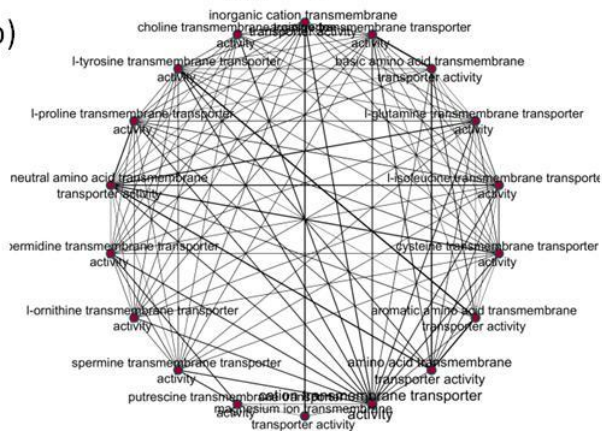[3] Department of Mathematics, MIT, Cambridge, MA, USA

## SUPPLEMENTARY INFORMATION
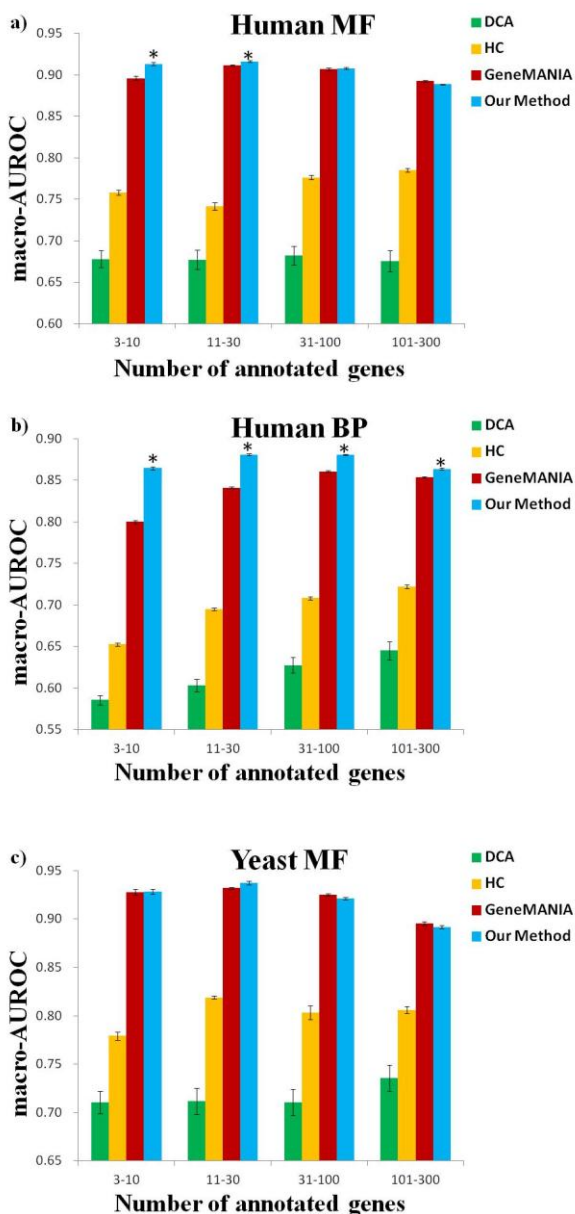
## 1 Two cluster structures discussed in Section 3.3



**Fig. S1.** (a) The cluster structure of several binding functions. (b) The cluster structure of GO labels related to transmembrane transporting.

## 2 Comparison of performance in terms of macro-AUROC

[†]These authors equally contribute to this work.

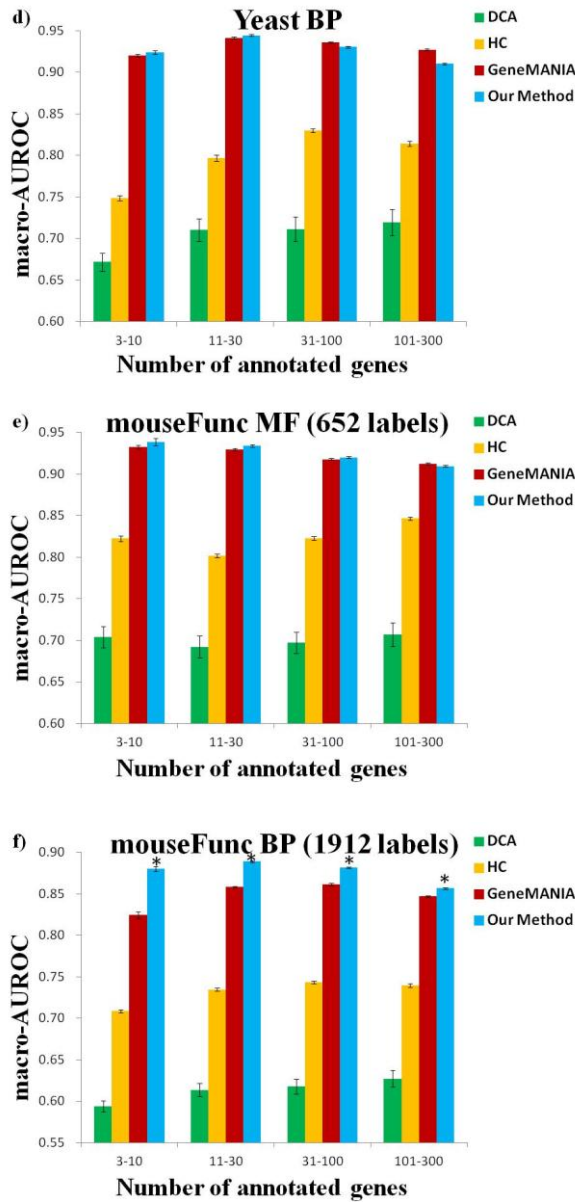[*]To whom correspondence should be addressed. *jianpeng@illinois.edu*

**Fig. S2.** Comparison of our approach with other methods in terms of macro-AUROC. * indicates that our approach is statistically significant in comparison with GeneMANIA. Performance is evaluated for different subsets of GO labels with varying sparsity levels as shown on the x-axis.

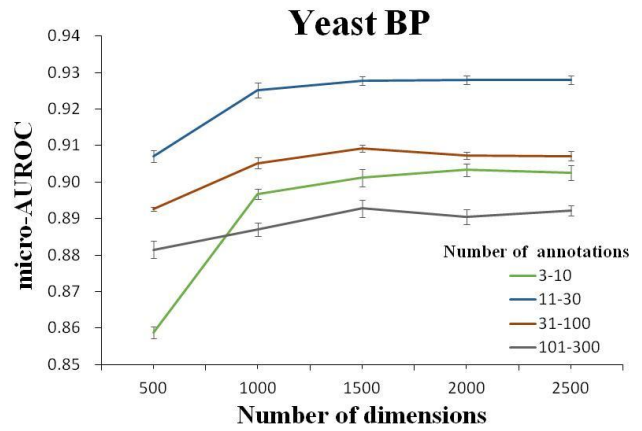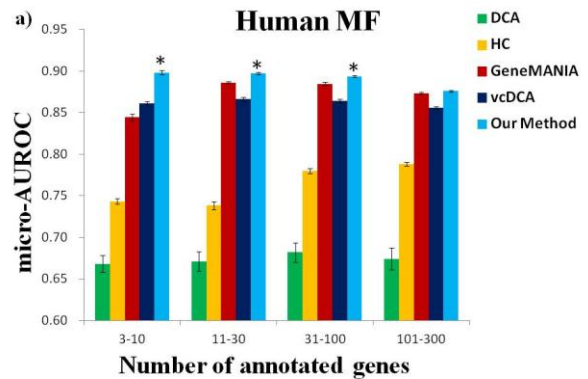# 3 Comparison of number of dimensions discussed in Section 3.7
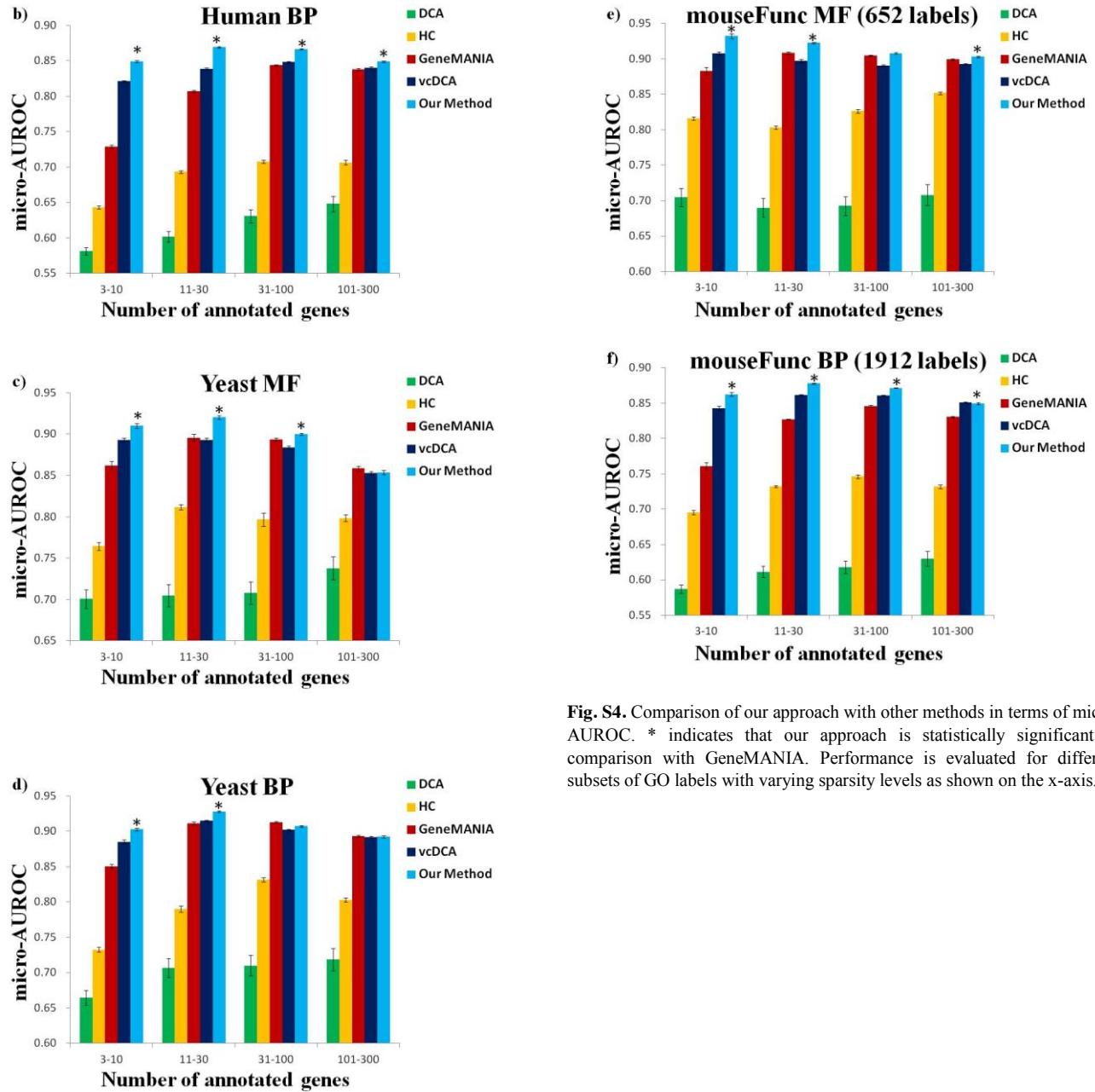


**Fig. S3.** Comparison of number of dimensions in terms of micro-AUROC of biological process in yeast.

# 4 Performance of clustering based on learned label vectors

Besides clustering GO labels based on the sparseness, we explored to cluster the GO labels based on learned label vectors. We denoted this method as vcDCA. We show the performance of this method in **Fig. S4** and **Fig. S5.**

**Fig. S4.** Comparison of our approach with other methods in terms of micro-AUROC. * indicates that our approach is statistically significant in comparison with GeneMANIA. Performance is evaluated for different subsets of GO labels with varying sparsity levels as shown on the x-axis.

**Fig. S5.** Comparison of our approach with other methods in terms of macro-AUROC. * indicates that our approach is statistically significant in comparison with GeneMANIA. Performance is evaluated for different subsets of GO labels with varying sparsity levels as shown on the x-axis.

## 5 Comparison of performance in terms of AP@10 and macro-APRUC

We show the comparison of performance in terms of AP@10 and macro-APRUC in **Table S1** and **Table S2.** In human, our method achieved 0.1849 AUPRC on BP labels with 101-300 annotations, which is higher than 0.1764 AUPRC for GeneMANIA. In mouse, our method achieved 0.4238 AP@10 on MF labels with 31-100 annotations, which is higher than 0.3706 AP@10 for GeneMANIA.

**Table S1.** Comparison of our approach with other methods in terms of AP@10. * indicates that our approach is statistically significant in comparison with GeneMANIA.

|  | #annotated genes | HC | GeneMANIA | clusDCA |
|---|---|---|---|---|
| Human MF | 3-10 | 0.1471 | 0.2123 | **0.2500** * |
|  | 11-30 | 0.1810 | 0.2338 | **0.2770** * |
|  | 31-100 | 0.2596 | 0.3502 | **0.3794** * |
|  | 101-300 | 0.3139 | 0.4883 | **0.5541** * |
| Human BP | 3-10 | 0.0566 | 0.0975 | **0.1224** * |
|  | 11-30 | 0.0769 | 0.1309 | **0.1448** * |
|  | 31-100 | 0.1305 | 0.2628 | **0.2845** * |
|  | 101-300 | 0.1802 | 0.4197 | **0.4938** * |

| | | | | |
|---|---|---|---|---|
| Yeast MF | 3-10 | 0.1851 | 0.3060 | **0.4152** * |
| | 11-30 | 0.2320 | 0.3499 | **0.3953** * |
| | 31-100 | 0.2921 | 0.5325 | **0.5559** * |
| | 101-300 | 0.3197 | 0.6160 | **0.7109** * |
| Yeast BP | 3-10 | 0.1658 | 0.2639 | **0.3453** * |
| | 11-30 | 0.2265 | 0.3482 | **0.3987** * |
| | 31-100 | 0.3434 | **0.5424** | 0.5292 |
| | 101-300 | 0.3875 | 0.7225 | **0.7249** |
| mouseFunc MF | 3-10 | 0.0961 | 0.2254 | **0.3514** * |
| | 11-30 | 0.1978 | 0.2931 | **0.3450** * |
| | 31-100 | 0.2877 | 0.3706 | **0.4238** * |
| | 101-300 | 0.3250 | 0.5278 | **0.6627** * |
| mouseFunc BP | 3-10 | 0.0559 | 0.1010 | **0.1523** * |
| | 11-30 | 0.1020 | 0.1380 | **0.1806** * |
| | 31-100 | 0.1365 | 0.2359 | **0.3019** * |
| | 101-300 | 0.1748 | 0.3437 | **0.4268** * |

**Table S2.** Comparison of our approach with other methods in terms of macro-AUPRC. * indicates that our approach is statistically significant in comparison with GeneMANIA.
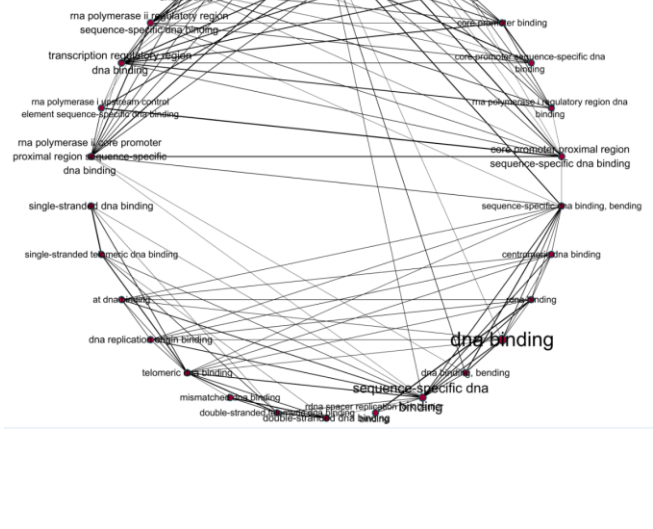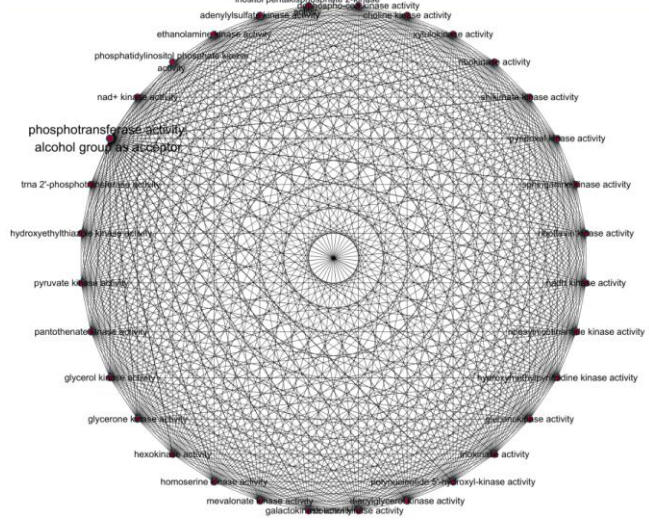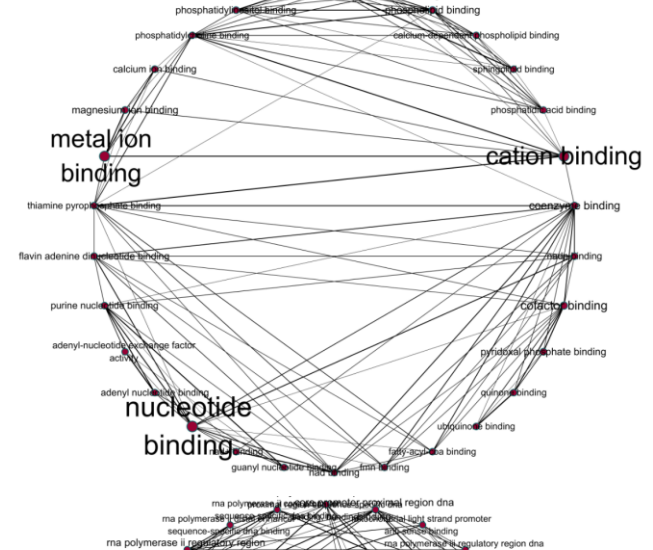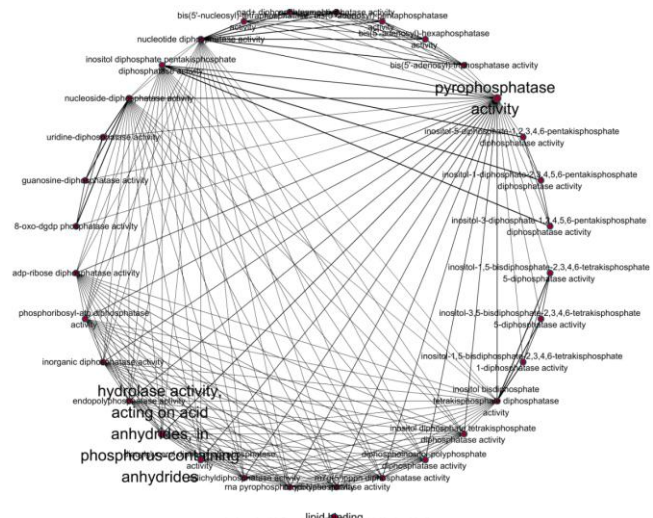
| | #annotated genes | HC | GeneMANIA | clusDCA |
|---|---|---|---|---|
| Human MF | 3-10 | 0.0153 | 0.0755 | **0.0818** * |
| | 11-30 | 0.0486 | 0.1344 | **0.1574** * |
| | 31-100 | 0.0677 | **0.1670** | 0.1625 |
| | 101-300 | 0.0637 | **0.2377** | 0.2151 |
| Human BP | 3-10 | 0.0098 | 0.0368 | **0.0429** * |
| | 11-30 | 0.0282 | 0.0673 | **0.0703** * |
| | 31-100 | 0.0384 | 0.1201 | **0.1228** |
| | 101-300 | 0.0475 | 0.1764 | **0.1849** * |
| Yeast MF | 3-10 | 0.0190 | 0.1075 | **0.1360** * |
| | 11-30 | 0.0776 | 0.1985 | **0.2334** * |
| | 31-100 | 0.0808 | **0.2751** | 0.2686 |
| | 101-300 | 0.0844 | 0.3348 | **0.3445** |
| Yeast BP | 3-10 | 0.0299 | 0.0954 | **0.1098** * |
| | 11-30 | 0.0857 | 0.1949 | **0.2217** * |
| | 31-100 | 0.1155 | **0.2939** | 0.2737 |
| | 101-300 | 0.1195 | **0.4321** | 0.3926 |
| mouseFunc MF | 3-10 | 0.0062 | 0.0965 | **0.1367** * |
| | 11-30 | 0.0451 | 0.1719 | **0.1950** * |
| | 31-100 | 0.0681 | 0.1866 | **0.1991** * |
| | 101-300 | 0.0651 | **0.2875** | 0.2814 |
| mouseFunc BP | 3-10 | 0.0074 | 0.0389 | **0.0516** * |
| | 11-30 | 0.0361 | 0.0680 | **0.0825** * |
| | 31-100 | 0.0394 | 0.0985 | **0.1088** * |
| | 101-300 | 0.0474 | 0.1350 | **0.1367** |

## 6 The complete list of clusters

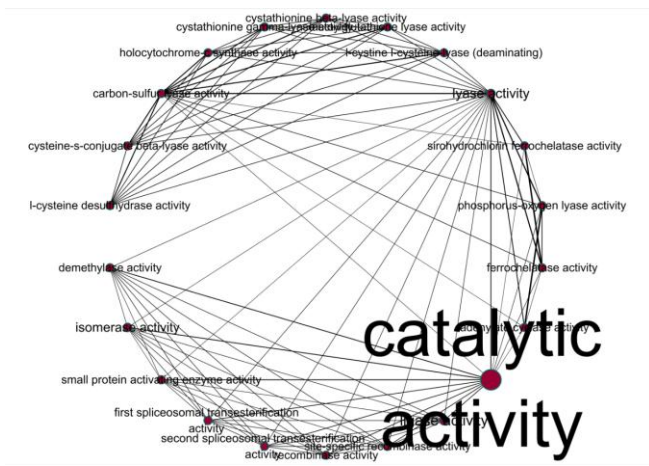We show the complete list of clusters here.

# Diffusion Component Analysis: Unraveling Functional Topology in Biological Networks

Hyunghoon Cho[1], Bonnie Berger[1,2,*] and Jian Peng[1,2,3,*]

[1] Computer Science and Artificial Intelligence Laboratory, MIT, Cambridge, MA, USA
[2] Department of Mathematics, MIT, Cambridge, MA, USA
[3] Department of Computer Science, University of Illinois at Urbana-Champaign, Champaign, IL, USA
* Corresponding authors: `bab@mit.edu, jianpeng@illinois.edu`

## 1  Introduction

Complex biological systems have been successfully modeled by biochemical and genetic interaction networks, typically gathered from high-throughput (HTP) data. These networks can be used to infer functional relationships between genes or proteins. Using t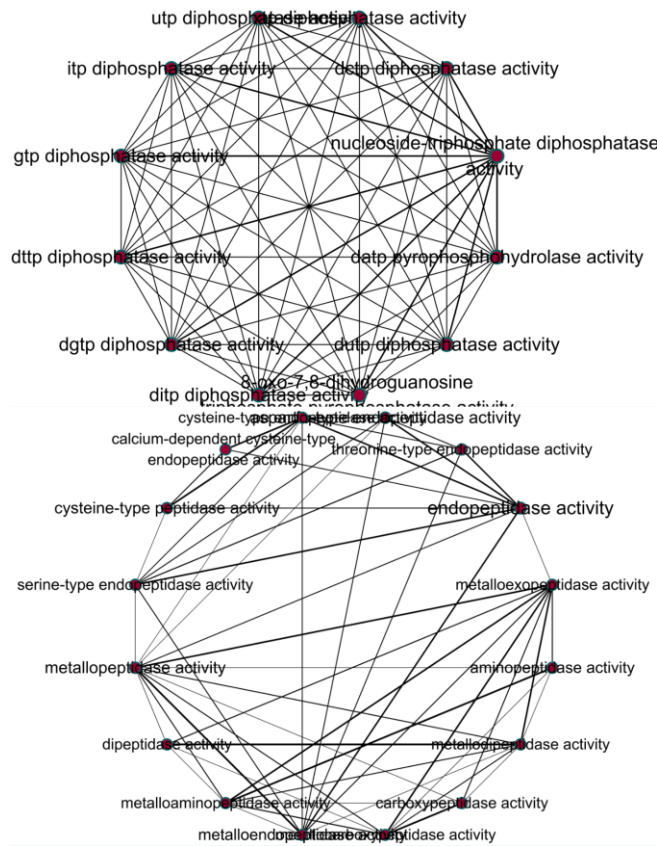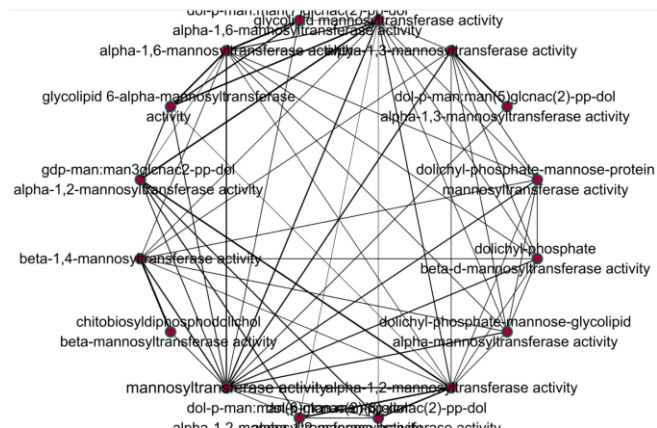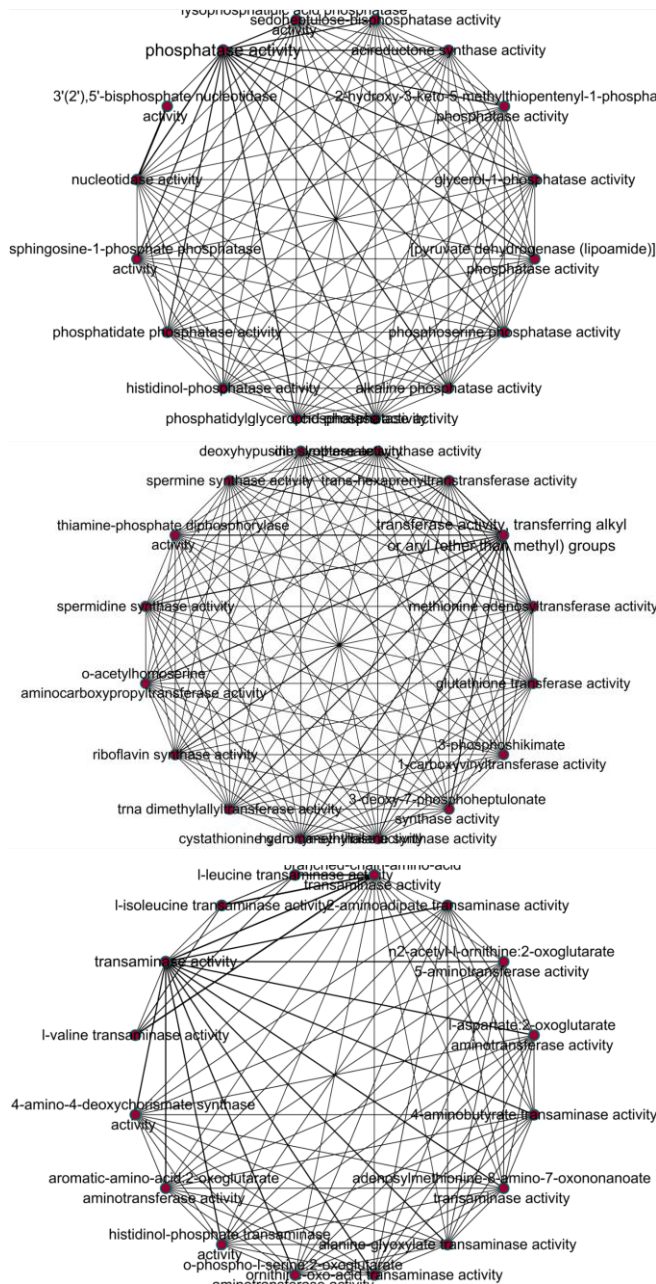he intuition that the topological role of a gene in a network relates to its biological function, local or diffusion-based "guilt-by-association" and graph-theoretic methods have had success in inferring gene functions [1, 2, 3]. Here we seek to improve function prediction by integrating diffusion-based methods with a novel dimensionality reduction technique to overcome the incomplete and noisy nature of network data.

A type of diffusion algorithm, also known as random walk with restart (RWR), has been extensively studied in the context of biological networks and effectively applied to protein function prediction (e.g., [1]). The key idea is to propagate information along the network, in order to exploit both direct and indirect linkages between genes. Typically, a distribution of topological similarity is computed for each gene, in relation to other genes in the network, so that researchers can select the most related genes in the resulting distribution or, rather, select genes that share the most similar distributions. Though successful, these approaches are susceptible to noise in the input networks due to the high dimensionality of the computed distributions.

## 2  Methods

We propose Diffusion Component Analysis (DCA), a novel analytical framework that combines diffusion-based methods and sophisticated dimensionality reduction to better extract topological network information in order to facilitate more accurate functional annotation of genes or proteins. The key idea behind DCA is to obtain informative, but low-dimensional features, which better encode the inherent topological properties of each node in the network. We first run a diffusion algorithm on a molecular network to obtain a distribution for each node that captures its relation to all other nodes in the network. We then approximate each of these distributions by constructing a multinomial logistic model, parameterized by low-dimensional feature vector(s), for each node. Feature vectors of all nodes are jointly learned by minimizing the Kullback-Leibler (KL) divergence (relative entropy) between the diffusion and parameterized-multinomial logistic distributions. A key differentiating factor of our novel dimensionality reduction from a more conventional approach, such as Principal Component Analysis (PCA), is the use of multinomial logistic models, which
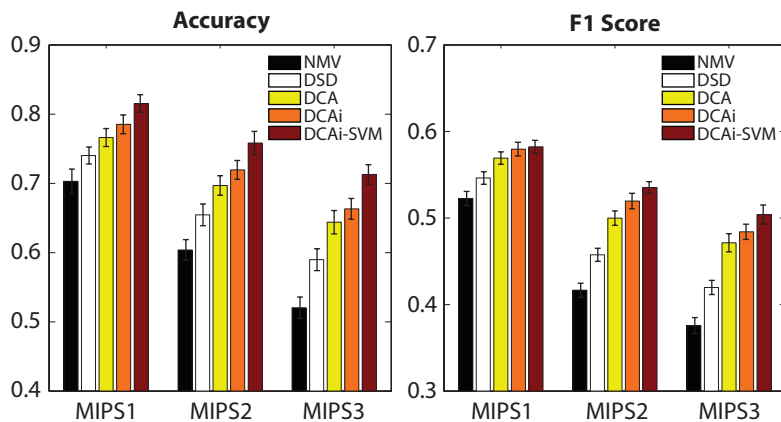
1

Figure 1: **Protein function prediction performance on yeast STRING networks in terms of both accuracy and F1 score–the harmonic mean of both precision and recall–with different levels of functional categories from MIPS.** Neighbor majority vote (NMV), Diffusion state distance (DSD), DCA with kNN (DCA), DCA combined with novel network integration with kNN (DCAi) or SVM (DCAi-SVM).

more naturally explain the input probability distributions from the diffusion. Moreover, DCA can be naturally extended to integrate multiple heterogeneous networks by performing diffusion on separate networks and jointly optimizing feature vectors. Given the low-dimensional vector representations of nodes, k-nearest neighbor (kNN) voting schemes or support vector machines (SVM) can be used for function prediction.

## 3   Results

We evaluated the ability of our DCA framework to uncover functional relationships in the interactome of yeast. By combining noise reduction via dimensionality reduction, improved integration of multiple heterogeneous networks (e.g., physical interaction, conserved co-expression), and the use of support vector machines, our DCA framework is able to achieve 71.29% accuracy with five-fold cross-validation on the STRING networks with third level functional annotations from MIPS, which is remarkably 12.31% higher than the previous state-of-the-art diffusion state distance (DSD) [1] method (Figure 1). We also observe improved performance over DSD in a different yeast PPI network, constructed from only physical interactions in the BioGRID database. In addition, we found that conventional approaches to dimensionality reduction, such as principal component analysis or non-negative matrix factorization, fail to achieve similar performance improvements. Our results demonstrate the potential of low-dimensional feature vectors learned by DCA to be plugged into other existing machine learning algorithms to decipher functional properties of and obtain novel insights into interactomes.

## References

1. Cao, Mengfei, et al. "New directions for diffusion-based network prediction of protein function: incorporating pathways with confidence." Bioinformatics 30.12 (2014): i219-i227.

2. Mostafavi, Sara, et al. "GeneMANIA: a real-time multiple association network integration algorithm for predicting gene function." Genome Biol 9.Suppl 1 (2008): S4.

3. Milenkovi, Tijana, and Nataa Prulj. "Uncovering biological network function via graphlet degree signatures." Cancer informatics 6 (2008): 257.