

Exploiting Ontology Graph for Predicting Sparsely Annotated Gene Function

Sheng Wang^{1,†}, Hyunghoon Cho^{2,†}, ChengXiang Zhai¹, 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

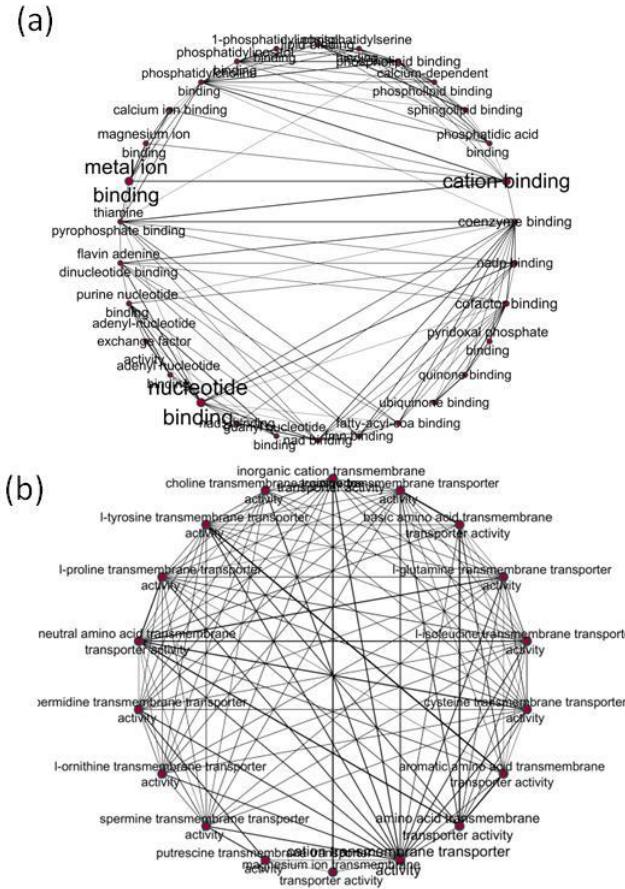
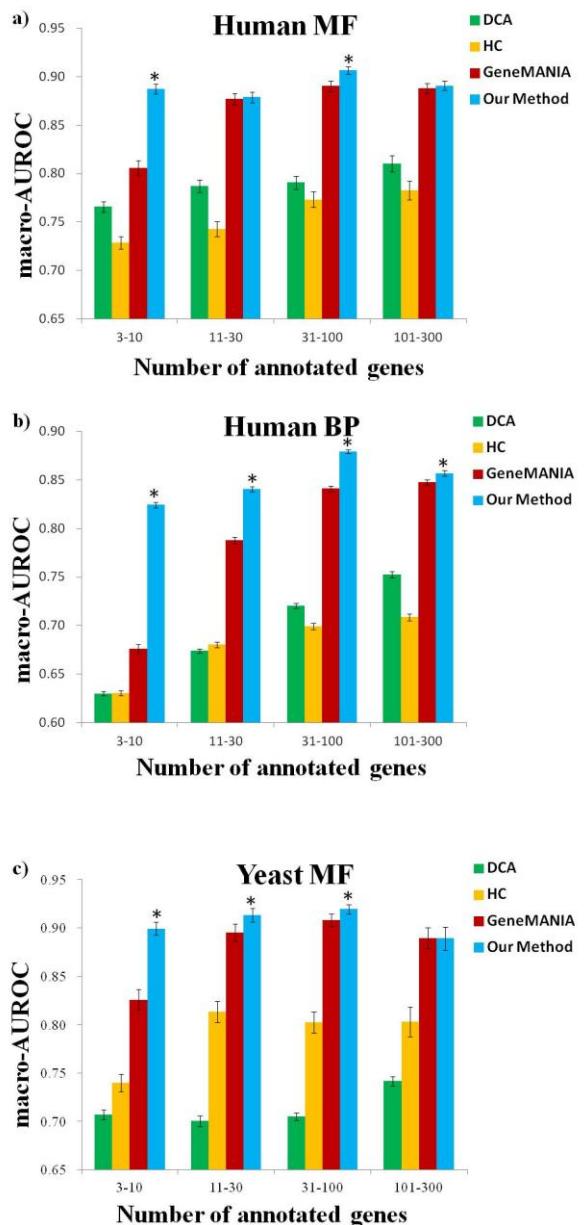


Figure S1: (a) The cluster structure of several binding functions. (b) The cluster structure of GO labels related to transmembrane transporting.

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

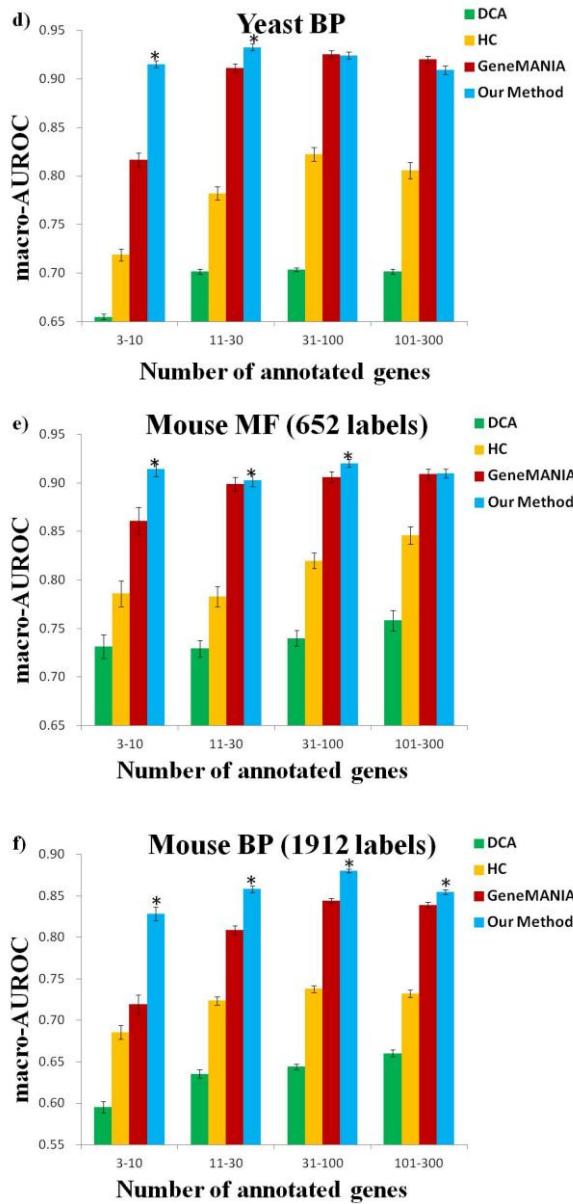


Figure S2: Comparison of our approach with other methods in terms of micro-AUROC. * indicates that our approach is statistically significant in comparison with GeneMANIA. Performance is evaluated for different subsets of GO labels with varying sparsity levels as shown on the x-axis.

3 Comparison of number of dimensions discussed in Section 3.7

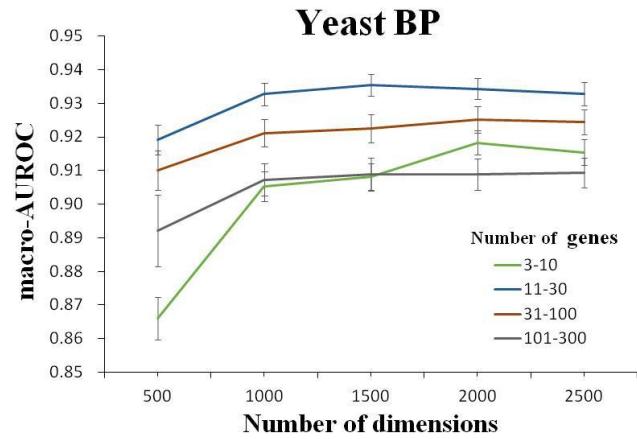
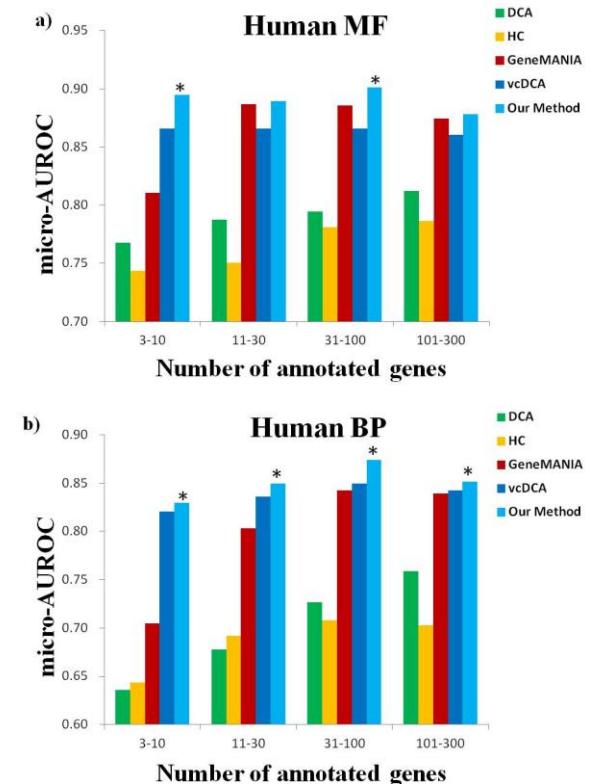


Figure S3: Comparison of number of dimensions in terms of macro-AUROC of biological process in yeast.

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 **Figure S4** and **Figure S5**.



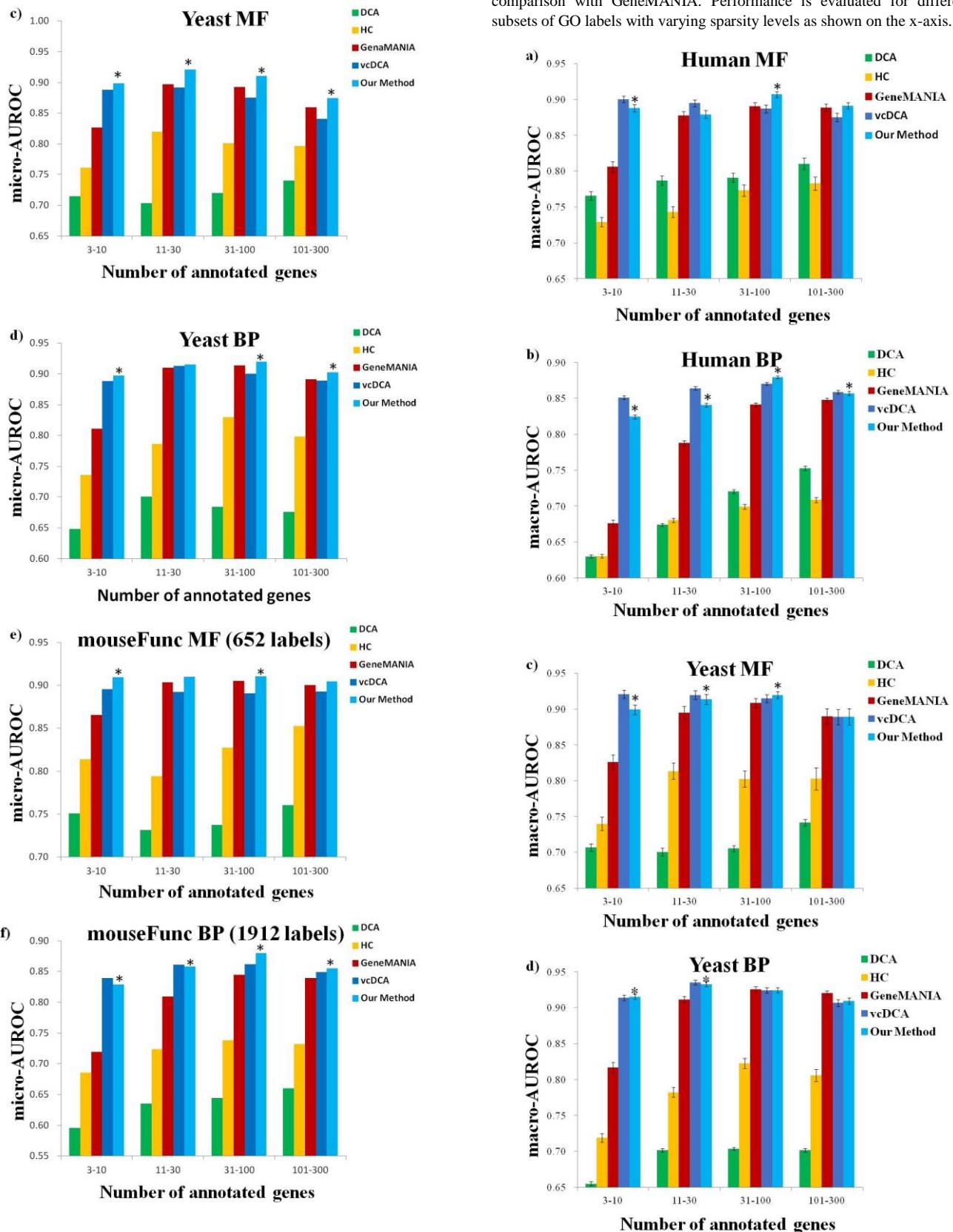


Figure S4: Comparison of our approach with other methods in terms of micro-AUROC. * indicates that our approach is statistically significant in

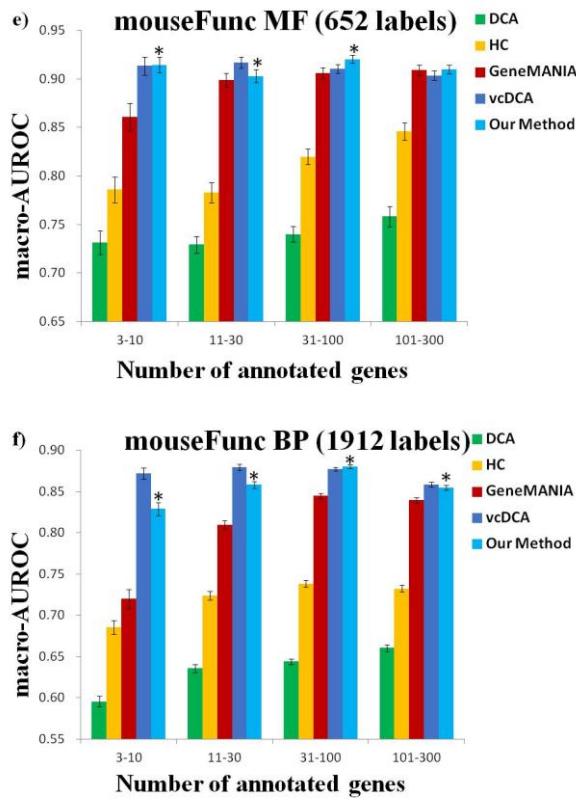


Figure S5: Comparison of our approach with other methods in terms of macro-AUROC. * indicates that our approach is statistically significant in comparison with GeneMANIA. Performance is evaluated for different subsets of GO labels with varying sparsity levels as shown on the x-axis.

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 0.1317 AUPRC on BP labels with 31-100 annotations, which is slightly higher than 0.1172 AUPRC for GeneMANIA. In yeast, our method achieved 0.3353 AUPRC on MF labels with 101-300 annotations, which is higher than 0.3198 AUPRC for GeneMANIA.

Table S1: Comparison of our approach with other methods in terms of AP@10.

	#annotated genes	HC	GeneMANIA	clusDCA
Human MF	3-10	0.0665	0.1265	0.1520
	11-30	0.1517	0.2227	0.2778
	31-100	0.2802	0.3307	0.3944
	101-300	0.2710	0.4642	0.5357
Human BP	3-10	0.0253	0.0541	0.0674
	11-30	0.0658	0.1169	0.1338
	31-100	0.1229	0.2457	0.3034
	101-300	0.1829	0.3860	0.4617

Yeast MF	3-10	0.0860	0.1187	0.2748
	11-30	0.2232	0.2521	0.3832
	31-100	0.3285	0.3911	0.5830
	101-300	0.3177	0.4451	0.6985
Yeast BP	3-10	0.0861	0.0973	0.2253
	11-30	0.2003	0.2312	0.3666
	31-100	0.3594	0.3837	0.5476
	101-300	0.4184	0.4953	0.6960
mouseFunc MF	3-10	0.0409	0.0993	0.2331
	11-30	0.1782	0.2130	0.3283
	31-100	0.3061	0.2637	0.4727
	101-300	0.3242	0.3629	0.6012
mouseFunc BP	3-10	0.0203	0.0406	0.0984
	11-30	0.0927	0.0958	0.1802
	31-100	0.1272	0.1678	0.3105
	101-300	0.1648	0.2504	0.3794

Table S2: Comparison of our approach with other methods in terms of macro-AUPRC.

	#annotated genes	HC	GeneMANIA	clusDCA
Human MF	3-10	0.0286	0.0882	0.0988
	11-30	0.0527	0.1430	0.1671
	31-100	0.0602	0.1678	0.1690
	101-300	0.0603	0.2367	0.2133
Human BP	3-10	0.0118	0.0368	0.0412
	11-30	0.0285	0.0662	0.0698
	31-100	0.0366	0.1172	0.1317
	101-300	0.0436	0.1764	0.1721
Yeast MF	3-10	0.0403	0.0946	0.1850
	11-30	0.0740	0.1768	0.2366
	31-100	0.0730	0.2466	0.2721
	101-300	0.0778	0.3198	0.3354
Yeast BP	3-10	0.0412	0.0760	0.1405
	11-30	0.0867	0.1634	0.2208
	31-100	0.1075	0.2610	0.2791
	101-300	0.1158	0.4010	0.3805
mouseFunc MF	3-10	0.0150	0.0804	0.1535
	11-30	0.0534	0.1519	0.1950
	31-100	0.0567	0.1664	0.2129
	101-300	0.0596	0.2738	0.2663
mouseFunc BP	3-10	0.0085	0.0313	0.0615
	11-30	0.0323	0.0564	0.0865
	31-100	0.0345	0.0860	0.1120
	101-300	0.0437	0.1270	0.1328

6 The complete list of clusters

We show the complete list of clusters here.

