

基于蛋白质网络功能模块的蛋白质功能预测*

卢宏超^{1,2,4)} 石秋艳³⁾ 石宝晨^{1,4)} 张治华^{1,4)}
赵屹²⁾ 唐素勤³⁾ 熊磊³⁾ 王强^{3)**} 陈润生^{1,2)**}

¹⁾中国科学院生物物理研究所, 北京 100101; ²⁾中国科学院计算技术研究所, 北京 100080;

³⁾广西师范大学物理与信息工程学院, 桂林 541004; ⁴⁾中国科学院研究生院, 北京 100049)

摘要 在破译了基因序列的后基因组时代, 随着系统生物学实验的快速发展, 产生了大量的蛋白质相互作用数据, 利用这些数据寻找功能模块及预测蛋白质功能在功能基因组研究中具有重要意义. 打破了传统的基于蛋白质间相似度的聚类模式, 直接从蛋白质功能团的角度出发, 考虑功能团间的一阶和二阶相互作用, 提出了模块化聚类方法 (MCM), 对实验数据进行聚类分析, 来预测模块内未知蛋白质的功能. 通过超几何分布 P 值法和增、删、改相互作用的方法对聚类结果进行预测能力分析和稳定性分析. 结果表明, 模块化聚类方法具有较高的预测准确度和覆盖率, 有很好的容错性和稳定性. 此外, 模块化聚类分析得到了一些具有高预测准确度的未知蛋白质的预测结果, 将会对生物实验有指导意义, 其算法对其他具有相似结构的网络也具有普遍意义.

关键词 蛋白质相互作用网络, 蛋白质功能预测, 聚类

学科分类号 Q811.4

在破译了基因序列的后基因组时代, 功能基因组学的一个目标是研究已破译基因的生物功能并控制它们. 生命体不能通过孤立的分子实现细胞功能, 而是必须通过分子之间的相互作用才能实现. 因而, 分析生命分子间的相互作用自然可能得到基因的生物功能更准确、详尽的信息, 而近年来多种实验方法产生的大规模蛋白质相互作用数据为揭示这些基因的功能提供了可能. 在这些数据中, 特别以酵母的蛋白质相互作用网络数据最为完备^[1-4]. 本文拟以酵母的蛋白质相互作用网络为研究对象, 采用新的模块化聚类方法, 对酵母的基因功能进行研究.

从大规模数据产生到现在, 对蛋白质网络的研究已经越来越得到大家的重视, 统计“投票”法^[5]、全局预测法^[6]、谱方法^[7]、概率方法^[8]、马尔可夫随机场方法^[9]、信息传递算法^[10]等都被用来寻找功能模块和预测蛋白质功能, 此外还有聚类方法^[11-16]. 目前, 聚类方法已成为寻找功能模块和预测蛋白质功能一类主要的方法, 其中采用最多的是层次聚类. Rives 等^[12]较早地对蛋白质相互作用网络进行了研究, 用聚类方法来寻找网络功能模块, 其基本假设就是近邻相互作用的蛋白质具有相似的

功能. Brun 等^[13]和 Samanta 等^[14]都利用具有相互作用蛋白质对的共同邻居定义它们之间的相似度, 进而采用聚类方法预测未知蛋白质的功能. Brun 等采用的是 Czekanowski-Dice 距离法, 而 Samanta 等则采用概率模型进行分析, 他们都得到了不错的结果. 我们研究组针对蛋白质网络的矩阵可视化的需要提出了基于链接矩阵的聚类算法^[16], 能够很好地反映网络的拓扑性质和生物特性, 并可用于蛋白质功能的预测.

现有的层次聚类算法都是用蛋白质相互作用网络的数据信息定义蛋白质间的距离或相似度, 并在此基础上进行聚类分析. 本文打破了这种设计模式, 直接从功能团的定义出发, 汲取 Brun 等和 Samanta 等方法的优点, 考虑功能团之间一阶和二阶相互作用, 定义功能模块间的相似度进行聚类分析, 进而

*国家自然科学基金资助项目(30500104 和 30570393).

** 通讯联系人.

陈润生. Tel: 010-64888543, E-mail: crs@sun5.ibp.ac.cn

王强. Tel: 0773-5838495, E-mail: qwang@mailbox.gxnu.edu.cn

收稿日期: 2005-11-04, 接受日期: 2005-12-30

预测模块内未知蛋白质的功能. 与以往的聚类方法相比较, 具有更好的预测效果. 该方法还可以应用于蛋白质代谢网、因特网、人际关系网等具有类似结构的网络, 对研究一般性网络具有普遍意义.

1 材料和方法

1.1 材料

本文主要对芽殖酵母的蛋白质相互作用网络进行分析, 其蛋白质相互作用的数据来源于大规模实验数据, 这些实验包括酵母双杂交实验, HMS-PCI 和 TAP 方法, 主要收集自 MIPS^[17] (<http://mips.gsf.de/>), PreBIND^[18] (<http://bind.ca/index2.phtml?site=prebind>), BIND^[19] (<http://bind.ca/>) 和 GRID^[20] (<http://biodata.mshri.on.ca/grid/servlet/Index>) 数据库. 在数据预处理阶段, 我们去除了自相互作用和冗余的相互作用, 对 HMS-PCI 和 TAP 方法产生的数据, 我们只考虑了“spoke-like”的数据^[21, 22], 即只对钓饵和被钓蛋白的直接链接建立相互作用. 本蛋白质相互作用数据集包括 4 537 个酵母蛋白的 13 344 对相互作用.

本文用于注释蛋白质功能的数据来自 MIPS 数据库. 这里的功能注释中, 每个蛋白质可能具有不止一种功能. 本文的功能注释条目有 12 397 条, 未知蛋白质有 999 个.

1.2 算法

蛋白质相互作用网络可以被表示成为一个无向图, 在这个图中, 每个节点表示一个蛋白质, 而每条边表示一个相互作用. 基于这个图可以定义两个功能团之间的相似度, 表示两个功能团在生物功能上的相似程度. 有了团之间的相似度的定义就可以进行层次聚类, 从每个点自成一团的初始状态出发, 采用聚集式的贪婪算法, 每次在所有团中找到最相似的两个, 把它们归并到一起, 组成新团, 在此基础上继续聚类, 直到所有团归并到一起, 在归并的过程中形成一棵聚类树.

两个团的蛋白质间具有直接相互作用称为具有一阶相互作用, 而两团与团外的某蛋白质都有直接相互作用称两团具有二阶相互作用, 功能团间的相似度定义就是基于两团间的一阶和二阶相互作用. 现将与两个功能团 C_1 和 C_2 都有相互作用的点归于集合 C_0 里, 如图 1 所示, 则团 C_1 和团 C_2 的相似度可以定义为:

$$S = \frac{m + f(S_1 + S_2)}{O_1 + O_2 - m}, \text{ 其中 } f = \frac{|C_0|}{|C_0| + |C_1| + |C_2|}$$

这里 m 是两个团 C_1 和团 C_2 间的一阶相互作用的个数, S_1 和 S_2 分别是团 C_1 和团 C_2 与 C_0 链接的相互作用个数, O_1 和 O_2 分别表示团 C_1 和团 C_2 的所有向外相互作用的总数, f 是调节二阶相互作用的因子, 随聚类团的增大而降低, $|C_1|$ 、 $|C_2|$ 和 $|C_0|$ 分别是对应团的大小.

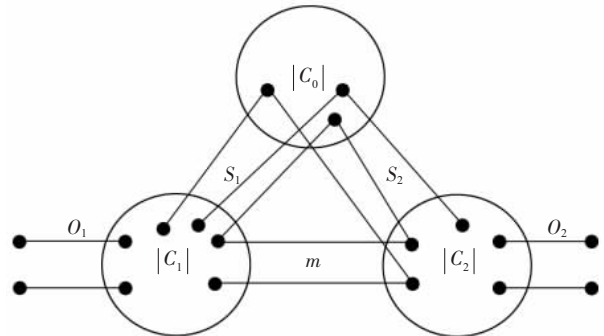


Fig. 1 The definition of the proximity of two functional clusters

C_1 and C_2 are the functional clusters, C_0 is the set of proteins which directly interact with two functional clusters.

经过实验比较, 这里的相似度定义没有像 Brun 等方法中采用类 Czekanowski-Dice 距离方式, 而采用类 Tanimono 相似度, 其预测能力稍好.

1.3 功能团的评测参数 P 值

对于聚类成的蛋白质团可能含有多种蛋白质功能, 用 P 值的方法可以给出每个功能相对于随机取团的几率, 进而把随机取团几率最小的功能赋给该团, 作为该团的主要功能, 而主要功能的 P 值就可以说明聚类团与生物功能的吻合程度.

功能团的 P 值是通过超几何聚集分布 (hypergeometric cumulative) 给出的^[23, 24], 对于蛋白质数为 n 的团, 含有某功能蛋白质 k 个, 设其所在的蛋白质组共有 G 个蛋白质, 该功能类共有 C 个蛋白质, 这样的团随机出现几率的 P 值是:

$$P = 1 - \sum_{i=0}^{k-1} \frac{\binom{C}{i} \binom{G-C}{n-i}}{\binom{G}{n}}$$

上面的 P 值用几率值表示了蛋白质团对于某个功能的富集程度, 如果 P 值越小, 越接近 0, 则说明蛋白质团能够随机出现这种功能的几率就越低, 当然可能更有生物学意义. 用蛋白质团最小 P 值的蛋白质功能给团定义主要功能, 进而可以评测理论

得出的蛋白质团的价值, 以衡量聚类算法的好坏. 由于这个值可能很小, 实际中我们往往用 P 值的负对数值来表示这个值.

2 结 果

我们把模块化聚类方法应用于芽殖酵母的蛋白质相互作用网络, 聚类的结果可以在我们的网站上下载 (<http://www.bioinfo.org.cn/clustering>), 该结果可以通过 TreeView [25] 软件和我们组开发的软件 PINC [16] 进行分析, 利用聚类的结果结合 MIPS 酵母蛋白质的功能注释信息, 我们可以进行进一步的功能模块分析和未知蛋白质功能预测.

为了考察我们聚类方法的预测能力, 我们采用了 Wu 等 [23] 对表达谱聚类分析文章中采用的分析方法, 对于聚类好的团利用 P 值把已知功能看成未知进行预测, 来考察预测的准确度和可预测功能覆盖

率的好坏. 具体做法是对一系列具有不同 P 值的团计算出其对应的预测准确度和覆盖率并作图, 如图 2 所示. 利用 P 值进行功能预测, 准确度和覆盖率是一对互相影响的因素, 高的覆盖率对应的必然是低的预测准确度, 而高的准确度必然导致低的覆盖率, 所以选择适当的 P 值才能保证较高的准确度和覆盖率. 利用 P 值进行预测, 影响准确度量度的因素还有功能团的大小. 一般来说同样的 P 值, 较小的团给出的预测其准确度会稍高一些, 如图 2a 所示, 我们利用 P 值采用了团的大小分别小于等于 50、100、150、200 和 250 的团对蛋白质的功能进行预测, 通过已知功能蛋白质给出了预测的准确度和覆盖率, 图中每个点表示一种 P 值的结果, P 值从小于 10^{-3} 一直到小于 10^{-13} 每个数量级一个点, 图中表明小于等于 150 的团能折衷地给出较好的预测准确度和覆盖率.

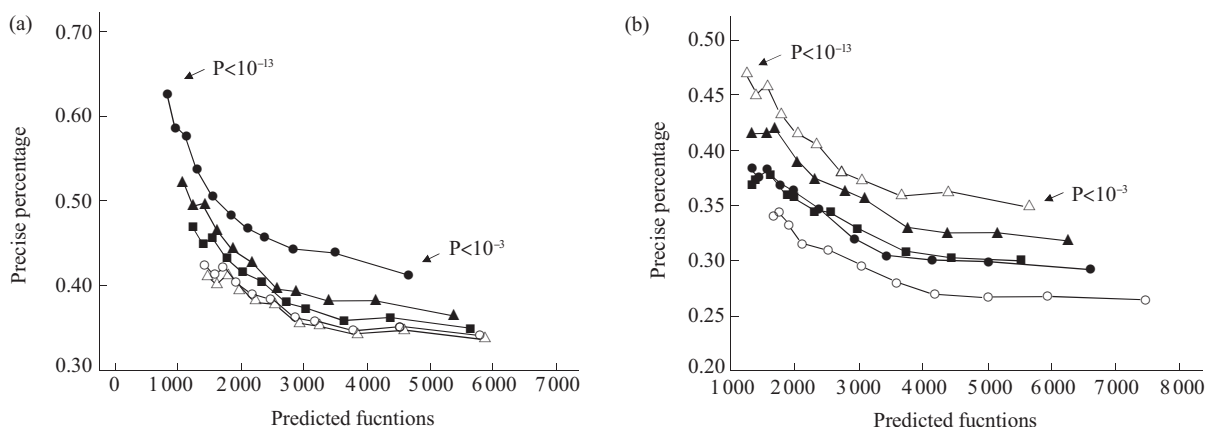


Fig. 2 The predicted precise percentage and coverage for the clustering methods

(a) For MCM clustering with different cut size. ●—●: Size ≤ 50 ; ▲—▲: Size ≤ 100 ; ■—■: Size ≤ 150 ; ○—○: Size ≤ 200 ; △—△: Size ≤ 250 . (b) Comparing with other clustering methods. ●—●: RIVES; ▲—▲: BRUN; ■—■: SAMANTA; ○—○: ADJW; △—△: MCM.

比较我们的模块化聚类方法和以前的其他层次聚类算法, 可以考察其预测能力. 把 Brun 等 [13]、Rives 等 [12]、Samanta 等 [15], ADJW 算法 [16] 和模块化聚类算法 (MCM), 都采用小于等于 150 的团对蛋白质功能进行预测, 图 2b 给出了每种方法不同 P 值下的预测准确度和覆盖率的曲线, 其作图方法同图 2a, 从中可以看出, 模块化方法具有较高的预测准确度和覆盖率.

利用 P 值可以对所有未知蛋白质的功能进行预测, 其预测的结果可以从我们的网站上得到, 这里

仅给出一些高准确度的预测结果, 可供生物实验参考. 考虑大小小于等于 50 而 P 值小于 10^{-13} 的团进行功能预测, 从图 2a 可以看出, 这样的团给出的预测准确度至少有 60%, 预测的结果见表 1.

表 1 列出了 25 个团对 84 个未知蛋白质的 42 种功能的预测结果, 给出了每个团对每个预测功能的 P 值来衡量预测的准确度、对应团该功能的蛋白质数目、未知蛋白质个数和其他蛋白质个数, 从表 1 中主要功能的富集程度, 我们可以看出预测可能具有高度的准确度.

Table 1 The predicted function for unknown function annotated from MIPS

<i>N</i>	<i>P</i>	<i>S</i>	<i>M</i>	<i>U</i>	<i>O</i>	Predicted function	Unknown proteins
1	22.0	44	23	1	20	RNA binding	YFR012W
	30.9	44	28	1	15	rRNA processing	
2	13.9	13	9	3	1	Modification by ubiquitination, deubiquitination	YJL149W,YLR352W,YOL087C
3	13.5	44	17	1	26	RNA binding	YMR125C
	70.0	44	43	1	0	Splicing	
4	13.3	34	11	3	20	Protein processing (proteolytic)	YMR155W,YDR126W,YFL006W
5	13.0	49	12	8	29	Structural protein	YLL049W,YKL061W,YLR031W,YBL031W, YPL056C,YFL068W,YBR144C,YJR134C
	13.3	49	20	8	21	Protein targeting, sorting and translocation	
	15.1	49	11	8	30	Nuclear membrane	
	19.6	49	18	8	23	RNA transport	
	23.8	49	20	8	21	Nuclear transport	
6	14.8	28	11	3	14	Small GTPase mediated signal transduction	YBL053W,YFR044C,YIL163C
7	16.5	40	22	5	13	Budding, cell polarity and filament formation	YBL053W,YFR044C,YIL163C,YDL129W, YOR081C
8	14.8	18	11	1	6	Actin cytoskeleton	YFR024C
9	15.1	27	10	5	12	Vesicle fusion	YPL095C,YJL151C,YMR253C,YLR064W, YHR105W
10	14.2	21	9	1	11	Extension/ Polymerization activity	YLR416C
11	20.4	43	22	4	17	DNA conformation modification (e.g. chromatin)	YER092W,YOR141C,YMR075W,YDR070C
12	16.2	23	15	1	7	DNA conformation modification (e.g. chromatin)	YOR225W
	17.6	23	11	1	11	Transcription initiation	
	18.7	23	10	1	12	G1 phase of mitotic cell cycle	
	21.3	23	15	1	7	Organization of chromosome structure	
	26.6	23	16	1	6	Modification by acetylation, deacetylation	
13	29.2	45	26	1	18	General transcription activities	YOR225W
	41.7	45	41	1	3	Transcriptional control	
14	14.0	28	9	4	15	MAPKKK cascade	YAR061W,YGR115C,YPR076W,YDR239C
15	20.2	47	16	14	17	rRNA synthesis	YAR064W,YBR089W,YJR162C,YBR157C, YDL183C,YDR067C,YDR199W,YDR286C, YDR340W,YER108C,YDR056C,YOL003C, YER181C,YJL182C
	23.2	47	16	14	17	tRNA synthesis	
	24.4	47	24	14	9	DNA binding	
16	13.7	40	8	4	28	Transcription termination	YDL218W,YHR100C,YOR227W,YPL009C
	24.8	40	14	4	22	3'-end Processing	
17	14.1	30	9	8	13	ori Recognition and priming complex formation	YDR344C,YKL047W,YBR047W,YLR331C, YLR269C,YCR022C,YCR050C,YFR057W
18	26.4	43	24	6	13	Mitochondrion	YOR205C,YGR150C,YCR072C,YHR059W, YOR333C,YOR331C
	24.0	43	24	6	13	Ribosomal proteins	
19	13.1	10	6	1	3	Peroxisomal transport	YDL139C
20	14.4	13	7	3	3	Regulation of glycolysis and gluconeogenesis	YBL049W,YMR157C,YNL300W
21	34.4	36	28	1	7	Ribosomal proteins	YGR033C
	39.5	36	29	1	6	Mitochondrion	
22	13.1	19	11	5	3	Proteasomal degradation (ubiquitin/proteasomal pathway)	YDR482C,YEL075C,YGL046W,YGL069C, YNL176C
	13.8	19	10	5	4	Modification by ubiquitination, deubiquitination	
23	15.8	25	11	7	7	Mitosis	YDR482C,YEL075C,YGL046W,YGL069C, YNL176C,YLR253W,YBR063C
24	15.5	21	11	2	8	ER to Golgi transport	YAL034C,YJR116W
25	13.1	7	6	1	0	ori Recognition and priming complex formation	YBL071C

N: The No. of the modules; *P*: *P* Value of the module for the function. ($-\lg P$); *S*: The size of module; *M*: Proteins with main function; *U*: Unknown proteins; *O*: Other proteins.

由于蛋白质网络具有较高的假阳性, 好的聚类算法应当具有很好的容错性, 我们对模块化聚类算法进行了稳定性分析. 我们人为地加入假阳性对聚类结果进行分析, 在保持蛋白质网络度分布的前提下分别随机增加(added)、删除(removed)和修改(rewired) 10%到 100%的相互作用, 按 10 个百分点间隔, 每个百分点随机分别对三种情况产生 10 组数据, 对这些数据分别进行聚类分析, 计算每个结果中 P 值小于 10^{-3} 的团的个数, 对每种情况的每个百分点计算统计其平均, 结果见图 3. 可见对于添加和删除 50%的相互作用, 仍有大于 70%的团有意义, 而修改 50%的相互作用, 仍大约有 40%的团有预测能力, 因此, 我们可以看到本聚类算法具有很好的容错性.

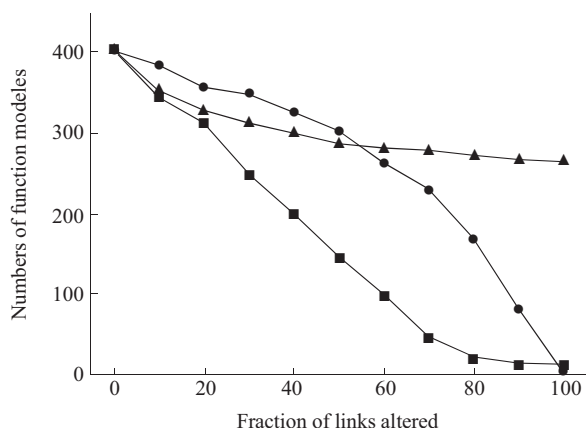


Fig. 3 Robustness test of the MCM method

To investigate the extent to which false positives affected the clustering tree, we randomly remove, add and rewire 10% ~100% of the interactions without changing the degree distribution in the PPI network. Each perturbation is repeated 10 times. ●—●: removed; ▲—▲: added; ■—■: rewired.

3 讨 论

在蛋白质相互作用网络功能研究中, 我们提出的模块化聚类方法, 汲取了前人在层次聚类算法中利用蛋白质相互作用邻居进行功能分析的好的生物前提, 打破了利用蛋白质之间距离和相似度进行聚类分析的传统思路, 提出了基于功能模块聚类的新思路. Harwell 等^[26]在 1999 年提出蛋白质功能模块化的思想, 具体表现在功能模块之间连接疏松而模块内部连接紧密, 这种聚类方法正是这种模块化思想的具体体现. 从聚类结果上看, 本聚类算法在预测未知蛋白质功能方面较以前的算法具有更好的预测能力, 并具有很好的容错能力和稳定性.

现在的聚类算法仍采用贪婪算法进行聚集式聚类, 该方法在聚类的初始阶段同 Brun 等的算法很相似, 从聚类的结果上我们也能看出(图 2), 对于较大的功能团, 由于聚类直接从功能团出发, 能给出更好的结果. 从计算复杂度上看, 该模块化聚类仅在聚类的递推公式上与传统层次聚类不同, 而复杂度没有增加. 这种基于功能团的定义可以自然地推广到基于优化的聚类方法上, 可能得到更加符合真实情况的功能团的分析结果. 由于本方法是基于图论的一般性方法, 对于其他的具有类似结构的一般网络如人际关系网、因特网和 Web 网等具有借鉴意义.

参 考 文 献

- Ito T, Chiba T, Ozawa R, *et al.* A comprehensive two-hybrid analysis to explore the yeast protein interactome. *Proc Natl Acad Sci USA*, 2001, **98** (8): 4569~4574
- Uetz P, Giot L, Cagney G, *et al.* A comprehensive analysis of protein-protein interactions in *Saccharomyces cerevisiae*. *Nature*, 2000, **403** (6770): 623~627
- Ho Y, Gruhler A, Heilbut A, *et al.* Systematic identification of protein complexes in *Saccharomyces cerevisiae* by mass spectrometry. *Nature*, 2002, **415** (6868): 180~183
- Gavin A C, Bosche M, Krause R, *et al.* Functional organization of the yeast proteome by systematic analysis of protein complexes. *Nature*, 2002, **415** (6868): 141~147
- Schwikowski B, Uetz P, Fields S. A network of protein-protein interactions in yeast. *Nat Biotechnol*, 2000, **18** (12): 1257~1261
- Vazquez A, Flammini A, Maritan A, *et al.* Global protein function prediction from protein-protein interaction networks. *Nat Biotechnol*, 2003, **21** (6): 697~700
- Bu D, Zhao Y, Cai L, *et al.* Topological structure analysis of the protein-protein interaction network in budding yeast. *Nucleic Acids Res*, 2003, **31** (9): 2443~2450
- Letovsky S, Kasif S. Predicting protein function from protein/protein interaction data: a probabilistic approach. *Bioinformatics*, 2003, **19** (Suppl 1): i197~204
- Deng M, Tu Z, Sun F, *et al.* Mapping Gene Ontology to proteins based on protein-protein interaction data. *Bioinformatics*, 2004, **20** (6): 895~902
- Leone M, Pagnani A. Predicting protein functions with message passing algorithms. *Bioinformatics*, 2005, **21** (2): 239~247
- Spirin V, Mirny L A. Protein complexes and functional modules in molecular networks. *Proc Natl Acad Sci USA*, 2003, **100** (21): 12123~12128
- Rives A W, Galitski T. Modular organization of cellular networks. *Proc Natl Acad Sci USA*, 2003, **100** (3): 1128~1133
- Brun C, Chevenet F, Martin D, *et al.* Functional classification of proteins for the prediction of cellular function from a protein-protein interaction network. *Genome Biol*, 2003, **5** (1): R6

- 14 Brun C, Herrmann C, Guenoche A. Clustering proteins from interaction networks for the prediction of cellular functions. *BMC Bioinformatics*, 2004, **5** (1): 95
- 15 Samanta M P, Liang S. Predicting protein functions from redundancies in large-scale protein interaction networks. *Proc Natl Acad Sci USA*, 2003, **100** (22): 12579~12583
- 16 Lu H, Zhu X, Liu H, *et al.* The interactome as a tree--an attempt to visualize the protein-protein interaction network in yeast. *Nucleic Acids Res*, 2004, **32** (16): 4804~4811
- 17 Mewes H W, Amid C, Arnold R, *et al.* MIPS: analysis and annotation of proteins from whole genomes. *Nucleic Acids Res*, 2004, **32** (Database issue): D41~44
- 18 Donaldson I, Martin J, De Bruijn B, *et al.* PreBIND and Textomy--mining the biomedical literature for protein-protein interactions using a support vector machine. *BMC Bioinformatics*, 2003, **4** (1): 11
- 19 Bader G D, Betel D, Hogue C W. BIND: the biomolecular interaction network database. *Nucleic Acids Res*, 2003, **31** (1): 248~250
- 20 Breitkreutz B J, Stark C, Tyers M. The GRID: the general repository for interaction datasets. *Genome Biol*, 2003, **4** (3): R23
- 21 Bader G D, Hogue C W. Analyzing yeast protein-protein interaction data obtained from different sources. *Nat Biotechnol*, 2002, **20** (10): 991~997
- 22 Von Mering C, Krause R, Snel B, *et al.* Comparative assessment of large-scale data sets of protein-protein interactions. *Nature*, 2002, **417** (6887): 399~403
- 23 Wu L F, Hughes T R, Davierwala a P, *et al.* Large-scale prediction of *Saccharomyces cerevisiae* gene function using overlapping transcriptional clusters. *Nat Genet*, 2002, **31** (3): 255~265
- 24 Tavazoie S, Hughes J D, Campbell M J, *et al.* Systematic determination of genetic network architecture. *Nat Genet*, 1999, **22** (3): 281~285
- 25 Eisen M B, Spellman P T, Brown P O, *et al.* Cluster analysis and display of genome-wide expression patterns. *Proc Natl Acad Sci USA*, 1998, **95** (25): 14863~14868
- 26 Hartwell L H, Hopfield J J, Leibler S, *et al.* From molecular to modular cell biology. *Nature*, 1999, **402** (6761 Suppl): C47~52

Predicting Protein Function Based on Modularized Protein Interaction Network*

LU Hong-Chao^{1,2,4}, SHI Qiu-Yan³, SHI Bao-Chen^{1,4}, ZHANG Zhi-Hua^{1,4}, ZHAO Yi², TANG Su-Qin³, XIONG Lei³, WANG Qiang^{3**}, CHEN Run-Sheng^{1,2**}

¹*Institute of Biophysics, The Chinese Academy of Sciences, Beijing 100101, China;*

²*Institute of Computing Technology, The Chinese Academy of Sciences, Beijing 100080, China;*

³*College of Physics and Information Engineering, Guangxi Normal University, Guilin 541004, China;*

⁴*Graduate University of The Chinese Academy of Sciences, Beijing 100049, China)*

Abstract In the post-genomics era in which gene sequences have been decoded, large-scale protein-protein interaction data are generated with the rapid development of system biology experiments. It is important in functional genomics to search for function modules and predict protein functions from the data. A new method called modularized clustering method (MCM), which are based on the direct and second-order interactions of modules, is applied to the latest high-throughput protein-protein network of yeast to predict the function of unknown proteins in the modules. *P* value of hypergeometric cumulative distribution of modules and the disturbance analysis on the data, including adding, removing and rewiring interactions, are employed to evaluate the prediction quality and robustness of the method. The results show that MCM has high prediction precise rate and coverage, and it is robust to high false-positive data and missing data. The predicted results of unknown proteins with high prediction precise rate can be instructive in biological analysis and the algorithm can be generalized to other networks with the similar structures.

Key words protein-protein interaction network, predicting function of proteins, clustering

*This work was supported by grants from The National Natural Science Foundation of China (30500104, 30570393).

**Corresponding author .

CHEN Run-Sheng. Tel: 86-10-64888543, E-mail: crs@sun5.ibp.ac.cn

WANG Qiang. Tel: 86-773-5838495, E-mail: qwang@mailbox.gxnu.edu.cn

Received: November 4, 2005 Accepted: December 30, 2005