

GOODIES: GO Based Data Mining Tool for Characteristic Attribute Interpretation on a Group of Biological Entities

Sung Geun Lee¹
sglee@istech21.com

Wan Seon Lee¹
konan@istech21.com

Yang Seok Kim^{1,2}
yskim@iste21.com

¹ Bioinformatics Unit, ISTECH Inc. #704, Hyundai Town Vill, 848-1 Janghang-dong, Ilsan-gu, Goyang city, Gyeonggi-do, 411-380, Korea

² Cancer Metastasis Research Center, Yonsei University College of Medicine, 134 Shinchon-dong, Seodaemun-gu, Seoul, 120-752, Korea

Keywords: Gene Ontology, GO tree, MaxPd, AverPd

1 Introduction

GOODIESTM is a Gene Ontology (GO) based data-mining tool with intuitive visualization on a GO tree. Its algorithm uses the graph structure of GO to interpret and classify aggregates of biological entities [3]. Given gene or protein lists, e.g. gene clusters obtained from DNA chip experiments, GOODIES takes the multiple functionalities of genes into account and computationally selects the optimal GO candidate terms for the most suitable biological interpretation of given lists.

2 Method and Results

The main usage of GOODIES can be as follows: biologically-oriented cluster analysis of DNA microarray data, automated functional annotation via clustering, and functional categorization of biological objects. First, in cluster analysis of DNA microarrays, biologists primarily want to know how well clusters of expression profiles are associated with known functional categories and cellular processes. GOODIES can perform such tasks in terms of GO that it can be complementary to statistical clustering methods. Secondly, the unknown function of genes can be putatively predicted through the clustering interpretation of GOODIES.

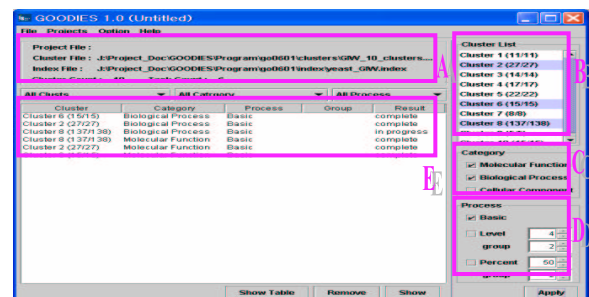
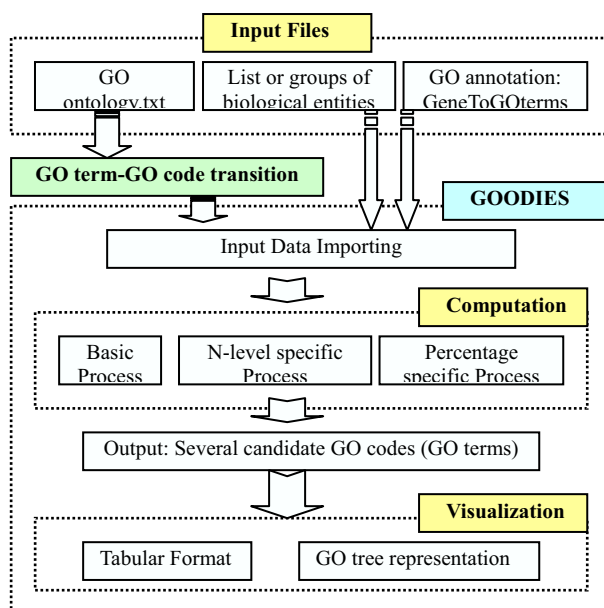


Figure 1: Schematic block diagram of GOODIES (left) and Sample analysis (right). Once GOODIES completes the matching process between the input groups and corresponding information from the GO annotation file, main window in A and B will be filled. After executing selected clusters(B), categories(C), and processes(D), GOODIES displays the results in E.

After the biological relationship among genes in each cluster is quantitatively estimated by *AverPd*, the clusters whose *AverPd* score is sufficiently low can be used for functional assignment of unknown

