

# Predicting 400 GO Functions of Proteins

Part of the Proteome Analyst Suite of Tools

Roman Eisner\*, Alona Fyshe, Russell Greiner, Paul Lu,  
Brandon Percy, Brett Poulin, Duane Szafron, David Wishart

\*eisner@cs.ualberta.ca



UNIVERSITY OF ALBERTA

[www.cs.ualberta.ca/~bioinfo/PA](http://www.cs.ualberta.ca/~bioinfo/PA)

## Introduction

- We predict over 400 molecular function categories from Gene Ontology (www.geneontology.org)
- We predict functions of proteins from sequence information
- Our technique is evaluated against experimentally verified data
- Contributions:

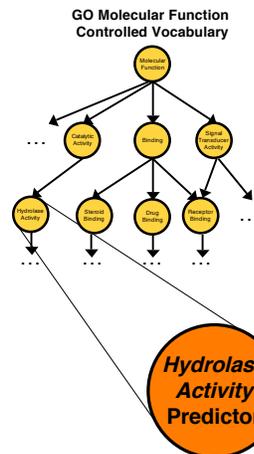
- More **accurate** than BLAST, especially on remotely related proteins
- Coverage of a **large ontology** with accurate predictions
- Exploiting the **hierarchy** to increase accuracy and minimize computational complexity

## Experimentally Consistent

- GO Term Predictors are trained, and evaluated against experimentally annotated proteins. No electronic annotations are used for evaluating our predictions.
- All accuracy statistics are obtained using 5-fold cross-validation.

## References

- The Gene Ontology Consortium [http://www.geneontology.org/]
- ProtFun 2.2 Server [http://www.cbs.dtu.dk/services/ProtFun/]
- Slim GO @ EBI [http://www.ebi.ac.uk/GOA/]
- Proteome Analyst Server [http://www.cs.ualberta.ca/~bioinfo/PA]
- GO Slim @ MGI [http://www.spatial.maine.edu/~mdolan/MGI\_GO\_Slim.html]



The Hydrolase Activity Predictor predicts whether a given protein is a hydrolase enzyme or not.

Figure 1

## ★ Creating Predictors for each GO Term

- Machine-learned
- Binary (predict yes/no for each function term)
- Created for every GO function with at least 20 experimentally verified proteins (Total of 406)
- A Weighted Combination of:
  - Probabilistic Suffix Trees (PSTs)
  - PFAM with Support Vector Machines
  - Proteome Analyst (See Proteome Analyst Poster)

## 1 Predicting the Function of Proteins

- Goal: To find all functions of a query protein. Input is a protein sequence (Fasta format)
- We increase the predictive accuracy on those proteins that are similar to experimentally verified ones (Table 1). Here, similar means BLAST E-value  $\leq 0.001$
- Our predictors work significantly better for query proteins which do not have a good BLAST result against the set of experimentally verified proteins (Table 2)

| Proteins with a good BLAST hit ( $\leq 1e-3$ )   |                   |                |
|--|-------------------|----------------|
|  | Overall Precision | Overall Recall |
| BLAST (E-value $\leq 0.001$ )  | 77%               | 78%            |
| <b>Our Method</b>  | <b>78%</b>        | <b>80%</b>     |
| How often does this occur?<br>60% of <i>D. melanogaster</i> proteins<br>62% of <i>S. cerevisiae</i> proteins |                   |                |

Table 1

| Proteins with no good BLAST hit  |                   |                |
|--|-------------------|----------------|
|  | Overall Precision | Overall Recall |
| BLAST (any hit accepted)   | 19%               | 20%            |
| <b>Our Method</b>  | <b>54%</b>        | <b>31%</b>     |
| How often does this occur?<br>40% of <i>D. melanogaster</i> proteins<br>38% of <i>S. cerevisiae</i> proteins |                   |                |

Table 2

## 2 Large Ontology

- 406 categories of Molecular Function
- Allows for very specific and general predictions of function

| GO Predictor                  | Size of Ontology |
|-------------------------------|------------------|
| Our Ontology                  | 406              |
| ProtFun <sup>2</sup>          | 14               |
| Slim-GO @ EBI <sup>3</sup>    | 30               |
| Proteome Analyst <sup>4</sup> | 12               |
| GO Slim @ MGI <sup>5</sup>    | 13               |

## 3 Exploiting the Hierarchy

- To find the functions of a query protein, first BLAST against experimental data. When BLAST provides a good match (Table 1), we use the hit's annotations as a guide to which term predictors should be computed. This reduces runtime, without penalty to accuracy.

- Currently working on reducing the computational runtime of predicting function for those proteins which do not return a good BLAST result.

- Predictions of predictors are propagated upwards in the hierarchy to maintain consistency.

- To create classifiers (Figure 1), positive and negative training examples for each term predictor are selected to maintain consistency with the hierarchy. This increases the accuracy of each term predictor.