EMBL
Heidelberg
Germany

# STRING

# Predicting protein networks from genomic context and external data

Lars Juhl Jensen (jensen@embl.de)
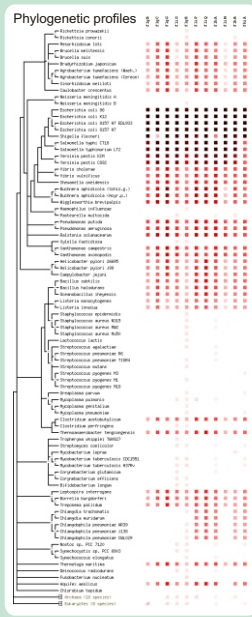
Christian von Mering (mering@embl.de)

Peer Bork (bork@embl.de)

The web-based database STRING (http://string.embl.de) is a large, pre-computed resource that allows any protein of interest to be placed into a high-confidence network of functionally associated protein partners. For 100+ genomes, the tool evaluates genomic context, protein interaction data, gene co-expression, co-mentioning of genes in literature, and certain database annotations.

## Genomic context

So-called genomic context methods allow discovery of functional relations from complete genome sequences. Two methods rely on functionally related genes often being arranged in operons (the gene neighborhood method) or even fused to a single gene (the gene fusion method). Finally, the genes of a functional module are typically either all present or all absent in each genome.

Genomic context methods tend to be most powerful for analyzing prokaryotic genomes. Part of the reason is that more prokaryotes than eukaryotes have been sequenced. The neighbourhood method, however, will never be powerful for eukaryotes, as it relies on detecting operons. In the near future, we intend to expand STRING with novel genomic context methods.


Phylogenetic profiles

- Calculate all-against-all pairwise alignments
- Identify fusion genes in each genome
- Score gene neighbors by intergenic distances

- Calculate best-hit profile for each protein
- Principal component analysis of profiles
- Calculate distances between PC profiles

## Large-scale experimental data


ARRAYPROSPECTOR

- Renormalize arrays by the Q-spline method
- Principal component analysis of arrays
- Construct kernel density predictor

ArrayProspector is an interactive web resource of gene associations predicted from all expression data in the Stanford Microarray Database (SMD). The underlying pre-computed database currently contains more than 200,000 high confidence gene associations in 12 different species. The resource allows every association to be visually inspected. STRING utilizes both the evidence scores and the visualization capabilities of ArrayProspector.

In addition to microarray expression data, STRING also incorporates other types of large-scale experimental data sets, e.g. from yeast two-hybrid screens, complex purifications, and chromatin immunoprecipitation experiments. For all such data, a quality score is associated with every individual binary interaction.

http://www.bork.embl.de/ArrayProspector

## Literature mining

Vast amounts of biological knowledge is buried in the scientific literature. STRING extracts protein associations from PubMed abstracts based on co-mentioning. We will soon also include more specific relations extracted by shallow parsing.



- Associate abstracts with species
- Identify gene names in each title/abstract
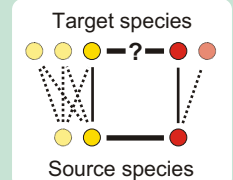- Count co-occurrences of gene names
- Calculate significance of each association

## Integration of diverse evidence

- Calibrate raw scores against KEGG maps
- Transfer evidence scores across species
- Combine evidence for each gene pair

The notion of "functional relations" is central to STRING. All raw evidence scores are calibrated against KEGG metabolic maps, which serves as our reference set. Calibrated scores for different evidence types are directly comparable and easily combined.


Target species
Source species

All types of evidence in STRING are transferred across species using a "fuzzy orthology" scheme.

## Related publications

C. von Mering, M. Huynen, D. Jaeggi, S. Schmidt, P. Bork and B. Snel
"STRING: A database of predicted functional associations between proteins"
*Nucleic Acids Research*, **31**, 258-261, 2003

J.O. Korbel, L.J. Jensen, C. von Mering and P. Bork
"Analysis of genomic context: Prediction of functional associations from bidirectionally transcribed gene pairs"
*Nature Biotechnology*, **22**, 911-917, 2004

L.J. Jensen, J. Lagarde, C. von Mering and P. Bork
"ArrayProspector: A web resource of functional associations inferred from microarray expression data"
*Nucleic Acids Research*, **32**, W445-W448, 2004

L.J. Jensen and P. Bork
"Quality analysis and integration of large-scale molecular data sets"
*Drug Discovery Today: TARGETS*, **3**, 51-56, 2004

P. Bork, L.J. Jensen, C. von Mering, A.K. Ramani, I. Lee and E.M. Marcotte
"Protein interaction networks from yeast to human"
*Current Opinions in Structural Biology*, **14**, 292-299, 2004

J. Saric, L.J. Jensen, R. Ouzounova, I. Rojas and P. Bork
"Extracting regulatory gene expression networks from PubMed"
in *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics*, 2004

## A network of functional modules

Through the cross-species integration of many, diverse evidence types, STRING provides a network of functional association for proteins from 100+ species. This entire network can be browsed online at http://string.embl.de

The network is highly ordered, and functional modules can be discovered by performing e.g. k means clustering of the network. The resulting functional modules are not isolated, but rather form a network of their own.

To the right is shown a small example network of three related functional modules from yeast: the ATP synthetase complex, the cytochrome C oxidase complex and the ubiquitinol-cytochrome C reductase complex.