# A tool for evaluating strategies for grouping of biological data

**KitEGA**

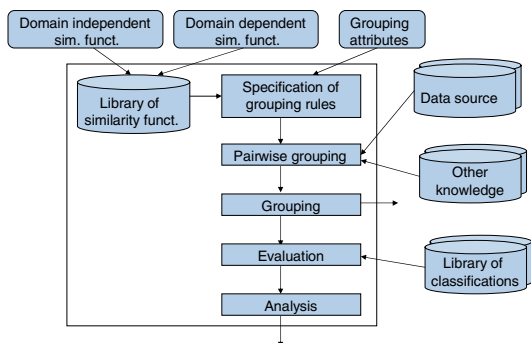Tool**Kit** for **E**valuation of **G**rouping **A**lgorithms

During the last decade an enormous amount of biological data has been generated and techniques and tools to analyze this data have been developed. Many of these tools use some form of grouping and are used in, for instance, data integration, data cleaning, prediction of protein functionality, and correlation of genes based on microarray data. Grouping of biological data is not a trivial task:

- a number of aspects influence the quality of the results: the data sources, the grouping attributes and the algorithms implementing the grouping procedure
- a variety of grouping algorithms is available, but it is often not clear which method performs best for which grouping tasks
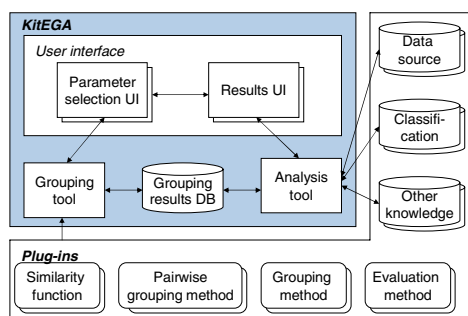- existing grouping algorithms may not be applied straightforward

Environments that support comparison and evaluation of different grouping strategies for different grouping tasks on different data sets are needed to:

- support study of the properties of biological data sources
- select the suitable grouping procedures
- get an insight in how the grouping procedures could be used in the best way
- lead to recommendations on how to improve the current procedures and develop new procedures

## Method for similarity-based grouping



## The KitEGA framework



## Example use

Grouping task. Grouping of proteins with respect to
1. biological function
2. class of isozymes they belong to

Formation of data sources:
- analyzed human proteins involved in glycolysis
- retrieved 190 data entries via Entrez
- used GO Consortium mappings ec2go and spkw2go to extend the set of available GO terms

The first prototype includes:

Library of similarity functions
- EditDist(v1,v2)
- SeqSim(v1,v2)
- SemSim(v1,v2)

Other knowledge
- GO ontology
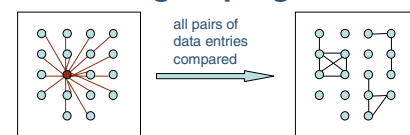
Classifications. Manual classification according to
- biological function
- classes of isozymes

Grouping metohds
- connected components (transitive similarity)
- cliques (overlapping groups)

## Specification of grouping rules



## Pairwise grouping



all pairs of data entries compared

## Grouping



## Evaluation

External quality measures, i.e. with respect to known classes of the grouped data

Number of entries: 92    Entropy: 1.0
Number of groups: 26    Purity: 1.0
Number of classes: 25   MutualInformation: 0.8810530832230519
                        FMeasure: 0.9939613526570048

## Analysis. Distribution of data entries in a test case



true positives
false positives
false negatives

## Analysis. Group vs class data entries



## Analysis. All test cases



Studied aspects, e.g. use of different data sources, grouping algorithms, and classifications, grouping on different attributes, impact of threshold.

## References

Jakoniene V, Lambrix P, `A Tool for Evaluating Strategies for Grouping of Biological Data', *Journal of Integrative Bioinformatics*, 4(3):83, 2007.

Jakoniene V, Rundqvist D, Lambrix P, `A method for similarity-based grouping of biological data', *Proceedings of the International Workshop on Data Integration in the Life Sciences - DILS06*, LNBI 4075, pp 136-151, 2006.

Contact: P. Lambrix
Department of Computer and Information Science
Linköping universitet, Sweden
kitega@ida.liu.se
http://www.ida.liu.se/~iislab/projects/KitEGA/

Conclusions for test cases obtained by using KitEGA:
- best suited grouping approaches for data source Glyc-Funct-AnnEc-onlyGO
  SemSim(GOcomb) for grouping on biological function
  SeqSim(Sequence) for grouping on classes of isozymes
- for the used grouping tasks spkw2go – too general, ec2go – specific enough