



# Protein World & Orthology Benchmarking

Tim Hulsen<sup>1\*</sup>, Peter Groenen<sup>2</sup>, Wilco Fleuren<sup>1,3</sup> and Martijn Huynen<sup>1</sup>

<sup>1</sup>CMBI, Toernooiveld 1, 6525 ED Nijmegen (NL) \*T.Hulsen@cmbi.kun.nl

<sup>2</sup>NV Organon, Molenstraat 110, 5342 CC Oss (NL)

<sup>3</sup>NBIC, Appelweg 16, 3818 NN Amersfoort (NL)



Organon

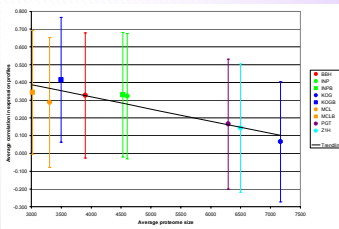
## Abstract

We tested the **quality** of **ortholog identification methods**, based on the assumption that orthologs have a **highly similar function**. The conservation of function between proteins in different species can be measured by using the **functional annotation** provided with the protein sequences but also with functional data like **co-expression** data, **protein interaction** data and **chromosomal position** data. Data were used in a **direct** way, looking at only one protein in each species, and in a **pairwise** way, looking at two proteins in each species having some kind of relationship (e.g. co-expressed, interacting, neighboring). The benchmarking methods were used to test a number of ortholog identification methods.

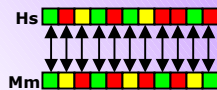
## 'Protein World' data set

For a fair comparison of all of the covered methods, the same data set was used at all time. This **'Protein World'** data set was created by **comparing** all of the currently known and predicted **proteins** (SpTrEMBL) through the **Smith-Waterman** algorithm with calculation of **Z-values**. This sensitive algorithm is more precise than the more common BLAST and FASTA algorithms, which were used to build other protein relation databases. The software package used to do the calculations is called **'Biofacet'**, formerly known as 'LASSAP' ('Large Scale Sequence compArison Package'). The data set is available through the dutch bioinformatics portal site <http://www.bioasp.nl>. As for both **human** and **mouse** good expression data and other functional data were available, we used the orthologous relationships between these two species for our study.

## Direct conservation of functional parameters

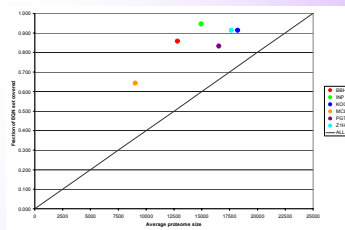


### Correlation of expression profiles



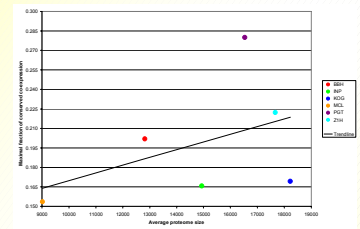
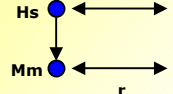
### Equal SwissProt name

ABCD\_HUMAN  
~  
ABCD\_MOUSE

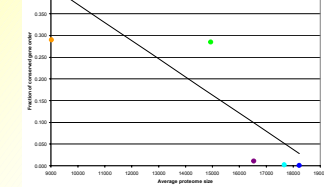


## Pairwise conservation of functional parameters

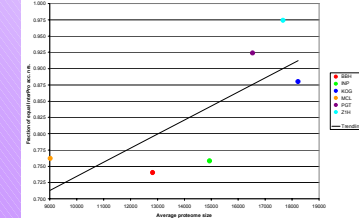
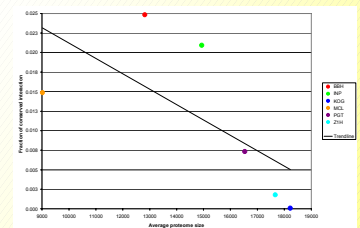
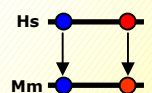
### Conservation of co-expression



### Conservation of interaction



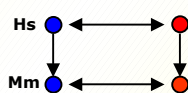
### Conservation of gene order



### Equal InterPro acc. nr.

5H6\_HUMAN (IPR000276)  
~  
5H7\_HUMAN (IPR000276)

### Conservation of interaction



## Compared ortholog identification methods

Best Bidirectional Hit (BBH), InParanoid (INP), KOG (KOG), OrthoMCL (MCL), Phylogenetic Tree (PGT) and Z>100 (Z1H)

## Overall scoring graph

Created by adding up all normalized benchmarking scores per ortholog identification method.

BBH and INP score best.

