# RESEARCHER'S DIGEST: USING GENE ANNOTATIONS TO CLASSIFY GENE LISTS INTO FUNCTIONAL GROUPS

**Juan Carlos Triviño, Juan Carlos Sánchez-Ferrero and <u>Juan Carlos Oliveros</u>**

(jctrivino@cnb.csic.es, jcsanchez@cnb.csic.es, oliveros@cnb.csic.es)

Service of Bioinformatics for Genomics and Proteomics (BioinfoGP), Centro Nacional de Biotecnología (CNB-CSIC). Darwin 3, Campus de Cantoblanco, 28049, Madrid, Spain.

## Introduction

DNA Microarrays are used in a routine way to measure the transcription level of thousands of genes in several experimental conditions. The main result of these experiments usually consists on one or several lists of genes whose transcription levels are induced (or repressed) compared against a reference, or share similar transcription patterns in different cellular states.

While there are many robust statistical methods to obtain these gene lists, today there is a growing interest on the development of new bioinformatic tools that would assist the researcher in the interpretation of their results going beyond the pure numerical classification.

Researcher's Digest is an on-line tool to classify gene lists, obtained from microarray experiments, into functional groups. To do that, the system makes use of text-mining techniques applied to the gene annotations.

## The program:

1. Classify the words in the gene annotations in biological categories (enzymes, accession codes, gene families, etc.).

2. Evaluate the degree of similarity between all pairs of sentences (the more relevant terms are shared between two annotations, the more similar are they).

3. Sort all genes starting from the first one in the list so the more similar gene is placed the next and so on.

4. Estimates a similarity threshold that best separates the gene groups.

## The results…

…are presented in a user friendly web page where the user can modify the similarity threshold used to separate the groups and change some visualization options.



Digest is one of the few bioinformatics tools that uses free-text as source of data for grouping genes and can compare up to ~1000 genes in few seconds. It is accessible at:

**http://bioinfogp.cnb.csic.es/tools/digest**