# Bacterial Feature Finder (BaFF) – A system for extracting features overrepresented in sets of prokaryotic organisms

## *Supplementary Material*

Javier López-Ibáñez, Laura T. Martín, Mónica Chagoyen & Florencio Pazos*

Computational Systems Biology Group. Systems Biology Program. Spanish National Centre for Biotechnology (CNB-CSIC). c/ Darwin, 3. 28049 Madrid. Spain. Phone: +34-915854669. Fax: +34-915854506.

pazos@cnb.csic.es

## 1.- Statistical Tests

Enrichment analysis statistics were computed using functions from the R *stats* package [1].

For quantitative features (continuous values), a two Sample, two-sided Kolmogorov-Smirnov test (R function: `ks.test()`) was used for checking if both sets of values (input set and background) have the same distribution. We report the p-value of rejecting the hypothesis that both distributions are different.

For qualitative features (discrete values), we compute the probability of obtaining a given annotation randomly from the background set using the cumulative hypergeometric density function (R function: `phyper()`).  In these cases a p-value correction for multiple testing is also performed, with Benjamini's False Discovery Rate (FDR) correction [2] (R function: `p.adjust()`).

[1] R Core Team (2015). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.

[2] Benjamini, Y. and Hochberg, Y. (1995) Controlling the false discovery rate: A practical and powerful approach to multiple testing. Journal of the Royal Statistical Society, Series B. 57: 289-300.

## 2.- Features and sources

| Feature | Source Database | URL |
|---|---|---|
| • Number of genes | EnsemblGenomes | https://img.jgi.doe.gov/cgi-bin/m/main.cgi<br><br>http://ensemblgenomes.org/info/genomes |
| • GC content | NCBI Genome | https://www.ncbi.nlm.nih.gov/genome/ |
| • Biosystems | NCBI biosystems | https://www.ncbi.nlm.nih.gov/biosystems |
| • COG Functional Classes (25) | Microbesonline | http://www.microbesonline.org/cgi-bin/genomeInfo.cgi |
| • Gram staining<br>• Sporulation<br>• Motility<br>• Shape<br>• Biotic Relationships<br>• Phenotype<br>• Oxigen requirement<br>• Temperature range<br>• Disease<br>• Host | Integrated Microbial genomes (IMG) | https://img.jgi.doe.gov/ |
| • Taxonomy | NCBI Taxonomy | https://www.ncbi.nlm.nih.gov/Taxonomy/ |

# 3.- Screenshots of the interface



Main interface



Paginated results of a database search

## Overrepresentation Analysis Results

Showing enrichment results for the findings of a previous search. Using the default background set. Found annotations for 692 out of 692 IDs..

You can show/hide different categories of features in the results table with the boxes below:

- ☑ Biosystems
- ☑ Biotic Relationships
- ☑ Motility
- ☑ Taxonomy
- ☑ Shape
- ☑ Temperature range
- ☑ Sporulation
- ☑ Oxigen requirement
- ☑ Disease
- ☑ Phenotype
- ☑ Host
- ☑ COG code

**Save Results**

| Feature | Category | p-value | FDR |
|---|---|---|---|
| %GC content | %GC content | 0.00e+0 | - |
| 2-Aminoethylphosphonate transport system | Biosystems | 1.33e-8 | 2.33e-8 |
| 2-Oxocarboxylic acid metabolism | Biosystems | 8.13e-4 | 1.19e-3 |
| *A:* RNA processing and modification | COG code | 0.00e+0 | - |
| Acyrthosiphon pisum | Host | 1.30e-7 | 3.65e-6 |
| Adenine ribonucleotide biosynthesis, IMP => ADP,ATP | Biosystems | 1.38e-6 | 2.26e-6 |
| Adhesin protein transport system | Biosystems | 7.64e-12 | 1.45e-11 |
| ADP-L-glycero-D-manno-heptose biosynthesis | Biosystems | 0.00e+0 | 0.00e+0 |
| AI-2 transport system | Biosystems | 0.00e+0 | 0.00e+0 |
| Alanine, aspartate and glutamate metabolism | Biosystems | 3.11e-3 | 4.40e-3 |
| alpha-Hemolysin/cyclolysin transport system | Biosystems | 2.07e-4 | 3.12e-4 |
| alpha-Linolenic acid metabolism | Biosystems | 0.00e+0 | 0.00e+0 |
| Aminobenzoate degradation | Biosystems | 4.74e-12 | 9.04e-12 |
| Aminoglycoside resistance, protease FtsH | Biosystems | 0.00e+0 | 0.00e+0 |
| Aminoglycoside resistance, protease HtpX | Biosystems | 0.00e+0 | 0.00e+0 |
| Arachidonic acid metabolism | Biosystems | 0.00e+0 | 0.00e+0 |
| ArcB-ArcA (anoxic redox control) two-component regulatory system | Biosystems | 0.00e+0 | 0.00e+0 |
| Arginine and proline metabolism | Biosystems | 6.28e-2 | 8.08e-2 |
| Arginine biosynthesis | Biosystems | 3.56e-5 | 5.48e-5 |
| Arginine transport system | Biosystems | 0.00e+0 | 0.00e+0 |
| Ascorbate and aldarate metabolism | Biosystems | 4.43e-13 | 8.65e-13 |
| Ascorbate degradation, ascorbate => D-xylulose-5P | Biosystems | 0.00e+0 | 0.00e+0 |

Results of the analysis of enriched features.