



# A procedure to recruit members to enlarge protein family databases - the building of UECOG (UniRef-Enriched COG Database) as a model

G.R. Fernandes<sup>1\*</sup>, D.V.C. Barbosa<sup>1\*</sup>, F. Prosdocimi<sup>1</sup>, I.A. Pena<sup>1</sup>,  
L. Santana-Santos<sup>1</sup>, O. Coelho Junior<sup>1</sup>, A. Barbosa-Silva<sup>1</sup>, H.M. Velloso<sup>1</sup>,  
M.A. Mudado<sup>2</sup>, D.A. Natale<sup>3</sup>, A.C. Faria-Campos<sup>4</sup>, S.V. A. Campos<sup>4</sup> and  
J.M. Ortega<sup>1</sup>

<sup>1</sup>Departamento de Bioquímica e Imunologia, Laboratório de Biodados,  
Instituto de Ciências Biológicas, Universidade Federal de Minas Gerais,  
Belo Horizonte, MG, Brasil

<sup>2</sup>Fundação Ezequiel Dias, Belo Horizonte, MG, Brasil

<sup>3</sup>Protein Information Resource, Georgetown University Medical Center,  
Washington, DC, USA

<sup>4</sup>Departamento de Ciência da Computação, ICEX,  
Universidade Federal de Minas Gerais, Belo Horizonte, MG, Brasil

\*These authors contributed equally to this study.

Corresponding author: J.M. Ortega

E-mail: miguel@icb.ufmg.br

Genet. Mol. Res. 7 (3): 910-924 (2008)

Received June 2, 2008

Accepted August 11, 2008

Published September 30, 2008

**ABSTRACT.** A procedure to recruit members to enlarge protein family databases is described here. The procedure makes use of UniRef50 clusters produced by UniProt. Current family entries are used to recruit additional members based on the UniRef50 clusters to which they belong. Only those additional UniRef50 members that are not fragments and whose length is within a restricted range relative to the original entry are recruited. The enriched dataset is then limited to contain only genomes from selected clades. We used the COG database - used for genome annotation and for studies of phylogenetics and gene evolution - as a model. To validate the method, a

UniRef-Enriched COG0151 (UECOG) was tested with distinct procedures to compare recruited members with the recruiters: PSI-BLAST, secondary structure overlap (SOV), Seed Linkage, COGnitor, shared domain content, and neighbor-joining single-linkage, and observed that the former four agree in their validations. Presently, the UniRef50-based recruitment procedure enriches the COG database for Archaea, Bacteria and its subgroups Actinobacteria, Firmicutes, Proteobacteria, and other bacteria by 2.2-, 8.0-, 7.0-, 8.8-, 8.7-, and 4.2-fold, respectively, in terms of sequences, and also considerably increased the number of species.

**Key words:** COG; Secondary database; UniRef; UniProt; UECOG