# RoundUp
## a repository of orthologs and corresponding evolutionary distances

I-Hsien Wu, Jian Pu, Todd Deluca, Thomas J. Monaghan, Dennis P. Wall
The Computational Biology Initiative, Department of Systems Biology, Harvard Medical School
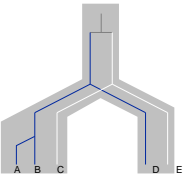
## Introduction

The ability to accurately detect orthologous or functionally equivalent proteins in different organisms is important to numerous biological research questions, including studies of variables influencing the rate of protein evolution [1-3], accurate genome annotation [4], and studies of proteins implicated in cancer [5]. In all cases, it is crucial that the orthologous pair be accurate (Figure 1). Moreover, the ability to retrieve orthologs from many genomes at the click of a button would help speed such research projects along by several orders of magnitude.

To meet these challenges, we have pre-computed orthologs for 251 genomes using the reciprocal smallest distance algorithm [6]. This algorithm is an improvement over approaches that rely on BLAST alone, since it uses global rather than local sequence alignments and evolutionary estimates of distance between sequences rather than blast probability scores, an approach that can often be positively misleading when trying to determine functional equivalence [7]. These pre-computed sets of orthologs are stored in a publicly accessible database, RoundUp, the most comprehensive of its kind [8, 9].

### Figure 1:

In the case of comparing the evolutionary rates of proteins in the absence of a normalizing molecular clock, such as the rate of synonymous substitutions, estimates of evolutionary rate must be based upon comparisons between sequences that are orthologs (sequences that diverged from each other at the species split), and not paralogs (sequences that diverged at another time). Only if all sequence comparisons share the same time of divergence are protein evolutionary distances expected to be proportional to relative evolutionary rates. For example, in this figure, the orthologous comparisons (A or B with D; C with E) would yield evolutionary distances indicative of the relative rates of protein evolution. By contrast, paralogs comparison of A or B with E would yield an evolutionary distance that would badly overestimate the evolutionary rate of these sequences. Procedures that rely on BLAST alone to detect ortholog pairs, such as reciprocal best BLAST hits (rbh) [3,12] are known to be faulty or incomplete [6,7]. To develop RoundUp, we used an improved method for detecting orthologs [6].

## Methods

### The reciprocal smallest distance algorithm

The method (Figure 2) employs BLAST [10] as a first step, starting with a subject genome, *J*, and a protein query sequence, *i*, belonging to genome *I*. A set of hits, *H*, exceeding a predefined significance threshold (we have generally used E-value < 1e-20) is obtained. Then, using Clustalw [13], each protein sequence in *H* is aligned separately with the original query sequence *i*. If the alignable region of the two sequences exceeds a threshold fraction of the alignment's total length (0.8 is our working cutoff), the program PAML [14] is used to obtain a maximum likelihood estimate of the number of amino acid substitutions separating the two protein sequences, given an empirical amino acid substitution rate matrix [11]. The model under which a maximum likelihood estimate is obtained may include variation in evolutionary rate among protein sites, and for more distant comparisons we have generally assumed a gamma distribution of rate across sites, with shape parameter a = 1.53 [15]. Of all sequences in *H* for which an evolutionary distance is estimated, only *j*, the sequence yielding the shortest distance, is retained. This sequence *j* is then used for a reciprocal BLAST against genome *I*, retrieving a set of high scoring hits, *L*. If any hit from *L* is the original query sequence, *i*, the distance between *i* and *j* is retrieved from the set of smallest distances calculated previously. The remaining hits from *L* are then separately aligned with *j* and maximum likelihood distance estimates are calculated for these pairs as described above. If the protein sequence from *L* producing the shortest distance to *j* is the original query sequence, *i*, it is assumed that a true orthologous pair has been found and their evolutionary distance is retained. This algorithm has been tested extensively and is known to produce a more comprehensive and accurate set of orthologs than alternative methods like reciprocal BLAST [6].

### The RoundUp data federation and Web Interface

Recovering the true set of orthologs between two lineages will depend on many parameters, including date of divergence between the lineages, rates of gene duplication in either lineage, intensity of molecular selection, and others. To account for such variables that can have a large impact on the size and content of a list of orthologs, and to provide users of RoundUp a certain degree of exploratory power, we adjusted two parameters in the reciprocal smallest distance algorithm – BLAST E-value and global pair-wise sequence divergence – when pre-calculating orthology between genomes. Specifically, in our pre-calculations, we used four increasingly stringent BLAST E-value thresholds, 1e-5, 1e-10, 1e-15, and 1e-20, and three increasingly stringent divergence thresholds, 0.8, 0.5, and 0.2. Therefore, for every pair of genomes, our RoundUp repository contains twelve ortholog lists representing all possible combinations of the two variables.

The pre-computed results are accessible through a web portal to a database federation. The federation, based on WSII technology [16] and our RoundUp data adaptor, allows SQL-based access to the orthology data. The federation is robust: updates and changes to the database repository as new proteomes become available and as older proteomes are updated will not interrupt access. NCBI and other standard warehouses of whole-genomic data are checked weekly for updates and additions, and the RoundUp pre-calculation procedure is run accordingly. Developed using Java technology, the web portal to this database federation is accessible from the Internet using any standard web browser. The web portal allows users without any knowledge of SQL to create arbitrarily complex SQL queries. Presently the portal allows a user to:

1. Discover all orthologs between a specified query genome and any number of other genomes in the database federation.
2. Discover any orthologs between a specified query genome and any number of other genomes in the database federation.
3. Discover all of the transitively closed groups of orthologs among a query genomes and any number of other genomes.
4. Focus results on genes of interest in a query genome by entering GI numbers for those genes.
5. Explore the effects of different parameter settings by choosing one of 12 possible parameter combinations, as well as the evolutionary distance threshold.
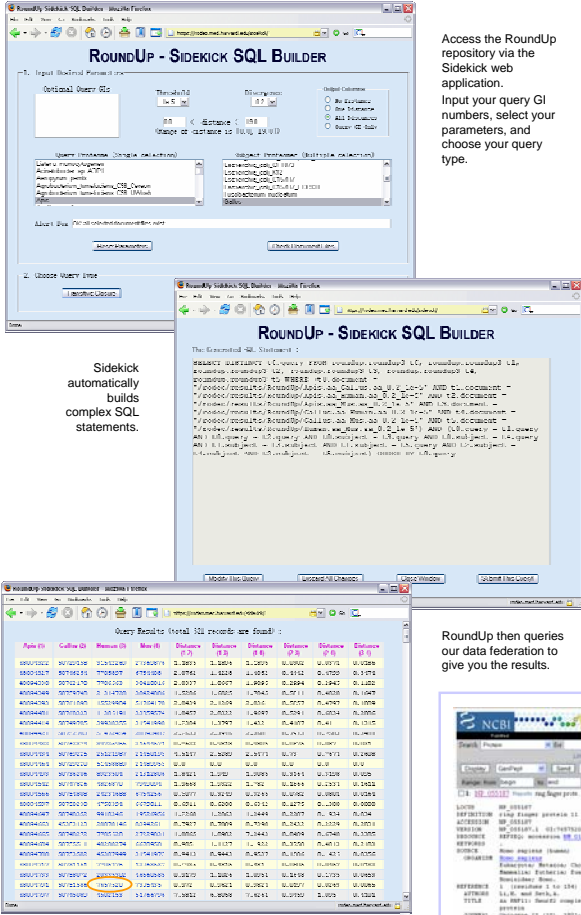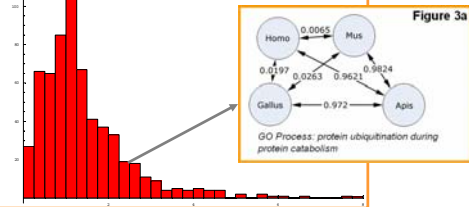6. View the evolutionary distance or orthologs by choosing to view none, one, or all evolutionary distances.



Access the RoundUp repository via the Sidekick web application.

Input your query GI numbers, select your parameters, and choose your query type.

Sidekick automatically builds complex SQL statements.

RoundUp then queries our data federation to give you the results.

### Figure 3:

A histogram of the 621 transitively closed results between Apis, Gallus, Homo, and Mus.



GO Process: protein ubiquitination during protein catabolism

## Present status and capabilities

RoundUp results are stored in relational form and accessible through a flexible web interface that allows the user to build arbitrarily complex and exploratory queries without needing to know how to construct SQL statements. Presently, the RoundUp repository contains orthologs from all possible pair-wise analyses of 251 genomes, comprising 376,500 ortholog files. This vast repository offers exciting opportunities to explore the conservation of genes, the evolution of genetic pathways, and the occurrences of gene gains and losses over evolutionary time scales and more.

A user of the RoundUp system can build queries to find all orthologs in all 251 genomes for a given search sequence that are either evolutionarily conserved or labile, according to some *a priori* chosen distance threshold. Alternatively, a user may quickly build clusters of orthologous genes for any number of genomes, i.e., build groups of orthologs that exhibit complete transitive closure.

As an example we executed the transitive closure query shown to the left. After 32 seconds the query returned the 621 transitively closed orthologs among Apis, Gallus, Homo, and Mus, with an average distance of 1.4. Figure 3a shows a single gene from this query. The gene is a ring finger protein, known to have ubiquitin-protein ligase activity and to be vital for proper functioning of protein ubiquitination during protein catabolism.
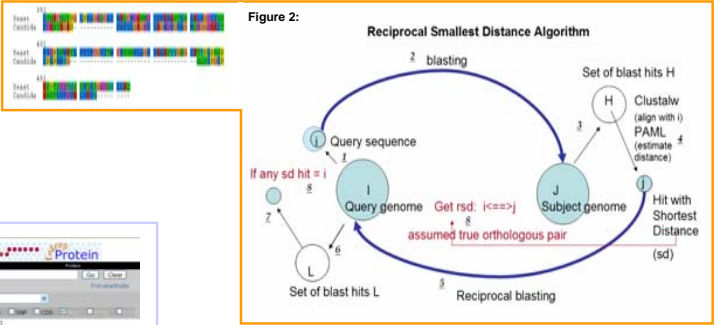
It will be interesting to investigate the functional annotations for each transitively closed gene in this and other lists from different combinations of organisms. Though many may turn out to be predictable for their functional importance, others may shed light on interesting layers of gene conservation not revealed by other methods of investigation.

### RoundUp is currently capable of executing queries such as:

What are the transitively closed orthologs among a variety of genomes?
What is the average distance of a set of orthologous sequences?
How does the set of orthologous sequences change when you relax the threshold parameters?

### Figure 2:



Reciprocal Smallest Distance Algorithm

## Future Objectives:

Integrate ortholog data with protein-protein interaction data, expression profile data, and data from Gene Ontology, etc., using the RoundUp database federation.
Extend and update orthology data as genomes are released and updated.
Provide more SQL query templates to answer interesting biological questions.
Extend RoundUp with functions over query result sets. E.g. statistical operations on distances or set operations to compare the effects of different parameters.

## References

[1]. Fraser, H.B., et al., *Evolutionary rate in the protein interaction network.* Science, 2002. **296**(5568): p. 750-2.
[2]. Fraser, H.B., D.P. Wall, and A.E. Hirsh, *A simple dependence between protein evolution rate and the number of protein-protein interactions.* BMC Evol Biol, 2003. **3**(1): p. 11.
[3]. Hirsh, A.E. and H.B. Fraser, *Protein dispensability and rate of evolution.* Nature, 2001. **411**(6841): p. 1046-9.
[4]. Waterston, R.H., et al., *Initial sequencing and comparative analysis of the mouse genome.* Nature, 2002. **420**(6915): p. 520-62.
[5]. Brown, J.R., et al., *Evolutionary relationships of Aurora kinases: implications for model organism studies and the development of anti-cancer drugs.* BMC Evol Biol, 2004. **4**(1): p. 39.
[6]. Wall, D.P., H.B. Fraser, and A.E. Hirsh, *Detecting putative orthologs.* Bioinformatics, 2003. **19**(13): p. 1710-1.
[7]. Koski, L.B. and G.B. Golding, *The closest BLAST hit is often not the nearest neighbor.* J Mol Evol, 2001. **52**(6): p. 540-2.
[8]. Remm, M., C.E. Storm, and E.L. Sonnhammer, *Automatic clustering of orthologs and in-paralogs from pairwise species comparisons.* J Mol Biol, 2001. **314**(5): p. 1041-52.
[9]. Tatusov, R.L., et al., *The COG database: a tool for genome-scale analysis of protein functions and evolution.* Nucleic Acids Res, 2000. **28**(1): p. 33-6.
[10]. Altschul, S.F., et al., *Basic local alignment search tool.* J Mol Biol, 1990. **215**(3): p. 403-10.
[11]. Jones, D.T., W.R. Taylor, and J.M. Thornton, *The rapid generation of mutation data matrices from protein sequences.* Comput Appl Biosci, 1992. **8**(3): p. 275-82.
[12]. Jordan,I.K., Rogozin,I.B., Wolf,Y.I. and Koonin,E.V. Essential genes are more evolutionarily conserved than are nonessential genes in bacteria. *Genome Res.*, **12**, 962–968.
[13]. Thompson,J.D., Plewniak,F., Thierry,J.C. and Poch,O. DbClustal: Rapid and reliable global multiple alignments of protein sequences detected by database searches. *Nucleic Acids Res.*, **28**, 2919–2926.
[14]. Yang,Z. (2000) Phylogenetic Analysis by Maximum Likelihood (PAML). University College London, London.
[15]. Nei,M., Xu,P. and Glazko,G. (2001) Estimation of divergence times from multiprotein sequences for a few mammalian species and several distantly related organisms. *Proc. Natl Acad. Sci. USA*, **98**, 2497–2502.
[16]. http://www.ibm.com/software/data/integration/federation.html

Department of Systems Biology at Harvard Medical School
Computational Biology Initiative