

# High resolution classification of orthogroups by recursive dynamic Markov clustering

Stephen R. Bond, Karl E. Keat, and Andreas D. Baxevasis\*

Computational and Statistical Genomics Branch, Division of Intramural Research, National Human Genome Research Institute, National Institutes of Health, 50 South Drive, Bethesda, MD, USA, 20892

\*Corresponding author: E-mail: andy@mail.nih.gov

Associate Editor:

## Abstract

Key words: orthogroup, clustering

## Introduction

When a gene evolves an important physiological role, purifying selection will often maintain that function through evolutionary time [REF]. As a result, orthology (i.e., homology via speciation) has become a well-accepted predictor of shared gene product function among species, and considerable effort has been made to develop computational methods to identify orthologs. The algorithms in current use fall into two distinct categories: Tree-based and graph-based clustering methods (Tekaia, 2016). Tree-based approaches (e.g., Ensembl Compara (Vilella *et al.*, 2009), LOFT (van der Heijden *et al.*, 2007), and SYNERGY (Wapinski *et al.*, 2007)) broadly rely on estimating a phylogenetic tree for a target gene family, and then reconciling the gene tree with a ‘known’ species tree to identify orthologous clades. Tree-based methods are very accurate under ideal conditions, although high quality species trees are often difficult to estimate.

Alternatively, pairwise similarity graph clustering methods can leverage graph theory to rapidly identify natural clusters of related sequences from genome scale datasets. Reciprocal best-hit methods were among the earliest developed for this purpose, but were restricted to assessing only two species at a time. Due to the non-transitive nature of orthology (i.e., paralogs in one species can be orthologous to a single gene in another species), it is more difficult (or impossible) to explicitly assign sequences into groups of pure orthologs. Instead, the term ‘orthogroup’ has come to represent a cluster of orthologs that may include closely related paralogs. InParanoid, EggNOG, and OMA are popular tools for assigning sequences to orthogroups using a ‘best-hit clique’ approach, where closed best-hit sub-graphs are identified in the dataset. While accurate within each sub-graph, these methods tend to be overly strict in their assignment; this causes an under-representation of actual orthologous relationships among many species.

Alternatively, Markov clustering (MCL) is very

© The Author 2016. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution. All rights reserved. For permissions, please email: journals.permissions@oup.com

**Table 1.** File format support provided by each BuddySuite module for reading (R) and writing (W).

Format	SeqBuddy	AlignBuddy	PhyloBuddy
Clustal	R & W <sup>†</sup>	R & W	None
EMBL <sup>‡</sup>	R & W	R <sup>†</sup> & W	None

<sup>†</sup>All sequences must be the same length

<sup>‡</sup>Supports rich sequence annotation

efficient at isolating more inclusive sub-graphs. OrthoMCL is one of the most popular MCL-based ortholog prediction methods, but it is prone to placing too many in-paralogs into orthogroups (i.e., it is less precise). In the current study we have increased the overall resolving power of de novo MCL-based orthogroup assignment with a number of novel enhancements, including refinement of the pairwise similarity metrics, using a supervised heuristic to dynamically select MCL parameters, recursively subdividing orthogroups, and testing putative orthogroups for best-hit cliques to maximize resolution.

## Methods

I used MAFFT (Katoh and Standley, 2013), because it’s awesome.

## Results

## Discussion

Spill some ink regarding in/out paralogs (Sonnhammer and Koonin, 2002; Tekaia, 2016).

## Conclusions

## Acknowledgments

This research was supported by the Intramural Research Program of the National Human Genome Research Institute, National Institutes of Health.

## References

- Katoh, K. and Standley, D. M. 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Molecular biology and evolution*, 30(4): 772–780.
- Sonnhammer, E. L. L. and Koonin, E. V. 2002. Orthology, paralogy and proposed classification for paralog subtypes. *Trends in genetics*, 18(12): 619–620.
- Tekaia, F. 2016. Inferring Orthologs: Open Questions and Perspectives. *Genomics insights*, 9: 17–28.
- van der Heijden, R. T., Snel, B., van Noort, V., and Huynen, M. A. 2007. Orthology prediction at scalable resolution by phylogenetic tree analysis. *BMC bioinformatics*, 8(1): 83.
- Vilella, A. J., Severin, J., Ureta-Vidal, A., Heng, L., Durbin, R., and Birney, E. 2009. EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees in vertebrates. *Genome research*, 19(2): 327–335.
- Wapinski, I., Pfeffer, A., Friedman, N., and Regev, A. 2007. Automatic genome-wide reconstruction of phylogenetic gene trees. *Bioinformatics*, 23(13): i549–58.