

PhylogeneticGraph
Program Documentation
Version 0.1

Ward C. Wheeler
Division of Invertebrate Zoology,
American Museum of Natural History,
200 Central Park West, New York, NY, 10024, USA;
wheeler@amnh.org

March 31, 2022

Contents

1	Introduction	3
2	Overview of Code base	3
3	Character reorganization and optimizing	3
3.1	Character weights	3
4	Command Parsing	3
5	Post-Order Graph Traversal	3
5.1	Trees	3
6	Pre-Order Graph Traversal	4
6.1	Trees	4
7	Adding New Character Types	4
7.1	Execution in Parallel	6
8	Bibliography	7

1 Introduction

This document contains descriptions of algorithms, procedures, data structures and other aspects of the source code documentation for the program PhylogeneticGraph (PhyG).

PhyG is a successor program to POY (Gladstein and Wheeler, 1997; Wheeler et al., 2005; Varón et al., 2008, 2010; Wheeler et al., 2013, 2015) <https://github.com/wardwheeler/POY5>, although a “complete” Haskell rewrite, optimized C (and even some assembler) was ported over from POY for pairwise alignment of small alphabet (j8) sequences. These functions are access via the Haskell FFI.

2 Overview of Code base

Source code structure.

3 Character reorganization and optimizing

Input data are passed through several fuincionts to:

- Rename taxa
- Exclude taxa
- Add missing data for taxa not present in all input files
- Check that input taxa and any input graphs contain teh same leaf set
- Data a reblocked if specified
- Static (Non-additive, additive, and matrix) characters are reorganized so that each class is put in a single (extensive) character (one for each type) type for each block.
- Non-additive characters are ‘bit-packed’ into new characters with state numbers $=2, \leq 4, \leq 5, \leq 8, \leq 64$, and > 64 . Invariant charcaters are filtered out.

3.1 Character weights

Static characters (Non-additive, additive, and matrix) with integer weights are reorganized by repeating th echaracter the number of times of its weight. This is to avoid alot of unnecessary ($\times 1$) operations. Non-integer weight characters are not reorganized or bit packed.

4 Command Parsing

5 Post-Order Graph Traversal

5.1 Trees

A decorated Graph (tree) is created for each character for each block for the graph. For exact characters, where no addition traversals are required, the specified or default outgroup sets the

direction of the graph. For non-exact (e.g. sequence) characters the best traversal rooting is stored for each character in each block although the cost of the graph is recalculated based on the best traversal (over all edges in the graph), the preliminary (post-order) states are not propagated back to the decorated graph (third field of phylogenetic graph). After the pre-order pass, the final states are propagated back. Vertices are not renumbered during the rerooting process, so indices remain unchanged.

Preliminary states (post-order) are determined for exact and non-exact characters as in [Wheeler \(2012\)](#).

6 Pre-Order Graph Traversal

6.1 Trees

Final state assignments of root vertices are set to the preliminary, post-order state. Final states are propagated back to the decorated graph (third field of phylogenetic graph). Vertices are not renumbered during the rerooting process, so indices remain unchanged.

Final states (pre-order) are determined for exact and non-exact characters as in [Wheeler \(2012\)](#). Currently final states for non-exact characters (e.g. sequence) are set as the median between the gapped preliminary state of the vertex and the final state of its parent (for a tree), ‘extra’ gaps in preliminary state are propagated to the gaped left and right descendant sequences, left, right, and parent final sequences should now line up and a 3-median can be calculated.

7 Adding New Character Types

Current character types include Additive, Non-Additive, Matrix, Slim Sequences, Wide Sequences, and Huge Sequences. Functions that branch on character types need to be updates and are found in:

- `GraphOptimization.Medians.hs`
 - `Median2Single`
 - `Median2SingleStaticIA`
 - `Union2Single`
 - `GetPrealignedUnion`
 - `getPreAligned2Median`
 - `median2SingleNonExact`
- `GraphOptimization.PreOrderFunctions.hs`
 - `updateCharacter`
 - `getCharacterDistFinal`
 - `setFinal`
 - `setPrelimToFinalCharacterData`
- `Commands.Transform.hs`

- transformCharacter
- Commands.CommandExecution.hs
 - makeCharLine
 - getCharacterLength
 - getCharCodeInfo
- Types.Types.hs
 - CharType
 - nonExactCharacterTypes
 - exactCharacterTypes
 - prealignedCharacterTypes
 - CharacterData
 - emptyCharacter
- Utilities.Utilities.hs
 - getCharacterInsertCost
 - makeBlockCharacterString
 - pairList2Fasta
 - splitBlockCharacters
 - getNumberNonExactCharacters
- Utilities.ThreeWayfunctions.hs
 - threeMedianFinal
- Support.Support.hs
 - subSampleStatic
 - makeSampledPairVect
- Input.Reorganize.hs
 - filterConst
 - getVariableChars
 - assignNewField
 - organizeBlockData'
- Input.FastAC.hs
 - Functions for sequence data processing on input
- Input.DataTransformation.hs

- These are for input—so not used by static approx
 - getMissingValue
 - getGeneralSequenceChar
 - getQualitativeCharacters—potentially depending on character features
 - createLeafCharacter
 - missingAligned
- Functions with “== NonAdd” etc will need extra cases for any new character type

7.1 Execution in Parallel

By default the program will execute multi-threaded based on the number processors available. By specifying the options ‘+RTS -NX -RTS’ where ‘X’ is the number of processors offered to the program. These are specified after the program as in (for 4 parallel threads):

PhyGraph +RTS -N4 -RTS other options...

Parallel code options are set using a parmap-type strategy throughout the code. This is usually specified by myParListChunkRDS from the PArallelUtilities module in PhyloLibs. The basic definitions of this functionality are found in ParallelUtilities.hs

Acknowledgments

The author would like to thank DARPA SIMPLEX N66001-15-C-4039, the Robert J. Kleberg Jr. and Helen C. Kleberg foundation grant “Mechanistic Analyses of Pancreatic Cancer Evolution”, and the American Museum of Natural History for financial support.

8 Bibliography

- Gladstein, D. S. and Wheeler, W. C. 1997. POY version 2.0. program and documentation available at <http://research.amnh.org/scicomp/projects/poy.php>. American Museum of Natural History, New York.
- Varón, A., Vinh, L. S., Bomash, I., and Wheeler, W. C. 2008. Poy 4.0. American Museum of Natural History. <http://research.amnh.org/scicomp/projects/poy.php>.
- Varón, A., Vinh, L. S., and Wheeler, W. C. 2010. POY version 4: Phylogenetic analysis using dynamic homologies. *Cladistics* 26:72–85.
- Wheeler, W. C. 2012. Systematics: A course of lectures. Wiley-Blackwell.
- Wheeler, W. C., Gladstein, D. S., and De Laet, J. 1996-2005. POY version 3.0. program and documentation available at <http://research.amnh.org/scicomp/projects/poy.php> (current version 3.0.11). documentation by D. Janies and W. C. Wheeler. commandline documentation by J. De Laet and W. C. Wheeler. American Museum of Natural History, New York.
- Wheeler, W. C., Lucaroni, N., Hong, L., Crowley, L. M., and Varón, A. 2013. POY version 5.0. American Museum of Natural History. <http://research.amnh.org/scicomp/projects/poy.php>.
- Wheeler, W. C., Lucaroni, N., Hong, L., Crowley, L. M., and Varón, A. 2015. POY version 5: Phylogenetic analysis using dynamic homologies under multiple optimality criteria. *Cladistics* 31. 189-196.