

A peer-reviewed version of this preprint was published in PeerJ on 12 May 2016.

[View the peer-reviewed version](https://peerj.com/articles/2038) (peerj.com/articles/2038), which is the preferred citable publication unless you specifically need to cite this preprint.

Stephens TG, Bhattacharya D, Ragan MA, Chan CX. (2016) PhySortR: a fast, flexible tool for sorting phylogenetic trees in R. PeerJ 4:e2038
<https://doi.org/10.7717/peerj.2038>

1 **PhySortR: a fast, flexible tool for sorting phylogenetic trees**
2 **in R**

3 Timothy G. Stephens¹, Debashish Bhattacharya², Mark A. Ragan¹ and Cheong Xin Chan¹

4 ¹ARC Centre of Excellence in Bioinformatics, and Institute for Molecular Bioscience, The
5 University of Queensland, Brisbane, QLD 4072, Australia

6 ²Department of Ecology, Evolution and Natural Resources, Rutgers University, New Brunswick,
7 NJ 08901, U.S.A.

8 Corresponding author:

9 Cheong Xin Chan¹

10 Institute for Molecular Bioscience, The University of Queensland, Brisbane, QLD 4072,
11 Australia

12 Email address: c.chan1@uq.edu.au

13 **Abstract**

14 A frequent bottleneck in interpreting phylogenomic output is the need to screen often thousands
15 of trees for features of interest, such as robust clades of specific taxa, as evidence of
16 monophyletic relationship and/or reticulated evolution. Here we present PhySortR, a fast,
17 flexible R package for sorting phylogenetic trees. Unlike existing utilities, PhySortR allows for
18 identification of both exclusive and non-exclusive clades uniting the target taxa, with
19 customisable options to assess clades within the context of the whole tree. PhySortR is a
20 command-line tool that is freely available, highly scalable, and easily automatable.

21 **Main Text**

22 Phylogenomics increasingly involves screening of thousands of phylogenetic trees using
23 specialised sorting algorithms for features of interest, e.g. strongly supported monophyletic
24 relationships of taxa in question (i.e. the “target” taxa). Currently available tree-sorting utilities,
25 e.g. PhyloSort (Moustafa and Bhattacharya 2008) and SICLE (Deblasio and Wisecaver 2013)
26 screen a set of trees (in Newick format) for the presence of clades that unite a set of user-defined
27 target taxa based on clade support (that exceeds a defined threshold), and sort these trees
28 accordingly. However, these tools do not consider the proportion of non-target taxa and overall
29 taxon composition in a tree during the sorting process.

30 Here we present PhySortR, a fast, flexible R package for screening and sorting phylogenetic
31 trees (Cardona et al. 2008). The package provides the quick and highly flexible *sortTrees*
32 function, allowing for screening (within a tree) for “Exclusive” clades that contain only the target
33 taxa and/or “Non-Exclusive” clades that include a defined portion of non-target taxa. The

34 algorithm and all available options of PhySortR are described in detail in Supplementary Text S1
35 and Supplementary Figures S1 and S2.

36 PhySortR allows the user to specify one or more target taxa using the *target.groups* argument,
37 providing that the taxa (tree-tip labels in the Newick files) are named consistently across all input
38 trees; this is a simple string-matching exercise, i.e. the terms specified here determine taxon-level
39 resolution of targets. The minimum support for a clade (*min.support*) can refer to bootstrap,
40 Bayesian posterior probability, or any other measure of support. Existing utilities such as
41 PhyloSort (Moustafa and Bhattacharya 2008), when assessing the relationship between two
42 target groups in a 100-taxon tree, would identify both (a) a robust two-member clade (one from
43 each target group), and (b) a robust 50-member clade containing multiple taxa from both groups,
44 as “Exclusive”, although (b) is likely the more-convincing evidence of a close association
45 between the targets and of greater biological significance. To address this issue, PhySortR allows
46 one to define *min.prop.target*, the minimum required proportion of target(s) present in a clade
47 relative to the total number of target(s) found in a tree.

48 Alternatively, one might wish to screen for a robust clade that contains the target groups and a
49 small proportion of “interrupting” non-target taxa, e.g. a 41-member clade consisting of 20 taxa
50 each from the targets Rhodophyta and Viridiplantae, as well as a stramenopile taxon (a non-
51 target such as a diatom). Whereas the diatom is “interrupting” the otherwise exclusive clade of
52 Rhodophyta + Viridiplantae, the association between the targets is still of interest and the diatom
53 presence might be readily explained by lateral gene transfer (LGT) due to plastid endosymbiosis.
54 Composite clades such as these are considered “Non-Exclusive” (Chan et al. 2011) and are not
55 identified by existing tree-sorting tools. The concept of exclusivity (Fig. 1A) versus non-

56 exclusivity (Fig. 1B) of clades in tree sorting has proven crucial in a number of genome-wide
57 studies that have investigated the impact of LGT on the evolution of diverse algal and protist
58 phyla (e.g., Chan et al. 2011; Curtis et al. 2012; Price et al. 2012; Bhattacharya et al. 2013).
59 PhySortR identifies both types of clades by default. The user has the option to define a “Non-
60 Exclusive” clade based on the proportion of target versus non-target taxa using the option
61 *clade.exclusivity*. At the default setting (*clade.exclusivity* = 0.9), the minimum proportion of
62 target taxa within a “Non-Exclusive” clade is 0.9, thus the maximum proportion of non-target
63 taxa allowed in the clade is 0.1 (i.e. 1 minus 0.9). This option accepts any value <1.0, and is only
64 applicable for sorting “Non-Exclusive” clades (see Supplementary Fig. S2); at 1.0 (no non-target
65 taxa allowed), the clade is considered “Exclusive”. PhySortR is an R implementation of the
66 algorithm deployed in Chan et al. (2011) that has been adopted in other phylogenomic studies
67 (Curtis et al. 2012; Price et al. 2012; Bhattacharya et al. 2013). The R platform is open source,
68 platform-independent, and broadly accessible to researchers, with continued support by the R
69 Core Team; functional modularity and the command-line interface enable batch automation and
70 workflow integration.

71 The runtime of PhySortR is dependent on the number of trees (N) to be sorted and the number of
72 taxa (X) within a tree. We benchmarked PhySortR using two simulated datasets to assess runtime
73 (t): one varying N at fixed $X = 100$ (Supplementary Data S1), another varying X at fixed $N =$
74 1000 (Supplementary Data S2). We required 20% of X within an “Exclusive” clade in a tree. We
75 observed that the runtime scales linearly with N (Figure 1C) and superlinearly with X (Figure
76 1D). In the extreme case, sorting through 10000 trees took <500 seconds (~8.3 minutes). We
77 observed negligible differences in t with negative controls (trees containing no identifiable
78 clades) as input, compared to the test set in Figure 1D.

79 PhySortR incorporates existing functionalities and data structures in the commonly used
80 phylogenetic packages *ape* (Paradis et al. 2004) and *phytools* (Revell 2012), allowing for
81 streamlined interoperability within the R environment. Whereas *ape* and *phytools* accept only
82 traditional Newick as input, PhySortR accepts tree files in both traditional and extended Newick
83 formats (Cardona et al. 2008). The function *convert.eNewick* in isolation can be used as a
84 general-purpose tool for converting extended Newick into traditional Newick format. PhySortR
85 is freely available from the Comprehensive R Archive Network (cran.r-project.org). See also
86 <https://cloudstor.aarnet.edu.au/plus/index.php/s/StQRSJBmcSZKd7y>.

87 **Supplementary Material**

88 Supplementary Text, Supplementary Figures S1-S2, Supplementary Data S1-S2.

89 **Acknowledgments**

90 This work was supported by the Australian Research Council Discovery Project (DP150101875)
91 grant awarded to M.A.R., C.X.C, and D.B. C.X.C. is supported by Great Barrier Reef
92 Foundation Bioinformatics Fellowship awarded to M.A.R. D.B. acknowledges support from the
93 National Science Foundation (1004213).

94 **References**

95 Bhattacharya D, Price DC, Chan CX, Qiu H, Rose N, Ball S, Weber APM, Arias MC, Henrissat
96 B, Coutinho PM, Krishnan A, Zauner S, Morath S, Hilliou F, Egizi A, Perrineau MM, Yoon HS.
97 2013. Genome of the red alga *Porphyridium purpureum*. *Nat. Commun.* 4:1941.
98 Cardona G, Rossello F, Valiente G. 2008. Extended Newick: it is time for a standard
99 representation of phylogenetic networks. *BMC Bioinformatics* 9:532.

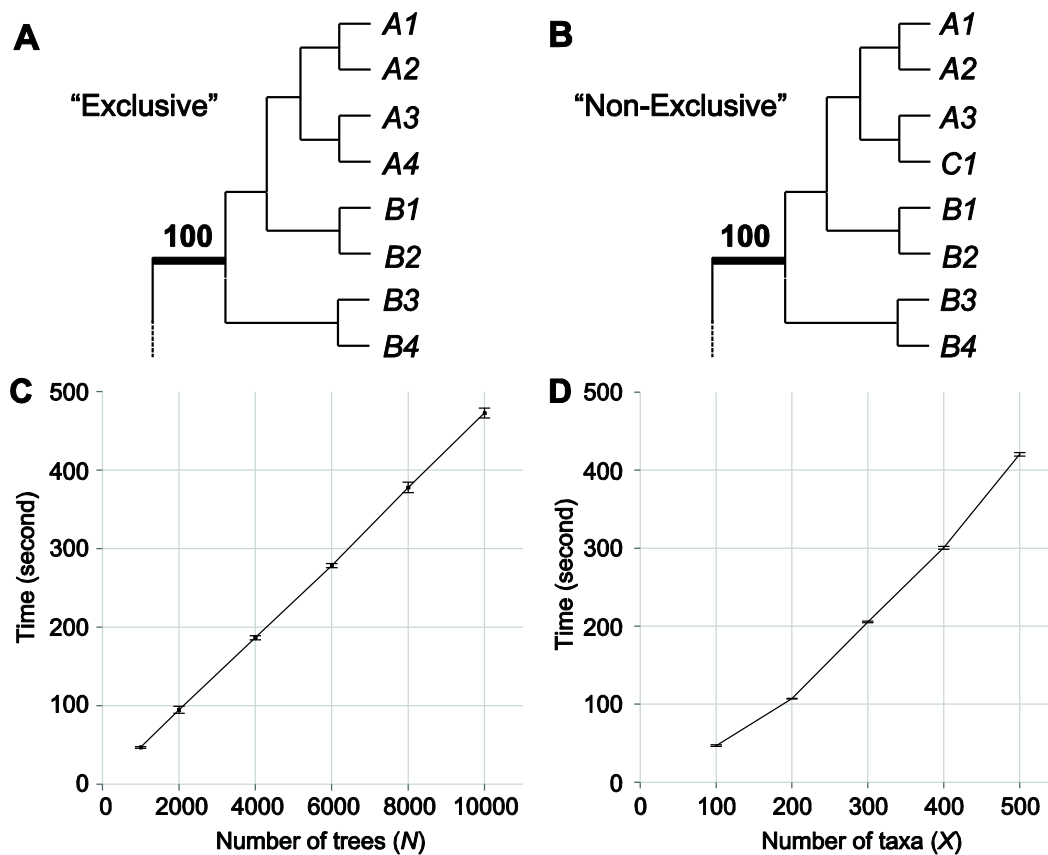
- 100 Chan CX, Yang EC, Banerjee T, Yoon HS, Martone PT, Estevez JM, Bhattacharya D. 2011. Red
101 and green algal monophyly and extensive gene sharing found in a rich repertoire of red algal
102 genes. *Curr. Biol.* 21:328-333.
- 103 Curtis BA, Tanifuji G, Burki F, Gruber A, Irimia M, Maruyama S, Arias MC, Ball SG, Gile GH,
104 Hirakawa Y, Hopkins JF, Kuo A, Rensing SA, Schmutz J, Symeonidi A, Elias M, Eveleigh
105 RJM, Herman EK, Klute MJ, Nakayama T, Obornik M, Reyes-Prieto A, Armbrust EV, Aves SJ,
106 Beiko RG, Coutinho P, Dacks JB, Durnford DG, Fast NM, Green BR, Grisdale CJ, Hempel F,
107 Henrissat B, Hoppner MP, Ishida KI, Kim E, Koreny LK, Kroth PG, Liu Y, Malik SB, Maier
108 UG, McRose D, Mock T, Neilson JAD, Onodera NT, Poole AM, Pritham EJ, Richards TA,
109 Rocap G, Roy SW, Sarai C, Schaack S, Shirato S, Slamovits CH, Spencer DF, Suzuki S, Worden
110 AZ, Zauner S, Barry K, Bell C, Bharti AK, Crow JA, Grimwood J, Kramer R, Lindquist E,
111 Lucas S, Salamov A, McFadden GI, Lane CE, Keeling PJ, Gray MW, Grigoriev IV, Archibald
112 JM. 2012. Algal genomes reveal evolutionary mosaicism and the fate of nucleomorphs. *Nature*
113 492:59-65.
- 114 Deblasio D, Wisecaver J. 2013. SICLE: A high-throughput tool for extracting evolutionary
115 relationships from phylogenetic trees. arXiv:1303.5785.
- 116 Moustafa A, Bhattacharya D. 2008. PhyloSort: a user-friendly phylogenetic sorting tool and its
117 application to estimating the cyanobacterial contribution to the nuclear genome of
118 *Chlamydomonas*. *BMC Evol. Biol.* 8:6.
- 119 Paradis E, Claude J, Strimmer K. 2004. APE: Analyses of Phylogenetics and Evolution in R
120 language. *Bioinformatics* 20:289-290.

121 Price DC, Chan CX, Yoon HS, Yang EC, Qiu H, Weber APM, Schwacke R, Gross J, Blouin
122 NA, Lane C, Reyes-Prieto A, Durnford DG, Neilson JAD, Lang BF, Burger G, Steiner JM,
123 Loffelhardt W, Meuser JE, Posewitz MC, Ball S, Arias MC, Henrissat B, Coutinho PM, Rensing
124 SA, Symeonidi A, Doddapaneni H, Green BR, Rajah VD, Boore J, Bhattacharya D. 2012.
125 *Cyanophora paradoxa* genome elucidates origin of photosynthesis in algae and plants. *Science*
126 335:843-847.

127 Revell LJ. 2012. phytools: an R package for phylogenetic comparative biology (and other
128 things). *Methods Ecol. Evol.* 3:217-223.

129

130

131 **Figure**

132

133 **Figure 1.** The concept of clade exclusivity and benchmarking results of PhySortR using
 134 simulated data. The schematic diagram of an "Exclusive" clade shown in (A) and consists solely
 135 of taxa from targets A and B, whereas the "Non-Exclusive" clade shown in (B) consists of
 136 targets A and B plus an "interrupting" taxon from group C; each clade has strong bootstrap
 137 support at 100%. The mean runtime (t) of PhySortR is shown for analysis across datasets (C)
 138 with different numbers of trees, N (Supplementary Data S1), and (D) with different numbers of
 139 taxa per tree, X (Supplementary Data S2). Values of t are mean across 100 replicates, error bars
 140 indicate the standard deviation of the mean.