

# PhyBin: binning trees by topology

A major goal of many evolutionary analyses is to determine the true evolutionary history of an organism. Molecular methods that rely on the phylogenetic signal generated by a few to a handful of loci can be used to approximate the evolution of the entire organism but fall short of providing a global, genome-wide, perspective on evolutionary processes. Indeed, individual genes in a genome may have different evolutionary histories. Therefore, it is informative to analyze the number and kind of phylogenetic topologies found within an orthologous set of genes across a genome. Here we present PhyBin: a flexible program for clustering gene trees based on topological structure. PhyBin can generate *bins* of topologies corresponding to exactly identical trees or can utilize Robinson-Fould's distance matrices to generate *clusters* of similar trees, using a user-defined threshold. Additionally, PhyBin allows the user to adjust for potential noise in the dataset (as may be produced when comparing very closely related organisms) by pre-processing trees to collapse very short branches or those nodes not meeting a defined bootstrap threshold. As a test case, we generated individual trees based on an orthologous gene set from 10 *Wolbachia* species across four different supergroups (A-D) and utilized PhyBin to categorize the complete set of topologies produced from this dataset. Using this approach, we were able to show that although a single topology generally dominated the analysis, confirming the separation of the supergroups, many genes supported alternative evolutionary histories. Because PhyBin's output provides the user with lists of gene trees in each topological cluster, it can be used to explore potential reasons for discrepancies between phylogenies including homoplasies, long-branch attraction, or horizontal gene transfer events.

1 **Title:** *PhyBin: binning trees by topology*

2 Authors: <sup>1</sup>Ryan R. Newton and <sup>2\*</sup>Irene L.G. Newton

3 \*corresponding author

4 Affiliations: 1 = School of Informatics and Computing, Indiana University, Bloomington, IN;

5 2=Department of Biology, Indiana University, Bloomington, IN

# 6 **Abstract:**

7 A major goal of many evolutionary analyses is to determine the true evolutionary history of an  
8 organism. Molecular methods that rely on the phylogenetic signal generated by a few to a handful  
9 of loci can be used to approximate the evolution of the entire organism but fall short of providing  
10 a global, genome-wide, perspective on evolutionary processes. Indeed, individual genes in a  
11 genome may have different evolutionary histories. Therefore, it is informative to analyze the  
12 number and kind of phylogenetic topologies found within an orthologous set of genes across a  
13 genome. Here we present PhyBin: a flexible program for clustering gene trees based on  
14 topological structure. PhyBin can generate *bins* of topologies corresponding to exactly identical  
15 trees or can utilize Robinson-Foulds distance matrices to generate *clusters* of similar trees, using  
16 a user-defined threshold. Additionally, PhyBin allows the user to adjust for potential noise in the  
17 dataset (as may be produced when comparing very closely related organisms) by pre-processing  
18 trees to collapse very short branches or those nodes not meeting a defined bootstrap threshold.  
19 As a test case, we generated individual trees based on an orthologous gene set from 10 *Wolbachia*  
20 species across four different supergroups (A-D) and utilized PhyBin to categorize the complete  
21 set of topologies produced from this dataset. Using this approach, we were able to show that  
22 although a single topology generally dominated the analysis, confirming the separation of the  
23 supergroups, many genes supported alternative evolutionary histories. Because PhyBin's output  
24 provides the user with lists of gene trees in each topological cluster, it can be used to explore  
25 potential reasons for discrepancies between phylogenies including homoplasies, long-branch  
26 attraction, or horizontal gene transfer events.

27 Availability: PhyBin is a standalone open-source program available:

28 <http://hackage.haskell.org/package/phybin>

# 29 **Introduction:**

30 The advent of genomic sequencing has produced a large amount of data available for phylogenetic  
31 analysis and many researchers have attempted to utilize the phylogenetic signal found across the  
32 bacterial genome to develop species trees ([Daubin, Gouy et al. 2001](#); [Sicheritz-Ponten and](#)  
33 [Andersson 2001](#); [Daubin, Moran et al. 2003](#); [Baptiste, Boucher et al. 2004](#); [Zhaxybayeva,](#)  
34 [Gogarten et al. 2006](#); [Ellegaard, Klasson et al. 2013](#)). What has become clear from these analyses  
35 is that significant fractions of bacterial genomes do not follow the evolutionary history of their  
36 resident genome ([Baptiste, Boucher et al. 2004](#)). These rogue genes are potentially undergoing  
37 evolutionary processes distinct from those felt by the rest of the resident genome or have arrived  
38 there via horizontal gene transfer events. In order, then, to understand the evolution of the  
39 genome, it would be useful to achieve an understanding of the evolution of each gene in the  
40 genome. Previous work by Sicheritz-Ponten and Andersson presented scripts combined the  
41 existing utilities BLAST, Clustalw, Paup 4.0\* to provide a complete pipeline from genome to tree-  
42 binning analysis ([Sicheritz-Ponten and Andersson 2001](#)). These kinds of complete solutions are

convenient but constrain the user to the specific utilities chosen by the authors for alignment and phylogeny generation.

Here we present PhyBin, a computer program aimed at binning precomputed sets of non-reticulated trees in Newick format, a file format produced by the majority of tree building software. PhyBin is a utility rather than a complete solution; it can serve as a component in many genomics pipelines, and provides a useful addition to the landscape of tools for dissecting and visualizing large numbers of trees. After the user applies their chosen ortholog prediction and tree-building algorithms, PhyBin offers a quick way to visualize and browse the different evolutionary histories, either binned by topology and sorted by bin size, or in the form of a full hierarchical clustering based on Robinson-Foulds distance: i.e. a *tree of trees*.

## Method and Implementation:

### *Generating orthologous sets and input trees*

Genomic sequences were downloaded from NCBI Microbial Genome Projects. The *Wobachia* species complex is made up of several major clades, called supergroups, designated by alphabetical letters ([Baldo and Werren 2007](#)). Accession numbers for the genomes analyzed here include: wUni and wVitA (submissions pending to genbank's ncbi), wBm (NC\_006833.1), wPip-Pel (NC\_010981.1), wHa (NC\_021089.1), wRi (NC\_012416.1), wMel (NC\_002978.6), wNo (NC\_021084.1), wAlbB (CAGB000000000.1), wBm (NC\_006833.1), wOo (NC\_018267.1). Orthologous gene sets were determined by Reciprocal Smallest Distance (RSD) algorithm ([Wall, Fraser et al. 2003](#)) with a  $10^3$  cutoff for significance threshold and alignment length threshold of 80%. Orthologs were then aligned using ClustalW ([Larkin, Blackshields et al. 2007](#)) and ML trees were generated using RAxML ([Stamatakis 2006](#)). The Newick format trees that resulted were used as input to PhyBin. The number of orthologous genes identified in this manner across all 10 taxa was 503.

## Description of the Program:

PhyBin is a standalone command-line program, portable across all major operating systems. It runs in batch-mode and is easily usable from scripts. PhyBin has two major modes: it can run very quickly and classify identical tree topologies into bins, or it can compute the distance ([Robinson and Foulds 1981](#)) between all pairs of trees and use that distance matrix to produce a configurable clustering of trees.

### *Fast Binning Mode*

The key algorithm PhyBin performs in this mode is tree normalization, computing a rooted, ordered *normal form* for all inputs (which are labeled, unrooted, unordered tree topologies). Previous work in this area has described a number of viable normal forms ([Chi, Yang et al. 2005](#)). Conversion to a normal form ensures that all equivalent unrooted trees are converted into the same rooted tree, with a canonical root chosen. After conversion, the rooted trees are much faster to compare for equality than the unrooted trees would be, which enables fast binning.

PhyBin chooses the following strategy: it attempts to order subtrees by *weight* (number of tree nodes) and select the root node which is most balanced by weight (not depth)---that is, which minimizes the maximum weight of any child of the root. Node labels are used only to "break ties" between equally weighted subtrees, or equally balanced roots. Because input trees in Newick format are typically labeled only on the *leaves* (taxa), PhyBin generates labels for intermediate nodes in the tree by creating a set of all the leaves contained in that subtree, given a root to determine up/down direction. This set can be represented as a bit-vector and is also a key ingredient of computing Robinson-Foulds distance, which relies on identifying all such subsets (i.e. bipartitions induced by the tree). With labels for all nodes, equally weighted subtrees are ordered by label, and ties between potential roots are broken by comparing the labels of their children.

Once input trees are normalized, testing for equality of two trees is as simple as comparing their representation in memory (a single, linear traversal). Normalization itself appears expensive due to the cost of labeling interior nodes with all leaves under them ( $O(N * I)$  for  $N$  taxa and  $I$  interior nodes), compounded by the fact that each intermediate node may have to consider each of its neighbors as a possible root and relabel itself  $b$  times in a tree of maximum branching factor  $b$ , yielding an  $O(N*I*b)$  asymptotic cost. However, in binning mode PhyBin runs much faster in the average case. One feature that enables PhyBin's efficiency is that it computes tree metadata---interior labels and "balanced" ratings---*lazily*, that is, on demand. Only when "tie breaking" is necessary between equally-weighted subtrees is an interior label computed at all. Likewise, only nodes "near the center" of the unrooted tree need to be considered for root status, those near the leaves need never be scored for balance.

After normalization, PhyBin performs binning, which amounts to inserting all normalized trees into a data structure *indexed* by tree topology. We define a total order over normalized trees (made possible by labels), and thereby represent the table of bins as a size-balanced binary tree supporting  $O(\log(n))$  insertion times. A hash-table would be an alternative, but the tree representation allows us to insert trees into the table without evaluating (forcing) unnecessary interior labels in the normal forms, whereas hashing requires traversing the entirety of each normalized tree to compute its hash. When execution completes, the contents of each bin are written out to disk, in addition to a visualization of a representative average tree for that topology, computed by averaging branch lengths of the bin members.

### *Pre-Processing Data*

PhyBin helps users extract a clean dataset and detect problems with the data, such as trees with mismatching numbers of taxa. In order to facilitate comparisons across trees with different taxon names (i.e. gene names), PhyBin can extract portions of designations or use a separate table of rules for mapping genes to taxa. In addition, PhyBin can restrict its analyses to a subset of taxon, ignoring others (--prune).

A problem with the simple binning approach is that it is fragile to minor differences in trees caused by noise (e.g. short length branches with high variability). This becomes increasingly problematic with large numbers of taxa, especially when closely related taxa (different strains) are compared. Fortunately, a simple pre-processing step that addresses this problem: PhyBin provides an option to collapse branches under two different conditions, a length threshold (for

example, a length threshold of 0.01 would collapse all branches less than 0.01, in their place inserting a star topology) or a bootstrap support threshold (such that nodes with less than that threshold would be collapsed and the branch lengths from the taxa to the parent node would be added).

#### *Full Clustering Mode using Robinson-Foulds Distance Matrix:*

PhyBin reimplements the HashRF algorithm for full all-to-all Robinson Foulds distance ([Sul and Williams 2007](#)), which is significantly faster than computing the distance matrix with repeated comparison of individual trees (e.g. PAUP ([Swofford and Sullivan 2009](#))). The HashRF algorithm is fast for today's data sizes (e.g. hundreds of taxa and thousands of trees), but scales much more poorly than the basic binning algorithm at significantly larger sizes.

Because ortholog sets across different genomic comparisons will produce trees with different taxon memberships (as a result of paralogs or gene losses), a user may consider decomposing their trees with other software solutions (such as treeKO, ([Marcet-Houben and Gabaldon 2011](#))). Further, PhyBin is also capable of directly comparing these trees with different numbers of taxa using the leaf pruning method implemented in STRAW ([Shaw, Ruan et al. 2013](#)). Specifically, in comparing trees with different taxa (--tolerant mode), the program first removes taxa that are not contained within each tree. If the taxon removed is in a polytomy, the parent and sister taxon are unchanged. However, in a binary node, taxon pruning would remove the intermediate node, retaining the branch lengths from the ancestor to the unpruned taxon. The --tolerant mode comes with a cost, however, as the more efficient HashRF algorithm cannot be used; instead Phybin falls back to the earlier PAUP-style algorithm.

A distance matrix alone is not directly useful for exploring the direct relationships between different gene trees. Thus, PhyBin uses the Robinson-Foulds distance matrix to compute a clustering of tree topologies, similar to the output of the simple binning mode, but able to identify trees that are merely *similar*, although not identical. A hierarchical clustering method is used. (If the user desires a different clustering method, they may use the distance matrix produced by PhyBin as input to a different processing pipeline.)

With the hierarchical clustering method, there remain several clustering options to configure. The choice of clustering options can dramatically alter bin membership (Supplementary Table 1), and running with several different options is a good way to get a sense for the range of possible outcomes. Specifically, the user may define the edit distance tolerated within clusters by providing a threshold, and may choose single, complete, or UPGMA linkage for clustering. Also if desired, rather than viewing a *flat* clustering of trees, the user may directly view a hierarchical clustering of the trees as a dendrogram. We believe PhyBin is the first program to date to provide this *tree-of-trees* output.

#### *Output Formats:*

PhyBin is meant to be used in scripts and by other programs. Every output produced by PhyBin goes into a separate, simple text file---for example, the consensus tree for each cluster and the

160 Robinson-Foulds distance matrix. Visualizations are produced separately and automatically in  
161 PDF files.

## 162 *Performance:*

163 There are very large differences in performance between existing programs for computing  
164 Robinson-Foulds distance matrices. The fundamental data-structures in this problem domain are  
165 sets and finite maps, for which there are many alternate representations (bit vectors, hash tables,  
166 balanced trees, etc), providing a large space of possible implementations to explore. The sharpest  
167 contrast is between those programs that directly compare individual pairs of trees (PAUP,  
168 DendroPy), vs. those that insert all tree's bipartitions into a global structure and summarize it as a  
169 separate phase (e.g. HashRF). The later approach achieves much better cache locality.

170 PhyBin is written in a very high level language, Haskell, which supports radical forms of  
171 optimization, including safe semi-automatic parallelism. PhyBin uses purely functional  
172 (immutable) data-structures for representing trees and their bipartitions; in particular it relies  
173 heavily on the balanced-tree implementations `Data.Map` and `Data.Set` from the standard  
174 library. Nevertheless, when computing a matrix for a 150-taxa, 100-tree test (Table 1), PhyBin is  
175 82 times faster than Philip (ANSI C) and 47.5 times faster than DendroPy (Python). However,  
176 PhyBin is still slower than HashRF by a factor of 2.8X-4.8X. HashRF was the first  
177 implementation that introduced high-performance techniques for RF matrices, and it introduced  
178 the algorithm on which PhyBin's implementation is based.

179 Unfortunately, the more widely used software (PAUP, DendroPy, Philip, etc), remains slow.  
180 HashRF, the currently available fast alternative, is delicate and must be used carefully (for  
181 example, an extra character of whitespace in the input file results in a segmentation fault with no  
182 error message in version 6.0.1). Additionally, because HashRF provides only the core RF-  
183 distance computation, other tools are required for a biologist to be able to derive any conclusions  
184 from the output.

185 As a final note on performance, PhyBin was straightforward to parallelize (using our "LVar"  
186 parallelism library) and achieves a 2.54X parallel speedup at four cores, and peaks at a 3.11X  
187 speedup at eight cores, making it a bit faster than HashRF on our target platform (Table 1).  
188 Future work will focus on reducing contention on shared data structures to improve scaling.

## 189 **Results and Discussion:**

190 We used PhyBin to identify how many phylogenies within the *Wolbachia* orthologous gene set  
191 support the supergroup divisions proposed by multi-locus sequence typing ([Baldo and Werren](#)  
192 [2007](#)). For comparative purposes in this analysis, a phylogeny for these 10 taxa was created using  
193 the concatenated, orthologous gene set (Figure 1A). In actuality, PhyBin does not require an  
194 expectation for tree topology and searches through tree space for distinct topological categories.  
195 As an illustration of PhyBin's ability to reduce the noise in a dataset produced by small branch  
196 lengths (i.e., closely related taxa), we used the program in *binning mode* on the set of *Wolbachia*  
197 orthologs under increasing branch length thresholds (Table 2). We chose a threshold of 0.01 for  
198 our dataset as the average branch length over the entire set of validated trees was 0.04 with  
199 minimum and maximum branch lengths of 0 and 2.31, respectively. Using this threshold, in



*binning* mode, the largest bin contains a topology that agrees with that of the published supergroups (133 members in largest bin, 175 total bins, Table 2, Figure 1B). However, 174 other potential topologies exist in the dataset with 129 alternative topologies supported by only a single ortholog tree (Table 2).

In order to better explore this tree set, we took advantage of PhyBin's ability to generate a distance matrix for all trees. By calculating the Robinson-Foulds (RF) distance between all trees, we can better assess the differences between clusters in the tree dataset. For example, by increasing the RF-distance threshold to 2 and using the average-neighbor clustering algorithm to group our trees, the number of clusters drops dramatically to only 77 with the largest cluster containing a majority (72%) of genes. Again, this topology agrees with the published supergroup data and our result from the binning approach (Figure 1C). Increasing the RF-distance threshold further provides increasing stringency in the detection of aberrant phylogenies – topologies not falling into the largest cluster at larger distance thresholds are likely to represent genes of interest in comparing evolutionary trajectories of these supergroups.

To test this hypothesis, we identified those *Wolbachia* genes that continue to display alternative evolutionary histories (that is, falling outside of the majority) even when clustering trees using increasingly large RF distances (Figure 2B, Table 3). As expected, a large number of distinct topologies are not inconsistent with the supergroup clades (65 distinct tree clusters do not support the major topology, using an RF-distance threshold of 1 and a branch length cutoff of 0.02, Table 3, Figure 2B). We further investigated the ortholog set supporting the dissolution of supergroup A (Table 4). Interestingly, a majority of these orthologs are predicted to be secreted (using the Effective database predictions of sec signal or eukaryotic domains ([Jehl, Arnold et al. 2011](#)), suggesting that perhaps interaction with the host would drive some of these orthologs in a different evolutionary direction compared to their resident genome. Another test of PhyBin's ability to detect orthologs under different evolutionary pressures would focus on the *Wolbachia* prophage, a mobile genetic element known to undergo horizontal transmission between strains ([Bordenstein and Wernegreen 2004](#); [Chafee, Funk et al. 2010](#); [Kent and Bordenstein 2010](#); [Kent, Salichos et al. 2011](#)). However, these phage orthologs do not occur across all of our 10 taxa included here and are therefore not suitable for testing support for the supergroups.

In conclusion, we PhyBin is a new software program that efficiently and quickly groups phylogenies either by strict topological congruence or by clustering using RF distance. We believe that this tool, due to its ease of use, its speed, and informative output, will be of interest to evolutionary biologists and bioinformaticians alike.

## Figure Legends:

**Figure 1.** In each of two modes (*full clustering* and *binning*) PhyBin is able to correctly recover the expected topology for the *Wolbachia pipientis* orthologs used herein. (A) Concatenated phylogeny based on 508 genes (using RAxML GTRGAMMA, bootstrap support based on 10,000 replicates). The four major supergroups are highlighted and denoted. (B) These same groups are recovered when PhyBin is run in either *binning* mode or (C) *full clustering* mode.

**Figure 2.** Robinson-Foulds distance matrices produced by PhyBin are also visualized as a dendrogram by the software. (A) A *tree of trees* for the *Wolbachia* ortholog set (508 trees), clustered using an edit distance of 0, where identical topologies (nodes – grey ovals) are shown connected by a red line. Length of the branches connecting each node is proportional to the RF distance. (B) This dendrogram is simplified by increasing the RF distance at which the

243 trees are clustered (shown  $RF = 3$ ). The top 10 clusters and their support different topologies are colored as indicated  
 244 in the legend (with largest bin size for each cluster cluster in parentheses).



## 245 References

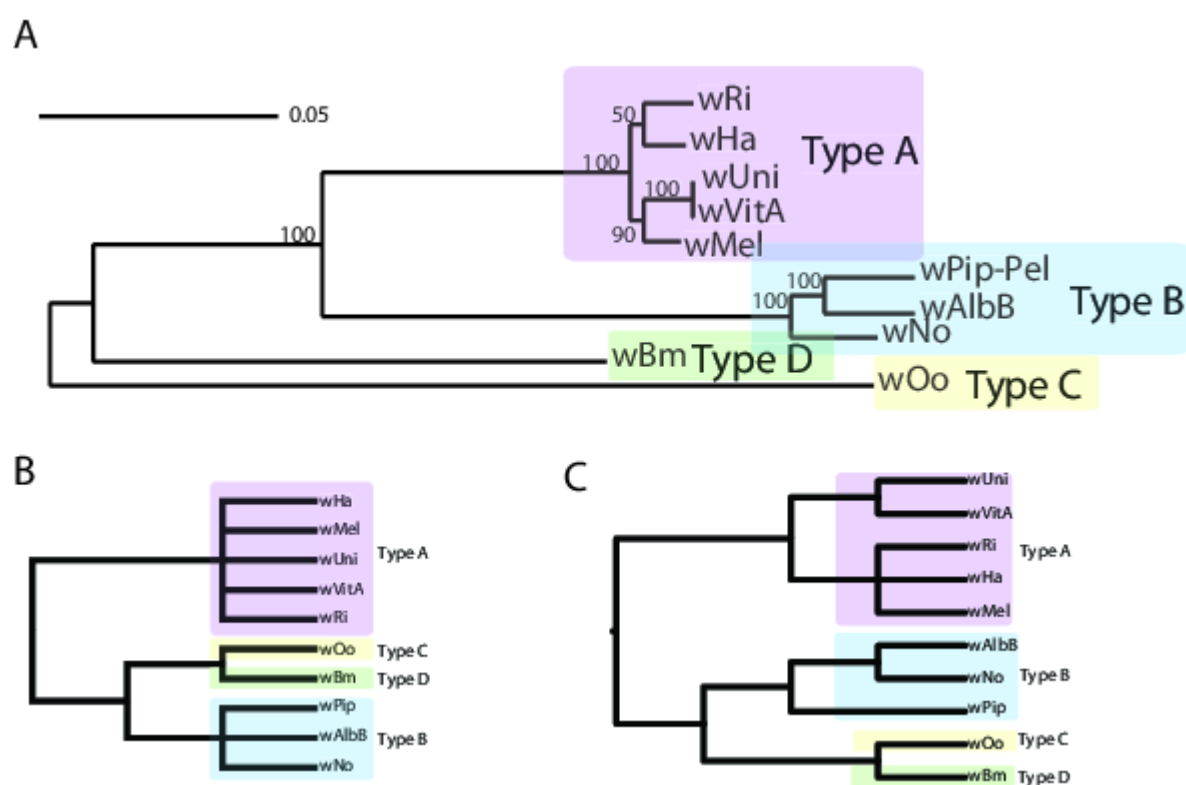
- 246 Baldo, L. and J. H. Werren (2007). "Revisiting Wolbachia supergroup typing based on WSP:  
247 spurious lineages and discordance with MLST." Curr Microbiol **55**(1): 81-87.
- 248 Baptiste, E., Y. Boucher, J. Leigh and W. F. Doolittle (2004). "Phylogenetic reconstruction and  
249 lateral gene transfer." Trends Microbiol **12**(9): 406-411.
- 250 Bordenstein, S. R. and J. J. Wernegreen (2004). "Bacteriophage flux in endosymbionts  
251 (Wolbachia): Infection frequency, lateral transfer, and recombination rates." Molecular  
252 Biology and Evolution **21**(10): 1981-1991.
- 253 Chafee, M. E., D. J. Funk, R. G. Harrison and S. R. Bordenstein (2010). "Lateral Phage Transfer  
254 in Obligate Intracellular Bacteria (Wolbachia): Verification from Natural Populations."  
255 Molecular Biology and Evolution **27**(3): 501-505.
- 256 Chi, Y., Y. R. Yang and R. R. Muntz (2005). "Canonical forms for labelled trees and their  
257 applications in frequent subtree mining." Knowledge and Information Systems **8**(2): 203-  
258 234.
- 259 Daubin, V., M. Gouy and G. Perriere (2001). "Bacterial molecular phylogeny using supertree  
260 approach." Genome Inform **12**: 155-164.
- 261 Daubin, V., N. A. Moran and H. Ochman (2003). "Phylogenetics and the cohesion of bacterial  
262 genomes." Science **301**(5634): 829-832.
- 263 Ellegaard, K. M., L. Klasson, K. Naslund, K. Bourtzis and S. G. E. Andersson (2013).  
264 "Comparative Genomics of Wolbachia and the Bacterial Species Concept." Plos Genetics  
265 **9**(4).
- 266 Jehl, M. A., R. Arnold and T. Rattei (2011). "Effective-a database of predicted secreted bacterial  
267 proteins." Nucleic Acids Research **39**: D591-D595.
- 268 Kent, B. N. and S. R. Bordenstein (2010). "Phage WO of Wolbachia: lambda of the endosymbiont  
269 world." Trends Microbiol **18**(4): 173-181.
- 270 Kent, B. N., L. Salichos, J. G. Gibbons, A. Rokas, I. L. G. Newton, M. E. Clark and S. R.  
271 Bordenstein (2011). "Complete Bacteriophage Transfer in a Bacterial Endosymbiont  
272 (Wolbachia) Determined by Targeted Genome Capture." Genome Biology and Evolution  
273 **3**: 209-218.
- 274 Larkin, M. A., G. Blackshields, N. P. Brown, R. Chenna, P. A. McGettigan, H. McWilliam, F.  
275 Valentin, I. M. Wallace, A. Wilm, R. Lopez, J. D. Thompson, T. J. Gibson and D. G.  
276 Higgins (2007). "Clustal W and Clustal X version 2.0." Bioinformatics **23**(21): 2947-  
277 2948.
- 278 Marcet-Houben, M. and T. Gabaldon (2011). "TreeKO: a duplication-aware algorithm for the  
279 comparison of phylogenetic trees." Nucleic Acids Res **39**(10): e66.
- 280 Robinson, D. F. and L. R. Foulds (1981). "Comparison of Phylogenetic Trees." Mathematical  
281 Biosciences **53**(1-2): 131-147.
- 282 Shaw, T. I., Z. Ruan, T. C. Glenn and L. Liu (2013). "STRAW: Species TRee Analysis Web  
283 server." Nucleic Acids Res **41**(Web Server issue): W238-241.
- 284 Sicheritz-Ponten, T. and S. G. E. Andersson (2001). "A phylogenomic approach to microbial  
285 evolution." Nucleic Acids Research **29**(2): 545-552.
- 286 Stamatakis, A. (2006). "RAxML-VI-HPC: Maximum likelihood-based phylogenetic analyses  
287 with thousands of taxa and mixed models." Bioinformatics **22**(21): 2688-2690.

- 288 Sul, S.-J. and T. Williams (2007). A randomized algorithm for comparing sets of phylogenetic  
289 trees. Asia-Pacific Bioinformatics Conference: 121-130.
- 290 Swofford, D. L. and J. Sullivan (2009). "Phylogeny inference based on parsimony and other  
291 methods using PAUP\*." Phylogenetic Handbook: A Practical Approach to Phylogenetic  
292 Analysis and Hypothesis Testing, 2nd Edition: 267-312.
- 293 Wall, D. P., H. B. Fraser and A. E. Hirsh (2003). "Detecting putative orthologs." Bioinformatics  
294 **19**(13): 1710-1711.
- 295 Zhaxybayeva, O., J. P. Gogarten, R. L. Charlebois, W. F. Doolittle and R. T. Papke (2006).  
296 "Phylogenetic analyses of cyanobacterial genomes: quantification of horizontal gene  
297 transfer events." Genome Res **16**(9): 1099-1108.

# Figure 1

Wolbachia supergroup trees produced by concatenation of a dataset of 508 orthologs or by PhyBin's binning and clustering algorithm.

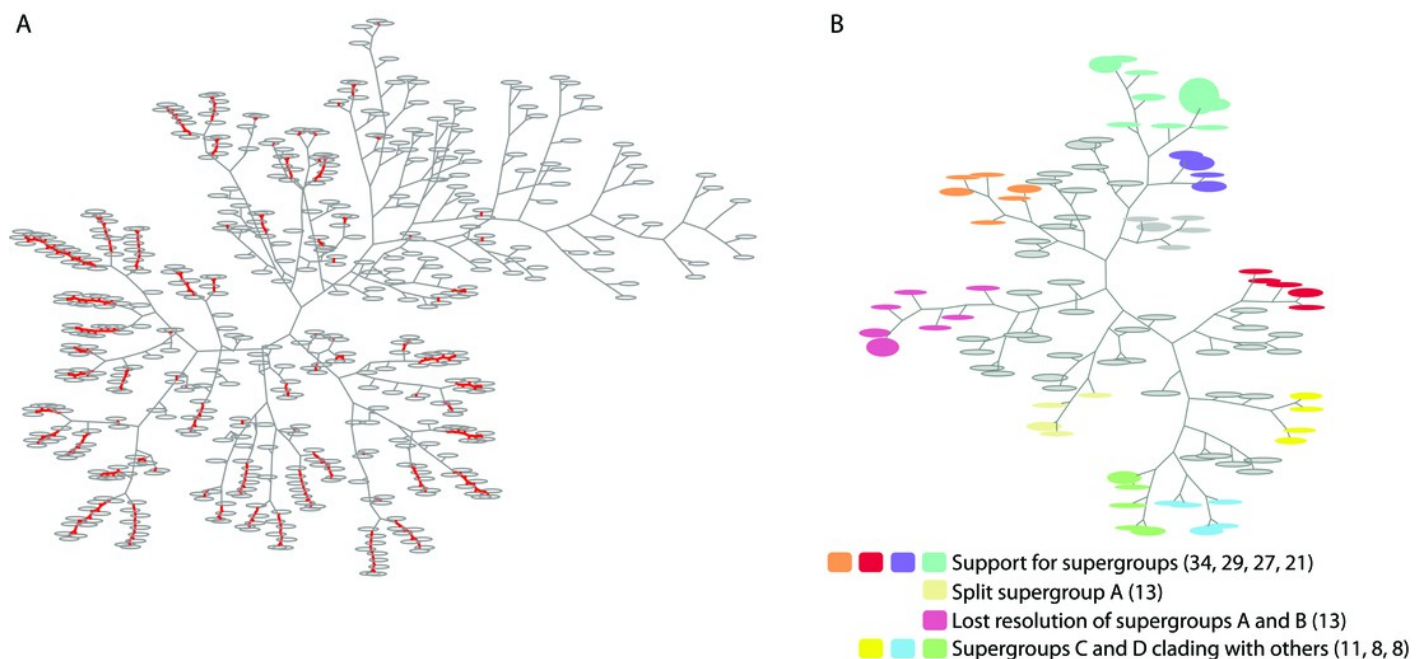
**Figure 1.** In each of two modes (*full clustering* and *binning*) PhyBin is able to correctly recover the expected topology for the *Wolbachia pipientis* orthologs used herein. (A) Concatenated phylogeny based on 508 genes (using RAXML GTRGAMMA, bootstrap support based on 10,000 replicates). The four major supergroups are highlighted and denoted. (B) These same groups are recovered when PhyBin is run in either *binning* mode or (C) *full clustering* mode.



# Figure 2

Two trees of trees for the *Wolbachia* ortholog set as visualized by PhyBin

**Figure 2.** Robinson-Foulds distance matrices produced by PhyBin are also visualized as a dendrogram by the software. (A) A *tree of trees* for the *Wolbachia* ortholog set (508 trees), clustered using an edit distance of 0, where identical topologies (nodes – grey ovals) are shown connected by a red line. Length of the branches connecting each node is proportional to the RF distance. (B) This dendrogram is simplified by increasing the RF distance at which the trees are clustered (shown RF = 3). The top 10 clusters and their support different topologies are colored as indicated in the legend (with largest bin size for each cluster cluster in parentheses).



# Table 1 (on next page)

Compute time for PhyBin compared to two other distance matrix calculation programs

**Table 1.** Timings for distance matrix computations only on 150-taxa benchmark included with HashRF. All times in seconds. PhyBin times are given with different numbers of threads in parentheses. All times were taken on a 4-socket, 32-core Intel Xeon E7-4830 server running at 2.13 GHz with RHEL 6. Phylip was compiled with gcc 4.4.7 and “-O2 “.

1 **Table 1.** Timings for distance matrix computations only on 150-taxa benchmark included with HashRF. All times in  
 2 seconds. PhyBin times are given with different numbers of threads in parentheses. All times were taken on a 4-  
 3 socket, 32-core Intel Xeon E7-4830 server running at 2.13 GHz with RHEL 6. Phylip was compiled with gcc 4.4.7  
 4 and “-O2”.

<b>Trees</b>	<b>PhyBin</b>	<b>HashRF</b>	<b>Phylip</b>	<b>DendroPy</b>
<b>100</b>	0.269	0.056	22.1	12.8
<b>1000</b>	4.7 (1), 3.0 (2), 1.9 (4), 1.4 (8)	1.7		

## Table 2<sub>(on next page)</sub>

The behavior of PhyBin on an example dataset from the *Wolbachia* genus using *binning* mode

**Table 2.** Using PhyBin in *binning mode* on the *Wolbachia* orthologous gene set (503 trees total) results in different size and number of bins depending on branch length threshold. The number of bins drops dramatically between a branch length threshold of 0 and 0.02, indicating a small amount of noise in the dataset due to the use of fairly similar taxa.



1 **Table 2.** Using PhyBin in *binning mode* on the *Wolbachia* orthologous gene set (503 trees total) results in different  
2 size and number of bins depending on branch length threshold. The number of bins drops dramatically between a  
3 branch length threshold of 0 and 0.02, indicating a small amount of noise in the dataset due to the use of fairly  
4 similar taxa.

Branch length threshold	Number of bins	Number of singletons	Size of largest bin
0	222	149	16
0.01	175	129	133
0.02	95	68	201
0.03	61	40	172
0.04	48	29	161

# Table 3(on next page)

The behavior of PhyBin on an example dataset from the *Wolbachia* genus using full clustering mode

**Table 3.** Using PhyBin in *full clustering mode* on the *Wolbachia* orthologous gene set (503 trees total) using average neighbor clustering produces a relatively small number of clusters, the largest comprised of a majority of orthologous genes.

2 **Table 3.** Using PhyBin in *full clustering mode* on the *Wolbachia* orthologous gene set (503 trees total) using average  
 3 neighbor clustering produces a relatively small number of clusters, the largest comprised of a majority of orthologous  
 4 genes.

<b>RF-distance threshold</b>	<b>Branch Length cutoff</b>	<b>Number of clusters</b>	<b>Number of singletons</b>	<b>Size of largest cluster</b>
<b>0</b>	n/a	222	149	16
<b>1</b>	n/a	140	67	34
<b>2</b>	n/a	77	29	56
<b>0</b>	0.01	175	129	133
<b>0</b>	0.02	95	68	201
<b>1</b>	0.02	66	35	246

## Table 4<sub>(on next page)</sub>

Wolbachia orthologs that do not conform to the dominant topology are highlighted by PhyBin

**Table 4.** List of *Wolbachia* orthologous gene sets not conforming to the dominant topology when PhyBin is run using *full clustering mode* (--UPGMA, --editdist=3). Protein products predicted to be secreted (based on screening using the Effective database ( Juhl, Arnold et al. 2011 ) are italicized.

1 **Table 4.** List of *Wolbachia* orthologous gene sets not conforming to the dominant topology when PhyBin is run using  
 2 *full clustering mode* (--UPGMA, --editdist=3). Protein products predicted to be secreted (based on screening using  
 3 the Effective database ([Jehl, Arnold et al. 2011](#))) are italicized.

Topology group	Orthologs (using wMel designations)
Support for splitting group A	Major facilitator family transporter (WD0470) <i>Diaminopimelate epimerase (WD1208)</i> GTP cyclohydrolase (WD0003) <i>Metalopeptidase (WD0059)</i> <i>Periplasmic divalent cation tolerance (WD0828)</i> <i>RodA (WD1108)</i>