

1 **Title**

2 ghost-tree: creating hybrid-gene phylogenetic trees for diversity analyses

3
4 **Authors**

5 Jennifer Fouquier¹, Jai Ram Rideout², Evan Bolyen², John Chase², Arron Shiffer^{2,3},
6 Daniel McDonald⁴, Rob Knight⁵, J Gregory Caporaso^{2,3} and Scott T. Kelley^{1,4*}

7
8 ¹ Graduate Program in Bioinformatics and Medical Informatics, San Diego State
9 University, San Diego, CA, USA.

10 ² Center for Microbial Genetics and Genomics, Northern Arizona University, Flagstaff,
11 AZ, USA.

12 ³ Department of Biological Sciences, Northern Arizona University, Flagstaff, AZ, USA.

13 ⁴ Department of Computer Science and BioFrontiers Institute, University of Colorado at
14 Boulder, CO, USA.

15 ⁵ Department of Pediatrics, and Department of Computer Science and Engineering,
16 University of California San Diego, San Diego, CA USA.

17 ⁶ Department of Biology, San Diego State University, San Diego, CA, USA.

18
19 **Corresponding Author**

20 Scott T. Kelley

21 Department of Biology

22 5500 Campanile Drive

23 San Diego State University

24 San Diego, CA, USA, 92182-4614

25 Phone: +1 619 206 8014

26 Email: skelley@mail.sdsu.edu

27
28
29 **Abstract**

30
31 ghost-tree is a bioinformatics tool that integrates sequence data from two genetic markers
32 into a single phylogenetic tree that can be used for diversity analyses. Our approach uses
33 one genetic marker whose sequences can be aligned across organisms spanning divergent
34 taxonomic groups (e.g., fungal families) as a “foundation” phylogeny. A second, more
35 rapidly evolving genetic marker is then used to build “extension” phylogenies for more
36 closely related organisms (e.g., fungal species or strains) that are then grafted on to the
37 foundation tree by mapping taxonomic names. We apply ghost-tree to graft fungal
38 extension phylogenies derived from ITS sequences onto a foundation phylogeny derived
39 from fungal 18S sequences. The result is a phylogenetic tree, compatible with the
40 commonly used UNITE fungal database, that supports phylogenetic diversity analysis
41 (e.g., UniFrac) of fungal communities profiled using ITS markers.

42
43 **Availability:** ghost-tree is pip-installable. All source code, documentation, and test code
44 are available under the BSD license at <https://github.com/JTFouquier/ghost-tree>.

45 Introduction

46

47 In recent years, next-generation sequencing and bioinformatics methods have rapidly
48 expanded our knowledge of the microbial world by enabling marker gene surveys of
49 bacterial, archaeal and eukaryotic microbial communities. Phylogenetic diversity metrics
50 such as *Phylogenetic Diversity* (PD) (Faith, 1992) and *UniFrac* (Lozupone and Knight,
51 2005) have improved resolution of community differences relative to their non-
52 phylogenetic analogs that were mostly developed for studying communities of macro-
53 organisms (e.g., Chao1 and Bray-Curtis dissimilarity). An ideal marker gene has highly
54 conserved regions that facilitate development of PCR primers and multiple sequence
55 alignment, and highly variable regions that support detailed taxonomic resolution.

56 For some taxonomic groups, such as the fungi, an ideal marker gene does not exist.
57 The small subunit ribosomal RNA (SSU) is highly conserved in the fungi, and therefore
58 sequences from different species are too similar for detailed taxonomic assignment. The
59 Internal Transcribed Spacer (ITS), a non-coding region found between the ribosomal
60 genes, has thus become popular for achieving high-resolution taxonomic profiles of
61 fungal communities. This marker “gene” is so non-conserved, however, that no
62 meaningful alignment is possible across distant fungal lineages, thus making
63 phylogenetic diversity analyses impossible.

64 Here we present ghost-tree, an open-source bioinformatics software tool for creating
65 phylogenetic trees using multiple genetic loci. Sequences from a more conserved locus
66 are aligned across distant taxonomic lineages to build a foundation tree, and sequences
67 from a less conserved locus are aligned for many groups of closely related taxa to create
68 extension trees that are then grafted onto the foundation tree. We apply “ghost-tree” to
69 build foundation trees from the fungal 18S rRNA and extension trees from fungal ITS, to
70 create a single “ghost-tree” that can be used in phylogenetic diversity analyses of fungal
71 communities. There are thus two outcomes of this project: a tree for use with popular
72 community analysis tools such as QIIME, and a software package for developing
73 phylogenetic trees for other sets of marker genes.

74

75 Methods

76

77 ghost-tree takes as input (1) the *Foundation Alignment* (for example, the Silva 18S
78 alignment) where sequences are annotated with taxonomy; (2) the *Extension Sequence*
79 *Collection* (for example, unaligned ITS sequences from the UNITE database); and (3) a
80 *taxonomy map*, which contains taxonomic annotations of the sequences in (2). The
81 Foundation Alignment is filtered to remove highly gapped and high entropy positions,
82 and FastTree is used to build a phylogenetic tree from the resulting filtered alignment.
83 This is the *Foundation Tree*. In parallel, the *Extension Sequence Collection* is clustered
84 with SUMACLUSt <http://metabarcoding.org/sumatra> resulting in an *operational*
85 *taxonomic unit (OTU) map* that groups sequences into OTUs by percent identity. For
86 each *Extension Sequence OTU* a consensus taxonomy is determined and OTUs with the
87 same consensus taxonomy are further grouped into a single OTU. The sequences in each
88 OTU are then aligned, and FastTree is applied to build an *Extension Tree*. The OTU's
89 consensus taxonomy is associated with the root of the *Extension Tree*. The taxa at the
90 root of the extension trees are then used to graft the extension tree on to the tip in the
91 foundation tree with the same taxonomy, resulting in the *Ghost Tree*. We applied ghost-
92 tree to build a phylogenetic tree from Silva (Ver. SSU 119.1) 18S sequences (the
93 foundation) (Pruesse *et al.*, 2007) and UNITE (Ver. 12_11_otus) ITS sequences (Köljalg
94 *et al.*, 2013). This tree is available in the GitHub repository, and this ghost-tree workflow
95 is illustrated in Supplementary Figure S1.

96 ghost-tree is hosted on GitHub under the BSD open source software license. It is
97 implemented in Python, using scikit-bio and Click, and adheres to the PEP8 Python style
98 guide. ghost-tree is subject to continuous integration testing using Travis CI which, on

99 each pull request, runs unit tests with nose, monitors code style using flake8, and
100 monitors test coverage with Coveralls.

101

102 **Results**

103

104 To evaluate whether ghost-tree supports improved resolution in studies of fungal
105 community analysis we compiled two real-world ITS sequence datasets: one collection of
106 human saliva samples (Ghannoum *et al.*, 2010) and one of surfaces in public restrooms
107 (Fouquier *et al.*, in review). These communities are expected to differ substantially from
108 one another (large effect size). We next simulated 10 samples based on each of the
109 samples from these studies using QIIME's `simsam.py` workflow. This generates sample
110 replicates that are phylogenetically similar to each other. We therefore consider the
111 grouping of each set of simulated samples to be a small effect size. We computed
112 distances between all samples using three approaches: using binary Jaccard, a qualitative
113 non-phylogenetic diversity metric where no tree is required; using unweighted UniFrac, a
114 qualitative phylogenetic metric of beta diversity with a tree generated using MUSCLE
115 and FastTree; and using unweighted UniFrac with a tree generated with ghost-tree.

116 Figure 1 contains principal coordinates analysis of these data based on binary Jaccard
117 distances (Fig. 1a), UniFrac/ITS distances (Fig. 1b) and UniFrac/ghost-tree distances
118 (Fig. 1c). UniFrac/ghost-tree resolves the small and the large effect sizes (ANOSIM
119 $R=0.38$ and $R=0.48$, respectively), while binary Jaccard ($R=0.02$ and $R=0.04$,
120 respectively) and unweighted UniFrac/ITS ($R=0.19$ and $R=0.08$, respectively) do not.

121

122 **Summary**

123

124 Widely used bacterial sequence databases such as GreenGenes (DeSantis *et al.*, 2006)
125 annotate 16S rRNA gene data, providing a reference tree and taxonomy for bacterial and
126 archaeal community analysis. ghost-tree facilitates development of phylogenetic trees that
127 can be used in a similar way for marker genes that are less conducive to phylogenetic
128 reconstruction at a global scale. We show that, as with bacterial community analysis,
129 incorporating phylogeny in diversity metrics is useful for resolving differences in fungal
130 communities. The Silva/UNITE-based ghost tree presented here integrate into a user's
131 exiting fungal analysis pipelines and the ghost-tree software package can be used to
132 develop phylogenetic trees for other marker gene sets that provide different taxonomic
133 resolution, or for bridging genome trees with amplicon trees.

134

135

136 **Acknowledgements**

137

138 This work was supported in part by a grant from the Alfred P. Sloan Foundation.

139

140

141

141 **Figure Legends**

142

143 **Figure 1.** Principal Coordinates comparing samples based on (a) Binary Jaccard
144 distances, (b) UniFrac distances where trees are computed by aligning ITS sequences
145 using MUSCLE, and (c) UniFrac distances where trees are computed using ghost-tree.
146 Blue points are simulated and real human saliva samples, and red points are simulated
147 and real restroom surface samples. Plots were made using EMPeror software (Vázquez-
148 Baeza *et al.*, 2013).

149 **Suppl. Figure S1.** ghost-tree workflow diagram.

150

151

151 **References**

152

153 DeSantis,T.Z., *et al.* (2006) Greengenes, a chimera-checked 16S rRNA gene database
154 and workbench compatible with ARB. *Appl. Environ. Microb.*, 72, 5069-5072.

155 Faith,D.P. (1992) Conservation Evaluation and Phylogenetic Diversity. *Biol. Conserv.*,
156 61, 1-10.

157 Ghannoum,M.A., *et al.* (2010) Characterization of the oral fungal microbiome
158 (mycobiome) in healthy individuals. *PLoS Pathog.*, 6, e1000713.

159 Kõljalg,U., *et al.* (2013) Towards a unified paradigm for sequence-based identification of
160 Fungi. *Mol. Ecol.*, 22, 5271-5277.

161 Lozupone,C. and Knight,R. (2005) UniFrac: a new phylogenetic method for comparing
162 microbial communities. *Appl. Environ. Microb.*, 71, 8228-8235.

163 Pruesse,E., *et al.* (2007) SILVA: a comprehensive online resource for quality checked
164 and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Res.*,
165 35, 7188-7196.

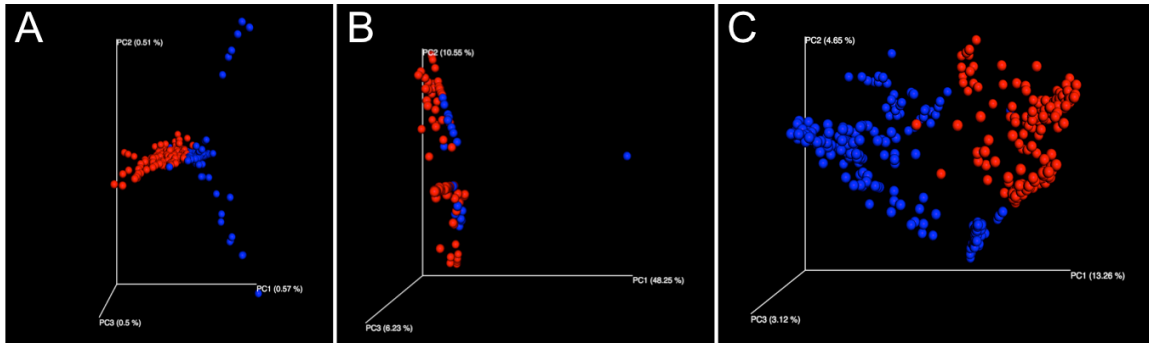
166 Quast,C., *et al.* (2013) The SILVA ribosomal RNA gene database project: improved data
167 processing and web-based tools. *Nucleic Acids Res.*, 41, D590-596.

168 Vázquez-Baeza,E., *et al.* (2013) EMPeror: a tool for visualizing high-throughput
169 microbial community data. *GigaScience.*, 2:16. doi:10.1186/2047-217X-2-16

170

171

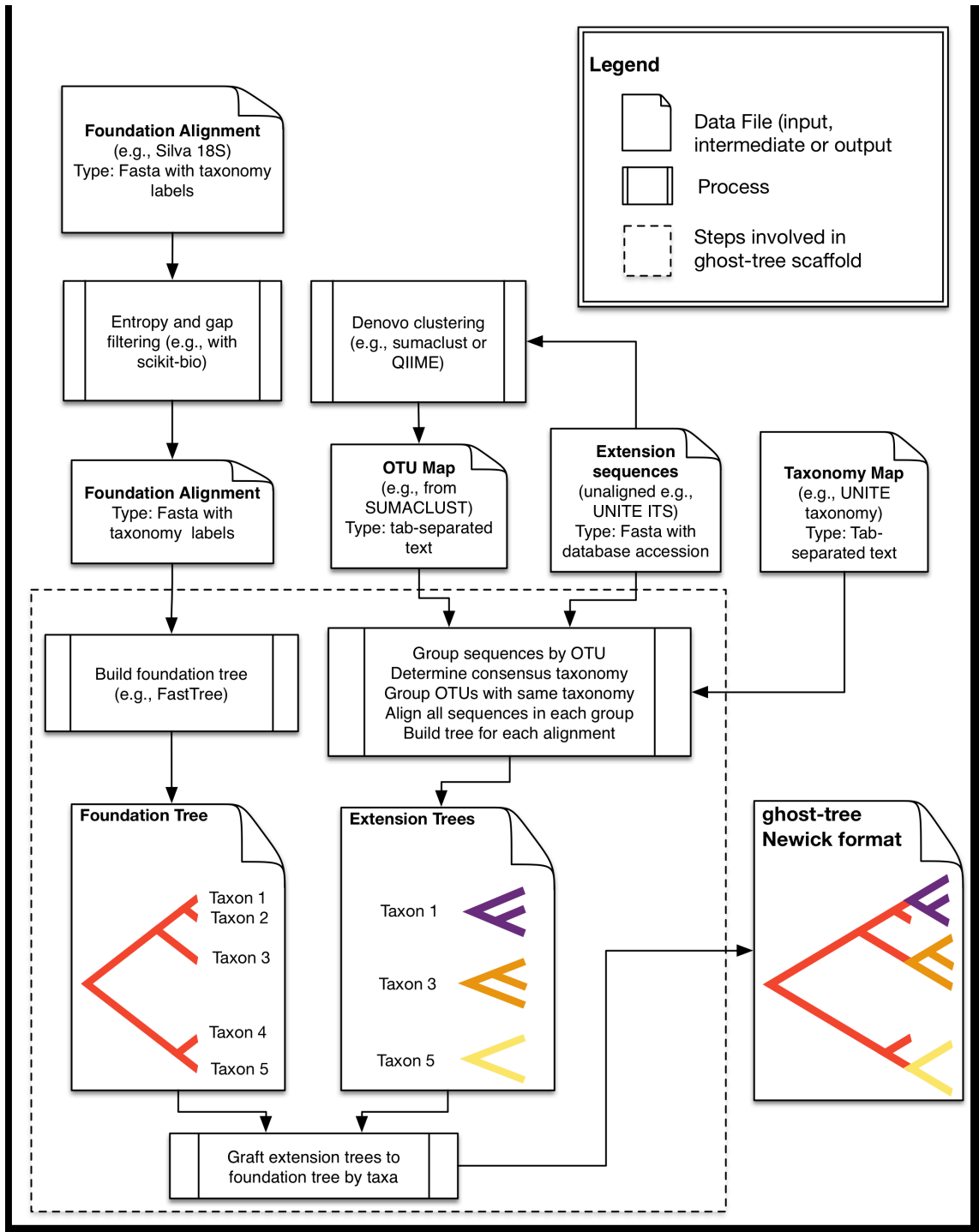
171 **Figure 1**
172



173
174

174
175
176

Supplemental Figure S1



177
178
179
180
181