

# APPLES: Distance-based Phylogenetic Placement for Scalable and Assembly-free Sample Identification

METIN BALABAN<sup>1</sup>, SHAHAB SARMASHGHI<sup>2</sup>, AND SIAVASH MIRARAB,<sup>2,\*</sup>

<sup>1</sup> *Bioinformatics and Systems Biology Graduate Program, UC San Diego, CA 92093, USA*

<sup>2</sup> *Department of Electrical and Computer Engineering, UC San Diego, CA 92093, USA*

*\*E-mail: smirarab@ucsd.edu.*

## ABSTRACT

1 Placing a new species on an existing phylogeny has increasing relevance to several  
2 applications. Placement can be used to update phylogenies in a scalable fashion and can  
3 help identify unknown query samples using (meta-)barcoding, skimming, or metagenomic  
4 data. Maximum likelihood (ML) methods of phylogenetic placement exist, but these  
5 methods are not scalable to reference trees with many thousands of leaves, limiting their  
6 ability to enjoy benefits of dense taxon sampling in modern reference libraries. They also  
7 rely on *assembled* sequences for the reference set and aligned sequences for the query.  
8 Thus, ML methods cannot analyze datasets where the reference consists of unassembled  
9 reads, a scenario relevant to emerging applications of genome-skimming for sample  
10 identification. We introduce APPLES, a distance-based method for phylogenetic  
11 placement. Compared to ML, APPLES is an order of magnitude faster and more memory  
12 efficient, and unlike ML, it is able to place on large backbone trees (tested for up to  
13 200,000 leaves). We show that using dense references improves accuracy substantially so  
14 that APPLES on dense trees is more accurate than ML on sparser trees, where it can run.  
15 Finally, APPLES can accurately identify samples without assembled reference or aligned  
16 queries using kmer-based distances, a scenario that ML cannot handle. APPLES is  
17 available publically at [github.com/balabanmetin/apples](https://github.com/balabanmetin/apples).

18 *Key words:* Phylogenetic placement, Distance-based methods, Genome-skimming.

19

20 Phylogenetic placement is the problem of finding the optimal position for a new  
21 *query* species on an existing *backbone* (or, reference) tree. Placement, as opposed to a  
22 *de novo* reconstruction of the full phylogeny, has two advantages. In some applications  
23 (discussed below), placement is all that is needed, and in terms of accuracy, it is as good  
24 as, and perhaps even better than, *de novo* reconstruction. Moreover, placement can be  
25 more scalable than *de novo* reconstruction when dealing with very large trees.

26 Earlier research on placement was motivated by scalability. For example, placement  
27 is used in greedy algorithms that start with an empty tree and add sequences sequentially  
28 (e.g., Felsenstein, 1981; Desper and Gascuel, 2002). Each placement requires polynomial  
29 (often linear) time with respect to the size of the backbone, and thus, these greedy  
30 algorithms are scalable (often requiring quadratic time). Despite computational challenges  
31 (Warnow, 2017), there has been much progress in the *de novo* reconstruction of ultra-large  
32 trees (e.g., thousands to millions of sequences) using both maximum likelihood (ML) (e.g.,  
33 Price *et al.*, 2010; Nguyen *et al.*, 2015) and the distance-based (e.g., Lefort *et al.*, 2015)  
34 approaches. However, these large-scale reconstructions require significant resources. As new  
35 sequences continually become available, placement can be used to update existing trees  
36 without repeating previous computations on full dataset.

37 More recently, placement has found a new application in sample identification:  
38 given one or more *query* sequences of unknown origins, detect the identity of the (set of)  
39 organism(s) that could have generated that sequence. These identifications can be made  
40 easily using sequence matching tools such as BLAST (Altschul *et al.*, 1990) when the  
41 query either exactly matches or is very close to a sequence in the reference library.  
42 However, when the sequence is novel (i.e., has lowered similarity to known sequences in the  
43 reference), this *closest* match approach is not sufficiently accurate (Koski and Golding,  
44 2001), leading some researchers to adopt a phylogenetic approach (Sunagawa *et al.*, 2013;

45 Nguyen *et al.*, 2014). Sample identification is essential to the study of mixed environmental  
46 samples, especially of the microbiome, both using 16S profiling (e.g., Gill *et al.*, 2006;  
47 Krause *et al.*, 2008) and metagenomics (e.g., von Mering *et al.*, 2007). It is also relevant to  
48 barcoding (Hebert *et al.*, 2003) and meta-barcoding (Clarke *et al.*, 2014; Bush *et al.*, 2017)  
49 and quantification of biodiversity (e.g., Findley *et al.*, 2013). Driven by applications to  
50 microbiome profiling, placement tools like pplacer (Matsen *et al.*, 2010) and EPA(-ng)  
51 (Berger *et al.*, 2011; Barbera *et al.*, 2018) have been developed. Researchers have also  
52 developed methods for aligning query sequence (e.g., Berger and Stamatakis, 2011;  
53 Mirarab *et al.*, 2012) and for downstream steps (e.g., Stark *et al.*, 2010; Matsen and Evans,  
54 2013). These publications have made a strong case that for sample identification,  
55 placement is sufficient (i.e., *de novo* is not needed). Moreover, some studies (e.g., Janssen  
56 *et al.*, 2018) have shown that when dealing with fragmentary reads typically found in  
57 microbiome samples, placement can be *more* accurate than *de novo* construction and can  
58 lead to improved associations of microbiome with clinical information.

59 Existing phylogenetic placement methods have focused on the ML inference of the  
60 best placement – a successful approach, which nevertheless, suffers from two shortcomings.  
61 On the one hand, ML can only be applied when the reference species are *assembled* into  
62 full-length sequences (e.g., an entire gene) and are *aligned*; however, in new applications  
63 that we will describe, assembling (and hence aligning) the reference set is not possible. On  
64 the other hand, ML, while somewhat scalable, is still computationally demanding,  
65 especially in memory usage, and cannot place on backbone trees with many thousands of  
66 leaves. As the density of reference substantially impacts the accuracy and resolution of  
67 placement, this inability to use ultra-large trees as backbone also limits accuracy. This  
68 limitation has motivated alternative methods using local sensitive hashing (Brown and  
69 Truszkowski, 2013) and divide-and-conquer (Mirarab *et al.*, 2012).

70 Assembly-free and alignment-free sample identification using genome-skimming  
71 (Dodsworth, 2015) can also benefit from phylogenetic placement. A genome-skim is a

72 shut-gun sample of the genome sequenced at low coverage (e.g., 1X) – so low that  
73 assembling the nuclear genome is not possible (though, mitochondrial or plastid genomes  
74 can often be assembled). Genome-skimming promises to replace traditional marker-based  
75 barcoding of biological samples (Coissac *et al.*, 2016) but limiting analyses to organelle  
76 genome can limit resolution. Sarmashghi *et al.* (2019) have recently shown that using  
77 shared  $k$ -mers, the distance between two unassembled genome-skims with low coverage can  
78 be accurately estimated. This approach, unlike assembling organelle genomes, uses data  
79 from the entire nuclear genome and hence promises to provide a higher resolution (e.g., at  
80 species or sub-species levels) while keeping the low sequencing cost. However, ML and  
81 other methods that require assembled sequences cannot analyze genome-skims, where both  
82 the reference and the query species are unassembled genome-wide bags of reads.

83 Distance-based approaches to phylogenetics are well-studied, but no existing tool  
84 can perform distance-based placement of a query sequence on a given backbone. The  
85 distance-based approach promises to solve both shortcomings of ML methods.  
86 Distance-based methods are computationally efficient and do not require assemblies. They  
87 only need distances (however computed). Thus, they can take as input assembly-free  
88 estimates of genomic distance estimated from low coverage genome-skims using Skmer  
89 (Sarmashghi *et al.*, 2019) or other alternatives (Haubold, 2014; Leimeister and  
90 Morgenstern, 2014; Leimeister *et al.*, 2017; Yi and Jin, 2013; Benoit *et al.*, 2016; Fan *et al.*,  
91 2015; Ondov *et al.*, 2016; Jain *et al.*, 2017). While alignment-based phylogenetics has been  
92 traditionally more accurate than alignment-free methods when both methods are possible,  
93 in these new scenarios, only alignment-free methods are applicable.

94 Here, we introduce a new method for distance-based phylogenetic placement called  
95 APPLES (Accurate Phylogenetic Placement using LEast Squares). APPLES uses dynamic  
96 programming to find the optimal distance-based placement of a sequence with running  
97 time and memory usage that scale linearly with the size of the backbone tree. We test  
98 APPLES in simulations and on real data, both for alignment-free and aligned scenarios.

MATERIALS AND METHODS

*Problem Statement.*

*Notations.* Let an unrooted tree  $T = (V, E)$  be a weighted connected acyclic undirected graph with leaves denoted by  $\mathcal{L} = \{1 \cdots n\}$ . We let  $T^*$  be the rooting of  $T$  on a leaf 1 obtained by directing all edges away from 1. For node  $u \in V$ , let  $p(u)$  denote its parent,  $c(u)$  denote its set of children,  $sib(u)$  denote its siblings, and  $g(u)$  denote the set of leaves at or below  $u$  (i.e., those that have  $u$  on their path to the root), all with respect to  $T^*$ . Also let  $l(u)$  denote the length of the edge  $(p(u), u)$ .

*Distances.* The tree  $T$  defines an  $n \times n$  matrix where each entry  $d_{ij}(T)$  corresponds to the path length between leaves  $i$  and  $j$ . We further generalize this definition so that  $d_{uv}(T^*)$  indicates the length of the undirected path between any two nodes of  $T^*$  (when clear, we simply write  $d_{uv}$ ). Given some input data, we can compute a matrix of all pairwise sequence distances  $\Delta$ , where the entry  $\delta_{ij}$  indicates the dissimilarity between species  $i$  and  $j$ . When the sequence distance  $\delta_{ij}$  is computed using (the correct) phylogenetic model, it will be a noisy but statistically consistent estimate of the tree distance  $d_{ij}(T)$  (Felsenstein, 2003). Given these “phylogenetically corrected” distances (e.g.  $\frac{3}{4} \ln(1 - \frac{4}{3}h)$  is the corrected hamming distance  $h$  under the Jukes and Cantor (1969) model), we can define optimization problems to recover the tree that best fits the distances. A natural choice is minimizing the (weighted) least square difference between tree and sequence distances:

$$Q^*(T) = \sum_{i=1}^n \sum_{j=1}^n w_{ij} (\delta_{ij} - d_{ij}(T))^2. \quad (1)$$

Here, weights (e.g.,  $w_{ij}$ ) are used to reduce the impact of large distances (expected to have high variance). A general weighting schema can be defined as  $w_{ij} = \delta_{ij}^{-k}$  for a constant value  $k \in \mathbb{N}$ . Standard choices of  $k$  include  $k = 0$  for the ordinary least squares (OLS) method of Cavalli-Sforza and Edwards 1967,  $k = 1$  due to Beyer *et al.* 1974 (BE), and  $k = 2$  due to Fitch and Margoliash 1967 (FM).

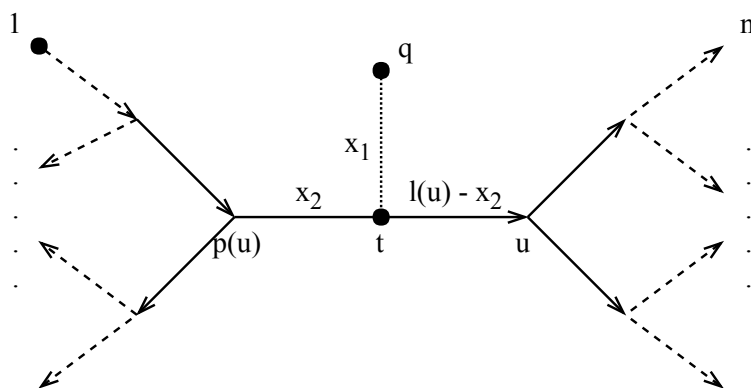


Fig. 1. Any placement of  $q$  can be characterized as a tree  $P(u, x_1, x_2)$ , shown here. The backbone tree  $T^*$  is an arborescence on leaves  $\mathcal{L} = \{1 \dots n\}$ , rooted at leaf 1. Query taxon  $q$  is added on the edge between  $u$  and  $p(u)$ , creating a node  $t$ . All placements on this edge are characterized by  $x_1$ , the length of the pendant branch, and  $x_2$ , the distance between  $t$  and  $p(u)$ .

112 Finding  $\arg \min_T Q^*(T)$  is NP-Complete (Day, 1987). However, decades of research  
 113 has produced heuristics like neighbor-joining (Saitou and Nei, 1987), alternative  
 114 formulations like (balanced) minimum evolution (Cavalli-Sforza and Edwards, 1967;  
 115 Desper and Gascuel, 2002), and several effective tools for solving the problem heuristically  
 116 (e.g., FastME by Lefort *et al.* 2015, DAMBE by Xia 2018, and Ninja by Wheeler 2009).

117 *Phylogenetic placement.* We let  $P(u, x_1, x_2)$  be the tree obtained by adding a  
 118 query taxon  $q$  on an edge  $(p(u), u)$ , creating three edges  $(t, q)$ ,  $(p(u), t)$ , and  $(t, u)$ , with  
 119 weights  $x_1$ ,  $x_2$ , and  $l(u) - x_2$ , respectively (Fig. 1). When clear, we simply write  $P$  and  
 120 note that  $P$  induces  $T$  both in topology and branch length. We now define the problem.

121 *Least Squares Phylogenetic Placement (LSPP).*

122 **Input:** A backbone tree  $T$  on  $\mathcal{L}$ , a query species  $q$ , and a vector  $\Delta_{q^*}$  with elements  $\delta_{qi}$   
 123 giving sequence distances between  $q$  and every species  $i \in \mathcal{L}$ ;

**Output:** The placement tree  $P$  that adds  $q$  on  $T$  and minimizes

$$Q(P) = \sum_{i=1}^n w_{qi} (\delta_{qi} - d_{qi}(P))^2 \quad (2)$$

124

### Linear Time Solution for LSPP

125

126

127

128

129

The number of possible placements of  $q$  is  $2n - 3$ . Therefore, LSPP can be solved by simply iterating over all the topologies, optimizing the score for that branch, and returning the placement with the minimum least square error. A naive algorithm can accomplish this in  $\Theta(n^2)$  running time by optimizing Eq. 2 for each of the  $2n - 3$  branches. However, using dynamic programming, the optimal solution can be found in linear time.

130

**THEOREM 1** The LSPP problem can be solved with  $\Theta(n)$  running time and memory.

131

132

133

134

135

The proof (given in Appendix A) follows easily from three lemmas that we next state. The algorithm starts with precomputing a fixed-size set of values for each nodes. For any node  $u$  and exponents  $a \in \mathbb{Z}$  and  $b \in \mathbb{N}^+$ , let  $S(a, b, u) = \sum_{i \in g(u)} \delta_{qi}^a d_{ui}^b$  and for  $b = 0$ , let  $S(a, 0, u) = S'(a, u) = \sum_{i \in g(u)} \delta_{qi}^a$ . Note that  $S'(0, u) = |g(u)|$ . Similarly, for  $u \in V \setminus \{1\}$ , let  $R(a, b, u) = \sum_{i \notin g(u)} \delta_{qi}^a d_{p(u)i}^b$  for  $b > 0$  and let  $R(a, 0, u) = R'(a, u) = \sum_{i \notin g(u)} \delta_{qi}^a$ .

136

137

**LEMMA 2** The set of all  $S(a, b, u)$  and  $R(a, b, u)$  values can be precomputed in  $\Theta(n)$  time with two tree traversals using the dynamic programming given by:

$$S(a, b, u) = \begin{cases} \delta_{qu}^a & u \in \mathcal{L} \setminus \{1\} \ \& \ b = 0 \\ 0 & u \in \mathcal{L} \setminus \{1\} \ \& \ b \neq 0 \\ \sum_{j=0}^b \sum_{v \in c(u)} l(v)^j \binom{b}{j} S(a, b-j, v) & u \notin \mathcal{L} \setminus \{1\} \end{cases} \quad (3)$$

$$R(a, b, u) = \begin{cases} \delta_{q1}^a & u = 1' = c(1) \ \& \ b = 0 \\ 0 & u = 1' = c(1) \ \& \ b \neq 0 \\ \sum_{j=0}^b \left( l(p(u))^j \binom{b}{j} R(a, b-j, p(u)) + \sum_{v \in sib(u)} l(v)^j \binom{b}{j} S(a, b-j, v) \right) & u \notin \{1, 1'\} \end{cases} \quad (4)$$

138

139

140

**LEMMA 3** Equation 2 can be rearranged (see Eq. S2 in Appendix A) such that computing  $Q(P)$  for a given  $P = P(u, x_1, x_2)$  requires a constant time computation using  $S(a, b, u)$  and  $R(a, b, u)$  values for  $-k \leq a \leq 2 - k$  and  $0 \leq b \leq 2$ .

141

142

143

Thus, after a linear time precomputation, we can compute the error for any given placement in constant time. It remains to show that for each node, the optimal placement on the branch above it (e.g.,  $x_1$  and  $x_2$ ) can be computed in constant time.

144 LEMMA 4 For a fixed node  $u \in V \setminus \{1\}$ , if  $(\hat{x}_1, \hat{x}_2) = \arg \min_{x_1, x_2} Q(P(u, x_1, x_2))$ , then

$$\begin{aligned} & \begin{bmatrix} R'(-k, u) + S'(-k, u) & R'(-k, u) - S'(-k, u) \\ R'(-k, u) - S'(-k, u) & R'(-k, u) + S'(-k, u) \end{bmatrix} \cdot \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \\ & \begin{bmatrix} R'(1-k, u) + S'(1-k, u) - l(u)S'(-k, u) - R(-k, 1, u) - S(-k, 1, u) \\ R'(1-k, u) - S'(1-k, u) + l(u)S'(-k, u) - R(-k, 1, u) + S(-k, 1, u) \end{bmatrix} \end{aligned} \quad (5)$$

145 and hence  $\hat{x}_1, \hat{x}_2$  can be computed in constant time.

146 *Non-negative branch lengths.* The solution to Equation 5 does not necessarily  
 147 conform to constraints  $0 \leq x_1$  and  $0 \leq x_2 \leq l(u)$ . However, the following lemma (proof in  
 148 Appendix A) allows us to easily impose the constraints by choosing optimal boundary  
 149 points when unrestricted solutions fall outside boundaries.

150 LEMMA 5 With respect to variables  $x_1$  and  $x_2$ ,  $Q(P(u, x_1, x_2))$  is a convex function.

151 *Minimum evolution* An alternative to directly using MLSE (Eq. 1) is the  
 152 minimum evolution (ME) principle (Cavalli-Sforza and Edwards, 1967; Rzhetsky and Nei,  
 153 1992). Our algorithm can also optimize the ME criterion: after computing  $x_1$  and  $x_2$  by  
 154 optimizing MLSE for each node  $u$ , we choose the placement with the minimum total  
 155 branch length. This is equivalent to using  $\arg \min_u x_1$ , since the value of  $x_2$  does not  
 156 contribute to total branch length. Other solution for ME placement exists (Desper and  
 157 Gascuel, 2002), a topic we return to in the Discussion section.

158 *Hybrid.* We have observed cases where ME is correct more often than MLSE, but when it  
 159 is wrong, unlike MLSE, it has a relatively high error. This observation led us to design a  
 160 hybrid approach. After computing  $x_1$  and  $x_2$  for all branches, we first select the top  $\log_2(n)$   
 161 edges with minimum  $Q(P(u, x_1, x_2))$  values (this requires  $\Theta(n \log \log n)$  time). Among this  
 162 set of edges, we place the query on the edge satisfying the ME criteria.



163

## Datasets

164

165

166

We benchmark accuracy and scalability of APPLES in two settings: sample identification using assembly-free genome-skims on real biological data and placement using aligned sequences on simulated data.

167

### *Real genome-skim datasets for the assembly-free scenario*

168

169

170

171

172

173

174

175

*Columbicola genome-skims.* We use a set of 61 genome-skims by Boyd *et al.* (2017), including 45 known lice species (some represented multiple times) and 7 undescribed species. We generate lower coverage skims of 0.1Gb or 0.5Gb by randomly subsampling the reads from the sequence read archives (SRA) provided by the original publication (NCBI BioProject PRJNA296666). We use BBTools (Bushnell, 2014) to filter subsampled reads for adapters and contaminants and remove duplicated reads. Since this dataset is not assembled, the coverage of the genome-skims is unknown; Skmer estimates the coverage to be between 0.2X and 1X for 0.1Gb samples (and 5 times that coverage with 0.5Gb).

176

177

178

179

180

181

182

*Anopheles and Drosophila datasets.* We also use two insect datasets used by Sarmashghi *et al.* (2019): a dataset of 22 *Anopheles* and a dataset of 21 *Drosophila* genomes (Table S1), both obtained from InsectBase (Yin *et al.*, 2016). For both datasets, genome-skims with 0.1Gb and 0.5Gb sequence were generated from the assemblies using the short-read simulator tool ART, with the read length  $l = 100$  and default error profile. Since species have different genome sizes, with 0.1Gb data, our subsampled genome-skims range in coverage from 0.35X to 1X for *Anopheles* and from 0.4X to 0.8X for *Drosophila*.

183

184

185

186

187

More recently, Miller *et al.* (2018) sequenced several *Drosophila* genomes, including 12 species shared with the InsectBase dataset. Sarmashghi *et al.* (2019) subsampled the SRAs from this second project to 0.1Gb or 0.5Gb and, after filtering contaminants, obtained artificial genome-skims. We can use these genome-skims as query and the genome-skims from the InsectBase dataset as the backbone. Since the reference and query

188 come from two projects, the query genome-skim can have a non-zero distance to the same  
189 species in the reference set, providing a realistic test of sample identification applications.

190 *Simulated datasets for the aligned sequence scenario*

191 *GTR.* We use a 101-taxon dataset available from Mirarab and Warnow 2015. Sequences  
192 were simulated under the General Time Reversible (GTR) plus the  $\Gamma$  model of site rate  
193 heterogeneity using INDELible (Fletcher and Yang, 2009) on gene trees that were  
194 simulated using SimPhy (Mallo *et al.*, 2016) under the coalescent model evolving on  
195 species trees generated under the Yule model. Note that the same model is used for  
196 inference under ML placement methods (i.e., no model misspecification). We took all 20  
197 replicates of this dataset with mutation rates between  $5 \times 10^{-8}$  and  $2 \times 10^{-7}$ , and for each  
198 replicate, randomly selected five estimated gene trees among those with  $\leq 20\%$  RF distance  
199 between estimated and true gene tree. Thus, we have a total of 100 backbone trees.

200 *RNASim.* Guo *et al.* 2009 designed a complex model of RNA evolution that does not  
201 make usual i.i.d assumptions of sequence evolution. Instead, it uses models of energy of the  
202 secondary structure to simulate RNA evolution by a mutation-selection population genetics  
203 model. This model is based on an inhomogeneous stochastic process without a global  
204 substitution matrix. The model complexity of RNASim allows us to test both ML and  
205 APPLES under a substantially misspecified model. An RNASim dataset of  $10^6$  sequences is  
206 available from Mirarab *et al.* 2015. We created several subsets of the full RNASim dataset.

207 *i)* Heterogeneous: We first randomly subsampled the full dataset to create 10  
208 datasets of size  $10^4$ . Then, we chose the largest clade of size at most 250 from each  
209 replicate; this gives us 10 backbone trees of mean size 249.

210 *ii)* Varied diameter: To evaluate the impact of the evolutionary diameter (i.e., the  
211 highest distance between any two leaves in the backbone), we also created datasets with  
212 low, medium, and high diameters. We sampled the largest five clades of size at most 250  
213 from each of the 10 replicates used for the heterogeneous dataset. Among these 50 clades,

214 we picked the bottom, middle, and top five clades in diameter, which had diameter in  
215  $[0.3, 0.4]$  (mean: 0.36),  $[0.5, 0.52]$  (mean: 0.51), and  $[0.65, 1.07]$  (mean: 0.82), respectively.

216 *iii*) Varied size: We randomly subsampled the tree of size  $10^6$  to create 5 replicates  
217 of datasets of size  $5 \times 10^2$ ,  $10^3$ ,  $5 \times 10^3$ ,  $10^4$ ,  $5 \times 10^4$ , and  $10^5$ , and 1 replicate (due to size)  
218 of size  $2 \times 10^5$ . For replicates that contain at least  $5 \times 10^3$  species, we removed sites that  
219 contain gaps in 95% or more of the sequences in the alignment.

## 220 *Methods*

221 *Alternative methods.* For aligned data, we compare APPLES to two ML methods:  
222 pplacer (Matsen *et al.*, 2010) and EPA-ng (Barbera *et al.*, 2018). Matsen *et al.* (2010)  
223 found pplacer to be substantially faster than EPA (Berger and Stamatakis, 2011) while  
224 their accuracy was similar. EPA-ng improves the scalability of EPA; thus, we compare to  
225 EPA-ng in analyses that concerned scalability (e.g., RNASim-Varied Size). We run pplacer  
226 and EPA-ng in their default mode using GTR+ $\Gamma$  model (the only option for pplacer). We  
227 also compare with a simple method referred to as CLOSEST that places the query as the  
228 sister to the species with the minimum distance to it. CLOSEST is meant to emulate the  
229 use of BLAST (if it could be used). For the assembly-free setting, existing phylogenetic  
230 placement methods cannot be used, and we compared only against CLOSEST.

231 *Distance calculation and models.* We modified FastME to compute distances only  
232 between query and backbone sequences, not among backbone sequences. This version,  
233 called FastME\* here, also ensures that when estimating model parameters, positions that  
234 have a gap in at least one of the two sequences are always ignored.

235 We compute phylogenetic distances under the parameter-free JC69 model, the  
236 six-parameter Tamura and Nei 1993 (TN93) model, and the 12-parameter general Markov  
237 model (Lockhart *et al.*, 1994). We compute distances independently for all pairs, and not  
238 simultaneously as suggested by Tamura *et al.* (2004). We also use the Gamma model of

239 sites rate heterogeneity for JC69 and TN93 using the standard approach (Waddell and  
240 Steel, 1997). Pairing Gamma with GTR is theoretically possible in the absence of noise;  
241 however, the method can run into problems on real data (Waddell and Steel, 1997). Thus,  
242 we do not include a GTR model directly. Instead, we use the log-det approach that can  
243 handle the most general (12-parameter) Markov model (Lockhart *et al.*, 1994); however,  
244 log-det cannot account for rate across sites heterogeneity (Waddell and Steel, 1997). The  $\alpha$   
245 parameter of the Gamma model cannot be computed from pairwise sequence comparisons  
246 (Steel, 2009); instead, we use the  $\alpha$  computed from the backbone tree. We used the  $\alpha$   
247 parameter computed by RAxML (Stamatakis, 2014) run on the backbone alignment and  
248 given the backbone tree.

249 In analyses on assembly-free datasets, we first compute genomic distances using  
250 Skmer (Sarmashghi *et al.*, 2019). We then correct these distances using the JC69 model,  
251 without the Gamma model of rate variation.

252 *Backbone trees.* For genome-skimming experiments, we estimated the backbone  
253 tree using FastME\* from the JC69 distance matrix computed from genome-skims using  
254 Skmer. For simulated datasets, we estimated the topology of the backbone tree by running  
255 RAxML (Stamatakis, 2014) on the true alignment using GTRGAMMA model and used  
256 this tree as the backbone for pplacer and EPA-ng. However, to handle large trees, we used  
257 FastTree-2 (Price *et al.*, 2010) to estimate the backbone tree for RNASim-varied size and  
258 re-estimated branch lengths on the fixed topology using RAxML. For the backbone of  
259 APPLES, we always used the same tree topology but re-estimated branch lengths using  
260 FastTree-2 under the JC69 model.

261 *APPLES parameters.* We have chosen default parameter settings for APPLES and refer  
262 to this version as APPLES\*. By default, we use FM weighting, the MLSE selection  
263 criterion, enforcement of non-negative branch lengths, and JC69 distances. When not  
264 specified otherwise, these default parameters are used.

### Evaluation Procedure

265

266 To evaluate the accuracy, we use a leave-one-out strategy. We remove each leaf  $i$   
267 from the backbone tree  $T$  and place it back on this  $T \setminus i$  tree to obtain the placement tree  
268  $P$ . However, on the RNAsim-varied size dataset, due to its large size, we only removed and  
269 added back 200 randomly chosen leaves per replicate.

270 *Delta error.* We measure the accuracy of the placement using delta error ( $\Delta e$ ): the  
271 number of branches of the true tree missing from  $P$  minus the number of branches of the  
272 true tree missing from  $T \setminus i$  (induced on the same leafset). Note that  $\Delta e \geq 0$  because  
273 adding  $i$  cannot decrease the number of missing branches in  $T \setminus i$ . Note that placing  $i$  to  
274 the same location as the backbone before leaving it out (e.g.,  $T$ ) can still have a non-zero  
275 delta error because the backbone tree is not the true tree. We refer to the placement of a  
276 leaf into its position in the backbone tree as the *de novo* placement.

277 On biological data, where the true tree is unknown, we use a reference tree  
278 (Fig. S1). For *Drosophila* and *Anopheles*, we use the tree available from the Open Tree Of  
279 Life (Hinchliff *et al.*, 2015) as the reference. For *Columbicola*, we use the ML concatenation  
280 tree published by Boyd *et al.* (2017) as the reference.

281

## RESULTS

282

### *Assembly-free Placement of Genome-skims*

283 On our three biological genome-skim datasets, APPLES\* successfully places the  
284 queries on the optimal position in most cases (97%, 95%, and 71% for *Columbicola*,  
285 *Anopheles*, and *Drosophila*, respectively) and is never off from the optimal position by  
286 more than one branch. Other versions of APPLES are less accurate than APPLES\*; e.g.,  
287 APPLES with ME can have up to five wrong branches (Table 1). On genome-skims, where  
288 assembly and alignment are not possible, existing placement tools cannot be used, and the  
289 only alternative is the CLOSEST method (emulating BLAST if assembly was possible).

	(a) <i>Columbicola</i>			(b) <i>Anopheles</i>			(c) <i>Drosophila</i>		
	%	$\Delta e$	$e_{max}$	%	$\Delta e$	$e_{max}$	%	$\Delta e$	$e_{max}$
<b>APPLES*</b>	97	0.03	1	95	0.05	1	71	0.29	1
<b>APPLES-ME</b>	84	0.28	5	95	0.05	1	67	0.42	2
<b>APPLES-HYBRID</b>	87	0.16	2	95	0.05	1	67	0.33	1
<b>CLOSEST</b>	54	1.15	7	91	0.09	1	57	0.62	3
<b>DE-NOVO</b>	98	0.02	1	95	0.05	1	71	0.29	1

Table 1. **Assembly-free placement of genome-skims.** We show the percentage of placements into optimal position (those that do not increase  $\Delta e$ ), average delta error ( $\Delta e$ ), and maximum delta error ( $e_{max}$ ) for APPLES, assignment to the CLOSEST species, and the placement to the position in the backbone (DE-NOVO) over the 61 (a), 22 (b), and 21 (c) placements. Results are shown for genome skims with 0.1Gbp of reads. Delta error is the increase in the missing branches between the reference tree and the backbone tree after placing each query.

290 CLOSEST finds the optimal placement only in 54% and 57% of times for *Columbicola* and  
291 *Drosophila*; moreover, it can be off from the best placement by up to seven branches for  
292 the *Columbicola* dataset. On the *Anopheles* dataset, where the reference tree is unresolved  
293 (Fig. S1), all methods perform similarly.

294 APPLES\* is less accurate on the *Drosophila* dataset than other datasets. However,  
295 here, simply placing each query on its position in the backbone tree would lead to identical  
296 results (Table 1). Thus, placements by APPLES\* are as good as the *de novo* construction,  
297 meaning that errors of APPLES\* are entirely due to the differences between our backbone  
298 tree and the reference tree. Moreover, these errors are not due to low coverage; increasing  
299 the genome-skim size 5x (to 0.5Gb) does not decrease error (Table S4).

300 On *Drosophila* dataset, we next tested a more realistic sample identification scenario  
301 using the 12 genome-skims from the separate study (and thus, non-zero distance to the  
302 corresponding species in the backbone tree). As desired, APPLES\* places all of 12 queries  
303 from the second study as sister to the corresponding species in the reference dataset.

### 304 *Alignment-based Placement*

305 We first compare the accuracy and scalability of APPLES\* to ML methods and  
306 then compare various settings of APPLES. For ML, we use pplacer (shown everywhere)  
307 and EPA-ng (shown only when we study scalability).

	Low			Medium			High			Heterogeneous		
	%	$\Delta e$	$e_{max}$	%	$\Delta e$	$e_{max}$	%	$\Delta e$	$e_{max}$	%	$\Delta e$	$e_{max}$
<b>APPLES*</b>	86	0.15	2	85	0.18	5	84	0.18	3	85	0.17	5
<b>CLOSEST</b>	59	0.88	13	60	0.88	13	60	0.85	14	60	0.87	14
<b>pplacer</b>	88	0.13	2	89	0.11	3	87	0.13	3	88	0.13	3

Table 2. The delta error for APPLES\*, CLOSEST match, and pplacer on the RNASim-varied diameter dataset (low, medium, or high) and the RNA-heterogeneous dataset. Measurements are shown over 1250 placements for each diameter size category, corresponding to 5 backbone trees and 250 placements per replicate.

### Comparison to Maximum Likelihood (ML)

*GTR dataset.* On this dataset, where it faces no model misspecification, pplacer has high accuracy. It finds the best placement in 84% of cases and is off by one edge in 15% (Fig. 2a); its mean delta error ( $\Delta e$ ) is only 0.17 edges. APPLES\* is also accurate, finding the best placement in 78% of cases and resulting in the mean  $\Delta e = 0.28$  edges. Thus, even though pplacer uses ML and faces no model misspecification and APPLES\* uses distances based on a simpler model, the accuracy of the two methods is within 0.1 edges on average. In contrast, CLOSEST has poor accuracy and is correct only 50% of the times, with the mean  $\Delta e$  of 1.0 edge.

*Model misspecification.* On the small RNASim data with subsampled clades of  $\approx 250$  species), both APPLES\* and pplacer face model misspecification. Here, the accuracy of APPLES\* is very close to ML using pplacer. On the heterogeneous subset (Fig. 2b and Table 2), pplacer and APPLES\* find the best placement in 88% and 86% of cases and have a mean delta error of 0.13 and 0.17 edges, respectively. Both methods are much more accurate than CLOSEST, which has a delta error of 0.87 edges on average.

*Impact of diameter.* When we control the tree diameter, APPLES\* and pplacer remain very close in accuracy (Fig. 2c). The changes in error are small and not monotonic as the diameters change (Table 2). The accuracies of the two methods at low and high diameters are similar. The two methods are most divergent in the medium diameter case, where pplacer has its lowest error ( $\Delta e = 0.11$ ) and APPLES\* has its highest error ( $\Delta e = 0.18$ ).

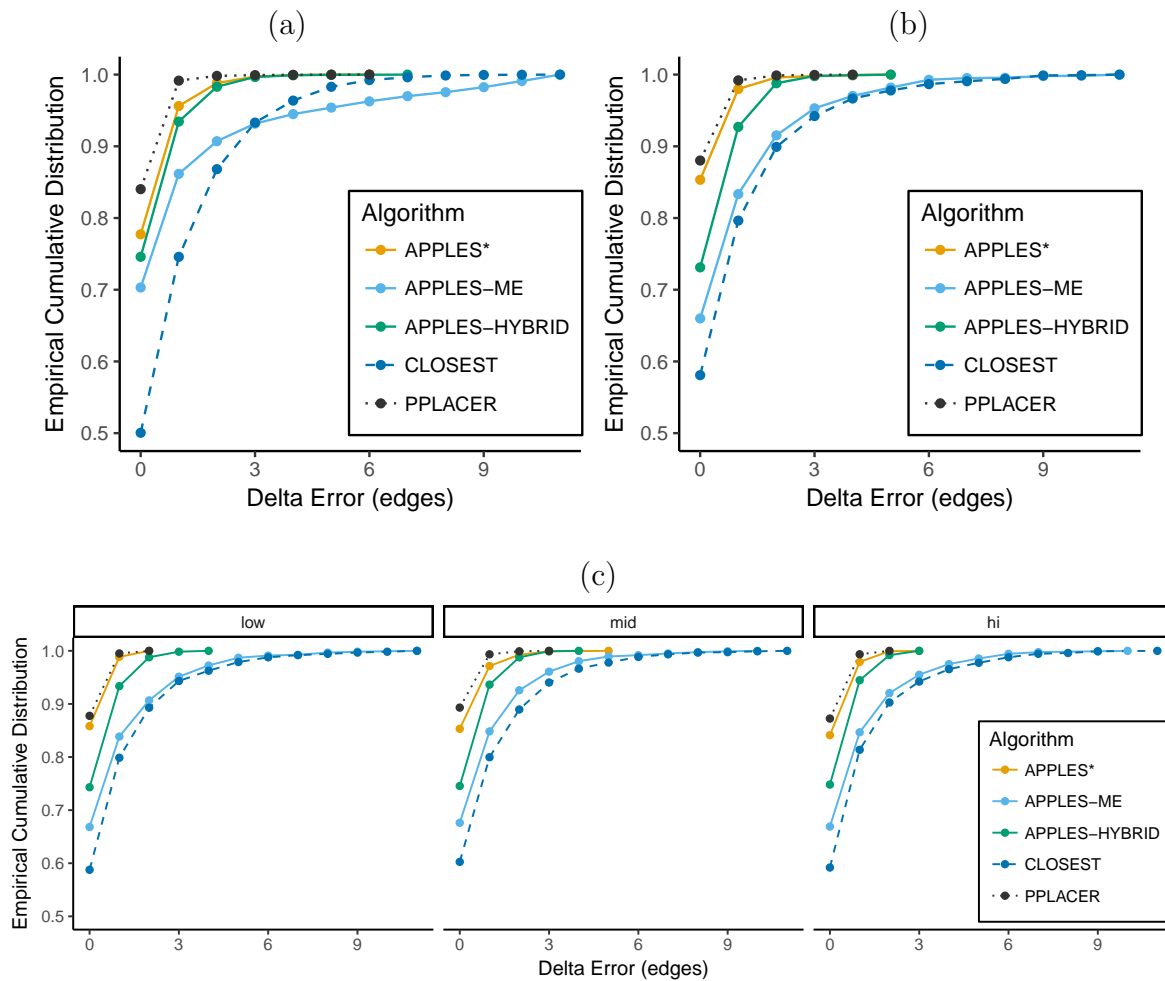


Fig. 2. **Accuracy on simulated data.** We show empirical cumulative distribution of the delta error, defined as the increase in the number of missing branches in the estimated tree compared to the true tree. We compare pplacer (dotted), CLOSEST match (dashed), and APPLES with FM weighting and JC69 distances and MLSE (APPLES\*), ME, or Hybrid optimization. (a) GTR dataset. (b) RNASim-Heterogeneous. (c) RNASim-varied diameter, shown in boxes: low, medium (mid), or high (hi). Distributions are over 10,000 (a), 2450 (b), and 3675 (c) points.

328 To summarize results on small RNASim dataset with model misspecification,  
 329 although APPLES\* uses a parameter-free model, its accuracy is extremely close to ML  
 330 using pplacer with the GTR+ $\Gamma$  model.

331 *Impact of taxon sampling.* The real advantage of APPLES\* over pplacer becomes clear for  
 332 placing on larger backbone trees (Fig. 3 and Table 3). For backbone sizes of 500 and 1000,  
 333 pplacer continues to be slightly more accurate than APPLES\* (mean  $\Delta e$  of pplacer is



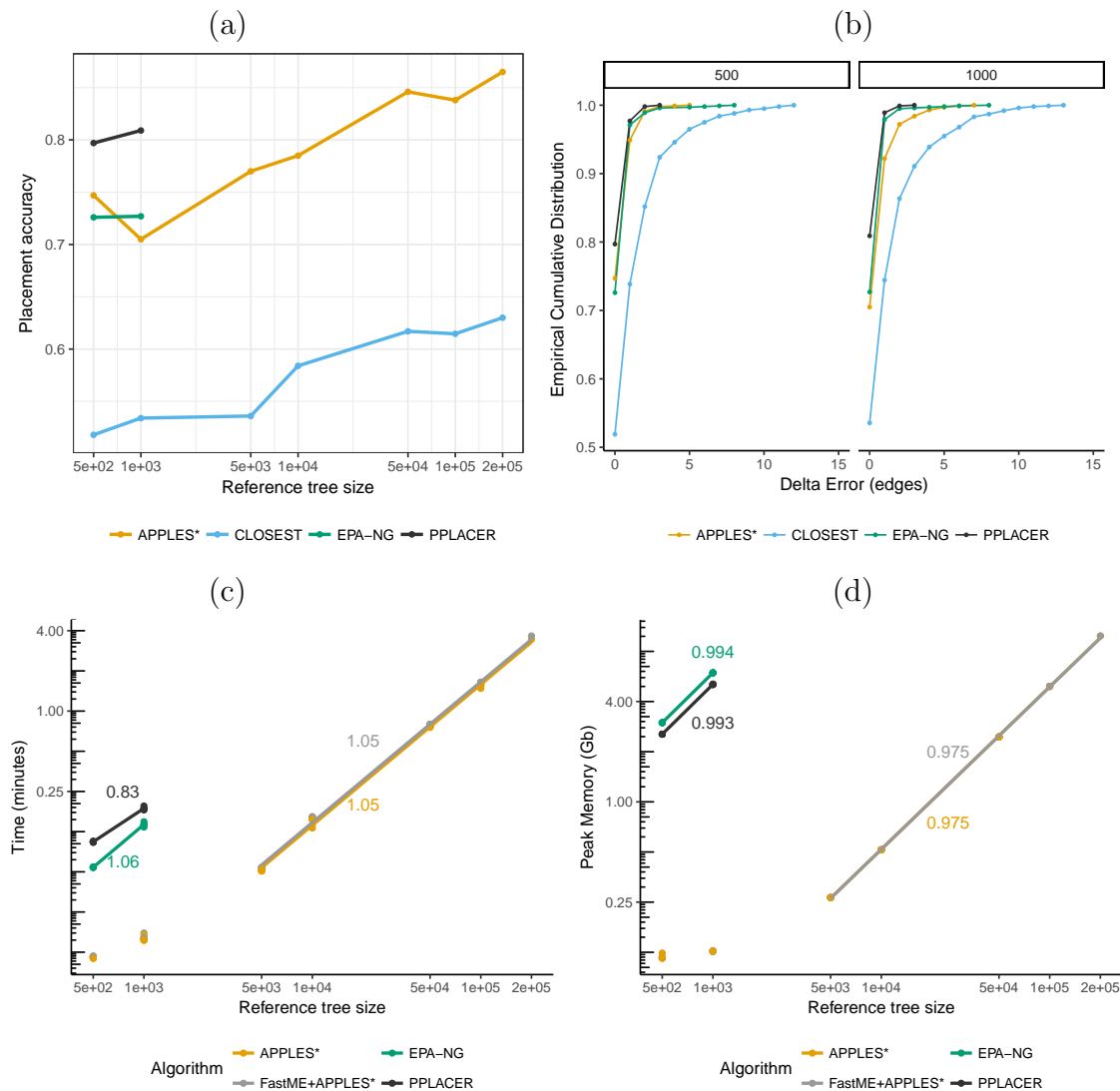


Fig. 3. **Results on RNASim-Variied size.** (a) Placement accuracy for various levels of taxon sampling, comparing pplacer, CLOSEST match, EPA-ng, and APPLES\* on RNASim dataset with backbone size ranging from 500 to 200,000. (b) The empirical cumulative distribution of the delta error on the same datasets, only shown for the backbone size 500 and 1000 where all methods can run. Distributions are over 1000 points. (c,d) Running time and peak memory usage of placement methods for a single placement. For APPLES\*, measurements are shown with and without the distance calculation step performed using FastME\*. On backbones of size 5000, pplacer managed to correctly run for only 551 out of 1000 placements, whereas EPA-ng managed to run for 200/1000 placements (Fig S2). Lines are fitted in the log-log scale; thus, the slope of the line (indicated on the figure) gives an empirical estimate of the polynomial degree of the asymptotic growth curve. All curves grow close to linearly (slopes  $\approx 1$ ). APPLES lines are fitted to  $\geq 5,000$  points because the first two values correspond to extremely low memory (100Mb) and are irrelevant to asymptotic behavior. All calculations are on 8-core, 2.6GHz Intel Xeon CPUs (Sandy Bridge) with 64GB of memory.

	$n = 500$		$n = 10^3$		$n = 5 \times 10^3$		$n = 10^4$		$n = 10^5$		$n = 2 \times 10^5$	
	%	$\Delta e$	%	$\Delta e$	%	$\Delta e$	%	$\Delta e$	%	$\Delta e$	%	$\Delta e$
<b>APPLES*</b>	75	0.32	71	0.43	77	0.37	79	0.33	84	0.25	87	0.25
<b>CLOSEST</b>	52	1.16	53	1.18	54	1.15	59	0.90	61	0.69	63	0.70
<b>EPA-ng</b>	73	0.33	73	0.31	fail (449)	n.p	n.p	n.p	n.p	n.p	n.p	n.p
<b>pplacer</b>	80	0.23	81	0.20	fail (800)	n.p	n.p	n.p	n.p	n.p	n.p	n.p

Table 3. Percentage of correct placements (shown as %) and the delta error ( $\Delta e$ ) on the RNASim datasets with various backbone size ( $n$ ). % and  $\Delta e$  is over 1000 placements (except  $n = 200,000$ , which is over 200 placements). Running pplacer and EPA-ng was not possible (n.p) for trees with at least 10,000 leaves and failed in some cases (number of fails shown) for 5,000 leaves.

334 better than APPLES\* by 0.09 and 0.23 edges, respectively). However, with backbones of  
 335 5000 leaves, pplacer fails to run on 449/1000 cases, producing infinity likelihood (perhaps  
 336 due to numerical issues) and has 41 times higher error than APPLES\* on the rest (Fig. S2).

337 Since pplacer could not scale to 5,000 leaves, we also tested the recent method,  
 338 EPA-ng (Barbera *et al.*, 2018). On datasets with up to 1000 leaves, EPA-ng was less  
 339 accurate than pplacer and close in accuracy to APPLES\* (Fig. 3ab). It also failed in  
 340 800/1000 replicates of the 5000-taxon backbone but had 4% less error than APPLES\* in  
 341 the minority of cases where it could run (Fig. S2).

342 For backbone trees with at least  $10^4$  leaves, pplacer and EPA-ng were not able to  
 343 run, and CLOSEST is not very accurate (finding the best placement in only 59% of cases).  
 344 However, APPLES\* continues to be accurate for all backbone sizes. As the backbone size  
 345 increases, the taxon sampling of the tree is improving (recall that these trees are all  
 346 random subsets of the same tree). With denser backbone trees, APPLES\* has increased  
 347 accuracy despite placing on larger trees (Fig. 3a, Table 3). For example, using a backbone  
 348 tree of  $2 \times 10^5$  leaves, APPLES\* is able to find the best placement of query sequences in  
 349 87% of cases, which is better than the accuracy of either APPLES\* or ML tools on any  
 350 backbone size. Thus, an increased taxon sampling helps accuracy, but ML tools are limited  
 351 in the size of the tree they can handle.

352 *Running time and memory.* As the backbone size increases, the running times of  
 353 all methods increase close to linearly with the size of the backbone tree (Fig. 3c). However,

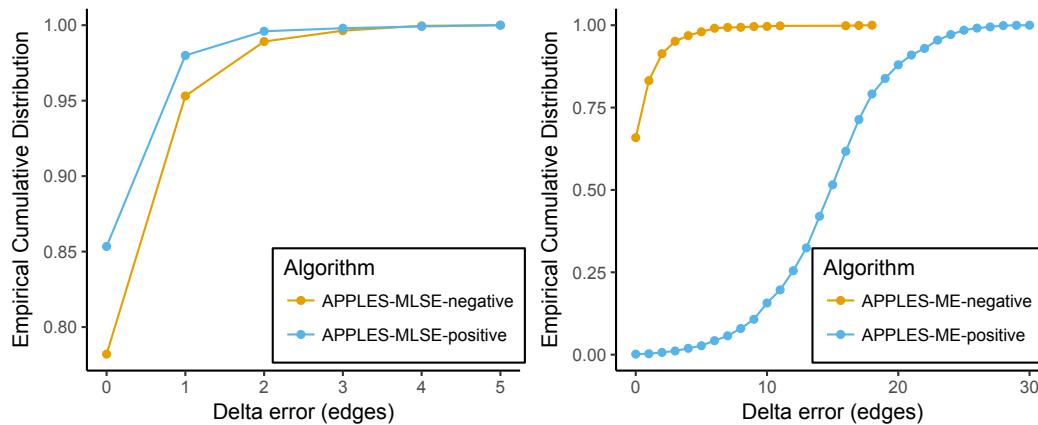


Fig. 4. **The effect of imposing positivity constraint on error.** We show the error of (a) APPLES-MLSE and (b) APPLES-ME run both with and without enforcement of non-negative branch lengths on RNASim heterogeneous dataset. Accuracy improves substantially for MLSE whereas it reduces drastically for ME.

354 APPLES is on average 15 times faster than pplacer and 12 times faster than EPA-ng on  
355 backbone trees with 5000 leaves in cases where those methods could run. Similarly, the  
356 memory of all methods increases linearly with the backbone size, but APPLES requires  
357 dramatically less memory (Fig. 3d). For example, for placing on a backbone with 5000  
358 leaves, pplacer requires 25GB of memory, EPA-ng requires 30GB whereas APPLES  
359 requires only 0.25GB. APPLES easily scales to a backbone of  $2 \times 10^5$  sequences, running in  
360 only 4 minutes and using 8GB of memory per query (including all precomputations in the  
361 dynamic programming). These numbers also include the time and memory needed to  
362 compute the distance between the query sequence and all the backbone sequences.

363 *Comparing parameters of APPLES.* We now compare different settings of  
364 APPLES. Comparing five models of sequence evolution, we see similar patterns of accuracy  
365 across all models despite their varying complexity, ranging from 0 to 12 parameters  
366 (Fig. S3). Since the JC69 model is parameter-free and results in similar accuracy to others,  
367 we have used it as the default. Next, we ask whether imposing the constraint to disallow  
368 negative branch lengths improves the accuracy. The answer depends on the optimization  
369 strategy. Forcing non-negative lengths marginally increases the accuracy for MLSE but  
370 dramatically reduces the accuracy for ME (Fig. 4). Thus, we always impose non-negative

371 constraints on MLSE but never for ME. Likewise, our Hybrid method includes the  
372 constraint for the first MLSE step but not for the following ME step (Fig. S4).

373 The next parameter to choose is the weighting scheme. Among the three methods  
374 available in APPLES, the best accuracy belongs to the FM scheme closely followed by the  
375 BE (Fig. S5). The OLS scheme, which does not penalize long distances, performs  
376 substantially worse than FM and BE. Thus, the most aggressive form of weighting (FM)  
377 results in the best accuracy. Fixing the weighting scheme to FM and comparing the three  
378 optimization strategies (MLSE, ME, and Hybrid), the MLSE approach has the best  
379 accuracy (Fig. 2), finding the correct placement 84% of the time (mean error: 0.18), and  
380 ME has the lowest accuracy, finding the best placement in only 67% of cases (mean error:  
381 0.70). The Hybrid approach is between the two (mean error: 0.34) and fails to outperform  
382 MLSE on this dataset. However, when we restrict the RNASim backbone trees to only 20  
383 leaves, we observe that Hybrid can have the best accuracy (Fig. S6).

384

## DISCUSSION

385 We introduced APPLES: a new method for adding query species onto large  
386 backbone trees using both unassembled genome-skims and aligned data. The accuracy of  
387 APPLES was very close to ML using pplacer in most settings where ML could run; the  
388 accuracy advantages of ML were particularly small for the more realistic simulation,  
389 RNASim, where both methods face model misspecification. As expected by the substantial  
390 evidence from the literature (Hillis *et al.*, 2003; Zwickl and Hillis, 2002), improved taxon  
391 sampling increased the accuracy of placement. Thus, overall, the best accuracy on  
392 RNASim dataset was obtained by APPLES\* run on the full reference dataset. This  
393 observation motivates the use of scalable methods such as APPLES\* instead of ML  
394 methods, which have to restrict their backbone to at most several thousand species. It is  
395 possible to follow up the APPLES\* placements with a round of ML placement on smaller  
396 trees, but the small differences in accuracy of pplacer and APPLES\* on smaller trees did

397 not give us compelling reasons to try such hybrid approaches.

398         Phylogenetic insertion using the ME criterion has been previously studied for the  
399 purpose of creating an algorithm for greedy minimum evolution (GME). Desper and  
400 Gascuel 2002 have designed a method that given the tree  $T$  can update it to get a tree  
401 with  $n + 1$  leaves in  $\Theta(n)$  after precomputation of a data-structure that gives the average  
402 *sequence* distances between all adjacent clusters in  $T$ . The formulation by Desper and  
403 Gascuel 2002 has a subtle but consequential difference from our ME placement. Their  
404 algorithm does not compute branch lengths for inserted sequence (e.g.,  $x_1$  and  $x_2$ ). It is  
405 able to compute the optimal placement topology *without* knowing branch lengths of the  
406 backbone tree. Instead, it relies on pairwise distances among backbone *sequences* ( $\Delta$ ),  
407 which are precomputed and saved in the data-structure mentioned before. In the context  
408 of the greedy algorithm for tree inference, in each iteration, the data structure can be  
409 updated in  $\Theta(n)$ , which does not impact the overall running time of the algorithm.  
410 However, if we were to start with a tree of  $n$  leaves, computing this structure from scratch  
411 would still require  $\Theta(n^2)$ . Thus, computing the placement for a new query would need  
412 quadratic time, unless if the  $\Theta(n^2)$  precomputation is allowed to be amortized over  $\Omega(n)$   
413 queries. Our formulation, in contrast, uses branch lengths of the backbone tree (which is  
414 assumed fixed) and thus never uses pairwise distances among the backbone sequences.  
415 Thus, using tree distances is what allows us to develop a linear time algorithm.

416         Our comparisons between versions of APPLES answered many questions but left  
417 others to future work. For example, we observed no advantage in using models more  
418 complex than JC69+ $\Gamma$  for distance calculation. However, these results may be due to our  
419 estimation of model parameters (e.g., base compositions) for each pair of sequences. More  
420 complex models may perform better if we instead estimate model parameters on the  
421 backbone alignment/tree and reuse the parameters for queries (or simultaneously among  
422 all queries and the reference sequences). Simultaneous estimation of distances has many  
423 advantages over using independent distances for the *de novo* case (Tamura *et al.*, 2004;

424 Xia, 2009); these results gives us hope that using simultaneous distances inside APPLES  
425 can further improve its accuracy.

426 In the aligned case, we were unable to test other methods. LSHPlace is theoretically  
427 fast, but we could not find an implementation of it. The distance-based insertion algorithm  
428 of FastME (Desper and Gascuel, 2002) is available only as part of a larger greedy  
429 algorithm but is not available as a stand-alone feature to place on a given tree. SEPP  
430 (Mirarab *et al.*, 2012) performs alignment and placement simultaneously (using alignment  
431 scores to help the placement); however, in our experiments, our goal was only to test the  
432 placement step and not the alignment. Thus, we used true alignments in all the simulation  
433 tests and left an exploration of the impact of alignment error on different methods to  
434 future work. On a related note, future work can incorporate APPLES inside SEPP to  
435 perform alignment and placement in a unified pipeline.

436 In our assembly-free test, we used Skmer to get distances because alternative  
437 alignment-free methods of estimating distance generally either require assemblies (e.g.,  
438 Haubold, 2014; Leimeister and Morgenstern, 2014; Leimeister *et al.*, 2017) or higher  
439 coverage than Skmer (e.g., Benoit *et al.*, 2016; Yi and Jin, 2013; Ondov *et al.*, 2016);  
440 however, combining APPLES with other alignment-free methods can be attempted in  
441 future (finding the best way of computing distances without assemblies was not our focus).  
442 Moreover, the Skmer paper has described a trick that can be used to compute log-det  
443 distances from genome-skims. Future studies should test whether using that trick and  
444 using GTR instead of JC69 improves accuracy.

445 Branch lengths of our backbone trees were computed using the same distance model  
446 as the one used for computing the distance of the query to backbone species. Using  
447 consistent models for the query and for the backbone branch lengths is essential for  
448 obtaining good accuracy (see Fig. S7 for evidence). Thus, in addition to having a large  
449 backbone tree at hand, we need to ensure that branch lengths are computed using the  
450 right model. Fortunately, FastTree-2 can compute both topologies and branch lengths on

## REFERENCES

23

451 large trees in a scalable fashion, without a need for quadratic time/memory computation  
452 of distance matrices (Price *et al.*, 2010).

453 APPLES was an order of magnitude or more faster and less memory-hungry than  
454 ML tools (pplacer and EPA-ng), but it has room for improvement. The python APPLES  
455 code is not optimized and can be dramatically improved. For example, APPLES can save  
456 precomputed values of Equations 3 and 4 for each backbone tree in a file, eliminating the  
457 need to recompute them. Also, online processing of the backbone alignment can  
458 dramatically reduce the memory usage for the distance calculation. Its current version uses  
459 Dendropy (Sukumaran and Holder, 2010), which is not optimized for large trees; switching  
460 to other platforms such as ETE (Huerta-Cepas *et al.*, 2010) can improve memory usage.  
461 Future implementations of APPLES will improve speed and memory by applying such  
462 optimizations.

## ACKNOWLEDGMENTS

463  
464 This work was supported by the National Science Foundation (NSF) grant  
465 IIS-1565862 and National Institutes of Health (NIH) subaward 5P30AI027767-28 to M.B.  
466 and S.M., and NSF grant NSF-1815485 to M.B., S.S., and S.M. Computations were  
467 performed on the San Diego Supercomputer Center (SDSC) through XSEDE allocations,  
468 which is supported by the NSF grant ACI-1053575.

## REFERENCES

469  
470 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. 1990. Basic local  
471 alignment search tool. *Journal of Molecular Biology*, 215(3): 403–410.  
472 Barbera, P., Kozlov, A. M., Czech, L., Morel, B., Darriba, D., Flouri, T., and Stamatakis,  
473 A. 2018. Epa-ng: massively parallel evolutionary placement of genetic sequences.  
474 *BioRxiv*, page 291658.

- 475 Benoit, G., Peterlongo, P., Mariadassou, M., Drezén, E., Schbath, S., Lavenier, D., and  
476 Lemaitre, C. 2016. Multiple comparative metagenomics using multiset k -mer counting.  
477 *PeerJ Computer Science*, 2: e94.
- 478 Berger, S. a. and Stamatakis, A. 2011. Aligning short Reads to Reference Alignments and  
479 Trees. *Bioinformatics*, 27(15): 2068–75.
- 480 Berger, S. A., Krompass, D., Stamatakis, A., D.K., Stamatakis, A., Krompass, D., and  
481 Stamatakis, A. 2011. Performance, Accuracy, and Web Server for Evolutionary  
482 Placement of Short Sequence Reads under Maximum Likelihood. *Systematic Biology*,  
483 60(3): 291–302.
- 484 Beyer, W. A., Stein, M. L., Smith, T. F., and Ulam, S. M. 1974. A molecular sequence  
485 metric and evolutionary trees. *Mathematical Biosciences*, 19(1-2): 9–25.
- 486 Boyd, B. M., Allen, J. M., Nguyen, N.-P., Sweet, A. D., Warnow, T., Shapiro, M. D., Villa,  
487 S. M., Bush, S. E., Clayton, D. H., and Johnson, K. P. 2017. Phylogenomics using  
488 target-restricted assembly resolves intrageneric relationships of parasitic lice  
489 (phthiraptera: Columbicola). *Systematic biology*, 66(6): 896–911.
- 490 Brown, D. G. and Truszkowski, J. 2013. LSHPlace: fast phylogenetic placement using  
491 locality-sensitive hashing. *Pacific Symposium on Biocomputing. Pacific Symposium on*  
492 *Biocomputing*, pages 310–9.
- 493 Bush, A., Sollmann, R., Wilting, A., Bohmann, K., Cole, B., Balzter, H., Martius, C.,  
494 Zlinszky, A., Calvignac-Spencer, S., Cobbold, C. A., Dawson, T. P., Emerson, B. C.,  
495 Ferrier, S., Gilbert, M. T. P., Herold, M., Jones, L., Leendertz, F. H., Matthews, L.,  
496 Millington, J. D. A., Olson, J. R., Ovaskainen, O., Raffaelli, D., Reeve, R., Rödel, M.-O.,  
497 Rodgers, T. W., Snape, S., Visseren-Hamakers, I., Vogler, A. P., White, P. C. L.,  
498 Wooster, M. J., and Yu, D. W. 2017. Connecting Earth observation to high-throughput  
499 biodiversity data. *Nature Ecology & Evolution*, 1(7): 41559–017.



- 500 Bushnell, B. 2014. Bbtools software package. URL <http://sourceforge.net/projects/bbmap>.
- 501 Cavalli-Sforza, L. L. and Edwards, A. W. 1967. Phylogenetic analysis. Models and  
502 estimation procedures. *American journal of human genetics*, 19(3 Pt 1): 233–57.
- 503 Clarke, L. J., Soubrier, J., Weyrich, L. S., and Cooper, A. 2014. Environmental  
504 metabarcodes for insects: In silico PCR reveals potential for taxonomic bias. *Molecular*  
505 *Ecology Resources*, 14(6): 1160–1170.
- 506 Coissac, E., Hollingsworth, P. M., Lavergne, S., and Taberlet, P. 2016. From barcodes to  
507 genomes: extending the concept of DNA barcoding. *Molecular Ecology*, 25(7):  
508 1423–1428.
- 509 Day, W. H. 1987. Computational complexity of inferring phylogenies from dissimilarity  
510 matrices. *Bulletin of Mathematical Biology*, 49(4): 461–467.
- 511 Desper, R. and Gascuel, O. 2002. Fast and Accurate Phylogeny Reconstruction Algorithms  
512 Based on the Minimum-Evolution Principle. *Journal of Computational Biology*, 9(5):  
513 687–705.
- 514 Dodsworth, S. 2015. Genome skimming for next-generation biodiversity analysis. *Trends*  
515 *in Plant Science*, 20(9): 525–527.
- 516 Fan, H., Ives, A. R., Surget-Groba, Y., and Cannon, C. H. 2015. An assembly and  
517 alignment-free method of phylogeny reconstruction from next-generation sequencing  
518 data. *BMC Genomics*, 16(1): 522.
- 519 Felsenstein, J. 1981. Evolutionary trees from DNA sequences: A maximum likelihood  
520 approach. *Journal of Molecular Evolution*, 17(6): 368–376.
- 521 Felsenstein, J. 2003. Inferring phylogenies. *Sunderland*.
- 522 Findley, K., Oh, J., Yang, J., Conlan, S., Deming, C., Meyer, J. A., Schoenfeld, D.,  
523 Nomicos, E., Park, M., Kong, H. H., and Segre, J. A. 2013. Topographic diversity of  
524 fungal and bacterial communities in human skin. *Nature*, 498(7454): 367–370.

- 525 Fitch, W. M. and Margoliash, E. 1967. Construction of phylogenetic trees. *Science*,  
526 155(3760): 279–284.
- 527 Fletcher, W. and Yang, Z. 2009. INDELible: A flexible simulator of biological sequence  
528 evolution. *Molecular Biology and Evolution*, 26(8): 1879–1888.
- 529 Gill, S. R., Pop, M., DeBoy, R. T., Eckburg, P. B., Turnbaugh, P. J., Samuel, B. S.,  
530 Gordon, J. I., Relman, D. A., Fraser-Liggett, C. M., and Nelson, K. E. 2006.  
531 Metagenomic Analysis of the Human Distal Gut Microbiome. *Science*, 312(5778):  
532 1355–1359.
- 533 Guo, S., Wang, L.-S., and Kim, J. 2009. Large-scale simulation of RNA macroevolution by  
534 an energy-dependent fitness model.
- 535 Haubold, B. 2014. Alignment-free phylogenetics and population genetics. *Briefings in*  
536 *Bioinformatics*, 15(3): 407–418.
- 537 Hebert, P. D. N., Cywinska, A., Ball, S. L., and deWaard, J. R. 2003. Biological  
538 identifications through DNA barcodes. *Proceedings of the Royal Society B: Biological*  
539 *Sciences*, 270(1512): 313–321.
- 540 Hillis, D. M., Pollock, D. D., McGuire, J. A., and Zwickl, D. J. 2003. Is sparse taxon  
541 sampling a problem for phylogenetic inference? *Systematic biology*, 52(1): 124–126.
- 542 Hinchliff, C. E., Smith, S. a., Allman, J. F., Burleigh, J. G., Chaudhary, R., Coghill, L. M.,  
543 Crandall, K. A., Deng, J., Drew, B. T., Gazis, R., Gude, K., Hibbett, D. S., Katz, L. a.,  
544 Laughinghouse, H. D., McTavish, E. J., Midford, P. E., Owen, C. L., Ree, R. H., Rees,  
545 J. a., Soltis, D. E., Williams, T. L., and Cranston, K. a. 2015. Synthesis of phylogeny  
546 and taxonomy into a comprehensive tree of life. *Proceedings of the National Academy of*  
547 *Sciences*, 112(41): 12764–12769.
- 548 Huerta-Cepas, J., Dopazo, J., and Gabaldón, T. 2010. ETE: a python Environment for  
549 Tree Exploration. *BMC Bioinformatics*, 11(1): 24.

- 550 Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T., and Aluru, S. 2017.  
551 High-throughput ANI Analysis of 90K Prokaryotic Genomes Reveals Clear Species  
552 Boundaries. *bioRxiv*.
- 553 Janssen, S., McDonald, D., Gonzalez, A., Navas-Molina, J. A., Jiang, L., Xu, Z. Z.,  
554 Winker, K., Kado, D. M., Orwoll, E., Manary, M., Mirarab, S., and Knight, R. 2018.  
555 Phylogenetic Placement of Exact Amplicon Sequences Improves Associations with  
556 Clinical Information. *mSystems*, 3(3): 00021–18.
- 557 Jukes, T. H. and Cantor, C. R. 1969. Evolution of protein molecules. In *In Mammalian*  
558 *protein metabolism, Vol. III (1969), pp. 21-132*, volume III, pages 21–132.
- 559 Koski, L. B. and Golding, G. B. 2001. The Closest BLAST Hit Is Often Not the Nearest  
560 Neighbor. *Journal of Molecular Evolution*, 52(6): 540–542.
- 561 Krause, L., Diaz, N. N., Goesmann, A., Kelley, S., Nattkemper, T. W., Rohwer, F.,  
562 Edwards, R. A., and Stoye, J. 2008. Phylogenetic classification of short environmental  
563 DNA fragments. *Nucleic acids research*, 36(7): 2230–9.
- 564 Lefort, V., Desper, R., and Gascuel, O. 2015. FastME 2.0: A comprehensive, accurate, and  
565 fast distance-based phylogeny inference program. *Molecular Biology and Evolution*,  
566 32(10): 2798–2800.
- 567 Leimeister, C.-A. and Morgenstern, B. 2014. kmacs: the k -mismatch average common  
568 substring approach to alignment-free sequence comparison. *Bioinformatics*, 30(14):  
569 2000–2008.
- 570 Leimeister, C.-A., Sohrabi-Jahromi, S., and Morgenstern, B. 2017. Fast and accurate  
571 phylogeny reconstruction using filtered spaced-word matches. *Bioinformatics*, page  
572 btw776.
- 573 Lockhart, P. J., Steel, M. A., Hendy, M. D., and Penny, D. 1994. Recovering evolutionary

- 574 trees under a more realistic model of sequence evolution. *Molecular Biology and*  
575 *Evolution*, 11(4): 605–612.
- 576 Mallo, D., De Oliveira Martins, L., and Posada, D. 2016. SimPhy : Phylogenomic  
577 Simulation of Gene, Locus, and Species Trees. *Systematic Biology*, 65(2): 334–344.
- 578 Matsen, F. A. and Evans, S. N. 2013. Edge Principal Components and Squash Clustering:  
579 Using the Special Structure of Phylogenetic Placement Data for Sample Comparison.  
580 *PLoS ONE*, 8(3).
- 581 Matsen, F. A., Kodner, R. B., and Armbrust, E. V. 2010. pplacer: linear time  
582 maximum-likelihood and Bayesian phylogenetic placement of sequences onto a fixed  
583 reference tree. *BMC bioinformatics*, 11(1): 538.
- 584 Miller, D. E., Staber, C., Zeitlinger, J., and Hawley, R. S. 2018. Highly Contiguous  
585 Genome Assemblies of 15 Drosophila Species Generated Using Nanopore Sequencing.  
586 *G3: Genes, Genomes, Genetics*, 8(10): 3131–3141.
- 587 Mirarab, S. and Warnow, T. 2015. ASTRAL-II: coalescent-based species tree estimation  
588 with many hundreds of taxa and thousands of genes. *Bioinformatics*, 31(12): i44–i52.
- 589 Mirarab, S., Nguyen, N., and Warnow, T. 2012. SEPP: SATé-Enabled Phylogenetic  
590 Placement. *Pacific Symposium On Biocomputing*, pages 247–58.
- 591 Mirarab, S., Nguyen, N., Guo, S., Wang, L.-S., Kim, J., and Warnow, T. 2015. PASTA:  
592 Ultra-Large Multiple Sequence Alignment for Nucleotide and Amino-Acid Sequences.  
593 *Journal of Computational Biology*, 22(05): 377–386.
- 594 Nguyen, L. T., Schmidt, H. A., Von Haeseler, A., and Minh, B. Q. 2015. IQ-TREE: A fast  
595 and effective stochastic algorithm for estimating maximum-likelihood phylogenies.  
596 *Molecular Biology and Evolution*, 32(1).
- 597 Nguyen, N. N.-p., Mirarab, S., Liu, B. B., Pop, M., and Warnow, T. 2014. TIPP:  
598 Taxonomic Identification and Phylogenetic Profiling. *Bioinformatics*, 30(24): 3548–3555.

- 599 Ondov, B. D., Treangen, T. J., Melsted, P., Mallonee, A. B., Bergman, N. H., Koren, S.,  
600 and Phillippy, A. M. 2016. Mash: fast genome and metagenome distance estimation  
601 using MinHash. *Genome Biology*, 17(1): 132.
- 602 Price, M. N., Dehal, P. S., and Arkin, A. P. 2010. FastTree-2 Approximately  
603 Maximum-Likelihood Trees for Large Alignments. *PLoS ONE*, 5(3): e9490.
- 604 Rzhetsky, A. and Nei, M. 1992. A simple method for estimating and testing  
605 minimum-evolution trees.
- 606 Saitou, N. and Nei, M. 1987. The neighbour-joining method: a new method for  
607 reconstructing phylogenetic trees. *Molecular biology and evolution*, 4(4): 406–425.
- 608 Sarmashghi, S., Bohmann, K., P. Gilbert, M. T., Bafna, V., and Mirarab, S. 2019. Skmer:  
609 assembly-free and alignment-free sample identification using genome skims. *Genome*  
610 *Biology*, 20(1): 34.
- 611 Stamatakis, A. 2014. RAxML version 8: A tool for phylogenetic analysis and post-analysis  
612 of large phylogenies. *Bioinformatics*, 30(9): 1312–1313.
- 613 Stark, M., Berger, S. A., Stamatakis, A., and von Mering, C. 2010. MLTreeMap—accurate  
614 Maximum Likelihood placement of environmental DNA sequences into taxonomic and  
615 functional reference phylogenies. *BMC genomics*, 11(1): 461.
- 616 Steel, M. 2009. A basic limitation on inferring phylogenies by pairwise sequence  
617 comparisons. *Journal of Theoretical Biology*, 256(3): 467–472.
- 618 Sukumaran, J. and Holder, M. T. 2010. DendroPy: a Python library for phylogenetic  
619 computing. *Bioinformatics*, 26(12): 1569–1571.
- 620 Sunagawa, S., Mende, D. R., Zeller, G., Izquierdo-Carrasco, F., Berger, S. a., Kultima,  
621 J. R., Coelho, L. P., Arumugam, M., Tap, J., Nielsen, H. B., Rasmussen, S., Brunak, S.,  
622 Pedersen, O., Guarner, F., de Vos, W. M., Wang, J., Li, J., Dore, J., Ehrlich, S. D.,

- 623 Stamatakis, a., and Bork, P. 2013. Metagenomic species profiling using universal  
624 phylogenetic marker genes. *Nature methods*, 10(12): 1196–1199.
- 625 Tamura, K. and Nei, M. 1993. Estimation of the Number of Nucleotide Substitutions in  
626 the Control Region of Mitochondrial-DNA in Humans and Chimpanzees. *Molecular*  
627 *biology and evolution*, 10(3): 512–526.
- 628 Tamura, K., Nei, M., and Kumar, S. 2004. Prospects for inferring very large phylogenies  
629 by using the neighbor-joining method. *Proceedings of the National Academy of Sciences*,  
630 101(30): 11030–11035.
- 631 von Mering, C., Hugenholtz, P., Raes, J., Tringe, S. G., Doerks, T., Jensen, L. J., Ward,  
632 N., and Bork, P. 2007. Quantitative Phylogenetic Assessment of Microbial Communities  
633 in Diverse Environments. *Science*, 315(5815): 1126–1130.
- 634 Waddell, P. J. and Steel, M. 1997. General Time-Reversible Distances with Unequal Rates  
635 across Sites: Mixing  $\Gamma$  and Inverse Gaussian Distributions with Invariant Sites.  
636 *Molecular Phylogenetics and Evolution*, 8(3): 398–414.
- 637 Warnow, T. 2017. *Computational phylogenetics: An introduction to designing methods for*  
638 *phylogeny estimation*. Cambridge University Press.
- 639 Wheeler, T. J. 2009. Large-scale neighbor-joining with NINJA. In *Algorithms in*  
640 *Bioinformatics*, pages 375–389. Springer.
- 641 Xia, X. 2009. Information-theoretic indices and an approximate significance test for testing  
642 the molecular clock hypothesis with genetic distances. *Molecular Phylogenetics and*  
643 *Evolution*, 52(3): 665–676.
- 644 Xia, X. 2018. DAMBE7: New and Improved Tools for Data Analysis in Molecular Biology  
645 and Evolution. *Molecular Biology and Evolution*, 35(6): 1550–1552.
- 646 Yi, H. and Jin, L. 2013. Co-phylog: an assembly-free phylogenomic approach for closely  
647 related organisms. *Nucleic Acids Research*, 41(7): e75–e75.

REFERENCES

31

- 648 Yin, C., Shen, G., Guo, D., Wang, S., Ma, X., Xiao, H., Liu, J., Zhang, Z., Liu, Y., Zhang,  
649 Y., Yu, K., Huang, S., and Li, F. 2016. InsectBase: a resource for insect genomes and  
650 transcriptomes. *Nucleic Acids Research*, 44(D1): D801–D807.
- 651 Zwickl, D. J. and Hillis, D. M. 2002. Increased taxon sampling greatly reduces  
652 phylogenetic error. *Systematic biology*, 51(4): 588–98.





APPENDIX A. PROOFS AND DERIVATIONS

Recall the following notations.

- For any node  $u$  and exponents  $a \in \mathbb{Z}$  and  $b \in \mathbb{N}^+$ , let

$$S(a, b, u) = \sum_{i \in g(u)} \delta_{qi}^a d_{ui}^b$$

$$R(a, b, u) = \sum_{i \notin g(u)} \delta_{qi}^a d_{p(u)i}^b \text{ defined for } u \in V \setminus \{1\}$$

- For  $b = 0$ , let  $S(a, 0, u) = \sum_{i \in g(u)} \delta_{qi}^a$  and let  $S'(a, u)$  be a shorthand for  $S(a, 0, u)$ .

$$\text{Similarly, let } R(a, 0, u) = R'(a, u) = \sum_{i \notin g(u)} \delta_{qi}^a.$$

*Proof of Lemma 2*

*Proof.* Recall the dynamic programming recursions of Equations 3 and 4:

$$S(a, b, u) = \sum_{j=0}^b \sum_{v \in c(u)} l(v)^j \binom{b}{j} S(a, b-j, v), \quad \text{for } u \notin \mathcal{L} \setminus \{1\}$$

$$R(a, b, u) = \sum_{j=0}^b \left( l(p(u))^j \binom{b}{j} R(a, b-j, p(u)) + \sum_{v \in sib(u)} l(v)^j \binom{b}{j} S(a, b-j, v) \right) \text{ for } u \notin \{1, 1'\}$$

Since  $u$  is not a leaf, for each leaf  $i \in g(u)$ , there exists a  $v \in c(u)$  such that the directed path from  $u$  to  $i$  passes through  $v$ . Therefore every leaf  $i$  can be grouped under its corresponding  $v$ .

$$\begin{aligned} S(a, b, u) &= \sum_{i \in g(u)} \delta_{qi}^a d_{ui}^b = \sum_{v \in c(u)} \sum_{i \in g(v)} \delta_{qi}^a (l(v) + d_{vi})^b = \sum_{j=0}^b \sum_{v \in c(u)} \sum_{i \in g(v)} \delta_{qi}^a d_{vi}^{b-j} l(v)^j \binom{b}{j} \\ &= \sum_{j=0}^b \sum_{v \in c(u)} l(v)^j \binom{b}{j} S(a, b-j, v) \end{aligned}$$

Similarly, given the condition  $u \neq 1$ , for each leaf  $i \notin g(u)$ , either (1) there exists  $v \in sib(u)$  such that the directed path from  $p(u)$  to  $i$  passes through  $v$ , or (2) undirected path between  $i$  and  $p(u)$  passes through  $p(p(u))$ .

$$\begin{aligned}
 R(a, b, u) &= \sum_{i \notin g(u)} \delta_{qi}^a d_{p(u)i}^b = \sum_{v \in \text{sib}(u)} \sum_{i \in g(v)} \delta_{qi}^a (l(v) + d_{vi})^b + \sum_{i \notin g(p(u))} \delta_{qi}^a (l(p(u)) + d_{p(u)i})^b \\
 &= \sum_{j=0}^b \sum_{v \in \text{sib}(u)} \sum_{i \in g(v)} \delta_{qi}^a d_{vi}^{b-j} l(v)^j \binom{b}{j} + \sum_{j=0}^b \sum_{i \notin g(p(u))} \delta_{qi}^a d_{p(u)i}^{b-j} l(p(u))^j \binom{b}{j} \\
 &= \sum_{j=0}^b \sum_{v \in \text{sib}(u)} l(v)^j \binom{b}{j} S(a, b-j, v) + \sum_{j=0}^b l(p(u))^j \binom{b}{j} R(a, b-j, p(u))
 \end{aligned}$$

665 Boundary conditions follow from definitions. For  $u \notin \mathcal{L} \setminus \{1\}$ , since  $d_{ii} = 0$ , we have  
 666  $S(a, b, u) = 0$  and it's trivial to see  $S'(a, u) = \delta_{qu}^a$ . For  $R(, , )$  recursions, the boundary case  
 667 happens at the unique child of the root, which we denote as  $1'$ . Based on the definition,  
 668 since the only  $i \notin g(1')$  is 1, and  $d_{p(1')1}^b = 0$ , we trivially have  $R(a, b, 1') = 0$ . For  $b = 0$ ,  
 669  $R'(a, 1') = \delta_{q1}^a$ .

670 A post-order traversal on  $T^*$  can compute  $S(a, b, u)$ , and a subsequent pre-order  
 671 traversal can compute  $R(a, b, u)$ , both in constant time and using constant memory per  
 672 node. Recall that  $a$  and  $b$  are both no more than  $k$ , which is a constant. Thus, time and  
 673 memory complexity of this dynamic programming is  $\Theta(bn)$ , which translates to  $\Theta(n)$  in  
 674 least squares setting, where  $b \leq 2$ . □

675 *Proof of Lemma 3.*

Recall  $w_{qi} = \delta_{qi}^{-k}$  and that Equation 2:

$$Q(P) = \sum_{i=1}^n w_{qi} (\delta_{qi} - d_{qi}(P))^2 = \sum_{i=1}^n \delta_{qi}^{-k} (\delta_{qi} - d_{qi}(P))^2$$

676 *Proof.* Equation 2 can be re-written as:

$$Q(P) = \sum_{i \in g(u)} \delta_{qi}^{-k} (\delta_{qi} - d_{ui}(P) - x_1 + x_2 - l(u))^2 + \sum_{i \notin g(u)} \delta_{qi}^{-k} (\delta_{qi} - d_{p(u)i}(P) - x_1 - x_2)^2 \tag{S1}$$

677 By simple rearrangement of the terms, we can rewrite Equation S1 as follows.

$$\begin{aligned}
 Q(P(u, x_1, x_2)) &= R'(2 - k, u) + S'(2 - k, u) + R(-k, 2, u) + S(-k, 2, u) \\
 &\quad + 2(x_1 + x_2)(R(-k, 1, u) - R'(1 - k, u)) \\
 &\quad + 2(x_1 + l(u) - x_2)(S(-k, 1, u) - S'(1 - k, u)) \tag{S2} \\
 &\quad + (x_1 + x_2)^2 R(-k, 1, u) + 2(x_1 + l(u) - x_2)^2 S(-k, 1, u) \\
 &\quad - 2R(1 - k, 1, u) - 2S(1 - k, 1, u)
 \end{aligned}$$

678 Note that computing  $Q(P(u, x_1, x_2))$  requires only  $S(\cdot, u)$  and  $R(\cdot, u)$  values and  $l(u)$ .

679 Thus, computing  $Q(P)$  requires only computing  $S(a, b, u)$  and  $R(a, b, u)$  values for

680  $-k \leq a \leq 2 - k$  and  $0 \leq b \leq 2$ .

681

□

682

*Proof of Lemma 4.*

Recall definitions

$$\begin{aligned}
 S(a, b, u) &= \sum_{i \in g(u)} \delta_{qi}^a d_{ui}^b \quad (\text{for } b > 0) & \text{and} & \quad S'(a, u) = S(a, 0, u) = \sum_{i \in g(u)} \delta_{qi}^a \\
 R(a, b, u) &= \sum_{i \notin g(u)} \delta_{qi}^a d_{p(u)i}^b \quad (\text{for } b > 0) & \text{and} & \quad R'(a, u) = R(a, 0, u) = \sum_{i \notin g(u)} \delta_{qi}^a
 \end{aligned}$$

and recall Eq. S1:

$$Q(P) = \sum_{i \in g(u)} \delta_{qi}^{-k} (\delta_{qi} - d_{ui}(P) - x_1 + x_2 - l(u))^2 + \sum_{i \notin g(u)} \delta_{qi}^{-k} (\delta_{qi} - d_{p(u)i}(P) - x_1 - x_2)^2.$$

683 *Proof.* We take the derivative of Eq. S1 with respect to  $x_1$  and set it equal to zero:

$$\begin{aligned}
 \frac{\partial Q(P)}{\partial x_1} &= -2 \sum_{i \in g(u)} \delta_{qi}^{-k} (\delta_{qi} - d_{ui}(P) - x_1 + x_2 - l(u)) - 2 \sum_{i \notin g(u)} \delta_{qi}^{-k} (\delta_{qi} - d_{p(u)i}(P) - x_1 - x_2) = 0 \\
 &\implies \left( \sum_{i \in g(u)} \delta_{qi}^{-k} + \sum_{i \notin g(u)} \delta_{qi}^{-k} \right) x_1 + \left( - \sum_{i \in g(u)} \delta_{qi}^{-k} + \sum_{i \notin g(u)} \delta_{qi}^{-k} \right) x_2 \\
 &\quad - \sum_{i \in g(u)} \delta_{qi}^{1-k} - \sum_{i \notin g(u)} \delta_{qi}^{1-k} + \sum_{i \in g(u)} \delta_{qi}^{-k} d_{ui}(P) + \sum_{i \notin g(u)} \delta_{qi}^{-k} d_{p(u)i}(P) + l(u) \sum_{i \in g(u)} \delta_{qi}^{-k} = 0 \\
 &\implies (S'(-k, u) + R'(-k, u)) x_1 + (-S'(-k, u) + R'(-k, u)) x_2 = \\
 &\quad S'(1-k, u) + R'(1-k, u) - S(-k, 1, u) - R(-k, 1, u) - l(u) S'(-k, u)
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 \frac{\partial Q(P)}{\partial x_2} &= 2 \sum_{i \in g(u)} \delta_{qi}^{-k} (\delta_{qi} - d_{ui}(P) - x_1 + x_2 - l(u)) - 2 \sum_{i \notin g(u)} \delta_{qi}^{-k} (\delta_{qi} - d_{p(u)i}(P) - x_1 - x_2) = 0 \\
 &\implies \left( - \sum_{i \in g(u)} \delta_{qi}^{-k} + \sum_{i \notin g(u)} \delta_{qi}^{-k} \right) x_1 + \left( \sum_{i \in g(u)} \delta_{qi}^{-k} + \sum_{i \notin g(u)} \delta_{qi}^{-k} \right) x_2 \\
 &\quad + \sum_{i \in g(u)} \delta_{qi}^{1-k} - \sum_{i \notin g(u)} \delta_{qi}^{1-k} - \sum_{i \in g(u)} \delta_{qi}^{-k} d_{ui}(P) + \sum_{i \notin g(u)} \delta_{qi}^{-k} d_{p(u)i}(P) - l(u) \sum_{i \in g(u)} \delta_{qi}^{-k} = 0 \\
 &\implies (-S'(-k, u) + R'(-k, u)) x_1 + (+S'(-k, u) + R'(-k, u)) x_2 = \\
 &\quad -S'(1-k, u) + R'(1-k, u) + S(-k, 1, u) - R(-k, 1, u) + l(u) S'(-k, u)
 \end{aligned}$$

684 These two linear equations have a unique solution for the pair  $x_1, x_2$  if and only if the  
 685 following matrix has the full rank:

$$H = \begin{bmatrix} R'(-k, u) + S'(-k, u) & R'(-k, u) - S'(-k, u) \\ R'(-k, u) - S'(-k, u) & R'(-k, u) + S'(-k, u) \end{bmatrix}.$$

686 Determinant of  $H$  is  $\det(H) = 4R'(-k, u)S'(-k, u)$ . Assuming that  $\delta_{qi} > 0$  for all  
 687  $i \in \mathcal{L}$ , both  $R'(-k, u) > 0$  and  $S'(-k, u) > 0$  hold. Therefore,  $H$  has the full rank.  
 688 However,  $\delta_{qi} = 0$  for  $q \neq i$  can be encountered on real data, especially for low divergence  
 689 times, low evolutionary rates, or short sequences. In this case, APPLES is designed to  
 690 place  $q$  on the pendant edge of  $i$  with  $x_1 = 0$  and  $x_2 = l(i)$ . In case there are multiple  
 691 leaves  $i$  that satisfy  $\delta_{qi} = 0$  for  $q \neq i$ , we pick one of them arbitrarily.  $\square$

*Proof of Theorem 1*

692

693 *Proof.* First, using two traversals of the tree, we compute all the  $S(a, b, u)$  and  $R(a, b, u)$   
694 values by Lemma 2. To find the optimal placement edge, we first optimize  $Q(P(u, x_1, x_2))$   
695 for all  $u \in V \setminus \{1\}$ . By Lemma 4, this task requires only constant time after the  
696 precomputations. Then, for each node, we compute  $Q(P(u, x_1, x_2))$  in constant time for the  
697 optimal  $u, x_1, x_2$  by Lemma 3. Thus, each node is processed in linear time and the whole  
698 optimization requires linear time. Note that the system of equations (shown in Lemma 4)  
699 will not have a solution iff  $\delta_{qi} \leq 0$  for some  $i$ ; if there is  $\delta_{qi} = 0$ , we make  $q$  sister to  $i$ ,  
700 breaking ties arbitrarily. □

*Proof of Lemma 5*

701

702 *Proof.* Eigenvalues of the Hessian matrix of  $Q(P(u, x_1, x_2))$  are  $2R'(-k, u)$  and  $2S'(-k, u)$ ,  
703 which are both non-negative since  $\delta_{qi} \geq 0$  for  $i \in \mathcal{L}$ . Thus, the Hessian matrix is positive  
704 semidefinite and therefore  $P(u, x_1, x_2)$  is a convex function of  $x_1$  and  $x_2$ . □

APPENDIX B. SUPPLEMENTARY FIGURES

REFERENCES

39

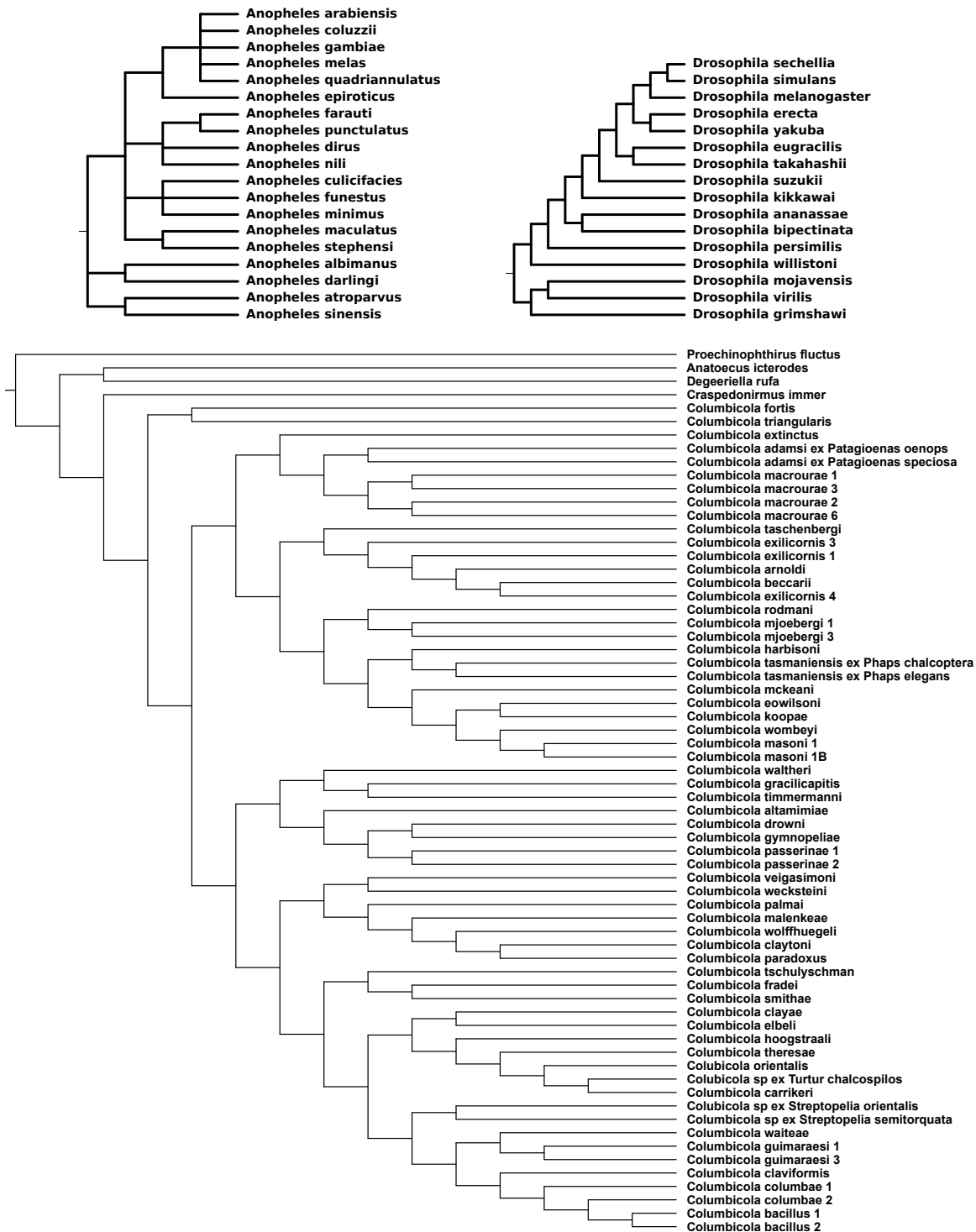


Fig. S1. The reference biological trees obtained from Open Tree of Life (*Drosophila* and *Anopheles*) and from Boyd *et al.* (2017) (*Columbicola*).

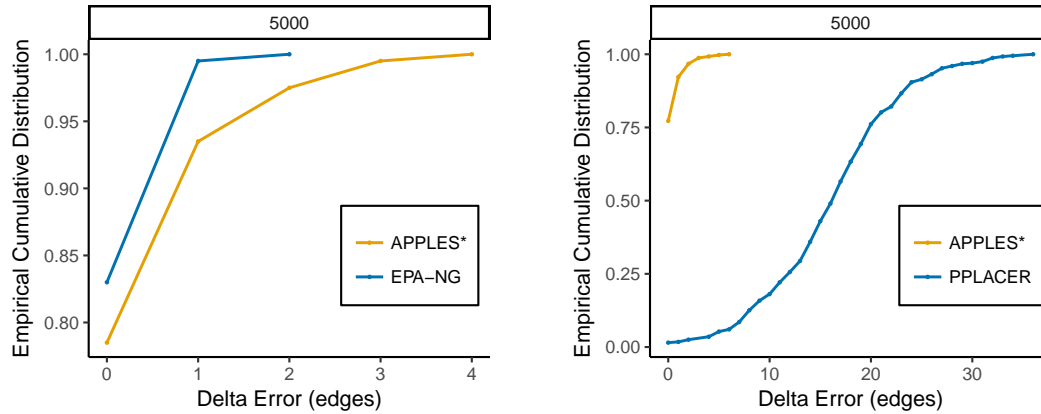


Fig. S2. **APPLES** versus **ML** tools on **5,000** backbone trees. The empirical cumulative distribution of the delta error is shown. We compare pplacer, EPA-ng, and APPLES\* on RNASim-varied backbone dataset with 5000 leaves. Distributions is over 551 cases where pplacer could run for the panel on the right and 200 points where EPA-NG could run for the panel on the left.

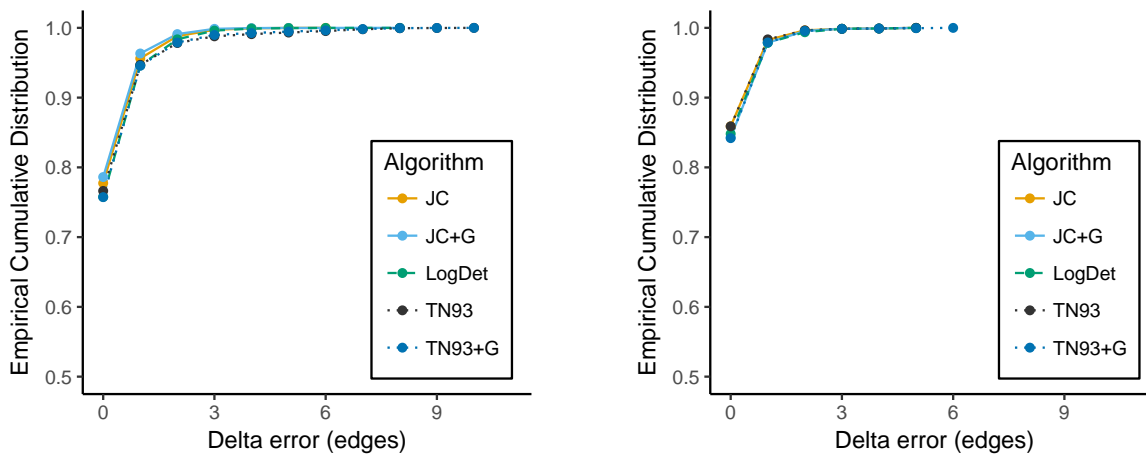


Fig. S3. **Comparing various models of DNA evolution.** For the GTR (a) and RNASim-heterogeneous (b) datasets, we show the delta error (edges) of APPLES\* run with five distance matrices calculated based on different models of DNA evolution. All model parameters are estimated per pair of sequences. The five models have similar accuracy.



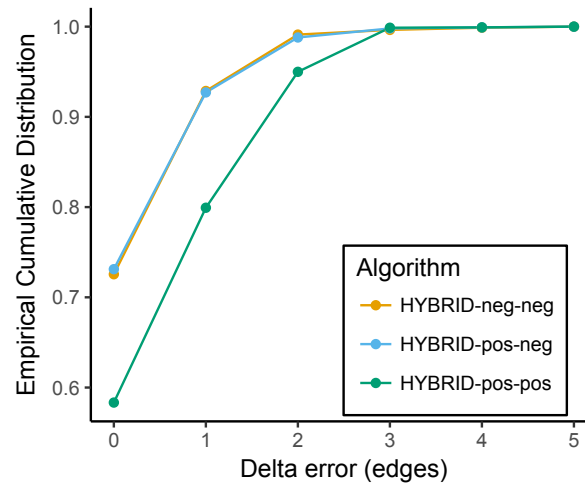


Fig. S4. **The effect of imposing positivity constraint on accuracy on Hybrid.** The HYBRID approach does not benefit from imposing positivity constraint on its second (ME) stage.

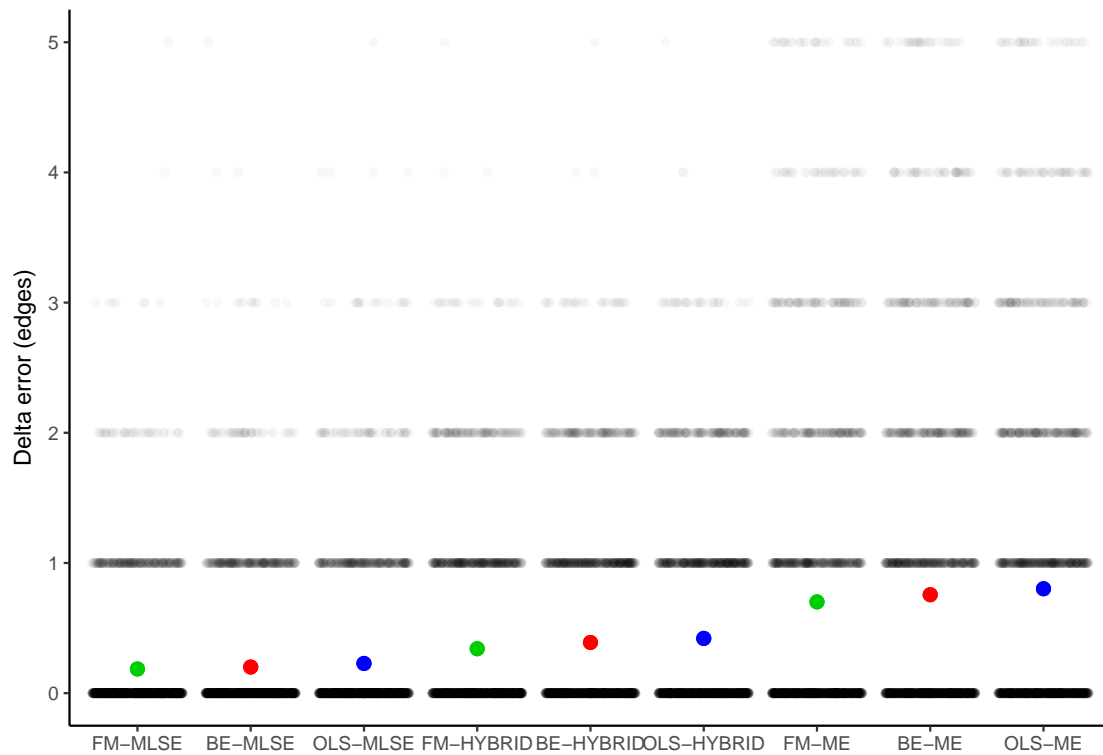


Fig. S5. **Comparing APPLES versions.** For the RNASim dataset (without controlling for the diameter), we show the delta error (edges) of APPLES run with three options for weighting: FM (green), BE (red), and OLS (blue), and three options for selection strategy (MLSE, ME, and Hybrid). For each method, the mean (colored circle) and standard errors (lines; too small to see) are shown over 2500 data points, each shown as dots. Some of the methods occasionally have error above 5 branches, but for better resolution, we cap the y-axis at 5.

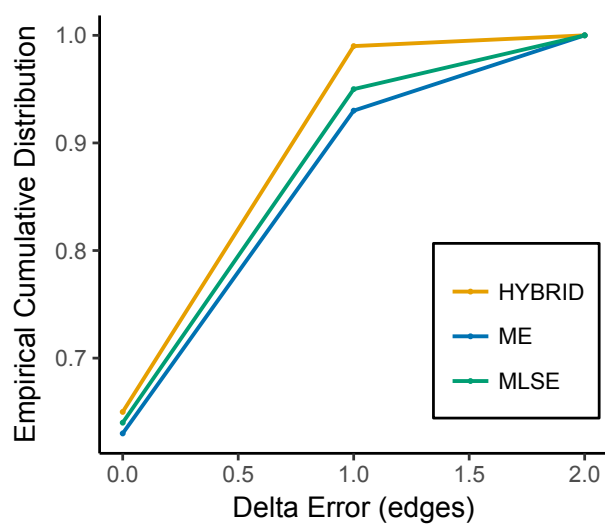


Fig. S6. **APPLES-HYBRID has higher accuracy on sparse RNAsim dataset** On the RNAsim dataset, we chose 20 sequences randomly from the larger RNAsim-heterogeneous dataset; here, APPLES-HYBRID has higher accuracy than APPLES\* (MLSE).

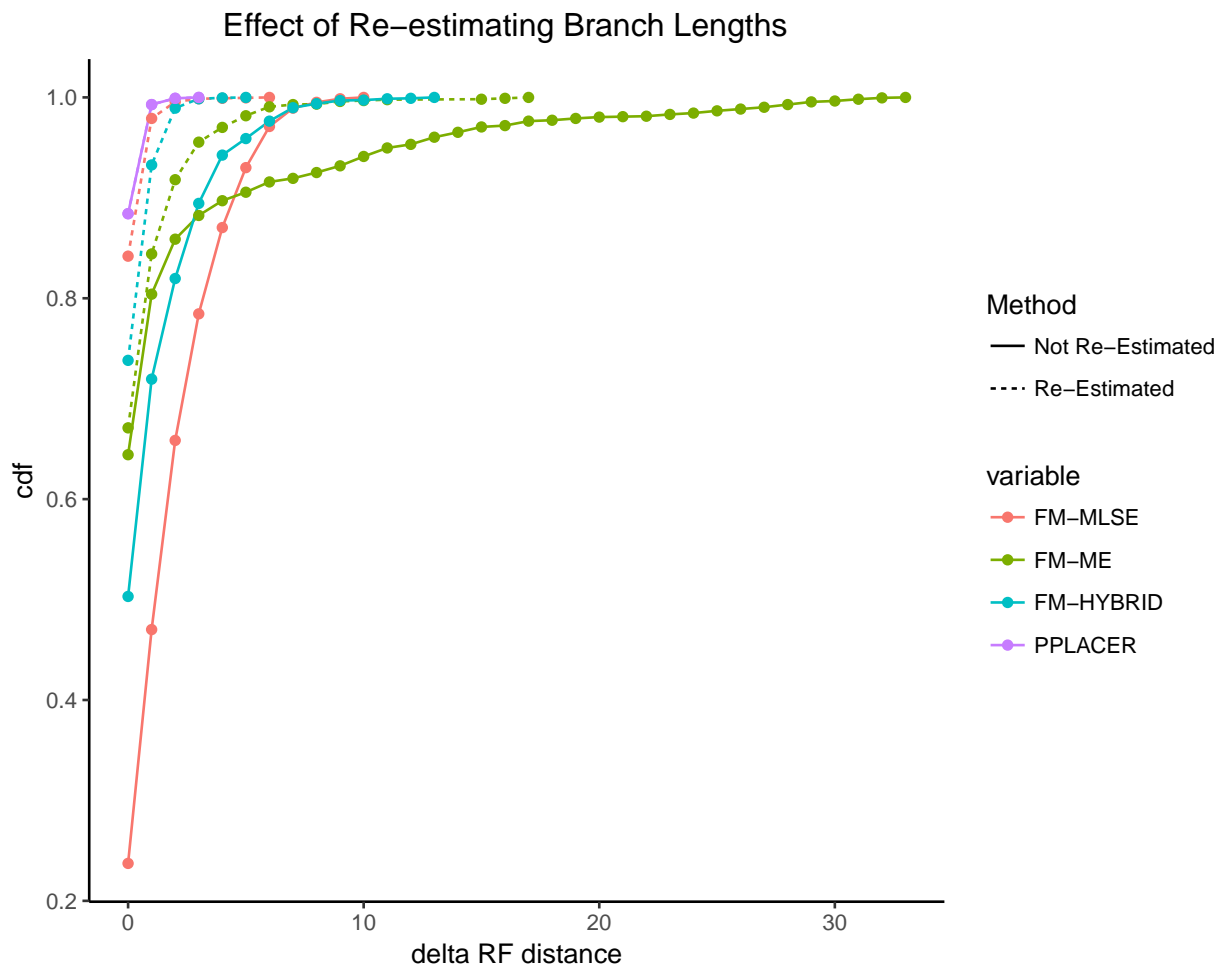


Fig. S7. **The effect of reestimating branch lengths of the backbone tree on accuracy.** We show the accuracy of pplacer and APPLES-FM with its three optimization criteria. APPLES is run both with (dotted) and without (solid) re-estimating branch lengths in the backbone tree using the same model (here, TN93+ $\Gamma$ ) used for computing distances of query sequences to backbone sequences. FastME\* is used to re-estimate branch lengths. Accuracy improves dramatically by recomputing backbone branch lengths using the same model. The case labeled “Not re-estimated” uses branch lengths produced using RAxML under the GTR+ $\Gamma$  model.

APPENDIX C. SUPPLEMENTARY TABLES

Table S1. GenBank accession numbers and URLs for the dataset of 22 Anopheles genomes

Species	GenBank assembly accession	URL
<i>Anopheles albimanus</i>	GCA_000349125.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_albimanus/Anopheles_albimanus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_albimanus/Anopheles_albimanus_genomic.fasta.gz</a>
<i>Anopheles arabiensis</i>	GCA_000349185.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_arabiensis/Anopheles_arabiensis_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_arabiensis/Anopheles_arabiensis_genomic.fasta.gz</a>
<i>Anopheles atroparvus</i>	GCA_000473505.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_atroparvus/Anopheles_atroparvus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_atroparvus/Anopheles_atroparvus_genomic.fasta.gz</a>
<i>Anopheles christyi</i>	GCA_000349165.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_christyi/Anopheles_christyi_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_christyi/Anopheles_christyi_genomic.fasta.gz</a>
<i>Anopheles coluzzii</i>	-	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_coluzzii/Anopheles_coluzzii_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_coluzzii/Anopheles_coluzzii_genomic.fasta.gz</a>
<i>Anopheles culicifacies</i>	GCA_000473375.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_culicifacies/Anopheles_culicifacies_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_culicifacies/Anopheles_culicifacies_genomic.fasta.gz</a>
<i>Anopheles darlingi</i>	GCA_000211455.3	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_darlingi/Anopheles_darlingi_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_darlingi/Anopheles_darlingi_genomic.fasta.gz</a>
<i>Anopheles dirus</i>	GCA_000349145.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_dirus/Anopheles_dirus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_dirus/Anopheles_dirus_genomic.fasta.gz</a>
<i>Anopheles epiroticus</i>	GCA_000349105.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_epiroticus/Anopheles_epiroticus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_epiroticus/Anopheles_epiroticus_genomic.fasta.gz</a>
<i>Anopheles farauti</i>	GCA_000956265.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_farauti/Anopheles_farauti_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_farauti/Anopheles_farauti_genomic.fasta.gz</a>
<i>Anopheles funestus</i>	GCA_000349085.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_funestus/Anopheles_funestus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_funestus/Anopheles_funestus_genomic.fasta.gz</a>
<i>Anopheles gambiae</i>	GCA_000150785.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_gambiae/Anopheles_gambiae_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_gambiae/Anopheles_gambiae_genomic.fasta.gz</a>
<i>Anopheles koliensis</i>	GCA_000956275.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_koliensis/Anopheles_koliensis_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_koliensis/Anopheles_koliensis_genomic.fasta.gz</a>
<i>Anopheles maculatus</i>	GCA_000473185.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_maculatus/Anopheles_maculatus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_maculatus/Anopheles_maculatus_genomic.fasta.gz</a>
<i>Anopheles melas</i>	GCA_000473525.2	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_melas/Anopheles_melas_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_melas/Anopheles_melas_genomic.fasta.gz</a>
<i>Anopheles merus</i>	GCA_000473845.2	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_merus/Anopheles_merus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_merus/Anopheles_merus_genomic.fasta.gz</a>
<i>Anopheles minimus</i>	GCA_000349025.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_minimus/Anopheles_minimus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_minimus/Anopheles_minimus_genomic.fasta.gz</a>
<i>Anopheles nili</i>	GCA_000439205.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_nili/Anopheles_nili_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_nili/Anopheles_nili_genomic.fasta.gz</a>
<i>Anopheles punctulatus</i>	GCA_000956255.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_punctulatus/Anopheles_punctulatus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_punctulatus/Anopheles_punctulatus_genomic.fasta.gz</a>
<i>Anopheles quadriannulatus</i>	GCA_000349065.1	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_quadriannulatus/Anopheles_quadriannulatus_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_quadriannulatus/Anopheles_quadriannulatus_genomic.fasta.gz</a>
<i>Anopheles sinensis</i>	GCA_000441895.2	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_sinensis/Anopheles_sinensis_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_sinensis/Anopheles_sinensis_genomic.fasta.gz</a>
<i>Anopheles stephensi</i>	GCA_000300775.2	<a href="http://www.insect-genome.com/data/genome_download/Anopheles_stephensi/Anopheles_stephensi_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Anopheles_stephensi/Anopheles_stephensi_genomic.fasta.gz</a>

Table S2. GenBank accession numbers and URLs for the dataset of 21 Drosophila genomes

Species	GenBank assembly accession	URL
<i>Drosophila ananassae</i>	GCA_000005115.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_ananassae/Drosophila_ananassae_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_ananassae/Drosophila_ananassae_genomic.fasta.gz</a>
<i>Drosophila biarmipes</i>	GCA_000233415.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_biarmipes/Drosophila_biarmipes_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_biarmipes/Drosophila_biarmipes_genomic.fasta.gz</a>
<i>Drosophila bipectinata</i>	GCA_000236285.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_bipectinata/Drosophila_bipectinata_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_bipectinata/Drosophila_bipectinata_genomic.fasta.gz</a>
<i>Drosophila elegans</i>	GCA_000224195.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_elegans/Drosophila_elegans_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_elegans/Drosophila_elegans_genomic.fasta.gz</a>
<i>Drosophila erecta</i>	GCA_000005135.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_erecta/Drosophila_erecta_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_erecta/Drosophila_erecta_genomic.fasta.gz</a>
<i>Drosophila eugracilis</i>	GCA_000236325.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_eugracilis/Drosophila_eugracilis_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_eugracilis/Drosophila_eugracilis_genomic.fasta.gz</a>
<i>Drosophila ficusphila</i>	GCA_000220665.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_ficusphila/Drosophila_ficusphila_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_ficusphila/Drosophila_ficusphila_genomic.fasta.gz</a>
<i>Drosophila grimshawi</i>	GCA_000005155.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_grimshawi/Drosophila_grimshawi_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_grimshawi/Drosophila_grimshawi_genomic.fasta.gz</a>
<i>Drosophila kikkawai</i>	GCA_000224215.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_kikkawai/Drosophila_kikkawai_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_kikkawai/Drosophila_kikkawai_genomic.fasta.gz</a>
<i>Drosophila melanogaster</i>	GCA_000778455.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_melanogaster/Drosophila_melanogaster_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_melanogaster/Drosophila_melanogaster_genomic.fasta.gz</a>
<i>Drosophila miranda</i>	GCA_000269505.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_miranda/Drosophila_miranda_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_miranda/Drosophila_miranda_genomic.fasta.gz</a>
<i>Drosophila mojavensis</i>	GCA_000005175.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_mojavensis/Drosophila_mojavensis_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_mojavensis/Drosophila_mojavensis_genomic.fasta.gz</a>
<i>Drosophila persimilis</i>	GCA_000005195.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_persimilis/Drosophila_persimilis_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_persimilis/Drosophila_persimilis_genomic.fasta.gz</a>
<i>Drosophila rhopaloa</i>	GCA_000236305.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_rhopaloea/Drosophila_rhopaloea_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_rhopaloea/Drosophila_rhopaloea_genomic.fasta.gz</a>
<i>Drosophila sechellia</i>	GCA_000005215.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_sechellia/Drosophila_sechellia_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_sechellia/Drosophila_sechellia_genomic.fasta.gz</a>
<i>Drosophila simulans</i>	GCA_000259055.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_simulans/Drosophila_simulans_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_simulans/Drosophila_simulans_genomic.fasta.gz</a>
<i>Drosophila suzukii</i>	GCA_000472105.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_suzukii/Drosophila_suzukii_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_suzukii/Drosophila_suzukii_genomic.fasta.gz</a>
<i>Drosophila takahashii</i>	GCA_000224235.2	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_takahashii/Drosophila_takahashii_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_takahashii/Drosophila_takahashii_genomic.fasta.gz</a>
<i>Drosophila virilis</i>	GCA_000005245.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_virilis/Drosophila_virilis_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_virilis/Drosophila_virilis_genomic.fasta.gz</a>
<i>Drosophila willistoni</i>	GCA_000005925.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_willistoni/Drosophila_willistoni_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_willistoni/Drosophila_willistoni_genomic.fasta.gz</a>
<i>Drosophila yakuba</i>	GCA_000005975.1	<a href="http://www.insect-genome.com/data/genome_download/Drosophila_yakuba/Drosophila_yakuba_genomic.fasta.gz">http://www.insect-genome.com/data/genome_download/Drosophila_yakuba/Drosophila_yakuba_genomic.fasta.gz</a>

Table S3. GenBank accession numbers of microbial species used in contamination removal.

Species	GenBank assembly accession
<i>Pasteurella langaensis</i>	GCA_003096995.1
<i>Providencia stuartii</i>	GCA_001558855.2
<i>Serratia marcescens</i>	GCA_000783915.2
<i>Shigella flexneri</i>	GCA_000006925.2
<i>Commensalibacter intestini</i>	GCA_002153535.1
<i>Acetobacter malorum</i>	GCA_002153605.1
<i>Acetobacter pomorum</i>	GCA_002456135.1
<i>Lactobacillus plantarum</i>	GCA_000203855.3
<i>Lactobacillus brevis</i>	GCA_003184305.1
<i>Enterococcus faecalis</i>	GCA_002208945.2
<i>Vagococcus teuberi</i>	GCA_001870205.1
<i>Wolbachia</i>	GCA_000022285.1

	0.1G			0.5G		
	%	$\Delta e$	$e_{max}$	%	$\Delta e$	$e_{max}$
<b>(a) <i>Columbicola</i></b>						
APPLES*	97	0.03	1	92	0.08	1
APPLES-ME	84	0.28	5	87	0.21	5
APPLES-HYBRID	87	0.16	2	87	0.16	2
CLOSEST	54	1.15	7	58	0.91	8
DE-NOVO	98	0.02	1	92	0.08	1
<b>(b) <i>Anopheles</i></b>						
APPLES*	95	0.05	1	95	0.05	1
APPLES-ME	95	0.05	1	95	0.05	1
APPLES-HYBRID	95	0.05	1	95	0.05	1
CLOSEST	91	0.09	1	95	0.05	1
DE-NOVO	95	0.05	1	95	0.05	1
<b>(c) <i>Drosophila</i></b>						
APPLES*	71	0.29	1	71	0.33	2
APPLES-ME	67	0.42	2	67	0.48	2
APPLES-HYBRID	67	0.33	1	67	0.38	2
CLOSEST	57	0.62	3	57	0.57	2
DE-NOVO	71	0.29	1	71	0.33	2

Table S4. **Assembly-free placement of genome-skims.** We show the percentage of correct placements (those that do not increase  $\Delta e$ ), average delta error ( $\Delta e$ ), and maximum delta error ( $e_{max}$ ) for APPLES, assignment to the CLOSEST species, and the placement to the position in the backbone (DE-NOVO) over the 61 (a), 22 (b), and 21 (c) placements. Results are shown for skims with 0.1 and 0.5Gbp of reads. Delta error is the increase in the number missing branches between the reference tree and the backbone tree after placing each query.

707

## APPENDIX D. COMMANDS

708

### *Sampling Clades*

709

For sampling clades of size at most 250 from a tree "tree.nwk", we used the

710

TreeCluster package, available at <https://github.com/niemasd/TreeCluster>.

```
#!/bin/bash
```

```
python TreeCluster/TreeCluster.py -i 250 -o clusters.txt -t tree.nwk
```

```
-m count_max_clade
```

711

### *Backbone tree estimation*

712

When multiple sequence alignment is available, we used the following RAxML

713

command to compute backbone tree for all datasets except RNAsim varied size dataset.

714

We used RAxML version 7.2.6

```
#!/bin/bash
```

```
raxmlHPC-PTHREADS -m GTRGAMMA -p 88 -n REF -s aln_dna.phy -T 6
```

715

For RNAsim varied size dataset, we used FastTreeMP version 2.1.10 for estimating

716

backbone topology. We run FastTreeMP with the following command:

```
#!/bin/bash
```

```
FastTreeMP -nosupport -gtr -gamma -nt -log tree.log < aln_dna.fa > tree.nwk
```

717

For alignment free datasets such as Drosophila dataset, we computed backbone tree

718

using FastME\* (based on FastME version 2.1.6.1) which is available at

719

<https://github.com/balabanmetin/FastME-personal-copy>. FastME\* is run with the

720

following command:

```
#!/bin/bash
```

```
fastme -i dist.mat -o tree.nwk -T 1
```

721 Note that we performed Jukes-Cantor correction on the distance matrix "dist.mat"  
722 before running FastME\*.

723 *Backbone tree branch length re-estimation*

724 When multiple sequence alignment is available, we used FastME\* to recompute  
725 backbone tree branch lengths for all datasets except RNAsim varied size dataset. We run  
726 FastME\* with the following command:

```
#!/bin/bash  
fastme -dJ -i aln_dna.phy -u RAxML_result.REF -o tree_me.nwk
```

727 For RNAsim varied size dataset, we used RAxML version 7.2.6 for re-estimating  
728 ML based branch lengths and used that tree for performing placements using pplacer.  
729 RAxML is run with the following command:

```
#!/bin/bash  
raxmlHPC-PTHREADS -f e -t tree.nwk -m GTRGAMMA -s aln_dna.phy -n REF -p 1984 -T 8
```

730 For the same dataset, we used FastTree again for re-estimating Minimum Evolution  
731 based branch lengths and used that tree for performing placements using APPLES.  
732 FastTree is run with the following command:

```
#!/bin/bash  
FastTreeMP -nosupport -nt -nome -noml -log tree.log  
-intree tree.nwk < aln_dna.fa > tree_me.nwk
```

733 *Performing placement*

734 We performed phylogenetic placement of a query using pplacer with the following  
735 commands:

```
#!/bin/bash  
nw_prune RAxML_result.REF query > backbone.nwk
```

```
pplacer -m GTR -s RAxML_info.REF -t backbone.nwk -o query.jplace aln_dna.fa
```