



Peer Community Journal

Section: Mathematical & Computational Biology

Research article

Published
2025-07-15

Cite as

Samson Weiner, Yutian Feng, J. Peter Gogarten and Mukul S. Bansal (2025) *A systematic assessment of phylogenomic approaches for microbial species tree reconstruction*, Peer Community Journal, 5: e72.

Correspondence

mukul.bansal@uconn.edu

Peer-review

Peer reviewed and
recommended by

PCI Mathematical &
Computational Biology,

<https://doi.org/10.24072/pci.mcb.100408>



This article is licensed
under the Creative Commons
Attribution 4.0 License.

A systematic assessment of phylogenomic approaches for microbial species tree reconstruction

Samson Weiner¹, Yutian Feng², J. Peter Gogarten^{ID,2,3}, and Mukul S. Bansal^{ID,1,3}

Volume 5 (2025), article e72

<https://doi.org/10.24072/pcjournal.579>

Abstract

A key challenge in microbial phylogenomics is that microbial gene families are often affected by extensive horizontal gene transfer (HGT). As a result, most existing methods for microbial phylogenomics can only make use of a small subset of the gene families present in the microbial genomes under consideration, potentially biasing their results and affecting their accuracy. To address this challenge, several methods have recently been developed for inferring microbial species trees from genome-scale datasets of gene families affected by evolutionary events such as HGT, gene duplication, and gene loss. In this work, we use extensive simulated and real biological datasets to systematically assess the accuracies of four recently developed methods for microbial phylogenomics, Species-Rax, ASTRAL-Pro 2, PhyloGTP, and AleRax, under a range of different conditions. Our analysis reveals important insights into the relative performance of these methods on datasets with different characteristics, identifies shared weaknesses when analyzing complex biological datasets, and demonstrates the importance of accounting for gene tree inference error/uncertainty for improved species tree reconstruction. Among other results, we find that (i) AleRax, the only method that explicitly accounts for gene tree inference error/uncertainty, shows the best species tree reconstruction accuracy among all tested methods, (ii) PhyloGTP (developed previously by the authors of this paper) shows the best overall accuracy among methods that do not account for gene tree error and uncertainty, (iii) ASTRAL-Pro 2 is less accurate than the other methods across nearly all tested conditions, and (iv) explicitly accounting for gene tree inference error/uncertainty can lead to substantial improvements in species tree reconstruction accuracy. Importantly, we also find that all methods, including AleRax and PhyloGTP, are susceptible to biases present in complex real biological datasets and can sometimes yield misleading phylogenies.

¹School of Computing, University of Connecticut, Storrs, Connecticut, USA, ²Department of Molecular and Cell Biology, University of Connecticut, Storrs, Connecticut, USA, ³The Institute for Systems Genomics, University of Connecticut, Storrs, Connecticut, USA

Peer Community Journal is a member of the
Centre Mersenne for Open Scientific Publishing
<http://www.centre-mersenne.org/>

e-ISSN 2804-3871



Introduction

The accurate inference of phylogenetic relationships between different microbes is an important problem in evolutionary biology. A key difficulty in estimating such phylogenies is the presence of extensive horizontal gene transfer (HGT) in microbial evolutionary histories (Lapierre et al., 2014). This can result in markedly different evolutionary histories for different gene families, obfuscating the underlying species-level or strain-level phylogeny. As a result, the traditional approach for reconstructing microbial phylogenies is to use only “well-behaved” gene families resistant to HGT. This includes the use of small-subunit ribosomal RNA genes (e.g., Olsen et al. 1994; Woese 1987) or of a concatenated alignment of a few core genes from the genomes of interest (e.g., Ciccarelli et al. 2006; Lang et al. 2013; Markowitz et al. 2014). Both these approaches, however, are known to be error-prone. For instance, ribosomal RNA genes are known to engage in horizontal transfer (Gogarten et al., 2002; Yap et al., 1999; Zhaxybayeva et al., 2009) and to yield histories that are inconsistent with those inferred using other core genes (Doolittle et al., 2003; Doolittle, 1999; Hilario and Gogarten, 1993; Hirt et al., 1999). Furthermore, ribosomal RNA genes often cannot be used when studying closely related species due to excessive sequence similarity. Similarly, concatenation based approaches, such as the widely used multilocus sequence analysis (MLSA) technique (Glaeser and Kampfer, 2015), essentially ignore HGT and aggregate the phylogenetic signal from several gene families with potentially distinct evolutionary histories (Gadagkar et al., 2005; McInerney et al., 2008). Indeed, the tree resulting from the concatenation might represent neither the organismal phylogeny nor any of the genes included in the concatenation (Lewis et al., 2016).

To overcome these limitations, several genome-scale methods have also been proposed for microbial phylogeny inference. These include methods such as Phylo SI that are based on gene order information (Sevillya et al., 2019; Shifman et al., 2014), supertree-based methods such as SPR supertrees (Whidden et al., 2014) and MRP (Beiko et al., 2005; Zhaxybayeva et al., 2009) that allow for the use of multiple orthologous gene families, and methods based on average nucleotide identity (ANI) of genomes (Gosselin et al., 2022; Henz et al., 2004; Konstantinidis and Tiedje, 2005). Such genome-scale methods are inherently preferable to methods that base phylogeny reconstruction on only a single gene or a small set of concatenated genes (McInerney et al., 2008). However, while these above methods all represent useful approaches for microbial phylogenomics, they are either targeted at analyzing closely related strains or species (gene order and ANI based methods), or are limited to using single-copy gene families or orthologous groups and do not model key evolutionary events affecting microbial gene family evolution (supertree based methods). Recently, truly genome-scale approaches for microbial phylogenomics, capable of using thousands of complete (multi-copy) gene families, have also been developed. Four of the most promising such methods are ASTRAL-Pro 2 (Zhang and Mirarab, 2022), SpeciesRax (Morel et al., 2022), PhyloGTP (Weiner et al., 2024), and AleRax (Morel et al., 2024). These methods all take as input a collection of unrooted gene trees, where each gene tree may contain zero, one, or multiple genes from a species/strain under consideration. ASTRAL-Pro 2 is based on quartets and seeks a species tree that maximizes a quartet based score (Zhang and Mirarab, 2022). While ASTRAL-Pro 2 does not directly model any specific evolutionary processes, such as HGT or gene duplication, responsible for gene tree discordance, it can handle complete (multi-copy) gene families and previous research suggests that its quartet based approach should be robust to HGT (Davidson et al., 2015). SpeciesRax uses an explicit Duplication-Transfer-Loss

(DTL) model of gene family evolution in microbes and seeks a species tree that maximizes the reconciliation likelihood of observing the input gene trees under that model (Morel et al., 2022). PhyloGTP, a method previously developed by the authors of the current paper, takes a similar overall approach as SpeciesRax but is based on the gene tree parsimony approach and uses a different heuristic search strategy. Specifically, PhyloGTP uses a parsimony-based DTL framework to account for HGT, gene duplication, and gene loss and uses local search heuristics to find a species tree with lowest total reconciliation cost with the input gene trees (Weiner et al., 2024). AleRax performs species tree inference under a probabilistic DTL model (Morel et al., 2024) and is more sophisticated than the other methods in that it can co-estimate gene trees along with the species tree. Unlike ASTRAL-Pro 2, SpeciesRax, and PhyloGTP, which all take as input a single, fixed gene tree per gene family, AleRax takes as input multiple MCMC gene tree samples for each gene family and can thus explicitly account for gene tree reconstruction error and uncertainty.

In this work, we use an extensive simulation study and two real biological datasets to evaluate the species tree reconstruction accuracies of ASTRAL-Pro 2, SpeciesRax, PhyloGTP, and AleRax. Our simulation study focuses on systematically evaluating the impact of number of input gene trees, realistic rates of duplication, HGT, and loss events, and input gene tree error rates on all methods, and on evaluating the impact of using multiple gene tree samples on AleRax. We find that AleRax, the only method that explicitly accounts for gene tree inference error/uncertainty, shows the best species tree reconstruction accuracy among all tested methods, while PhyloGTP shows the best overall accuracy among methods that do not explicitly account for gene tree error and uncertainty. AleRax shows similar accuracy as PhyloGTP when using true (error-free) simulated gene trees, but yields a substantial improvement in reconstruction accuracy compared to PhyloGTP when using estimated (error-prone) simulated gene trees. Between PhyloGTP and SpeciesRax, we find that PhyloGTP can substantially outperform SpeciesRax when the number of input gene trees is small or when DTL rates are high, while SpeciesRax often outperforms PhyloGTP on datasets with low DTL rates. ASTRAL-Pro 2 shows worse accuracy than all other methods across nearly all tested conditions. Overall, the average reconstruction accuracies (defined formally later) of ASTRAL-Pro 2, SpeciesRax, PhyloGTP, and AleRax across our core simulated datasets with estimated gene trees are 81.2%, 86.9%, 88.8% and 91.3%, respectively. We find, however, that the improved accuracies of AleRax and PhyloGTP over the other methods come at the expense of substantially longer running times. We also investigate how AleRax's ability to handle gene tree error/uncertainty contributes to its species tree reconstruction accuracy. We find that when AleRax is provided as input only a single estimated gene tree per gene family (as with the other methods), its accuracy becomes comparable to that of PhyloGTP. This suggests that explicit handling of gene tree error/uncertainty can lead to an approximately 20% reduction in species tree reconstruction error.

We also use the four methods to analyze two real microbial datasets; a more complex 174-taxon Archaeal dataset exhibiting extreme divergence and compositional biases, and a less complex dataset of 44 Frankiales exhibiting low divergence. We find that all four methods perform well on the less complex dataset, recovering identical relationships among the major clades. On the more complex dataset, PhyloGTP, SpeciesRax and ASTRAL-Pro 2 result in some incorrect placements, but appear to perform better than AleRax which produces a tree that is markedly different than any highly supported previously calculated Archaeal tree. Overall, this suggests

that all tested methods, including AleRax and PhyloGTP, are potentially susceptible to biases present in complex datasets.

Overall, our results indicate that AleRax and PhyloGTP may be the two best methods currently available for microbial phylogenomics, though their improved accuracies come at the cost of significantly longer running times. Our results also suggest that phylogenomics methods can benefit substantially from explicitly account for gene tree error and uncertainty. At the same time, our results show that even the best existing methods for microbial phylogenomics may not produce accurate results for certain complex microbial datasets and that their results should be interpreted with caution. A preliminary version of this manuscript, which focused on describing and evaluating PhyloGTP, appeared in the proceedings of the RECOMB Comparative Genomics 2024 conference (Weiner et al., 2024). The current manuscript focuses more heavily on a systematic experimental assessment of the four methods and expands upon the preliminary version by (i) including the recently published method AleRax in the experimental evaluation, (ii) studying the impact of event cost assignments on PhyloGTP's accuracy, (iii) providing descriptions of all four methods evaluated, (iv) evaluating the methods on additional datasets with different ratios of evolutionary events, (v) evaluating memory requirements of all methods, (vi) providing an updated and expanded assessment of the methods on the two real datasets, and (vii) presenting a more extensive discussion of the experimental evaluation and our findings.

Materials and Methods

Description of evaluated methods

We provide brief descriptions of the four methods considered in this work and state their specific objective functions.

Basic definitions and preliminaries. Let T be a leaf-labeled tree with node, edge, and leaf sets denoted by $V(T)$, $E(T)$, and $Le(T)$. If T is rooted, we denote its root by $rt(T)$. For any node $v \in V(T)$, where T is a rooted tree, the (maximal) subtree rooted at v is denoted T_v . Unless otherwise specified, all trees are binary and unrooted.

We use the term *species tree* for the tree depicting evolutionary relationships for the taxa (e.g., species, strains, etc.) under consideration. Given a gene family from the taxa under consideration, a *gene tree* is a tree that depicts the evolutionary relationships of the genes in the gene family. We assume that each edge in a gene tree has an associated branch length (representing substitutions per site), though not all methods make use of branch lengths. Note that a gene tree may have zero, one, or multiple genes from the same taxon.

We assume that the taxon set under consideration is denoted by Ω and that the species tree, denoted S , depicts the evolutionary relationships for taxa in Ω , i.e., $Le(S) = \Omega$. We use \mathcal{G} to denote a collection of gene trees $\{G_1, \dots, G_k\}$, where each G_i , $1 \leq i \leq k$, describes the evolutionary history of a gene family present in the taxon set Ω . We implicitly assume that $Le(S) = \bigcup_{i=1}^k Le(G_i)$. ASTRAL-Pro 2, SpeciesRax, and PhyloGTP assume that each input gene tree (i.e., each gene tree in \mathcal{G}) corresponds to a different gene family, while AleRax requires multiple gene tree samples per gene family, as explained below. All methods considered in this work assume that the input gene trees are unrooted.

The methods SpeciesRax, AleRax, and PhyloGTP utilize DTL reconciliation to assess the fit of input gene trees with candidate species trees. DTL reconciliation provides a framework for reconciling the differences between a gene tree and the corresponding rooted species tree by

invoking gene duplication, HGT, and gene loss events. The method ASTRAL-Pro 2 uses quartet trees to assess the fit between the input gene trees and candidate species trees. A quartet tree is an unrooted tree on four leaves, and the number (or fraction) of quartet trees shared between a gene tree and an unrooted species tree can serve as a measure of similarity between the two trees. Further details on each method appear below.

ASTRAL-Pro 2. This quartet-based method builds upon the widely used ASTRAL method (Mirarab et al., 2014). ASTRAL is designed to work with single-copy gene trees constructed from orthologous sequences and seeks a species tree that maximizes the *quartet score* with the input gene trees. ASTRAL-Pro 2 uses a related but different similarity measure called *per-locus quartet score*, designed to avoid over-counting of quartets in multi-copy gene trees. Thus, given a collection of unrooted gene trees \mathcal{G} as input, ASTRAL-Pro 2 seeks an unrooted species tree S that maximizes the per-locus quartet score with the input gene trees. Since the underlying computational problem is NP-hard (Lafond and Scornavacca, 2019), ASTRAL-Pro 2 implements a heuristic for the problem, using dynamic programming to efficiently find an optimal species tree within a restricted search space. We note that ASTRAL-Pro 2 does not make use of gene tree branch lengths. Further details on the method appear in Zhang and Mirarab 2022.

SpeciesRax. This method uses the “undatedDTL” probabilistic DTL reconciliation framework of GeneRax (Morel et al., 2020) to estimate the species tree and model parameters (rates of duplication, HGT, and loss) given a collection of input gene trees. Specifically, SpeciesRax takes a collection of unrooted gene trees \mathcal{G} as input and seeks a rooted species tree S and model parameters Θ that maximize the reconciliation likelihood $L(S, \Theta | \mathcal{G})$. SpeciesRax uses a distance-based method, *MiniNJ*, to estimate a starting species tree and then executes a local search heuristic to further optimize this starting tree. We note that SpeciesRax utilizes gene tree branch lengths and infers a *rooted* species tree. Further technical details on SpeciesRax appear in Morel et al. 2022.

AleRax. AleRax is similar to SpeciesRax in that it also uses a probabilistic model of DTL reconciliation and uses the same search strategy as SpeciesRax (miniNJ followed by local search) for finding a maximum likelihood species tree and model parameters given the input gene trees. However, unlike SpeciesRax, AleRax accounts for gene tree inference error and uncertainty by taking as input multiple gene tree samples for each gene family. Specifically, AleRax uses the ALE algorithm (Szollosi et al., 2013) to integrate over gene tree uncertainty by approximating the probability of observing a gene family sequence alignment given a rooted species tree. Thus, AleRax seeks to find a rooted species tree S and model parameters that maximize the likelihood $L(S | A)$, where A denotes the collection of sequence alignments for all gene families represented in \mathcal{G} . We note that AleRax does not directly take sequence alignments as input and instead uses the multiple (typically 1000) input gene tree samples for each gene family sequence alignment to estimate the fit of an alignment with a species tree. Like SpeciesRax, AleRax utilizes gene tree branch lengths and infers a *rooted* species tree. Further details on this method appear in Morel et al. 2024.

PhyloGTP. Unlike SpeciesRax and AleRax, PhyloGTP uses the parsimony-based DTL reconciliation model of Bansal et al. (2012), David and Alm (2011), and Tofigh et al. (2011). Under this model, each event type has an associated (user-defined) cost and the objective is to find a reconciliation of minimum total cost. This model allows for an unrooted gene tree to be optimally

reconciled with a rooted species tree within $O(mn)$ time, where m and n denote the number of leaves in the gene tree and species tree, respectively (Bansal et al., 2012).

In the following, we denote the event costs for gene duplications, HGTs, and gene losses by P_d , P_t , and P_l , respectively. Given a gene tree $G \in \mathcal{G}$, species tree S , and event costs P_d , P_t , and P_l , we denote by $\mathcal{R}_{P_d, P_t, P_l}(G, S)$ the reconciliation cost of an optimal DTL reconciliation of G and S under the event costs P_d , P_t , and P_l . Given a species tree S , a collection of gene trees $\mathcal{G} = \{G_1, \dots, G_k\}$, and event costs P_d , P_t , and P_l , we define the *total DTL reconciliation cost* of \mathcal{G} with S to be the sum of the DTL reconciliation costs of each $G \in \mathcal{G}$ with S , i.e., $\sum_{i=1}^k \mathcal{R}_{P_d, P_t, P_l}(G_i, S)$.

Given as input a collection of gene trees, PhyloGTP seeks a species tree that minimizes the total DTL reconciliation cost against the collection of input gene trees. More formally, we can define the *Most Parsimonious Species Tree (MPST)* problem as follows: Given a collection of gene trees \mathcal{G} and event costs P_d , P_t , and P_l , find a species tree S that minimizes the total DTL reconciliation cost with \mathcal{G} .

We note that under this formulation only the topology of the gene tree is used and branch lengths are ignored. The MPST problem can be shown to be NP-hard, W[2]-hard, and inapproximable to within log factor through a reduction from the NP-hard gene duplication problem (Bansal and Shamir, 2011; Ma et al., 2000). The gene duplication problem is a special case of MPST problem defined in this manuscript and seeks a species tree minimizing just the total number of gene duplications. Details of the reduction are straightforward and therefore omitted. PhyloGTP uses a local search heuristic to solve the MPST problem. For completeness, further methodological details on PhyloGTP appear in the supplement (Weiner et al., 2025).

Experimental setup

We use both simulated and real biological datasets to carefully assess the reconstruction accuracy of the four methods. ASTRAL-Pro 2, SpeciesRax, and PhyloGTP take as input a single unrooted maximum-likelihood gene tree per gene family, while the recommended input for AleRax is 1000 unrooted gene trees, sampled from the posterior using Mr.Bayes (Ronquist et al., 2012), per gene family. All methods were run using their default/recommended parameter settings. For AleRax, we present results both for the recommended number (1000) of gene trees per gene family, allowing it to account for gene tree inference error/uncertainty, as well as when only a single gene tree is used per gene family, effectively disabling its ability to account for gene tree error/uncertainty.

Evaluating reconstruction accuracy. To evaluate the species tree reconstruction accuracies of the different methods, we compare the species tree estimated by each method with the corresponding ground truth species tree. To perform this comparison we utilize the widely used (unrooted) normalized Robinson-Foulds distance (NRFD) (Robinson and Foulds, 1981) between the reconstructed and ground truth species trees. For any reconstructed species tree, the NRFD reports the fraction of non-trivial splits in that species tree that do not appear in the corresponding ground truth species tree. For ease of interpretation, we report results in terms of *percentage accuracy*, defined to be the percentage of non-trivial splits in the reconstructed species tree that also appear in the ground truth species tree. Thus, percent accuracy is simply $(1 - \text{NRFD}) \times 100$. Thus, for example, a percentage accuracy of 87% is equivalent to an NRFD of 0.13.

Description of simulated datasets

We used simulated datasets with known ground truth species trees to assess the impact of three key parameters on reconstruction accuracy: Number of input gene trees, rates of gene duplication, HGT, and gene loss (or DTL rates for short), and estimation error in the input gene trees.

Our core simulated datasets were created using a four-step pipeline: (1) simulation of a ground-truth species tree and corresponding true gene trees (one per gene family) with varying DTL rates, (2) simulation of sequence alignments of different lengths for each true gene tree, (3) inference of estimated maximum likelihood gene trees from the sequence alignments, and (4) inference of 1000 estimated gene trees, sampled from the posterior using Mr.Bayes (Ronquist et al., 2012), for each true gene tree (i.e., per gene family) from the sequence alignments. In the first step, we used SaGePhy (Kundu and Bansal, 2019) to first simulate ground-truth species trees, each with exactly 50 leaves (taxa) and a height (root to tip distance) of 1, under a probabilistic birth-death framework. We then used these species trees to simulate multiple gene trees under the probabilistic duplication-transfer-loss model implemented in SaGePhy. This resulted in 9 different datasets of simulated true gene trees, each corresponding to a different number of true input gene trees (10, 100, or 1000), and a different DTL rate (low, medium, or high; see Table 1). Each dataset comprised of 10 replicates. The chosen DTL rates are based on the relative rates and frequencies of gene duplication and HGT events in real microbial gene families from species sampled broadly across the tree of life (Bansal et al., 2015; David and Alm, 2011). In each case, the gene loss rate is assigned to be equal to the gene duplication rate plus the additive HGT rate, so as to balance the number of gene gains with the number of gene losses (Table 1). Basic statistics on these simulated true gene trees, including average sizes and numbers of gene duplication and HGT events, are provided in Table 2. We note that numbers of inferred gene duplications and HGTs are larger for estimated gene trees, where reconstruction error manifests itself as closely matching rates inferred for real microbial gene families based on estimated gene trees (Bansal et al., 2015; David and Alm, 2011).

In the second step, we used AliSim (Ly-Trong et al., 2022) to simulate DNA sequence alignments along each true gene tree under the General Time-Reversible (GTR) model (using default AliSim GTR model settings) with three different sequence lengths: 400, 100, and 50 bp. In the third step, maximum-likelihood gene trees were inferred using IQ-TREE 2 (Minh et al., 2020) from the simulated sequence alignments under the Jukes-Cantor (JC) model. We use the simpler JC model when estimating gene trees, instead of the GTR model used to generate the sequences, since this better captures the biases and limitations of applying standard substitution models to real biological sequences when inferring biological gene trees. Thus, from each dataset of true gene trees, we derive 3 additional datasets of estimated gene trees corresponding to the three sequence lengths. The purpose of the second and third steps above is to generate error-prone gene trees that reflect the reconstruction/estimation error present in real gene trees. We found that the estimated gene trees had average normalized Robinson-Foulds (RF) distances (Robinson and Foulds, 1981) of 0.08, 0.22, and 0.35 for sequence lengths 400, 100, and 50 bp, respectively, to the corresponding true gene trees. Since some of the methods also make use of gene tree branch lengths, we additionally measured branch length inference error in the estimated medium DTL gene trees. In particular, since the topologies of the true and estimated gene trees can be different, we compared the leaf branch lengths and the path lengths between each pair of leaves using

two metrics. First, we used the mean absolute percentage error (MAPE) to measure the error itself, and second, we used the Pearson correlation coefficient (PCC) to measure the correlation between the estimated and inferred lengths. These branch length error statistics are summarized in Supplementary Table S3. For medium DTL datasets and sequence lengths of 400, 100, and 50 bp, path length MAPEs were 6.9%, 15.8%, and 24.1%, respectively. Overall, estimated branch lengths show strong correlation with true lengths. For example, even with 50 bp sequences, we observe PCC over leaf branch lengths and path lengths of 0.824 and 0.833, respectively.

In the fourth and final step, we used Mr.Bayes 3.2.7 (Ronquist et al., 2012) to generate the 1000 posterior gene tree samples per gene family required by AleRax. Consistent with the previous step, these trees were inferred by applying Mr.Bayes to the simulated sequence alignments and using the JC model. Following Morel et al. (2024), we ran each MCMC chain for 100,000 generations and sampled every 100 generations, using the maximum-likelihood gene trees generated in the previous (third) step above as starting trees. (Note that burn-in is not needed since we start the MCMC chain from the maximum-likelihood tree; this is consistent with how Mr.Bayes is used by Morel et al. (2024).) As before, from each dataset of true gene trees, this creates 3 additional datasets of estimated gene trees, with 1000 estimated gene trees per gene family, corresponding to the three sequence lengths.

Table 1 – Key parameters used to generate the core simulated datasets. The table lists the main parameters and their values explored in the simulation study for the core datasets. All 36 (= 3 × 3 × 4) combinations of these three parameters were evaluated at 10 replicates each. DTL rates are specified in the form (*d*, *t*, *l*), where *d* is the gene duplication rate, *t* is the HGT rate (split evenly between additive and replacing HGTs), and *l* is the gene loss rate. The number of species was fixed at 50 for these datasets.

Parameter	Values
Number of gene families	10, 100, 1000
DTL rates	low = (0.3, 0.6, 0.6) med = (0.6, 0.12, 0.12) high = (0.12, 0.24, 0.24)
Sequence length (nucleotides)	400, 100, 50, and true gene trees

Table 1 summarizes the specific ranges of parameter values we explored for the number of gene families, DTL rates, and sequence lengths in the core simulated datasets described above. We evaluated all combinations of these parameter values, resulting in a total of 36 core simulated datasets, with each dataset comprising of 10 replicates created using that specific assignment of parameter values. In addition to these core simulated datasets, we also created corresponding simulated datasets with different relative rates of HGT and gene duplication (as described later in the Results section; see Supplementary Table S1 for specific DTL parameters used), and created datasets with 10 and 100 taxa for the runtime and memory usage analysis. The specific commands used to generate the simulated datasets are available in the supplement (Weiner et al., 2025).

Table 2 – Basic statistics for simulated true gene trees in the core dataset. Average number of leaves, duplications, and HGTs, and losses in the simulated low, medium, and high DTL true gene trees in the core simulated datasets. For each DTL rate, the number of losses is roughly equal to the number of duplications plus half the number of HGTs. Results were averaged over all 10 replicates of the 100 gene tree datasets.

DTL rate	Leaves	Duplications	HGTs
Low	53.618	3.408	6.586
Med	55.121	6.15	11.125
High	59.718	10.077	18.37

Description of biological datasets

We assembled two previously used biological datasets of different size, composition, and complexity to assess the accuracy and consistency of species trees inferred by AleRax, PhyloGTP, SpeciesRax and ASTRAL-Pro 2, and traditional non-DTL cognizant methods such as MLSA and tANI (Gosselin et al., 2022) (Table 3). To examine the effect of extreme divergence and genome complexity variation on species tree inference, we used a dataset composed of 176 Archaea, which was drawn from Feng et al. 2021. The Archaea included in the dataset span 2-3 kingdoms (or superphylums), and radically different lifestyles (from extremophiles inhabiting Antarctic lakes to mammal gut constituents). Because the pan-genome of an entire domain would be immeasurably large and computationally infeasible to accurately infer, we have reduced the number of gene families in this dataset to 282 core genes, which are shared by all members. This also allows direct comparison of the species trees inferred by the four methods to previously calculated phylogenies by Feng et al. (2021) which used the same loci. It should be noted that the 282 gene families used in this analysis have been expanded to include all homologs (paralogs, xenologs, etc.) found in each genome, while only orthologs were used by Feng et al. (2021).

To examine the impact of low sequence divergence on species tree inference, we used a dataset of 44 Frankiales genomes, drawn from Gosselin et al. 2022. These included taxa are all closely related members of the order Frankiales, and as such the entire pan-genome (8,862 gene families with at least 4 sequences) was used for inference. The order Frankiales are composed of nitrogen-fixing symbionts of pioneer flora, and although they demonstrate variation in GC content and genome size these factors were previously shown to not bias phylogenetic inference (Gosselin et al., 2022).

Archaea dataset assembly. Annotated genomes of 176 Archaea used in Feng et al. 2021 were collected. The 282 core gene loci described in Feng et al. 2021 were used as amino acid query sequences to search every collected genome, using blastp (Camacho et al., 2009) with default parameters (-evalue was changed to 1e-10). All significant sequence for every loci across all genomes were collected (provided they met a length threshold of 50% in reference to the average gene family sequence size to filter partial sequences). Each gene family was then aligned using mafft-linsi (Katoh and Standley, 2013) with default parameters. These alignments were used for inferring maximum-likelihood gene trees in IQ-Tree 2 (Minh et al., 2020), where the best substitution model for each gene family was determined using Bayesian Inference Criterion (Kalyaanamoorthy et al., 2017). The resulting maximum-likelihood gene trees were used as input for ASTRAL-Pro 2, SpeciesRax, and PhyloGTP. To generate the input gene trees for AleRax, we used Mr.Bayes to sample 1000 posterior gene trees under the LG amino-acid model (Le and

Gascuel, 2008) running MCMC chain for 100,000 generations, sampling every 100 generations, and starting each chain with the corresponding maximum-likelihood gene tree for each of the 282 gene families.

Frankiales dataset assembly. Annotated proteomes of the 44 Frankiales used in Gosselin et al. 2022 were collected. Protein sequences were clustered into gene families and using the OrthoFinder2.4 pipeline (Emms and Kelly, 2019) with default parameters (the search algorithm was changed to blast). Briefly, all-vs-all blastp (evalue of 1e-3) was used to find the best hits between input species. The set of query-matches were then clustered into gene families using the MCL algorithm, and the subsequent gene families were aligned using mafft-linsi with default parameters. Resulting alignments were used to create maximum-likelihood gene trees using FastTree (Price et al., 2010) using the JTT model and default parameters. To create the input gene trees for AleRax, we used Mr.Bayes to sample 1000 posterior gene trees for each of the 8,862 gene families using the same parameter settings as before.

Table 3 – Summary of the two biological datasets.

Dataset	Number of gene families	Potential biases	Previous methods used to infer species tree
176 Archaea (domain)	282	Extreme divergence, long branch attraction, compositional bias	tANI, MLSA, single gene
44 Frankiales (order)	8,862	Low divergence, contamination, genome size difference	tANI, MLSA

Results

Results on simulated data

Accuracy on true (error-free) gene trees. We first evaluate the accuracy of the species tree reconstruction methods when given true (error-free) gene trees as input (effectively skipping steps 2, 3 and 4 of the simulation pipeline). While error-free gene trees do not capture the complexities of real data, this analysis helps us understand how the different methods perform in a controlled, ideal setting. Figure 1 shows the results for low, medium, and high DTL rates with varying numbers of gene families for 50-taxon datasets. Unsurprisingly, we find that both DTL rates and number of input gene families are highly impactful parameters. The performance of all four methods worsens as DTL rates increase, and improves as the numbers of input gene families increase. As the figure shows, AleRax shows the highest overall accuracy on these datasets, with PhyloGTP showing comparable but slightly worse accuracy than AleRax. Between PhyloGTP and SpeciesRax, we find that PhyloGTP shows higher accuracy when the number of gene families is small (100 or fewer), particularly when DTL rates are medium or high. For the remaining datasets, both PhyloGTP and SpeciesRax show nearly identical accuracies. Notably, AleRax, PhyloGTP, and SpeciesRax substantially outperform ASTRAL-Pro 2, especially on the medium and high DTL datasets. In particular, we find that ASTRAL-Pro 2 is highly susceptible to high DTL rates, and that it also shows poor performance when the number of input gene families is small.

Interestingly, the accuracy of Astral-pro 2 improves rapidly as the number of gene families increases, with the method performing equivalently to the other methods on the low and medium DTL datasets when the input consists of 1000 gene families.

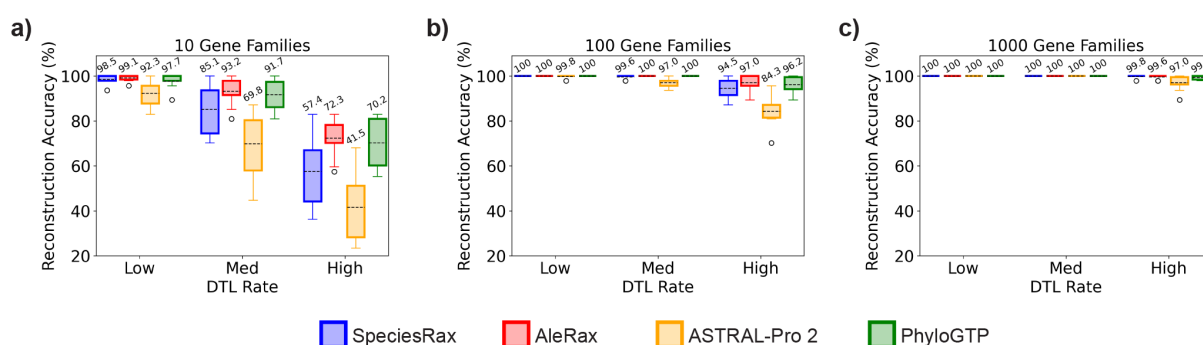


Figure 1 – Accuracy on true gene trees. Tree reconstruction accuracies are shown for SpeciesRax, AleRax, ASTRAL-Pro 2, and PhyloGTP when applied to error-free or ‘true’ gene trees. Results are shown for increasing numbers of input gene families (10, 100, and 1000) and for low, medium, and high DTL rates. The number of taxa (i.e., number of leaves in the species tree) is fixed at 50. Higher percentages (y-axis) imply greater accuracy. The number above each box is the mean value across 10 replicate runs, and the dotted line within each box represents the median value.

Accuracy on estimated (error-prone) gene trees. We next assess the accuracy of reconstructed species trees when the input consists of estimated (error-prone) gene trees. Figure 2 shows the results of this analysis for all 27 combinations of number of input gene families, DTL rates, and sequence lengths (or gene tree estimation error rates). As expected, the accuracy of all three methods is substantially affected by the quality of the estimated gene trees, with higher accuracies achieved using gene trees estimated from longer sequences. We also find that an increased number of input gene trees can partly make up for error in the input gene trees. AleRax, the only method that explicitly accounts for gene tree inference error and uncertainty, shows the best overall performance, achieving an average reconstruction accuracy of 91.3% when averaged across all 27 datasets with estimated gene trees. PhyloGTP shows the next best accuracy, with an average reconstruction accuracy of 88.8%, and SpeciesRax and ASTRAL-Pro 2 show average accuracies of 86.9% and 81.2%, respectively. As the figure shows, the magnitude of improvement offered by AleRax over the other methods increases as the quality of the input gene trees decreases (i.e., with decreasing sequence length). This is not surprising and points to the significant impact of AleRax’s ability to handle gene tree error and uncertainty. Interestingly, PhyloGTP outperforms AleRax on 6 of the 9 datasets that use the highest quality estimated gene trees (400 base pair sequences; plots in first column of Figure 2). We also find that all methods still outperform ASTRAL-Pro 2 across most datasets and that ASTRAL-Pro 2 continues to be more susceptible to high DTL rates than the other methods. As before, the performance of ASTRAL-Pro 2 improves rapidly with increasing number of input gene trees, even sometimes outperforming all other methods when DTL rates are low or medium. This suggests that ASTRAL-Pro 2 may be well-suited for microbial phylogenomics on datasets with lots of gene trees and relatively low prevalence of HGT. Comparing PhyloGTP with SpeciesRax, we find that both methods have similar performance overall, with PhyloGTP and SpeciesRax showing average percent accuracies of 88.8% and 86.9%, respectively, when averaged across all 27 datasets. However, PhyloGTP shows substantially higher accuracy than SpeciesRax on datasets with high DTL rates,

as well as on datasets with 10 input gene trees. This suggests that PhyloGTP may be especially useful for analyzing datasets with high levels of HGT or with a small number of gene trees.

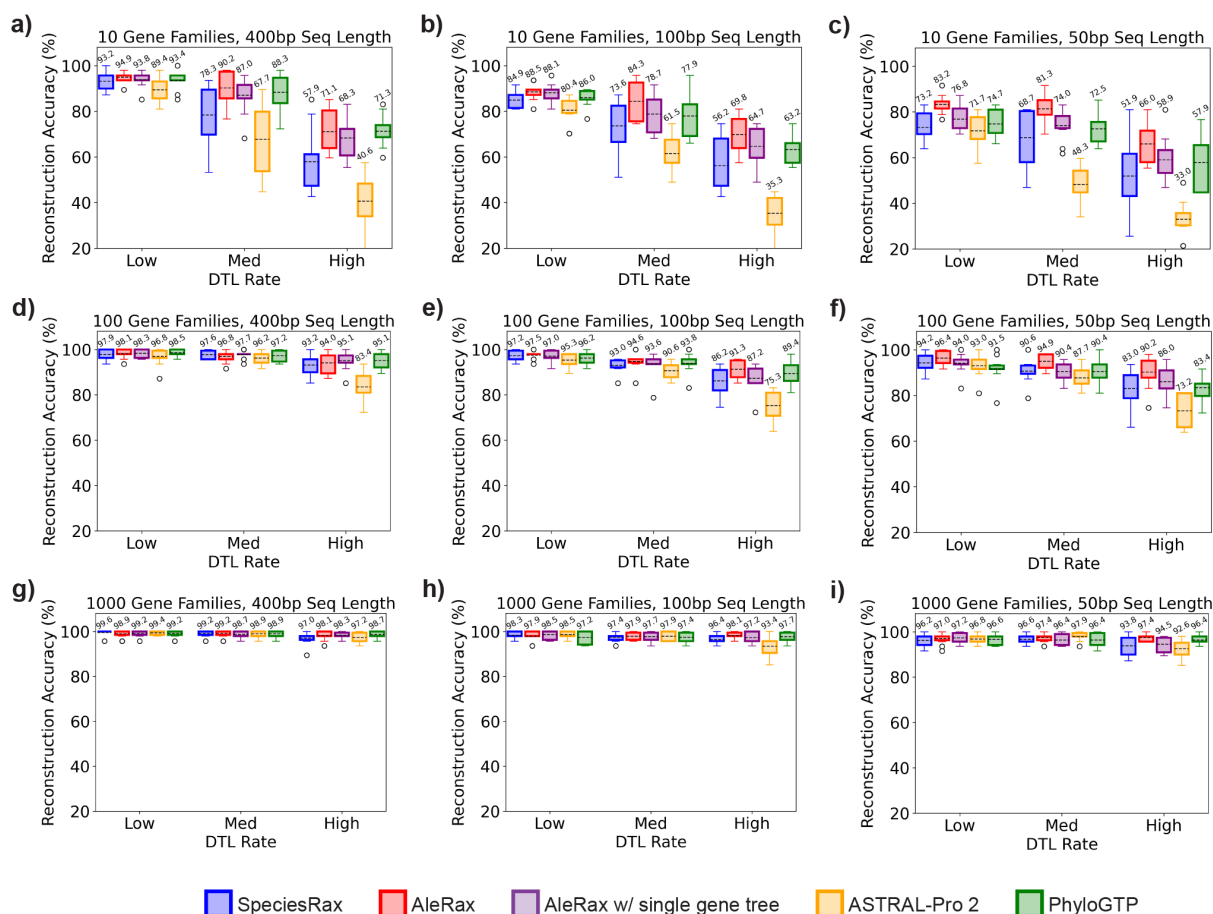


Figure 2 – Accuracy on estimated gene trees. Tree reconstruction accuracies are shown for SpeciesRax, AleRax, ASTRAL-Pro 2, and PhyloGTP when applied to estimated gene trees. Results are shown for all 27 combinations of number of input gene families, sequence lengths (shorter sequence lengths imply greater gene tree estimation error), and DTL rates. The first, second, and third rows correspond to datasets with 10, 100, and 1000 gene families, respectively, and the first, second, and third columns correspond to 400, 100, and 50 base pair sequence lengths, respectively. The number of taxa (i.e., number of leaves in the species tree) is fixed at 50. “AleRax” (red) refers to the default execution of AleRax with 1000 gene tree samples per gene family (i.e., with gene tree error-correction), and “AleRax w/single gene tree” (purple) refers to the modified execution where only a single gene tree per gene family is provided as input (i.e., no gene tree error-correction). Higher percentages imply greater accuracy. The number above each box is the mean value across 10 replicate runs, and the dotted line within each box represents the median value.

Impact of gene tree error-correction on AleRax’s accuracy. Our results show that AleRax can significantly outperform the other methods on datasets with error-prone gene trees. To better understand how AleRax’s ability to handle gene tree error/uncertainty contributes to its species tree reconstruction accuracy, we also apply AleRax to the estimated (error-prone) gene tree datasets with only a single gene tree per gene family provided as input. Providing only a single gene tree per gene family, instead of the default of 1000, effectively prevents AleRax from being able to account for gene tree error/uncertainty. Figure 2 shows the results of this analysis and reveals several interesting insights (see results for “AleRax w/single gene tree” in that figure). First,

we find that the accuracy of this restricted version of AleRax remains quite high and comparable to that of PhyloGTP overall. Second, gene tree error-correction becomes much more impactful on datasets with gene trees of lower quality (100 and 50 base pair sequence lengths). In fact, on the datasets with 400 base pair gene trees, the restricted version of AleRax often slightly outperforms regular AleRax. And third, AleRax's ability to handle gene tree error/uncertainty is responsible for about a 20% average reduction in reconstruction error on these datasets. These results show that explicit handling of gene tree error and uncertainty can lead to substantially improved species tree reconstruction accuracy, especially on datasets with low-quality gene trees.

Robustness of PhyloGTP to event costs. PhyloGTP relies on a parsimony-based DTL reconciliation algorithm (Bansal et al., 2012) to assess the "fit" of input gene trees with candidate species trees. This reconciliation framework relies on user-specified costs for duplication (D), HGT (T), and loss (L) events to compute optimal reconciliations. By default, PhyloGTP uses D-T-L costs of 2-4-1. To assess the robustness of PhyloGTP to different event costs, we apply PhyloGTP variants using five different D-T-L event costs, 2-2-1, 2-3-1, 2-4-1 (default), 2-5-1, and 1-3-1, to the simulated datasets. Supplementary Figures S1 and S2 shows the results of this analysis for true and estimated gene trees, respectively. As the figures show, the performance of PhyloGTP remains robust to the specific event costs used. We also find that no single variant outperforms the others across all, or even most, datasets, and that each of the five variants emerges as the top performer across at least one of the simulated datasets.

Robustness of results to relative rates of HGT and gene duplication events. Since biological microbial datasets can have different relative rates of HGT and gene duplication events, we additionally evaluated the methods on simulated datasets with a different ratio of DTL events. In particular, we used higher rates of HGT than in the core datasets ($1.5\times$) and near-zero rates of gene duplication. These event rates are based on an analysis of over 7,500 gene families from 103 *Aeromonas* strains representing 28 different species (Rangel et al., 2019). As before the loss rate was set to be equal to the gene duplication rate plus the additive HGT rate. As with the core simulated dataset, we used low, medium, and high rates of DTL, different numbers of gene families, and different sequence lengths to obtain 36 alternative simulated datasets, each with 10 replicates. Supplementary Table S1 shows the specific DTL rates and other parameter settings used to generate these alternative datasets.

Results on these alternative datasets are consistent with those reported above for the core simulated datasets. For example, on the alternative datasets with true (error-free) gene trees, we again find that AleRax and PhyloGTP show greatest accuracy and that all methods substantially outperform ASTRAL-Pro 2 on most datasets (Supplementary Figure S3). Likewise, on the alternative datasets with estimated gene trees, AleRax and PhyloGTP continue to be the two best methods, with AleRax, PhyloGTP, SpeciesRax, and ASTRAL-Pro 2 showing average accuracies of 91.5%, 89.0%, 86.2%, and 81.7% across the 27 datasets, respectively (Supplementary Figure S4).

Runtimes and memory usage. We compare the runtimes of the four methods when varying the number of taxa (10, 50, and 100) over low, medium, and high DTL rates. In addition, we also evaluate the impact of the number of input gene trees (100 and 1000) using the 50-taxon dataset. These runtimes are shown in Table 4. All methods have parallel implementations and were allocated 12 cores on a 2.1 GHz Intel Xeon processor with 64 GB of RAM. We find that ASTRAL-Pro 2 is, by far, the fastest method, requiring only about 5 seconds on the high-DTL 50-taxon 1000

Table 4 – Impact of number of taxa and gene trees on running time. Runtimes in seconds are shown for the three methods for datasets with 10, 50, and 100 taxa and low medium, and high rates of DTL. For the 10- and 100-taxon datasets, the number of input gene trees is 100. For 50-taxon datasets, results are shown for both 100 and 1000 gene trees. The runtimes are based on simulated true input gene trees and are averaged over 10 replicate runs. Each method was allocated 12 cores on a 2.1 GHz Intel Xeon processor with 64 GB of RAM.

Dataset size	DTL rate	SpeciesRax	AleRax	ASTRAL-Pro 2	PhyloGTP
10 taxa, 100 gene trees	low	1.45	8.02	0.08	4.02
	med	1.36	8.89	0.08	4.45
	high	1.35	10.6	0.09	4.82
50 taxa, 100 gene trees	low	5.69	503.8	1.11	1,299.56
	med	6.25	947.45	1.14	1,374.92
	high	8.9	1,276.77	1.43	2,015.03
50 taxa, 1000 gene trees	low	50.34	7,903.59	5.38	10,011.79
	med	52.33	9,443.45	5.29	11,433.19
	high	59.95	11,376.92	5.55	13,137.93
100 taxa, 100 gene trees	low	22.05	10,659.92	3.61	19,871.48
	med	28.90	14,110.78	4.84	32,606.87
	high	47.19	22,091.67	7.15	38,259.04

gene tree datasets and less than 10 seconds on the high-DTL 100-taxon 100 gene tree datasets. SpeciesRax is also extremely fast, requiring only about 60 seconds and 50 seconds, respectively, on those datasets. AleRax and PhyloGTP are much slower than the other two methods, with AleRax requiring over 3 hours and 6 hours, and PhyloGTP requiring about 3.5 hours and 10.5 hours, respectively, on those same datasets. Thus, the improved accuracy provided by AleRax and PhyloGTP comes at the expense of significantly longer running times. It is surprising that PhyloGTP, despite being parsimony based, has the longest runtimes. This is partly due to PhyloGTP’s use of a more extensive SPR-based local search heuristic, while AleRax and SpeciesRax both use a simpler NNI-based local search heuristic. We note that AleRax also requires additional Bayesian analysis runs to generate its input posterior gene tree samples, which add an additional computational burden not accounted for in the reported runtimes. Further research on AleRax and PhyloGTP may lead to faster running times without negatively impacting their accuracies. For example, using fewer posterior gene tree samples (say 100 instead of 1000) per gene family may be sufficient for most analyses. Likewise, it may be possible to speed up the algorithms and local search heuristics implemented in the current prototype version of PhyloGTP.

We also compare the computational memory requirements of the four methods by measuring the peak memory usage during execution. To evaluate the impact of dataset size, memory was profiled on the high DTL rate datasets with 1) 50 taxa and 1000 gene trees and 2) 100 taxa and 100 gene trees. Peak memory usage statistics appear in Supplementary Table S2. We find that PhyloGTP has the lowest memory footprint for the 1000 gene tree dataset while ASTRAL-Pro 2 has the lowest memory footprint for the 100 taxa dataset. Overall, PhyloGTP, ASTRAL-Pro 2, and

AleRax have very modest memory requirements in the low hundreds of MBs. SpeciesRax uses substantially more memory than the other methods, but still less than 1 GB on both datasets.

Results on biological data

Archaeal dataset. A myriad of controversies surround the phylogeny of Archaea. These controversies include the monophyly of the DPANN superphylum (Aouad et al., 2018; Brochier-Armanet et al., 2011; Feng et al., 2021; Narasingarao et al., 2012; Raymann et al., 2014), the placement of extreme halophiles (Aouad et al., 2019; Feng et al., 2021; Narasingarao et al., 2012; Sorokin et al., 2019), and the root of the Archaea (Raymann et al., 2015). These differences in phylogenetic inference are driven by many factors including, but not limited to compositional bias, long branch attraction, extremely small genomes, numerous HGT events, and biased sampling of metagenome-assembled genomes. Thus, it is interesting to evaluate the performance of the four studied methods in the face of these factors. Using 282 unrooted input gene trees, all four methods inferred Archaeal species trees with inaccuracies with respect to commonly accepted placements of groups in previous analyses. These inaccuracies should be interpreted in the context that for several Archaeal clades (mostly halophiles) there is no consistent, consensus position that is universally accepted amongst Archaeologists. For example, the monophyly of the DPANN superphylum is considered by some to be an artifact (driven by long branch attraction or biased genome sampling) (Aouad et al., 2018; Feng et al., 2021; Zhaxybayeva et al., 2013).

There are no extreme topological differences between the SpeciesRax, ASTRAL-Pro 2 and PhyloGTP Archaea tree reconstructions. Those three methods fail to recover a monophyletic Euryarchaea kingdom (Figures 3A and 4A,B) although these resulting topologies (with the Methanomada and Thermococcales on the branches leading to the TACK group) are consistent with trees in previous attempts to find an alternative root of the Archaea (Raymann et al., 2015). One notable difference is that compositional attraction may have played a larger role in PhyloGTP and ASTRAL-Pro 2, particularly with halophiles. The Haloarchaea were attracted to the Methanonatronarchaea and were left out of their accepted position within the Methanotecta (Figure 3A) in the PhyloGTP tree. The Methanonatronarchaea are typically seen as basal to the Methanotecta but have been attracted closer to the other methanogens in the ASTRAL-Pro 2 tree (Figure 4A). Although the Haloarchaea were correctly placed in the SpeciesRax tree (Figure 4B), they are on an extremely long branch. Incorrect placements of the halophiles Nanohaloarchaea, Haloarchaea and Methanonatronarchaea are often attributed to compositional bias (Feng et al., 2021). These halophiles prefer acidic amino acid residues (such as aspartate and glutamate), on account of their survival strategies in hypersaline environments, and these acidified proteomes attract the placement of these groups together in phylogenetic reconstructions.

In contrast, the AleRax species tree exhibits a markedly different topology (Figure 3B) compared to the above methods. This was the only method that does not recover a monophyletic DPANN group: the Nanohaloarchaea, Parvarchaeota, and Aenigmarchaeota form a clade basal to the Euryarchaea. AleRax fails to recover the correct position of the Haloarchaea, and instead places the group basal to the Methanotecta super-class. Additionally, several Methanomada leaves (*Methanopyrus* sp. Kol6 and *Methanobacteriota archaeaon*) failed to associate with the larger Methanomada clade and were placed in very different places along the tree.

These four different tree topologies for the same input data reveal interesting contrasts in the reconstruction methods. PhyloGTP, SpeciesRax, and ASTRAL-Pro 2 are susceptible to the

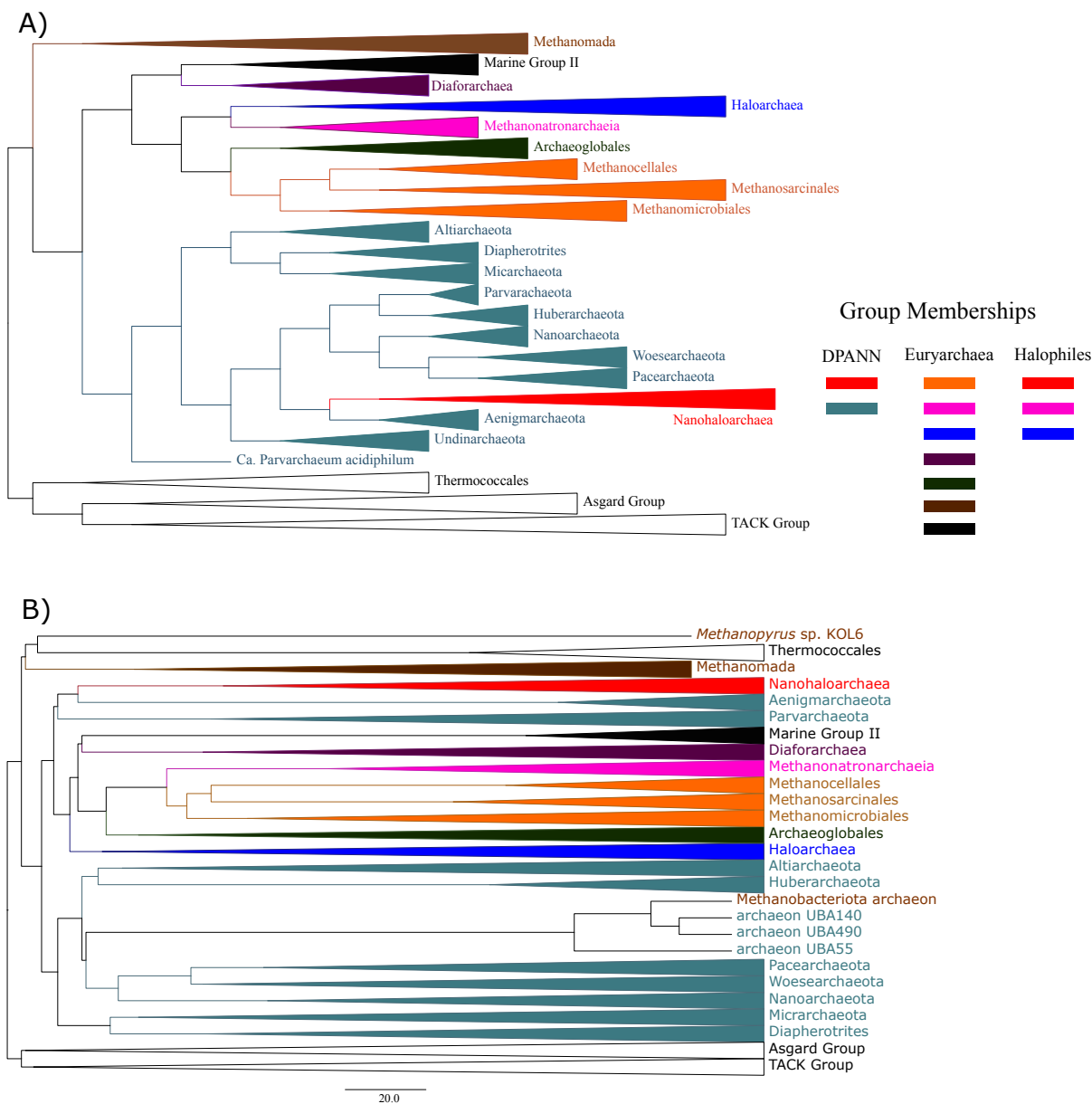


Figure 3 – Archaeal species tree reconstructions by PhyloGTP and AleRaX. Individual taxa on both trees have been collapsed into clades and are colored corresponding to higher level classifications (clades with the same color are part of the same class or phylum). The legend shows previously reported Kingdom memberships of these collapsed clades, and also the halophiles which may group together as a result of compositional bias. Part A) Unrooted Archaeal tree inferred by PhyloGTP, shown as a cladogram since PhyloGTP does not infer branch lengths. Part B) Unrooted Archaeal tree inferred by AleRaX.

presence of problematic groups (such as the extreme halophiles) and other biases in complex datasets, potentially limiting their accuracy in some cases. Still, the trees inferred by these three methods are more consistent with previous estimates of the Archaeal tree, demonstrating that these methods can produce a mostly accurate Archaeal tree even in the face of the many biases present in the dataset. In contrast, the tree calculated by AleRax does not resemble any highly supported previously calculated Archaeal tree. This can be explained by two possible causes: 1) AleRax is not suitable for domain level comparisons, where divergence and numerous DTL events have saturated over the extreme time scale (at least 3Ga (Martinez-Gutierrez et al., 2023)). 2) The

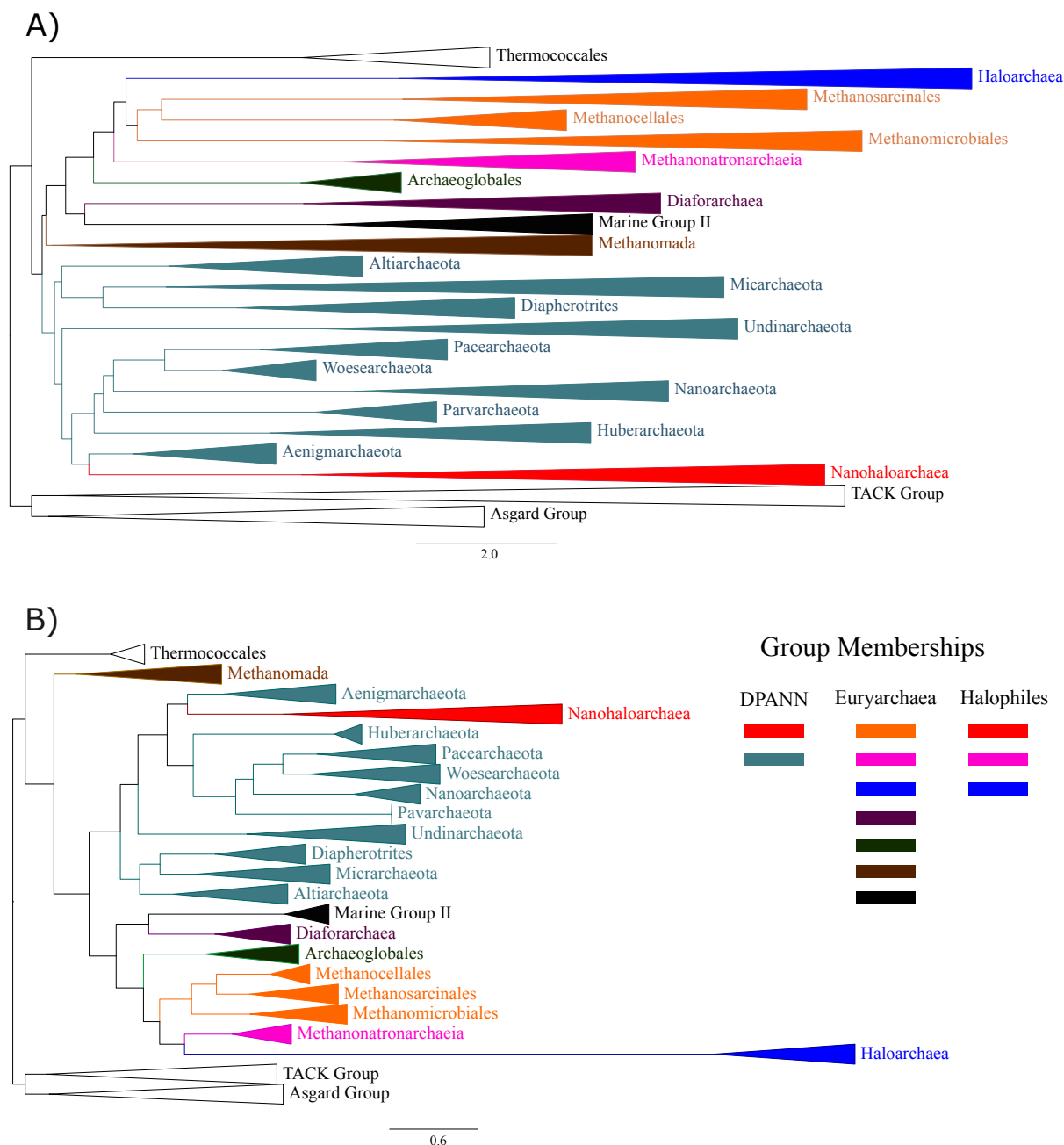


Figure 4 – Archaeal species tree reconstructions by ASTRAL-Pro 2 and SpeciesRax. Individual taxa on both trees have been collapsed into clades and are colored corresponding to higher level classifications (clades with the same color are part of the same class or phylum). The legend shows previously reported Kingdom memberships of these collapsed clades, and also the halophiles which may group together as a result of compositional bias. Part A) Unrooted Archaeal tree inferred by ASTRAL-Pro 2. Part B) Unrooted Archaeal tree inferred by SpeciesRax.

tree produced by AleRax (Figure 3B) reflects the true evolution of the group. Very few previous attempts to consider DTL events in Archaeal phylogenomics exist (Davín et al., 2018; Williams et al., 2017), and were not done at this scale. While these previous studies recovered a monophyletic Euryarchaea and DPANN (in contrast to AleRax), there is not enough information to discount either topology. The topology inferred by AleRax could therefore be viewed as lending credence to the idea that the DPANN are not monophyletic.

Frankiales dataset. In the case of the Frankiales, reconstructions with the four methods yield identical relationships between the major clades (Figures 5 and 6). This suggests that all four methods have comparable efficacy when the dataset analyzed is less complex and less divergent. Since this analysis used the entire pan-genome of the Frankiales, a possible concern is that small gene families (such as those that are only found in 4-8 genomes) may negatively impact the methods. To assess the impact of small gene families on species tree reconstruction, a subset of 1,702 genes families present in at least 20 genomes and in the smallest Frankia genome (*Frankia* sp. DG2) was used for inference using PhyloGTP and SpeciesRax. The trees produced from this subset recovered the same topologies for major clades as those in the full complement, indicating that the smaller gene families are not a problem for these methods.

In comparison to previous trees inferred on the same genomes using previous methods, such as those shown in Gosselin et al. 2022, there are a few rearrangements of early branching clades in the backbone of the Frankiales. In phylogenies inferred using tANI and MLSA sequence methods, Group 1 (Figure 5) is basal to the rest of the Frankiales. In the trees inferred by the four methods, Group 3 is basal to the other Frankiales, with Group 1 as a later branching basal group. In addition to the movement of these clades, *Frankia* sp. NRRLB16219 and *Frankia* sp. CgIS1 have swapped positions, where *Frankia* sp. CgIS1 has moved from Group 2 to Group 5. These rearrangements may be attributed to the additional genomic data used to reconstruct the four genome-scale trees. Only 24 loci were used in Gosselin et al. 2022, and the inclusion of thousands of additional gene families have painted a slightly different picture of evolution throughout the Frankiales. This suggests that truly genome-scale methods like AleRax, PhyloGTP, ASTRAL-Pro 2, and SpeciesRax could lead to more accurate phylogenomic inference on real datasets compared to other methods. These results also suggest that all four methods can perform well when analyzing less divergent datasets with large numbers of gene families.

Discussion

In this work, we systematically evaluated four recently developed methods, ASTRAL-Pro 2, SpeciesRax, PhyloGTP, and AleRax, for microbial phylogenomics. These methods can all use thousands of complete (multi-copy) gene families, thereby enabling truly genome-scale microbial phylogenomic species tree inference. Our simulation study identifies AleRax, the only method that explicitly accounts for gene tree inference error/uncertainty, as the best species tree reconstruction accuracy among all tested methods. PhyloGTP shows the best overall accuracy among methods that do not explicitly account for gene tree error and uncertainty, performing particularly well on datasets with high DTL rates or a small number of gene families. Experiments using AleRax with and without gene tree error-correction show that error-correction can lead to an approximately 20% reduction in species tree reconstruction error, especially when the input gene trees are of poor quality.

We also find that AleRax, PhyloGTP, and SpeciesRax almost always outperform ASTRAL-Pro 2, a highly scalable but HGT-naïve method. However, our experiments also show that ASTRAL-Pro 2 matches the accuracies of the best methods on datasets with low or medium rates of DTL and a large number of input gene families (1000 in our tests). This suggests that ASTRAL-Pro 2, as well as other closely related quartet based methods such as DISCO-ASTRAL (Willson et al., 2021), could be potentially useful for analyzing such datasets, especially given their exceptional

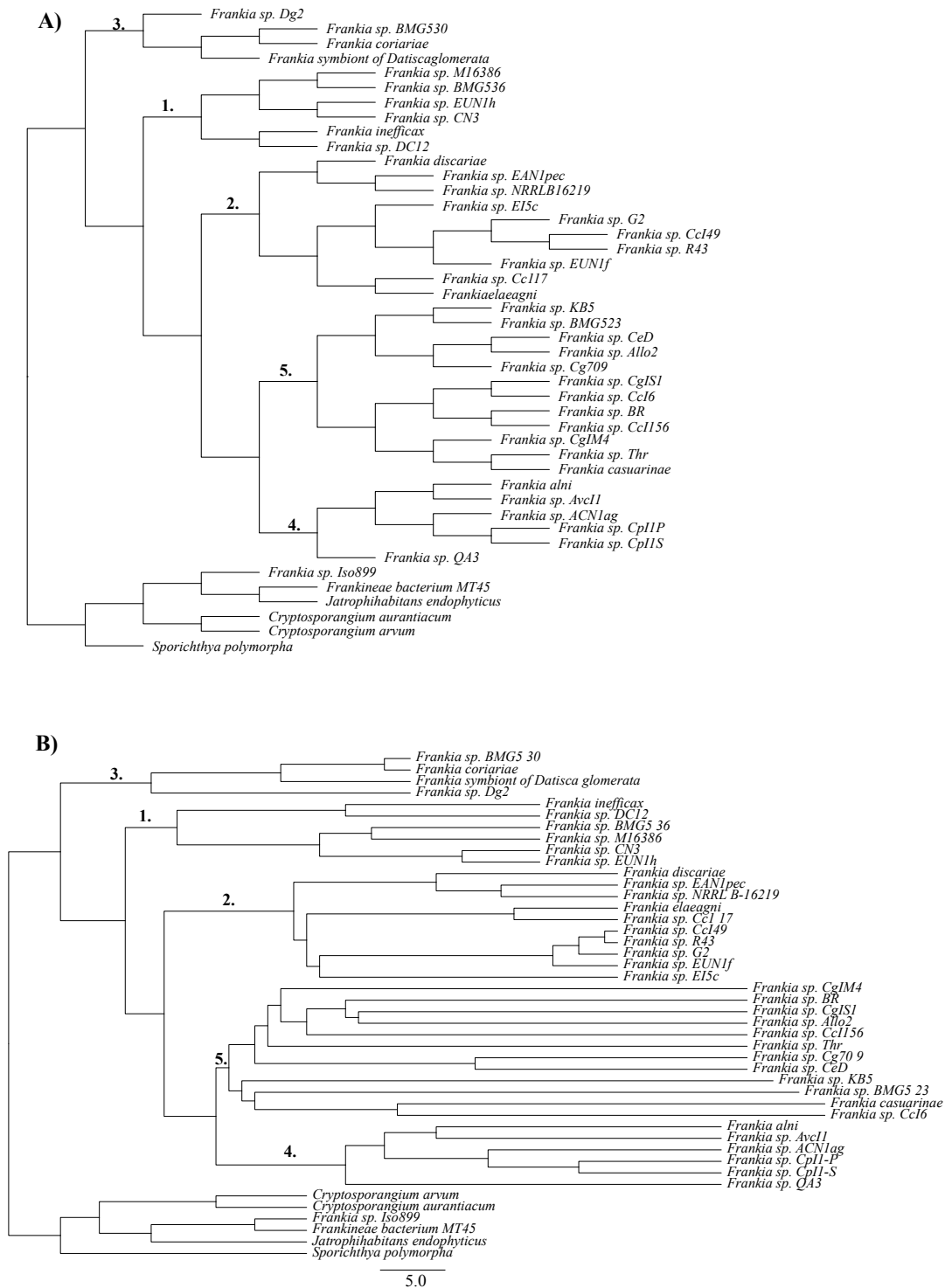


Figure 5 – Frankiales species tree reconstructions by PhyloGTP and AleRaX. Clades on both trees are categorized and enumerated based on the group designations described in Gosselin et al. 2022. Note that both trees show identical relationships among the labeled clades, but not necessarily within those clades. Part A) Unrooted Frankiales cladogram inferred by PhyloGTP. Part B) Unrooted Frankiales tree inferred by AleRaX.

speed. Another potential advantage of ASTRAL-Pro 2 is that it is agnostic to the underlying evolutionary processes and may therefore be more robust to the effect of evolutionary processes such as incomplete lineage sorting on the datasets being analyzed. On the other hand, the robustness of ASTRAL-Pro 2 to HGT may only hold under simple stochastic models of HGT (Davidson et al.,

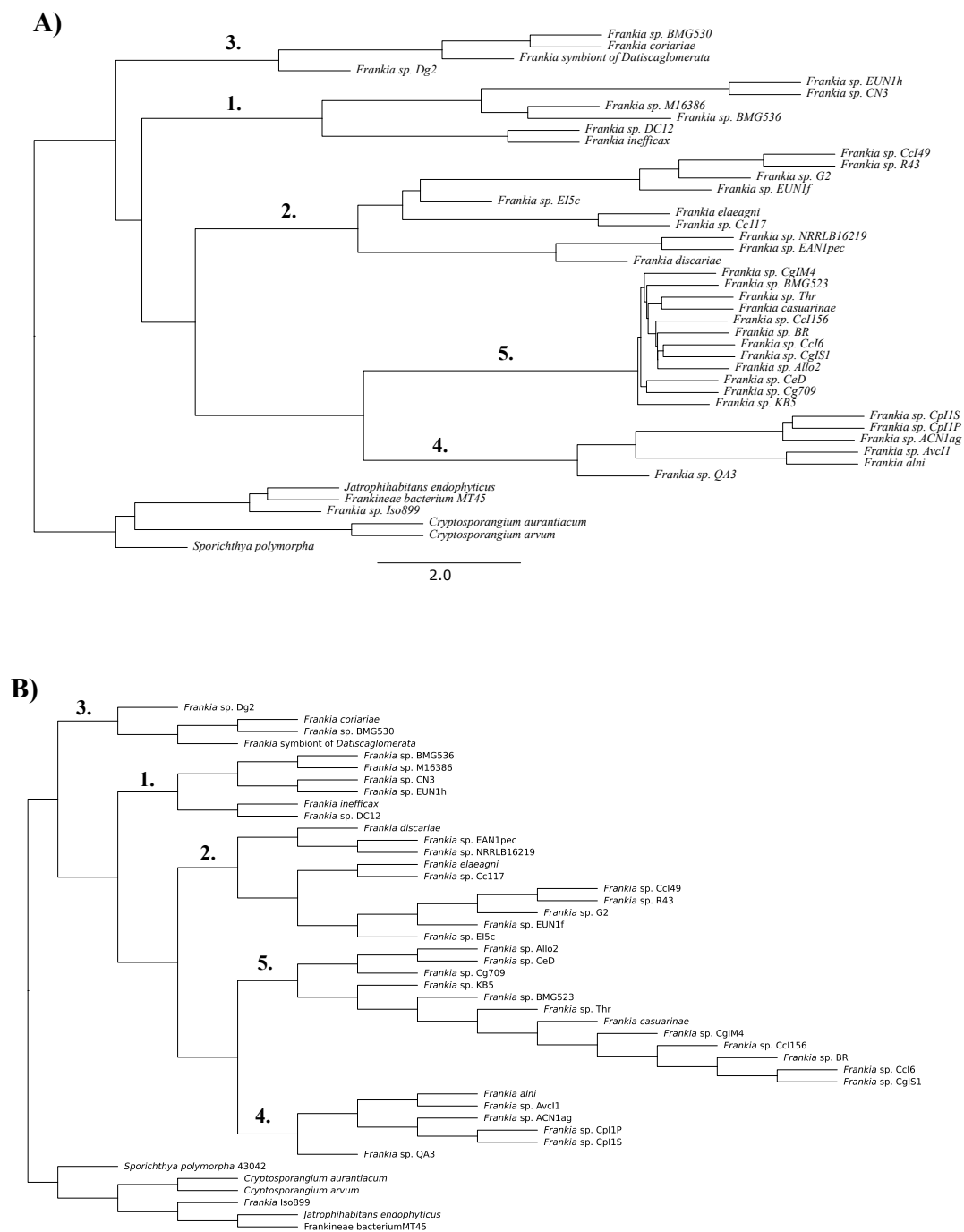


Figure 6 – Frankiales species tree reconstructions by ASTRAL-Pro 2 and SpeciesRax. Clades on both trees are categorized and enumerated based on the group designations described in Gosselin et al. 2022. Note that both trees show identical relationships among the labeled clades, but not necessarily within those clades. Part A) Unrooted Frankiales tree inferred by ASTRAL-Pro 2. Part B) Unrooted Frankiales cladogram inferred by SpeciesRax.

2015), of the kind employed in our simulation study. ASTRAL-Pro 2 may therefore be more susceptible than other methods to non-random patterns of HGT, such as when HGT frequency depends on the phylogenetic distance between donor and recipient or when some donor-recipient pairs are more likely to engage in HGT than others.

Application of these methods to the two biological datasets of different complexities provides additional valuable insights. We find that all four methods perform well on the less complex Frankiales dataset, but see mixed results on the more complex Archaeal dataset. This disparity was an intentional element of our study design, as these datasets represent opposite ends of a DTL complexity spectrum. The Frankiales dataset, with its higher phylogenetic resolution serves as a transitional case between simulation and highly complex empirical scenarios. Conversely, the Archaeal dataset, with its extensive evolutionary time and documented phylogenetic controversies, represents a stress test for these methods under confounding conditions.

PhyloGTP, SpeciesRax and ASTRAL-Pro 2 produce Archaeal trees that are mostly consistent with current estimates but also have some clearly incorrect placements. On the other hand, AleRax produces a tree that is markedly different than any highly supported previously calculated Archaeal tree. This may be because AleRax is unable to perform well on this complex, highly divergent dataset, or because the AleRax tree more accurately reflects the true evolutionary history of this group. The difficulty in resolving Archaeal phylogenies is well-documented in the literature and stems from a manifold of factors including HGT, compositional biases and long phylogenetic distances. Rather than viewing these inconsistencies as methodological failures, we interpret them as informative indicators of each method's behavior when confronted with increased evolutionary complexity. Such benchmarking against empirical datasets of varying complexity provides a more realistic assessment of method performance than simulations alone.

Overall, these results suggest that all tested methods are potentially susceptible to compositional and other biases present in complex datasets, and that the results of AleRax, in particular, may need to be interpreted with caution. Our findings underscore the importance of method selection based on expected dataset complexity and highlight the need for careful evaluation of results when analyzing domain-level phylogenies with extensive evolutionary histories. The disparity between simulation and empirical outcomes also suggests opportunities for further methodological refinement.

This work identifies several directions for future research on microbial phylogenomics. First, given our findings on the Archaeal dataset, it would be useful to develop simulation frameworks that better incorporate the specific challenges present in complex empirical datasets. Such advances would help bridge the current gap between theoretical performance and practical application. Second, and related to the above point, it may be informative to perform an expanded simulation study that assesses the impact of additional evolutionary parameters not assessed in the current study. For example, one could assess the impact of heterotachy along species tree branches, distance-dependent (non-uniform) HGT rates, presence of incomplete lineage sorting and gene conversion, compositional biases, etc. Third, the two most accurate methods, AleRax and PhyloGTP, are also the slowest, by far, and could therefore benefit from further methodological and algorithmic development and optimization. And fourth, our results indicate that most methods could benefit by implementing gene tree error-correction or using other strategies to account for gene tree error and uncertainty.

Acknowledgements

We thank Kashvi Parashar, a high-school student who worked in MSB's lab during Summer 2024, for helping assess the biological realism of our simulated datasets by comparing their rates of duplication, HGT and loss to those found in real biological datasets.

Preprint version 4 of this article has been peer-reviewed and recommended by Peer Community In Mathematical and Computational Biology (<https://doi.org/10.24072/pci.mcb.100408>; Scornavacca 2025).

Fundings

This work was supported in part by a University of Connecticut Research Excellence Program award to JPG and MSB.

Conflict of interest disclosure

The authors declare that they comply with the PCI rule of having no financial conflicts of interest in relation to the content of the article.

Data, script, code, and supplementary information availability

Software implementations of all methods evaluated in this work are freely available through their respective web pages or repositories.

The two real datasets are derived from Feng et al. 2021 and Gosselin et al. 2022 and the corresponding alignments and scripts are available online via Zenodo (<https://doi.org/10.5281/zenodo.14213473>; Bansal et al. 2024).

Supplementary information (including supplementary text, tables, and figures) and scripts used to generate simulated datasets are available online via Zenodo (<https://doi.org/10.5281/zenodo.15811297>; Weiner et al. 2025).

References

- Aouad M, Borrel G, Brochier-Armanet C, Gribaldo S (2019). Evolutionary placement of Methanona-tronarchaea. *Nature microbiology* **4**, 558–559. <https://doi.org/10.1038/s41564-019-0359-z>.
- Aouad M, Taib N, Oudart A, Lecocq M, Gouy M, Brochier-Armanet C (2018). Extreme halophilic archaea derive from two distinct methanogen Class II lineages. *Molecular phylogenetics and evolution* **127**, 46–54. <https://doi.org/10.1016/j.ympev.2018.04.011>.
- Bansal M, Gogarten JP, Feng Y, Weiner S (2024). *Archaea and Frankia gene family alignments and associated scripts*. Zenodo. <https://doi.org/10.5281/zenodo.14213473>.
- Bansal MS, Alm EJ, Kellis M (2012). Efficient algorithms for the reconciliation problem with gene duplication, horizontal transfer and loss. *Bioinformatics* **28**, 283–291. <https://doi.org/10.1093/bioinformatics/bts225>.
- Bansal MS, Shamir R (2011). A Note on the Fixed Parameter Tractability of the Gene-Duplication Problem. *IEEE/ACM Trans. Comput. Biology Bioinform.* **8**, 848–850. <https://doi.org/10.1109/TCBB.2010.74>.

- Bansal MS, Wu YC, Alm EJ, Kellis M (2015). Improved gene tree error correction in the presence of horizontal gene transfer. *Bioinformatics* **31**, 1211–1218. <https://doi.org/10.1093/bioinformatics/btu806>.
- Beiko RG, Harlow TJ, Ragan MA (2005). Highways of gene sharing in prokaryotes. *Proceedings of the National Academy of Sciences of the United States of America* **102**, 14332–14337. <https://doi.org/10.1073/pnas.0504068102>.
- Brochier-Armanet C, Forterre P, Gribaldo S (2011). Phylogeny and evolution of the Archaea: one hundred genomes later. *Current opinion in microbiology* **14**, 274–281. <https://doi.org/10.1016/j.mib.2011.04.015>.
- Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL (2009). BLAST+: architecture and applications. *BMC bioinformatics* **10**, 1–9. <https://doi.org/10.1186/1471-2105-10-421>.
- Ciccarelli FD, Doerks T, von Mering C, Creevey CJ, Snel B, Bork P (2006). Toward Automatic Reconstruction of a Highly Resolved Tree of Life. *Science* **311**, 1283–1287. <https://doi.org/10.1126/science.1123061>.
- David LA, Alm EJ (2011). Rapid evolutionary innovation during an Archaeal genetic expansion. *Nature* **469**, 93–96. <https://doi.org/10.1038/nature09649>.
- Davidson R, Vachaspati P, Mirarab S, Warnow T (2015). Phylogenomic species tree estimation in the presence of incomplete lineage sorting and horizontal gene transfer. *BMC Genomics* **16** (Suppl 10), S1. <https://doi.org/10.1186/1471-2164-16-S10-S1>.
- Davín AA, Tannier E, Williams TA, Boussau B, Daubin V, Szöllősi GJ (2018). Gene transfers can date the tree of life. *Nature ecology & evolution* **2**, 904–909. <https://doi.org/10.1038/s41559-018-0525-3>.
- Doolittle WF, Boucher Y, Nesbo CL, Douady CJ, Andersson JO, Roger AJ (2003). How big is the iceberg of which organellar genes in nuclear genomes are but the tip? *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences* **358**, 39–58. <https://doi.org/10.1098/rstb.2002.1185>.
- Doolittle WF (1999). Phylogenetic Classification and the Universal Tree. *Science* **284**, 2124–2128. <https://doi.org/10.1126/science.284.5423.2124>.
- Emms DM, Kelly S (2019). OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome biology* **20**, 1–14. <https://doi.org/10.1186/s13059-019-1832-y>.
- Feng Y, Neri U, Gosselin S, Louyakis AS, Papke RT, Gophna U, Gogarten JP (2021). The evolutionary origins of extreme halophilic archaeal lineages. *Genome biology and evolution* **13**, evab166. <https://doi.org/10.1093/gbe/evab166>.
- Gadagkar SR, Rosenberg MS, Kumar S (2005). Inferring species phylogenies from multiple genes: Concatenated sequence tree versus consensus gene tree. *Journal of Experimental Zoology Part B: Molecular and Developmental Evolution* **304B**, 64–74. <https://doi.org/10.1002/jez.b.21026>.
- Glaeser SP, Kampfer P (2015). Multilocus sequence analysis (MLSA) in prokaryotic taxonomy. *Systematic and Applied Microbiology* **38**, 237–245. <https://doi.org/10.1016/j.syapm.2015.03.007>.
- Gogarten JP, Doolittle WF, Lawrence JG (2002). Prokaryotic Evolution in Light of Gene Transfer. *Molecular Biology and Evolution* **19**, 2226–2238. <https://doi.org/10.1093/oxfordjournals.molbev.a004046>.

- Gosselin S, Fullmer MS, Feng Y, Gogarten JP (2022). Improving phylogenies based on average nucleotide identity, incorporating saturation correction and nonparametric bootstrap support. *Systematic Biology* **71**, 396–409. <https://doi.org/10.1093/sysbio/syab060>.
- Henz SR, Huson DH, Auch AF, Nieselt-Struwe K, Schuster SC (2004). Whole-genome prokaryotic phylogeny. *Bioinformatics* **21**, 2329–2335. <https://doi.org/10.1093/bioinformatics/bth324>.
- Hilario E, Gogarten JP (1993). Horizontal transfer of {ATPase} genes – the tree of life becomes a net of life. *Biosystems* **31**, 111–119. [https://doi.org/10.1016/0303-2647\(93\)90038-e](https://doi.org/10.1016/0303-2647(93)90038-e).
- Hirt RP, Logsdon JM, Healy B, Dorey MW, Doolittle WF, Embley TM (1999). Microsporidia are related to Fungi: Evidence from the largest subunit of RNA polymerase II and other proteins. *Proceedings of the National Academy of Sciences* **96**, 580–585. <https://doi.org/10.1073/pnas.96.2.580>.
- Kalyanamoorthy S, Minh BQ, Wong TK, Von Haeseler A, Jermini LS (2017). ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature methods* **14**, 587–589. <https://doi.org/10.1038/nmeth.4285>.
- Katoh K, Standley DM (2013). MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Molecular biology and evolution* **30**, 772–780. <https://doi.org/10.1093/molbev/mst010>.
- Konstantinidis KT, Tiedje JM (2005). Genomic insights that advance the species definition for prokaryotes. *Proceedings of the National Academy of Sciences* **102**, 2567–2572. <https://doi.org/10.1073/pnas.0409727102>.
- Kundu S, Bansal MS (2019). SaGePhy: an improved phylogenetic simulation framework for gene and subgene evolution. *Bioinformatics* **35**. <https://doi.org/10.1093/bioinformatics/btz081>.
- Lafond M, Scornavacca C (2019). On the Weighted Quartet Consensus problem. *Theoretical Computer Science* **769**, 1–17. <https://doi.org/10.1016/j.tcs.2018.10.005>.
- Lang JM, Darling AE, Eisen JA (2013). Phylogeny of Bacterial and Archaeal Genomes Using Conserved Genes: Supertrees and Supermatrices. *PLoS ONE* **8**, e62510. <https://doi.org/10.1371/journal.pone.0062510>.
- Lapierre P, Lasek-Nesselquist E, Gogarten JP (2014). The impact of HGT on phylogenomic reconstruction methods. *Briefings in Bioinformatics* **15**, 79–90. <https://doi.org/10.1093/bib/bbs050>.
- Le SQ, Gascuel O (2008). An Improved General Amino Acid Replacement Matrix. *Molecular Biology and Evolution* **25**, 1307–1320. <https://doi.org/10.1093/molbev/msn067>.
- Lewis PO, Chen MH, Kuo L, Lewis LA, Fucikova K, Neupane S, Wang YB, Shi D (2016). Estimating Bayesian Phylogenetic Information Content. *Systematic Biology* **65**, 1009–1023. <https://doi.org/10.1093/sysbio/syw042>.
- Ma B, Li M, Zhang L (2000). From Gene Trees to Species Trees. *SIAM J. Comput.* **30**, 729–752. <https://doi.org/10.1137/s0097539798343362>.
- Markowitz VM, Chen IMA, Palaniappan K, Chu K, Szeto E, Pillay M, Ratner A, Huang J, Woyke T, Huntemann M, Anderson I, Billis K, Varghese N, Mavromatis K, Pati A, Ivanova NN, Kyrpides NC (2014). IMG 4 version of the integrated microbial genomes comparative analysis system. *Nucleic Acids Research* **42**, D560–D567. <https://doi.org/10.1093/nar/gkt963>.

- Martinez-Gutierrez CA, Uyeda JC, Aylward FO (2023). A timeline of bacterial and archaeal diversification in the ocean. *eLife* **12**. Ed. by John McCutcheon and George H Perry, RP88268. <https://doi.org/10.7554/eLife.88268.3>.
- McInerney JO, Cotton JA, Pisani D (2008). The prokaryotic tree of life: past, present... and future? *Trends in Ecology & Evolution* **23**, 276–281. <https://doi.org/10.1016/j.tree.2008.01.008>.
- Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, Lanfear R (2020). IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Molecular Biology and Evolution* **37**, 1530–1534. <https://doi.org/10.1093/molbev/msaa015>.
- Mirarab S, Reaz R, Bayzid MS, Zimmermann T, Swenson MS, Warnow T (2014). ASTRAL: genome-scale coalescent-based species tree estimation. *Bioinformatics* **30**, i541–i548. <https://doi.org/10.1093/bioinformatics/btu462>.
- Morel B, Kozlov AM, Stamatakis A, Szöllősi GJ (2020). GeneRax: A Tool for Species-Tree-Aware Maximum Likelihood-Based Gene Family Tree Inference under Gene Duplication, Transfer, and Loss. *Molecular Biology and Evolution* **37**, 2763–2774. <https://doi.org/10.1093/molbev/msaa141>.
- Morel B, Schade P, Lutteropp S, Williams TA, Szöllősi GJ, Stamatakis A (2022). SpeciesRax: A Tool for Maximum Likelihood Species Tree Inference from Gene Family Trees under Duplication, Transfer, and Loss. *Molecular Biology and Evolution* **39**, msab365. <https://doi.org/10.1093/molbev/msab365>.
- Morel B, Williams TA, Stamatakis A, Szollosi GJ (2024). AleRax: a tool for gene and species tree co-estimation and reconciliation under a probabilistic model of gene duplication, transfer, and loss. *Bioinformatics* **40**, btae162. <https://doi.org/10.1093/bioinformatics/btae162>.
- Narasimharao P, Podell S, Ugalde JA, Brochier-Armanet C, Emerson JB, Brocks JJ, Heidelberg KB, Banfield JF, Allen EE (2012). De novo metagenomic assembly reveals abundant novel major lineage of Archaea in hypersaline microbial communities. *The ISME journal* **6**, 81–93. <https://doi.org/10.1038/ismej.2011.78>.
- Olsen GJ, Woese CR, Overbeek R (1994). The winds of (evolutionary) change: breathing new life into microbiology. *Journal of Bacteriology* **176**, 1–6. <https://doi.org/10.1128/jb.176.1.1-6.1994>.
- Price MN, Dehal PS, Arkin AP (2010). FastTree 2—approximately maximum-likelihood trees for large alignments. *PloS one* **5**, e9490. <https://doi.org/10.1371/journal.pone.0009490>.
- Rangel LT, Marden J, Colston S, Setubal JC, Graf J, Gogarten JP (2019). Identification and characterization of putative *Aeromonas* spp. T3SS effectors". *PLOS ONE* **14**, 1–20. <https://doi.org/10.1371/journal.pone.0214035>.
- Raymann K, Brochier-Armanet C, Gribaldo S (2015). The two-domain tree of life is linked to a new root for the Archaea. *Proceedings of the National Academy of Sciences* **112**, 6670–6675. <https://doi.org/10.1073/pnas.1420858112>.
- Raymann K, Forterre P, Brochier-Armanet C, Gribaldo S (2014). Global phylogenomic analysis disentangles the complex evolutionary history of DNA replication in archaea. *Genome biology and evolution* **6**, 192–212. <https://doi.org/10.1093/gbe/evu004>.
- Robinson D, Foulds L (1981). Comparison of phylogenetic trees. *Mathematical Biosciences* **53**, 131–147. [https://doi.org/10.1016/0025-5564\(81\)90043-2](https://doi.org/10.1016/0025-5564(81)90043-2).

- Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Höhna S, Larget B, Liu L, Suchard MA, Huelsenbeck JP (2012). MrBayes 3.2: Efficient Bayesian Phylogenetic Inference and Model Choice Across a Large Model Space. *Systematic Biology* **61**, 539–542. <https://doi.org/10.1093/sysbio/sys029>.
- Scornavacca C (2025). Inferring species trees under extensive horizontal gene transfer: insights from simulated and empirical data. *Peer Community in Mathematical and Computational Biology*. <https://doi.org/10.24072/pci.mcb.100408>.
- Sevillya G, Doerr D, Lerner Y, Stoye J, Steel M, Snir S (2019). Horizontal Gene Transfer Phylogenetics: A Random Walk Approach. *Molecular Biology and Evolution* **37**, 1470–1479. <https://doi.org/10.1093/molbev/msz302>.
- Shifman A, Ninyo N, Gophna U, Snir S (2014). Phylo SI: a new genome-wide approach for prokaryotic phylogeny. *Nucleic Acids Research* **42**, 2391–2404. <https://doi.org/10.1093/nar/gkt1138>.
- Sorokin DY, Makarova KS, Abbas B, Ferrer M, Golyshin PN, Galinski EA, Ciorda S, Mena MC, Merkel AY, Wolf YI, et al. (2019). Reply to ‘Evolutionary placement of Methanonatronarchaeia’. *Nature Microbiology* **4**, 560–561. <https://doi.org/10.1038/s41564-019-0358-0>.
- Szollósi GJ, Tannier E, Lartillot N, Daubin V (2013). Lateral Gene Transfer from the Dead. *Systematic Biology* **62**, 386–397. <https://doi.org/10.1093/sysbio/syt003>.
- Tofigh A, Hallett MT, Lagergren J (2011). Simultaneous Identification of Duplications and Lateral Gene Transfers. *IEEE/ACM Trans. Comput. Biology Bioinform.* **8**, 517–535. <https://doi.org/10.1109/TCBB.2010.14>.
- Ly-Trong N, Naser-Khdour S, Lanfear R, Minh BQ (2022). AliSim: A Fast and Versatile Phylogenetic Sequence Simulator for the Genomic Era. *Molecular Biology and Evolution* **39**, msac092. <https://doi.org/10.1093/molbev/msac092>.
- Weiner S, Feng Y, Gogarten JP, Bansal MS (2024). Assessing the Potential of Gene Tree Parsimony for Microbial Phylogenomics. In: *Comparative Genomics*. Ed. by Celine Scornavacca and Maribel Hernández-Rosales. Cham: Springer Nature Switzerland, pp. 129–149. https://doi.org/10.1007/978-3-031-58072-7_7.
- Weiner S, Feng Y, Gogarten JP, Bansal M (2025). Supplementary Material: A systematic assessment of phylogenomic approaches for microbial species tree reconstruction. <https://doi.org/10.5281/zenodo.15811297>.
- Whidden C, Zeh N, Beiko RG (2014). Supertrees Based on the Subtree Prune-and-Regraft Distance. *Systematic Biology* **63**, 566–581. <https://doi.org/10.1093/sysbio/syu023>.
- Williams TA, Szollósi GJ, Spang A, Foster PG, Heaps SE, Boussau B, Ettema TJG, Embley TM (2017). Integrative modeling of gene and genome evolution roots the archaeal tree of life. *Proceedings of the National Academy of Sciences* **114**, E4602–E4611. <https://doi.org/10.1073/pnas.1618463114>.
- Willson J, Roddur MS, Liu B, Zaharias P, Warnow T (2021). DISCO: Species Tree Inference using Multicopy Gene Family Tree Decomposition. *Systematic Biology* **71**, 610–629. <https://doi.org/10.1093/sysbio/syab070>.
- Woese CR (1987). Bacterial evolution. *Microbiological Reviews* **51**, 221–271. <https://doi.org/10.1128/mr.51.2.221-271.1987>.
- Yap WH, Zhang Z, Wang Y (1999). Distinct Types of rRNA Operons Exist in the Genome of the Actinomycete *Thermomonospora chromogena* and Evidence for Horizontal Transfer of an

- Entire rRNA Operon. *Journal of Bacteriology* **181**, 5201–5209. <https://doi.org/10.1128/jb.181.17.5201-5209.1999>.
- Zhang C, Mirarab S (2022). ASTRAL-Pro 2: ultrafast species tree reconstruction from multi-copy gene family trees. *Bioinformatics* **38**, 4949–4950. <https://doi.org/10.1093/bioinformatics/btac620>.
- Zhaxybayeva O, Doolittle WF, Papke RT, Gogarten JP (2009). Intertwined Evolutionary Histories of Marine Synechococcus and Prochlorococcus marinus. *Genome Biology and Evolution* **1**, 325–339. <https://doi.org/10.1093/gbe/evp032>.
- Zhaxybayeva O, Stepanauskas R, Mohan NR, Papke RT (2013). Cell sorting analysis of geographically separated hypersaline environments. *Extremophiles* **17**, 265–275. <https://doi.org/10.1007/s00792-013-0514-z>.