

Syst. Biol. 00(0):1–7, 2016
 © The Author(s) 2016. Published by Oxford University Press, on behalf of the Society of Systematic Biologists.
 This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>),
 which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.
 DOI:10.1093/sysbio/syw064

Species Tree Inference under the Multispecies Coalescent on Data with Paralogs is Accurate

ZHI YAN¹, PENG DU¹, MATTHEW W. HAHN², AND LUAY NAKHLEH^{1,3,*}

¹Department of Computer Science, Rice University, 6100 Main Street, Houston, TX 77005, USA; ²Department of Biology and Department of Computer Science, Indiana University, 1001 East Third Street, Bloomington, IN 47405, USA; ³Department of BioSciences, Rice University, 6100 Main Street, Houston, TX 77005, USA;

*Correspondence to be sent to: Department of Computer Science, Rice University, 6100 Main Street, Houston, TX 77005, USA;
 E-mail: nakhleh@rice.edu.

Received ...; reviews returned ...; accepted ...

Abstract.—The multispecies coalescent (MSC) has emerged as a powerful and desirable framework for species tree inference in phylogenomic studies. Under this framework, the data for each locus is assumed to consist of orthologous, single-copy genes, and heterogeneity across loci is assumed to be due to incomplete lineage sorting (ILS). These assumptions have led biologists that use ILS-aware inference methods, whether based directly on the MSC or proven to be statistically consistent under it (collectively referred to here as MSC-based methods), to exclude all loci that are present in more than a single copy in any of the studied genomes. Furthermore, such analyses entail orthology assignment to avoid the potential of hidden paralogy in the data. The question we seek to answer in this study is: What happens if one runs such species tree inference methods on data where paralogy is present, in addition to or without ILS being present? Through simulation studies and analyses of two biological data sets, we show that running such methods on data with paralogs provide very accurate results, either by treating all gene copies within a family as alleles from multiple individuals or by randomly selecting one copy per species. Our results have significant implications for the use of MSC-based phylogenomic analyses, demonstrating that they do not have to be restricted to single-copy loci, thus greatly increasing the amount of data that can be used. [Multispecies coalescent; incomplete lineage sorting; gene duplication and loss; orthology; paralogy.]

1 Species phylogeny inference from genome-wide data
 2 entails accounting for the fact that the evolutionary
 3 histories of different loci can disagree with each other,
 4 as well as with the phylogeny of the species. The reasons
 5 for this incongruence include biological causes such as
 6 incomplete lineage sorting (ILS) and horizontal gene
 7 transfer (broadly interpreted to include all biological
 8 processes involving gene flow), as well as technical causes
 9 such as the misidentification of paralogs as orthologs
 10 (“hidden paralogy”; (Maddison, 1997)).

11 Inference of species phylogenies in light of
 12 incongruence relies on a model of how gene trees
 13 evolve within the branches of the species phylogeny (in
 14 addition to probabilistic models of sequence evolution
 15 on the gene trees). The multispecies coalescent (MSC)
 16 (Degnan and Rosenberg, 2009; Hudson, 1983; Rannala
 17 and Yang, 2003; Takahata, 1989) has emerged as
 18 the most commonly employed model of such gene
 19 genealogies. Indeed, in the last two decades, a wide
 20 array of methods and computer programs have been
 21 developed for species tree inference under the MSC; see
 22 (Knowles and Kubatko, 2011; Liu *et al.*, 2009, 2015;
 23 Nakhleh, 2013) for recent reviews and surveys of these
 24 methods.

25 MSC-based inference of species trees assumes that
 26 the data consist of only orthologous sequences obtained
 27 from multiple loci across the genomes of the species
 28 under study. Indeed, most phylogenetic methods require
 29 the identification of orthologs. Therefore, before such
 30 inference methods are applied to a phylogenomic data
 31 set, paralogs must be identified and removed from the
 32 data. One way to accomplish this is by using orthology
 33 detection tools within larger gene families, and another
 34 is to restrict the data to only loci that are present in

35 exactly (and, in some cases, at most) one copy in each of
 36 the genomes. Neither of these two approaches guarantees
 37 that the resulting data set includes only orthologous
 38 sequences (Koonin, 2005). Furthermore, restricting the
 39 data to single-copy genes, which is by far the most
 40 common practice in the community, means much of the
 41 data must be excluded from the analysis. In particular,
 42 as more species are sampled, the frequency of genes that
 43 are present in single-copy across all species is expected
 44 to decrease.

45 MSC-based inference relies on the distribution of gene
 46 trees as the signal for inference, and it is likely that the
 47 MSC induces a distribution that differs from that arising
 48 from gene duplication and loss (GDL). An obvious way
 49 to handle data sets where ILS and gene duplication/loss
 50 could have simultaneously acted on gene families is to
 51 employ models of gene evolution that go beyond the
 52 MSC in order to incorporate GDL as well. Indeed,
 53 such models are beginning to emerge (Boussau *et al.*,
 54 2013; Rasmussen and Kellis, 2012). However, the more
 55 complex the models of gene family evolution, the more
 56 computationally prohibitive statistical inference under
 57 these models becomes (Du and Nakhleh, 2018), rendering
 58 their applicability infeasible except for very small data
 59 sets in terms of the number of species and gene families.

60 Given that much progress in terms of accuracy and
 61 computational efficiency has been made on MSC-based
 62 species tree inference methods, we ask in this paper the
 63 following question: Is MSC-based species tree inference
 64 robust to the presence of paralogs in the data? If the
 65 answer to this question is positive, then the reach of
 66 MSC-based inference methods is significantly extended
 67 and the exclusion of loci from phylogenomic data sets
 68 is deemed unnecessary, thus providing more signal for

the inference task. To answer this question, we study the performance of three species tree inference methods, all of which use gene trees as the input data: The maximum pseudo-likelihood method of [Yu and Nakhleh \(2015\)](#) as implemented by the function `InferNetwork_MPL` in PhyloNet ([Wen et al., 2018](#)), ASTRAL-III ([Zhang et al., 2018](#)), and NJ_{st} ([Liu and Yu, 2011](#)). It is important to point out that we ran `InferNetwork_MPL` restricting the search to the tree space (i.e., not allowing reticulations), which amounts to running the species tree inference method MP-EST ([Liu et al., 2010](#)).

For the sake of conclusions that we draw from this study, it may be helpful to highlight the difference between `InferNetwork_MPL` on the one hand and ASTRAL and NJ_{st} on the other hand. The former optimizes a pseudo-likelihood function that is derived based on the assumptions of the MSC. This function is very different, for example, from a likelihood function that is based on a birth-death model of gene family evolution ([Arvestad et al., 2009](#)). Therefore, its accuracy in inferring species trees from data with paralogs reflects directly on the performance of MSC-based methods on such data. The latter two methods make no direct use of the MSC, but have been shown to be statistically consistent under the MSC. Therefore, their accuracy on data with paralogs reflects the suitability of these methods, rather than the MSC itself, for analyzing such data.

As these methods assume only orthologs are used for each locus—and to facilitate their use on data with paralogs—we either treat multiple copies of a gene in the same genome as multiple alleles from different individuals, or randomly sample a single copy per genome. Of course, the former violates the mathematical assumptions of the MSC, and the latter does not guarantee that the single-copy genes are orthologs (in fact, some would not be with very high probability). Our results show that inferences made by all methods are very accurate, and are mostly identical to the accuracy of their inferences when using only genes that are present in exactly one copy in each of the species. Particularly striking is the finding that these methods infer very accurate species trees when all gene tree incongruence is due to GDL, and ILS is not a factor. We find that gene tree estimation error affects the methods' performances at a similar, or even higher, level than ILS.

METHODS

Species tree inference methods

For species tree inference, we use three methods that are statistically consistent under the MSC:

- The maximum pseudo-likelihood inference function `InferNetwork_MPL` in PhyloNet, which implements the method of [Yu and Nakhleh \(2015\)](#) and amounts to running MP-EST ([Liu et al., 2010](#)) when restricted to trees with no reticulations.

- ASTRAL-III ([Zhang et al., 2018](#)), Version 5.6.3.

- NJ_{st} ([Liu and Yu, 2011](#)).

While the maximum likelihood method of [Yu et al. \(2014\)](#) as implemented by the `InferNetwork_ML` function in PhyloNet ([Wen et al., 2018](#)) is relevant here, it is much more computationally demanding than maximum pseudo-likelihood, and we chose not to run it.

Given a collection of trees corresponding to gene families (one tree per gene family), we generated three types of input to each of the methods:

- **ONLY:** The input consists of trees of *only* gene families that are present in exactly one copy in each of the species.

- **ONE:** The input consists of trees of *all* gene families, but where a single copy per species per gene family is selected at random and the remaining copies are removed. If a gene family has no copies at all for some species, then the resulting tree of that gene family also has no copies for that species.

- **ALL:** The input consists of trees of *all* gene families, but where all copies of a gene in a species are treated as multiple alleles from different individuals within the species. Similar to ONE, if a gene family has no copies at all for some species, then the resulting tree of that gene family also has no copies for that species.

ONLY corresponds to the practice that is followed in almost all phylogenomic studies, as it reduces, though not necessarily completely eliminate, the potential of hidden paralogy—a term that was introduced by [Daubin and Gouy \(2001\)](#); [Gribaldo and Philippe \(2002\)](#) to denote cases where single-copy paralogs within a gene family are mistaken for orthologs. Unless GDL is not involved, ONE is almost always guaranteed to have hidden paralogs in the input. By construction, ALL has all orthologs and paralogs in the input, but all are effectively labeled as orthologs, with the (wrong) implied assumption that multiple individuals are sampled per species.

Simulation setup

For model species trees, we used the trees on 16 fungi species and 12 fly species reported in ([Rasmussen and Kellis, 2012](#)) and shown in Fig. 1. The 16 fungi species are: *Candida albicans* (Calb), *Candida tropicalis* (Ctro), *Candida parapsilosis* (Cpar), *Lodderomyces elongisporus* (Lelo), *Candida guilliermondii* (Cgui), *Debaryomyces hansenii* (Dhan), *Candida lusitaniae* (Clus), *Saccharomyces cerevisiae* (Scer), *Saccharomyces paradoxus* (Spar), *Saccharomyces mikatae* (Smik), *Saccharomyces bayanus* (Sbay), *Candida glabrata* (Cgla), *Saccharomyces castellii* (Scas), *Kluyveromyces lactis*

(Klac), *Ashbya gossypii* (Agos), and *Kluyveromyces waltii* (Kwal). The 12 fly species are: *Drosophila melanogaster* (Dmel), *Drosophila simulans* (Dsim), *Drosophila sechellia* (Dsec), *Drosophila erecta* (Dere), *Drosophila yakuba* (Dyak), *Drosophila ananassae* (Dana), *Drosophila pseudoobscura* (Dpse), *Drosophila persimilis* (Dper), *Drosophila willistoni* (Dwil), *Drosophila mojavensis* (Dmoj), *Drosophila virilis* (Dvir), and *Drosophila grimshawi* (Dgri).

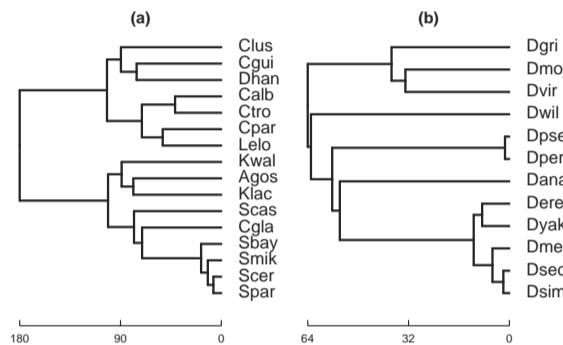


FIGURE 1. The species trees reported in (Rasmussen and Kellis, 2012) for the yeast and fly species, which are also used as the true species trees in the simulations. (a) The species tree of 16 fungi species. (b) The species tree of 12 fly species. The species tree topologies and their branch lengths in units of million years are reported on <http://compbio.mit.edu/dlcoal/>.

For the fungal simulated data sets, we used four different duplication and loss rates (assuming duplication and loss rates are equal): 0 (to investigate the performance when ILS, but not GDL, acted on the gene families), 1×10^{-10} , 2×10^{-10} , and 5×10^{-10} per generation where 1×10^{-10} is similar to the duplication rate of 7.32×10^{-11} and loss rate of 8.59×10^{-11} used by Rasmussen and Kellis (2011). We used two effective population sizes of 10^7 and 5×10^7 , where the former was also used by Rasmussen and Kellis (2012) as the true population size. We assumed 0.9 year per generation as in (Rasmussen and Kellis, 2012) and used 4×10^{-10} as the nucleotide mutation rate per generation per site, similar to the rates of 3.3×10^{-10} and 3.8×10^{-10} used by Zhang and Wu (2017) and Lang and Murray (2008), respectively.

For the fly simulated data sets, we used four different duplication and loss rates (assuming duplication and loss rates are equal): 0, 1×10^{-10} , 2×10^{-10} , and 5×10^{-10} per generation. A GDL rate of 1.2×10^{-10} was used in (Rasmussen and Kellis, 2012; Zhang and Wu, 2017) and reported by Hahn et al. (2007). We used two effective population sizes of 10^6 and 5×10^6 , similar to the values used in (Rasmussen and Kellis, 2012) and estimated value of 1.15×10^6 reported in (Pollard et al., 2006; Sawyer and Hartl, 1992). We assumed 10 generations per year as in (Rasmussen and Kellis, 2012; Zhang and Wu, 2017) and used 3×10^{-9} as the mutation rate per generation per site, similar to the rate of 5×10^{-9} found in (Schrider et al., 2013).

To generate gene trees while allowing for ILS and GDL, we used SimPhy (Mallo et al., 2015) with the parameters specified above (assuming all species are diploid). SimPhy uses the three-tree model developed in (Rasmussen and Kellis, 2012) to simulate data. In this model, a *locus tree* is simulated within the branches of the species tree. All incongruence between the locus tree and the species tree is due to GDL. Then, a gene tree is simulated within the branches of the locus tree, where all incongruence between the locus tree and the gene tree is due to ILS. The resulting gene tree differs from the species tree due to a combination of ILS and GDL. Using the locus trees as input to an inference method amounts to using data where all incongruence is solely due to GDL (but not ILS). Setting the rates of GDL to 0.0 amounts to generating gene trees where all incongruence is solely due to ILS. 10,000 gene families (each containing a locus tree and its corresponding gene tree) were simulated in this fashion as one data set for each combination of GDL rate and population size. Ten such data sets, each with 10,000 gene families, were generated.

To study the effect of using data sets of varying sizes, for each of the 10 data sets, we randomly sampled 10, 50, 100, and 250 gene families from the 10,000 gene families under the ALL, ONE, and ONLY scenarios. In case the number of available gene families that fits ONLY is smaller than the desired size, that number of gene families is used (e.g., when only 6 gene family trees are available when data sets of size 10 are desired, the 6 trees are used as input). Finally, for each data set of trees (true gene trees or true locus trees, that is, trees without estimation error) of a given size, we fed the data set as input to InferNetwork_MPL, ASTRAL, and NJ_{st} and computed the Robinson-Foulds distance (Robinson and Foulds, 1981), normalized by the number of internal branches in the (unrooted) species tree to obtain a value between 0 and 1, between the true and inferred species trees.

To study the effect of error in the gene tree estimates, we simulated the evolution of sequences of length 500 on all gene trees under the HKY model, using Seq-gen (Rambaut and Grassly, 1997). We then inferred gene trees from the simulated sequence data using IQTREE (Nguyen et al., 2014). InferNetwork_MPL assumes that the input gene trees are rooted. In this study, we rooted the gene tree estimates by minimizing deep coalescences (Maddison, 1997; Than and Nakhleh, 2009); that is, we rooted each gene tree in a way that minimizes the number of extra lineages when reconciled with the true species tree. Furthermore, to study the effect of error in the locus tree estimates, we treated the true locus tree as a gene tree and simulated the evolution of sequences of length 500 on all locus trees under the HKY model, using Seq-gen, and inferred locus trees from the simulated sequence data using IQTREE. It is important to note that in practice only gene trees, but not locus trees, are inferable, as the locus tree is an artifact of the three-tree model (Rasmussen and Kellis, 2012) and not a biological entity. However, conducting analysis using inferred locus trees

1 gives a picture of the performance when all incongruence
2 is due to GDL and gene tree error only.

3 Biological data

4 We randomly selected 250 gene trees estimated
5 with PhyML (Guindon and Gascuel, 2003) reported
6 in <http://compbio.mit.edu/dlcoal/> for the 16 fungi
7 genomes, as well as 250 gene trees estimated
8 with RAxML (Stamatakis, 2006) reported in
9 <http://www.treefam.org> for the 12 fly genomes for
10 the ALL, ONE and ONLY settings. We then estimated
11 the species trees using ASTRAL, NJ_{st}, and maximum
12 pseudo-likelihood with these gene trees as input.

13 We again rooted each gene tree in the empirical data
14 with respect to the species tree of Fig. 1 so as to minimize
15 deep coalescences (Maddison, 1997; Than and Nakhleh,
16 2009) using the method of Yu *et al.* (2011a, b), as
17 implemented by the function ProcessGT in PhyloNet
18 (Wen *et al.*, 2018). contrast

19 RESULTS

20 Characteristics of the simulated data

21 Before we describe the inference results, we discuss the
22 characteristics of the simulated data. First, we inspect
23 the effects of gene duplication and loss on the number
24 of gene copies per species in each gene family. Fig. 2(a-
25 b) and Fig. S1(a-b) show data on the sizes (numbers of
26 copies) of gene families in the 16-taxon and 12-taxon data
27 sets, respectively, under the various settings of effective
28 population sizes and duplication and loss rates.

29 Clearly, the higher the GDL rates, the larger the
30 variance in size of gene families. The figure also shows
31 that the average size of a gene family is roughly equal
32 to the number of species, with the largest gene families
33 having 51 copies for the 16-taxon data sets, and 58
34 copies for the 12-taxon data sets. We then counted the
35 average (over the 10 data sets per setting) number of
36 gene families for each setting that have exactly one copy
per species. The results are shown in Table 1. The table

TABLE 1. The average number of gene families that fit the ONLY setting out of the 10,000 gene families for 10 data sets.

GDL rate $\frac{N_e}{\rightarrow}$	16-taxon data		12-taxon data	
	10^7	5×10^7	10^6	5×10^6
1×10^{-10}	7619	7585	4591	4584
2×10^{-10}	5794	5787	2197	2176
5×10^{-10}	2554	2538	268	266

37 shows that as the GDL rates increase, the number of
38 single-copy gene families decreases.

39 We then set out to assess the extent of incongruence
40 in the gene trees due to GDL and ILS. For every pair
41 of true species tree and true locus tree, we computed
42 the number of extra lineages (Maddison, 1997) using the
43 DeepCoalCount_tree command in PhyloNet (Than and
44 Nakhleh, 2009; Wen *et al.*, 2018) as a proxy for the
45 amount of incongruence in the data. Here, we treated all
46 gene copies from the same species as different individuals.

47 Zero extra lineages mean there is no incongruence
48 between the two trees, and the higher the value, the more
49 incongruence there is. In particular, no incongruence
50 means that all gene copies from the same species are
51 monophyletic in the locus tree, and when restricted to
52 a single arbitrary copy per species, the locus tree and
53 species tree have identical topologies.

54 Fig. 2(c-d) and Fig. S1(c-d) show data on the number
55 of extra lineages in the 16-taxon and 12-taxon data
56 sets, respectively, under the various settings of effective
57 population sizes and duplication and loss rates. It is
58 important to note that all incongruence in this case
59 is exclusively due to GDL (ILS is not a factor in the
60 results in these two panels). The panels do not have
61 results for the GDL rate of 0x, because in such cases
62 there is no incongruence at all between the locus tree
63 and the species tree, and thus there are zero extra
64 lineages. The results show that, unsurprisingly, there is
65 much more incongruence for the ALL scenario than the
66 ONE scenario. For the ONLY scenario, there is very
67 little incongruence in the data, and such incongruence
68 is relatively more noticeable in the 12-taxon data set, as
69 all nonzero levels of GDL rates have nonzero number
70 of extra lineages. In contrast, for the 16-taxon data
71 set only the largest GDL rate causes incongruence
72 between locus tree and species tree. In both data sets,
73 the incongruence indicates the phenomenon of hidden
74 paralogy: single-copy genes are paralogs, and their gene
75 trees therefore do not always agree with the species tree
76 in terms of topology.

77 Finally, we computed the number of extra lineages
78 when reconciling the true gene trees with the true locus
79 trees. Here, incongruence is exclusively due to ILS.
80 Fig. 2(e-f) and Fig. S1(e-f) show data on the number
81 of extra lineages in the 16-taxon and 12-taxon data
82 sets, respectively, under the various settings of effective
83 population sizes and duplication and loss rates. When
84 the gene tree topology is identical to the locus tree
85 topology, the number of extra lineages is 0, and the larger
86 the number of extra lineages, the more ILS in the data.
87 The figure shows that, as expected, the amount of ILS
88 is larger for larger population sizes. Furthermore, there
89 is much more ILS in the 16-taxon data set than in the
90 12-taxon data set. One other trend to observe is that, on
91 average, the amount of incongruence due to ILS increases
92 with the increase in the GDL rate. This is a reflection
93 of the fact that for higher GDL rates, the locus trees
94 are larger (more leaves and internal branches) and this
95 naturally results in more ILS.

96 Results on Simulated Data

97 We are now in position to describe the inference
98 results. We show figures for the 16-taxon data sets in
99 the main text, while figures for the 12-taxon data sets
100 Figs. S8 to S11 are all in the Supplementary Materials.
101 The results for the 12-taxon data sets are consistently
102 better in terms of accuracy, so we chose to focus here on
103 the less-optimal results.

104

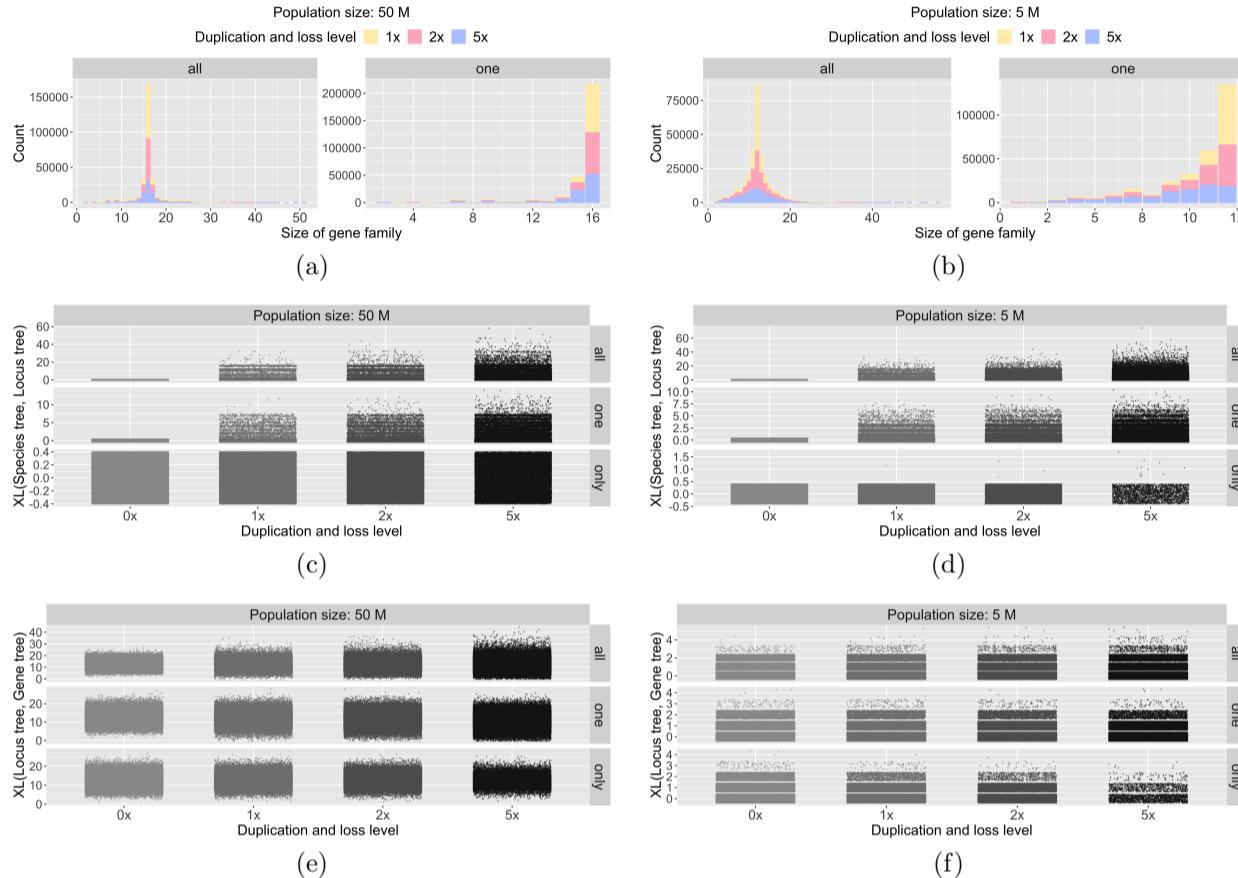


FIGURE 2. Characteristics of the simulated data under different settings of the duplication/loss rates and effective population sizes using all 10 replicates. The duplication/loss rates are denoted by the rate multiplier (0x, 1x, 2x and 5x), where 1x is the rate used by (Rasmussen and Kellis, 2011, 2012; Zhang and Wu, 2017). (a-b) Distribution of the total number of gene copies in individual gene families in the 16-taxon and 12-taxon data sets, respectively. (c-d) Scatter plots of XL(Species tree, Locus tree), the number of extra lineages when reconciling the true locus trees with the true species tree, for the 16-taxon and 12-taxon data sets, respectively. (e-f) Scatter plots of XL(Locus tree, Gene tree), the number of extra lineages when reconciling the true gene trees with the true locus tree, for the 16-taxon and 12-taxon data sets, respectively.

1 We first ran the three inference methods ASTRAL,
 2 **InferNetwork_MPL**, and **NJ_{st}**, on the true gene trees
 3 for all three input scenarios: ALL, ONE, and ONLY.
 4 In this case, gene tree estimation error is not a cause
 5 of gene tree incongruence. Instead, all incongruence is
 6 due to a combination of ILS and GDL. Results on the
 7 full 16-taxon tree are shown in Fig. 3 and Fig. S4. Note
 8 that in all cases, using input data with GDL levels of 0
 9 amounts to inferring a species tree from gene trees whose
 10 incongruence is solely due to ILS.

11 There are several observations based on the results.
 12 First, the accuracy of the inferred 16-taxon trees is
 13 much lower in general than that of the inferred 12-
 14 taxon trees. In particular, for the 12-taxon data sets,
 15 the species trees are perfectly estimated in almost all
 16 cases (Fig. S3), whereas the species tree estimation
 17 error is high, especially for the larger population sizes,
 18 for the 16-taxon data sets. As shown in Fig. 2 and
 19 Fig. S1, both data sets have similar gene family sizes,
 20 but differ significantly in terms of the amount of ILS in
 21 the data, with the 12-taxon data sets having very little

22 ILS. Therefore, the straightforward explanation for the
 23 observed differences in the RF distances between the 16-
 24 and 12-taxon data sets is the higher level of ILS in the
 25 former. Given that the level of incongruence due to GDL
 26 is similar between the 16-taxon and 12-taxon data sets
 27 (Fig. 2(c-d) and Fig. S1(c-d)), these results point to the
 28 larger role that ILS plays in the methods' performances
 29 than GDL does.

30 Second, in the case of the 16-taxon data, the
 31 performance of all three methods improves as the number
 32 of gene families used as input to the method increases.
 33 Note also that the largest dataset used here consists
 34 of only 250 gene trees, which is much smaller than
 35 the number available in most phylogenomic data sets.
 36 While there is very little difference observed in the
 37 performance among the three methods on the 16-taxon
 38 data, ASTRAL and **NJ_{st}** are more similar to each other
 39 in terms of performance than either of them is to
 40 inference under maximum pseudo-likelihood. This makes
 41 sense as both ASTRAL and **NJ_{st}** are summary methods
 42 that make inference based on statistics derived from the

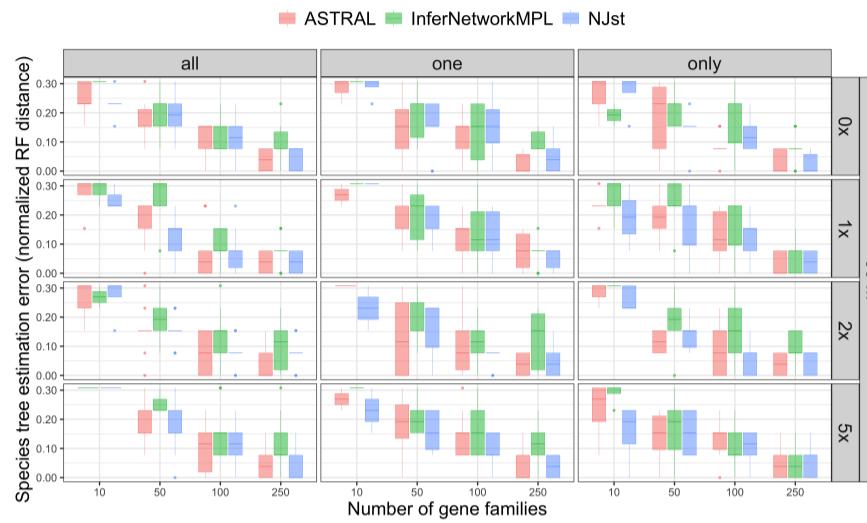


FIGURE 3. The normalized RF distances between the true species tree and the ones inferred by ASTRAL, InferNetwork_MPL, and NJ_{st} from true gene trees of the 16-taxon simulated yeast gene families under population size 5.0×10^7 and different GDL rates. The duplication/loss rates are denoted by the rate multiplier (0x, 1x, 2x and 5x), where 1x is the rate used by (Rasmussen and Kellis, 2011, 2012; Zhang and Wu, 2017). Each row corresponds to a combination of population size and GDL rates. The X axis in each panel represents the number of gene families used and the Y axis represents the normalized RF distance.

1 input gene trees, whereas maximum pseudo-likelihood
2 uses calculations based on the multispecies coalescent
3 directly.

4 Third, the level of ILS for a population size of 50M is
5 higher than for a population size of 10M, and this results
6 in lower accuracy of inferred species trees by all three
7 methods in the former case. This behavior is expected
8 for any MSC-based method, regardless of whether GDL
9 is acting.

10 Last, but not least, we observe very little difference
11 in the accuracy of inferred species trees across the three
12 input scenarios: ALL, ONE, and ONLY. This implies
13 that the presence of paralogs in the data, no matter how
14 they are treated, does not have much of an effect on the
15 performance of the three methods. In the case of the 12-
16 taxon data, all three methods make perfect inferences in
17 almost all cases (see Supplementary Materials Figs. S8
18 and S10).

19 These results raise the important question: Does GDL
20 have any effect on the performance of these three
21 methods? To answer this question, we ran ASTRAL,
22 InferNetwork_MPL, and NJ_{st} on the locus trees as input
23 to infer species trees. By the three-tree model, this
24 amounts to feeding these methods “gene trees” whose
25 incongruence is solely due to GDL; that is, ILS plays
26 no role in incongruence here. It is important to point
27 out here that locus trees are mathematical constructs of
28 the three-tree model; in practice, inferring a locus tree is
29 not possible, unless the data has no ILS. We conducted
30 this experiment to study the performance of methods
31 when GDL, but not ILS, causes all incongruence. Results
32 on the full 16-taxon data sets are shown in Fig. 4 and
33 Fig. S5. As the results now show, all three methods infer
34 the true species tree on almost all data sets regardless of
35 the parameter settings and the input scenario. In other

36 words, when these three methods—which have been
37 developed based on the multispecies coalescent directly
38 (InferNetwork_MPL) or “inspired” by it (ASTRAL and
39 NJ_{st})—are applied to data that have no ILS but do have
40 paralogs in them, they have almost perfect accuracy in
41 terms of the species tree topology they infer, under the
42 conditions of our simulations. Combined with the results
43 summarized in Fig. 3 and Fig. S4, these results show,
44 perhaps surprisingly, that methods developed to handle
45 ILS but not GDL do much better in handling GDL than
46 they do in handling ILS.

47 In practice, gene trees are unknown and inferred from
48 the sequence data. To simulate more realistic data,
49 we inferred gene trees and locus trees from simulated
50 sequence data and fed these tree estimates as input to
51 the three methods. In this case, gene tree error is a
52 factor in the observed incongruences. Fig. S2 shows the
53 extent of error of in the estimated gene and locus trees
54 for the 16-taxon data. The gene tree estimation error is
55 measured by the normalized RF distance between the
56 true gene tree and the reconstructed gene tree. For the
57 12-taxon data set, the average gene tree estimation error
58 ranges from 0.581 to 0.623, whereas the average locus
59 tree estimation error is slightly lower, ranging from 0.57
60 to 0.612 (Fig. S3). For the 16-taxon data set, the average
61 gene tree estimation error ranges between 0.101 to 0.106
62 while the average locus tree estimation error ranges from
63 0.0792 to 0.0848. In other words, there is much less gene
64 tree estimation error in the 16-taxon data sets than in
65 the 12-taxon data sets.

66 Results of species tree inference using the full 16-
67 taxon data set based on estimated gene trees are shown
68 in Fig. 5 and Fig. S6 and based on the locus tree
69 estimates are shown in Fig. 6 and Fig. S7. These results
70 should be contrasted with Fig. 3, Fig. S4, Fig. 4 and

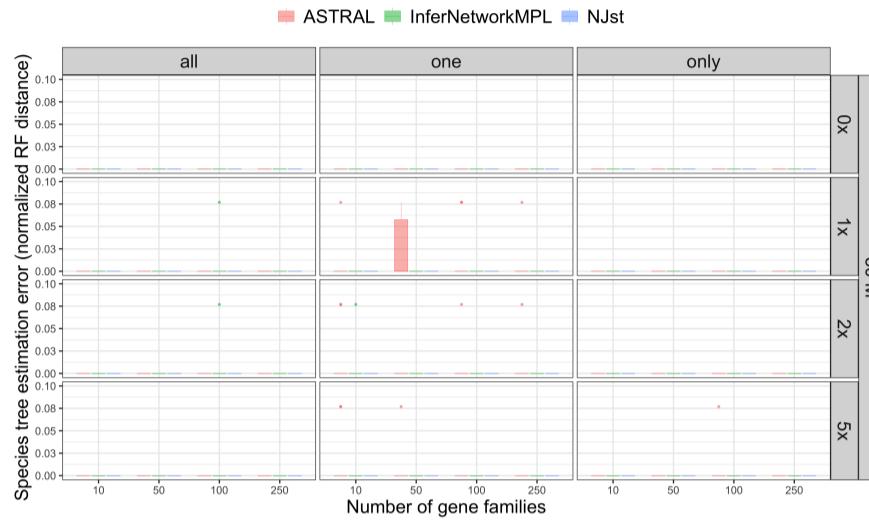


FIGURE 4. The normalized RF distances between the true species tree and the ones inferred by ASTRAL, InferNetwork.MPL, and NJ_{st} from true locus trees of the 16-taxon simulated yeast gene families under population size 5.0×10^7 and different GDL rates. The duplication/loss rates are denoted by the rate multiplier (0x, 1x, 2x and 5x), where 1x is the rate used by (Rasmussen and Kellis, 2011, 2012; Zhang and Wu, 2017). Each row corresponds to a combination of population size and GDL rates. The X axis in each panel represents the number of gene families used and the Y axis represents the normalized RF distance.

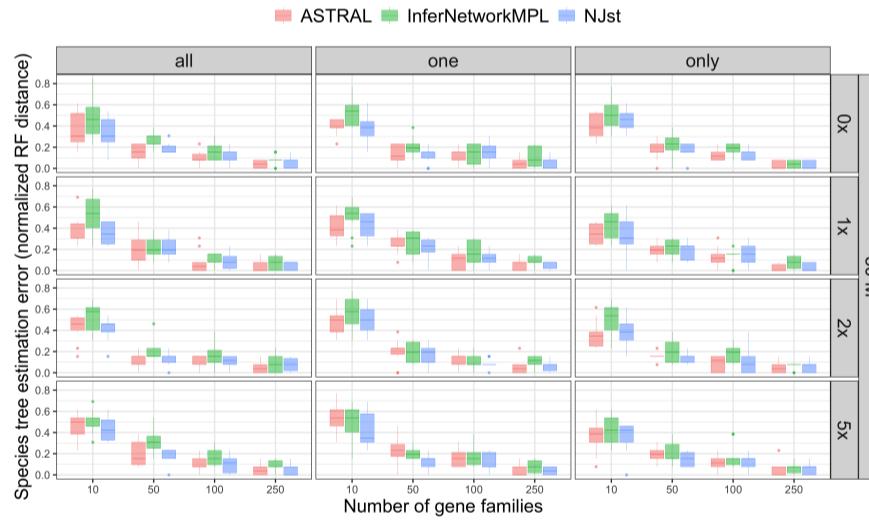


FIGURE 5. The normalized RF distances between the true species tree and the ones inferred by ASTRAL, InferNetwork.MPL, and NJ_{st} from estimated gene trees of the 16-taxon simulated yeast gene families under population size 5.0×10^7 and different GDL rates. The duplication/loss rates are denoted by the rate multiplier (0x, 1x, 2x and 5x), where 1x is the rate used by (Rasmussen and Kellis, 2011, 2012; Zhang and Wu, 2017). Each row corresponds to a combination of population size and GDL rates. The X axis in each panel represents the number of gene families used and the Y axis represents the normalized RF distance.

1 Fig. S5, respectively, to understand the effect of gene tree
2 estimation error on the accuracy of the three species tree
3 inference methods.

4 In the case of inferences using gene tree estimates
5 where ILS, GDL, and gene tree estimation error are
6 involved, the error rates of all three species tree inference
7 methods went up, as expected, due to gene tree error
8 (Fig. 5 and Fig. S6), but only slightly. The accuracy
9 of the species trees improves as the number of gene
10 families increases. As we discussed above, the error in
11 gene tree estimates in the 16-taxon data sets is very low.

12 Since gene tree estimation error in the 12-taxon data
13 sets is much higher (higher substitution rates resulting
14 in noisier sequence data), we now observe a larger impact
15 on this error on the performance of methods on the
16 12-taxon data sets. While the methods had an almost
17 perfect accuracy on true gene trees, species tree estimates
18 now have as high as 50% error when 10 gene family trees
19 are used, and close to 25% error when 250 gene family
20 trees are used (Fig. S10). These results illustrate the big
21 impact gene tree estimation error has on these methods.
22 In the case of the 12-taxon data sets, the impact of gene

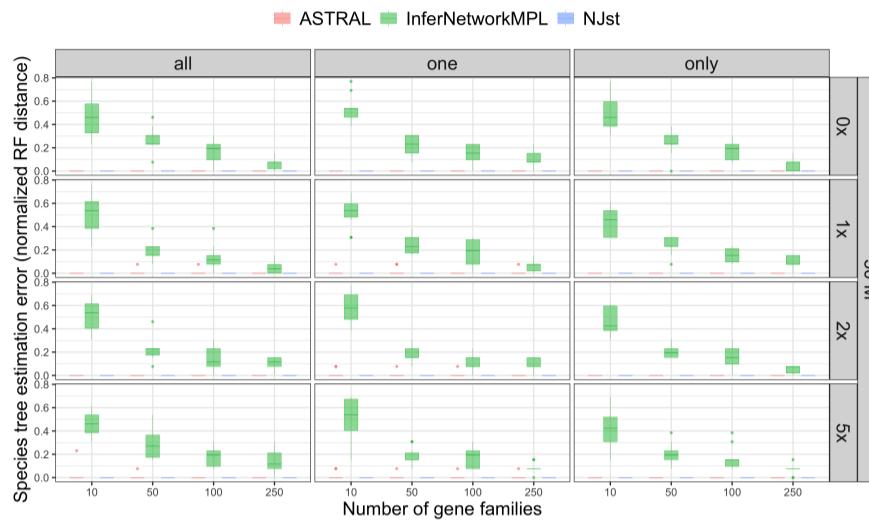


FIGURE 6. The normalized RF distances between the true species tree and the ones inferred by ASTRAL, InferNetwork_MPL, and NJ_{st} from estimated locus trees of the 16-taxon simulated yeast gene families under population size 5.0×10^7 different GDL rates. Each row corresponds to a combination of population size and GDL rates. The X axis in each panel represents the number of gene families used and the Y axis represents the normalized RF distance.

1 tree estimation error significantly outweighs that of ILS
2 and/or GDL.

3 Fig. 6 demonstrates how GDL and gene tree inference
4 error (but no ILS) impact species tree inference. As
5 Fig. 4 and Fig. S5 show almost perfect performance
6 of species tree inference from true locus trees (i.e.,
7 GDL and no ILS), all the error observed in Fig. 6 and
8 Fig. S7 can be ascribed to error in the gene trees. The
9 results demonstrate that in the absence of ILS, both
10 ASTRAL and NJ_{st} are robust to gene tree estimation
11 error, whereas InferNetwork_MPL is very sensitive to
12 gene tree estimation error. In the case of the 12-taxon
13 data sets, where locus tree estimation errors are much
14 higher, the three species tree inference methods have
15 comparable accuracies (Fig. S11). It is worth mentioning
16 that as InferNetwork_MPL requires rooted gene trees,
17 poorer performance of this method could be due to the
18 incorrect rooting of the estimated gene trees. We rooted
19 the gene tree estimates by minimizing the number of
20 extra lineages when reconciling the gene/locus trees with
21 the true species tree (according to the formulation given
22 in (Yu *et al.*, 2011a, b)). This further highlights the
23 challenge with identifying the root of gene trees when
24 GDL occurs.

25 All of these results combined point to a very small
26 impact of GDL on the performance of the three studied
27 species tree inference methods, regardless of how the
28 paralogs are handled. Across all data sets it was evident
29 that gene tree estimation error has noticeable impact on
30 the methods' performances.

31 *Results on Biological Data*

32 We ran the three methods, InferNetwork_MPL,
33 ASTRAL, and NJ_{st}, on two biological data sets. Again,

each data set consists of only 250 gene trees, although
many more trees are available for each.

34
35 For the 16 fungi genome empirical data set, NJ_{st}
36 inferred an identical tree to that of Fig. 1(a) under
37 all three input scenarios. Assuming this tree is the
38 “true” tree, NJ_{st} has performed remarkably well. For the
39 other two methods, for all three sampling schemes, the
40 inferred tree differs from the tree shown in Fig. 1(a).
41 In particular, the positions of *Kluyveromyces waltii* and
42 *Kluyveromyces laticis* have been switched, as have the
43 positions of *Candida glabrata* and *Saccharomyces castellii*
44 (Fig. 7(a)).

45 For the 12 fly genome empirical data set,
46 ASTRAL(ONLY) inferred the exact same tree as
47 the species tree shown in Fig. 1(b). The trees inferred
48 by ASTRAL(ALL) and ASTRAL(ONE) differed from
49 the true tree of Fig. 1(b) in terms of the placement of
50 *Drosophila melanogaster*, as shown in Fig. 7(b) and Fig.
51 7(c), respectively. InferNetwork_MPL and NJ_{st} inferred
52 an identical tree to that of Fig. 1(b) under all three
53 input scenarios.

54 DISCUSSION

55 As phylogenomic data sets grow, our ability to use
56 them within the bounds of current analysis paradigms
57 shrinks. One of the main problems is the decreasing
58 number of gene families that are single-copy as the
59 number of sampled species increases (Emms and Kelly,
60 2018). Because most current phylogenetic methods
61 assume that only single-copy orthologs are being used,
62 this restriction means that such methods cannot be used
63 for data sets with even several dozen taxa without severe
64 downsampling or other *ad hoc* solutions (e.g., (Thomas
65 *et al.*, 2020)). Here, we set out to ask whether MSC-based
66

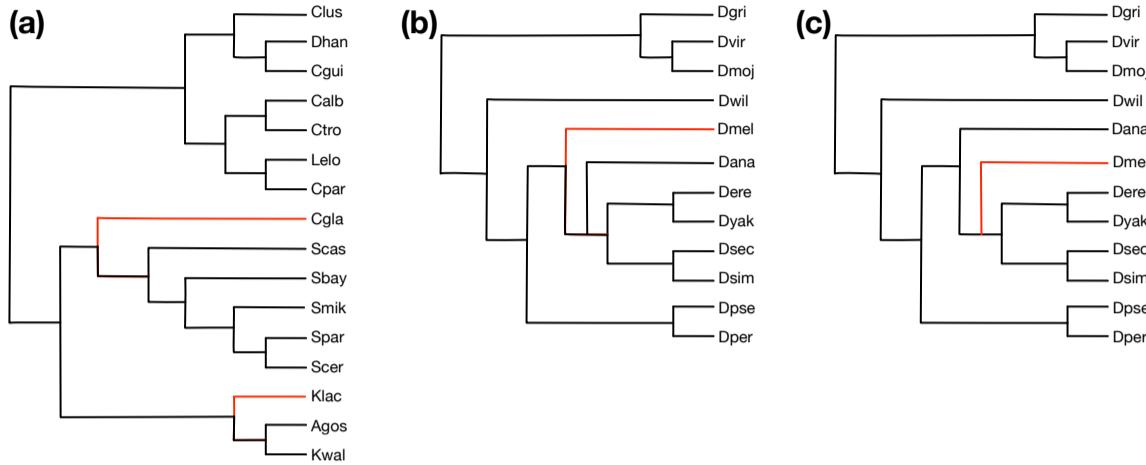


FIGURE 7. Inferred yeast and fly species trees. (a) The yeast species tree inferred by `InferNetwork_MPL` and ASTRAL. (b) The fly species tree inferred by ASTRAL(ALL). (c) The fly species tree inferred by ASTRAL(ONE). Differences between the inferred species trees and the trees of Fig. 1 are highlighted in red.

1 phylogenomic methods can be applied to data containing
2 both orthologs and paralogs.

3 On simulated data sets where only ILS acted on gene
4 families, and GDL was not a factor, all methods had the
5 expected performance: accurate species tree estimates
6 that improved as the number of gene family trees
7 used increases. In the case where the level of ILS was
8 very low (the 12-taxon data), the methods had perfect
9 performance under almost all conditions, regardless of
10 the number of gene trees used.

11 In the cases where both ILS and GDL acted on
12 gene families, the performance of the three methods
13 was hardly affected. These results imply that running
14 these methods on data with paralogs—either by treating
15 them as alleles from different individuals or by sampling
16 a single copy at random—is a safe practice when the
17 desired outcome is the species tree topology. This is
18 especially important for data sets with a large number
19 of species or high GDL rates.

20 When the methods were run on the locus tree data,
21 where ILS did not play a role and the data consisted
22 of many gene families with multiple copies, the methods
23 obtained very accurate species trees. This was true even
24 when the data consisted of only 10 gene family trees. This
25 further demonstrates that GDL has very little effect on
26 the performance of the three methods.

27 While at first it may be surprising that these
28 methods have done very well in terms of accuracy, the
29 majority of signal in any input gene tree reflects species
30 relationships. Gene duplication—if random across the
31 species tree—simply adds noise to the data, while
32 at the same time often doubling the amount of
33 information on the relationships among species carrying
34 an extra gene-copy. Similarly, gene loss does not
35 positively mislead these methods, leading to accurate
36 reconstructions of the species tree. Nevertheless, upon
37 close inspection, some of these results are not intuitive,
38 especially for the maximum pseudo-likelihood inference.
39 `InferNetwork_MPL` makes direct use of the MSC, whose

40 assumptions are clearly violated in all data sets except
41 when the GDL rates are set to 0, whereas ASTRAL
42 and NJ_{st} are summary methods that make no direct
43 use of the MSC. Consequently, one would have expected
44 that `InferNetwork_MPL` would be very sensitive to the
45 presence of paralogs in the data, while ASTRAL and
46 NJ_{st} less so. However, we did not observe any such
47 behavior of the methods.

48 In practice, gene trees are estimated from sequence
49 data and can be erroneous. Error in the gene tree
50 estimates, rather than ILS, could explain much of the
51 heterogeneity observed in phylogenomic analyses,
52 especially at deeper nodes in a species tree (Scornavacca
53 and Galtier, 2017). We showed the gene tree
54 estimation error can indeed impact species tree inference
55 significantly, and the level of its impact is similar to that of
56 ILS, if not larger.

57 In analyses of two biological data sets where a species
58 tree has been inferred using hundreds or thousands of
59 loci, we found high accuracy of methods using paralogs
60 with only 250 gene trees. All methods accurately inferred
61 the published fly species tree, except for ASTRAL(ALL)
62 and ASTRAL(ONE), which only differed in terms of
63 placing one species. For the yeast species tree, NJ_{st}
64 inferred the published species tree, while maximum
65 pseudo-likelihood and ASTRAL inferred trees that
66 differed from the true one in only two groupings of taxa
67 (around two very short branches in the species tree).

68 All these results point to a clear message: Under the
69 conditions of our simulations and on the two biological
70 data sets we used, running MSC-based or MSC-inspired
71 species tree inference on gene trees with paralogs yields
72 very accurate results. This conclusion is powerful for
73 at least two reasons. First, it implies that orthology
74 assignment and paralogy removal are not necessary
75 for running MSC-based species tree inference; simply
76 treating all copies as different individuals or randomly
77 selecting a single copy would yield very accurate species
78 tree topologies. Second, in many practical cases, too few

single-copy gene families are available to ensure good performance of species tree inference from those data alone. In these cases, our results suggest a ready source of more phylogenetic signal. Summary methods that do not explicitly use the MSC model (i.e., ASTRAL and NJ_{st}) are expected to be more robust in the presence of GDL than methods that explicitly use the model.

While this study focused on the accuracy of the inferred species tree topology, treating paralogs as orthologs from multiple individuals would have impact on the estimated branch lengths of the species tree. In particular, under the ALL setting, there could be much more incongruence due to the large number of lineages, and, consequently, MSC-based methods would estimate the branch lengths to be shorter than they truly are. For this reason, branch lengths inferred by such methods should not be used. Instead, an alternative approach is needed.

The results of our analyses point to the following potential approach for inferring accurate species trees (topologies and branch lengths) by utilizing as much of the phylogenomic data as possible:

1. Use all available gene family trees as input, either treating all copies from the same species within a family as multiple individuals from the same species, or while randomly selecting one copy per species.
2. Feed all gene trees to an MSC-based method and obtain a species tree topology.
3. Using a smaller subset of truly single-copy genes, and fixing the species tree topology obtained from Step (2), optimize the branch lengths of the species tree.

For Steps (1) and (2), one option is to infer a species tree based on ALL, and repeat the random sampling of single copies to obtain multiple versions of the input under ONE, and then score all the inferred species trees under some criterion that combines the MSC with a model of gene duplication/loss. This overcomes the issue of fixing a single species tree as input to Step (3), and avoids searching the species tree space while evaluating a likelihood function that is very complex and computationally very demanding to compute. As an alternative to using only single-copy orthologs in Step (3), one could also use a statistical model that combines the MSC and GDL models (e.g., (Boussau *et al.*, 2013; Rasmussen and Kellis, 2012)). Such methods allow for paralogy detection and orthology assignment, conditional on the fixed species tree (or species trees), by using a more detailed evolutionary model and the full signal in the sequence data. For example, the orthology assignment could be “integrated out” or sampled, depending on the desired outcomes of the analysis.

CONCLUSIONS

In this paper we set out to study how MSC-based species tree inference would perform on data with paralogs. The motivation for exploring this question was two-fold. First, as MSC-based species tree inference has become commonplace and as practitioners are almost never certain that their data contain no paralogs, it is important to understand the effect of such hidden paralogy on the quality of the inference. Second, as larger phylogenomic data sets become available, insistence on single-copy genes would mean throwing away most of the data and potentially keeping a number of loci that is severely inadequate for the data requirement for MSC-based inference methods to perform well. We investigated the question through a combination of simulations and biological data analyses. Our results show that MSC-based inference is robust to the presence of paralogs in the data, at least under the simulation conditions and on the empirical data sets we investigated.

Our results highlight the issue that MSC-based inference could result in very accurate species trees even when ILS is not a factor or not the only factor. This finding implies that orthology detection and restricting data to single-copy genes as a requirement for employing MSC-based inference can be mostly eliminated, thus making use of as much of the data as possible. In particular, for very large data sets (in terms of the number of species), eliminating all but single-copy genes might leave too few loci for the species tree to be inferred accurately. Our findings show that this data exclusion could be an unnecessary practice. It is important to note however, that our results do not apply to concatenated analyses, and in such cases the presence of paralogs may indeed have a large, negative effect (Brown and Thomson, 2016). Species tree inference from a concatenation of the sequences with gene families is challenging in the presence of paralogs for at least two reasons. First, when gene families have different numbers of copies across species, the concatenated alignment will have very large gaps. Second, correct orthology detection is still required, so that orthologous gene copies are placed in correct correspondence across the multiple genomes in the concatenated alignment. This issue is very important to examine so as to avoid aligning non-orthologous sequences in the concatenated data set.

In our simulations, we generated gene families under a neutral model and with GDL rates that were the same across all families. It is well known that the functional implications of gene duplication and the ways in which they are fixed and maintained in the genome result in much more complex scenarios than those captured in our simulations (Hahn, 2009; Innan and Kondrashov, 2010). However, analyses of the two biological data sets yield results with very similar trends to those observed in our simulations. Furthermore, the assumptions underlying the MSC model are also undoubtedly violated in most empirical data sets. Another direction for future research entails characterizing the conditions under which species

55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112

1 tree inference under the MSC is robust to the presence
 2 of paralogs in the data.

3 Finally, while we did not discuss or incorporate gene
 4 flow in our study, it is possible that all three processes—
 5 ILS, GDL, and gene flow—are simultaneously involved
 6 in some phylogenomic analyses. Studies of the robustness
 7 of MSC-based species tree inference under some models
 8 of gene flow exist (Davidson *et al.*, 2015; Roch and
 9 Snir, 2012; Solís-Lemus *et al.*, 2016; Steel *et al.*, 2013;
 10 Zhu *et al.*, 2016), but, to the best of our knowledge,
 11 such studies under scenarios that incorporate all the
 12 aforementioned processes do not exist yet. It is important
 13 to highlight, as well, that great strides have been made
 14 in developing methods for phylogenetic network inference
 15 in the presence of ILS (Elworth *et al.*, 2019), but none
 16 currently incorporate gene duplication and loss. We
 17 believe methods along the lines described in the previous
 18 section could be promising for accurate and scalable
 19 phylogenomic inferences without sacrificing much of the
 20 data.

21 SUPPLEMENTARY MATERIAL

22 Supplementary material, including data files and
 23 online-only appendices, can be found in the Dryad data
 24 repository at .

25 FUNDING

26 Funding was provided by the National Science
 27 Foundation (DBI-1355998, CCF-1514177, CCF-1800723,
 28 and DMS-1547433) to L.N. and (DBI-1564611 and DEB-
 29 1936187) to M.W.H.

30 REFERENCES

- 31 Arvestad, L., Lagergren, J., and Sennblad, B. 2009. The gene
 32 evolution model and computing its associated probabilities.
Journal of the ACM (JACM), 56(2): 7.
- 33 Boussau, B., Szöllösi, G. J., Duret, L., Gouy, M., Tannier, E., and
 34 Daubin, V. 2013. Genome-scale coestimation of species and gene
 35 trees. *Genome research*, 23(2): 323–330.
- 36 Brown, J. M. and Thomson, R. C. 2016. Bayes factors unmask
 37 highly variable information content, bias, and extreme influence
 38 in phylogenomic analyses. *Systematic Biology*, 66(4): 517–530.
- 39 Daubin, V. and Gouy, M. 2001. Bacterial molecular phylogeny
 40 using supertree approach. *Genome Informatics*, 12: 155–164.
- 41 Davidson, R., Vachaspati, P., Mirarab, S., and Warnow, T.
 42 2015. Phylogenomic species tree estimation in the presence of
 43 incomplete lineage sorting and horizontal gene transfer. *BMC
 44 Genomics*, 16(10): S1.
- 45 Degnan, J. H. and Rosenberg, N. A. 2009. Gene tree discordance,
 46 phylogenetic inference and the multispecies coalescent. *Trends
 47 in Ecology & Evolution*, 24(6): 332–340.
- 48 Du, P. and Nakhleh, L. 2018. Species tree and reconciliation
 49 estimation under a duplication-loss-coalescence model.
*Proceedings of the 9th ACM Conference on Bioinformatics,
 50 Computational Biology, and Health Informatics*, pages 376–385.
- 51 Elworth, R. L., Ogilvie, H. A., Zhu, J., and Nakhleh, L. 2019.
 52 Advances in computational methods for phylogenetic networks
 53 in the presence of hybridization. In T. Warnow, editor,
Bioinformatics and Phylogenetics, pages 317–360. Springer.
- 54 Emms, D. and Kelly, S. 2018. STAG: Species tree inference from
 55 all genes. *bioRxiv*, page 267914.
- 56 Gribaldo, S. and Philippe, H. 2002. Ancient phylogenetic
 57 relationships. *Theoretical Population Biology*, 61(4): 391–408.
- 58 Guindon, S. and Gascuel, O. 2003. A simple, fast, and accurate
 59 algorithm to estimate large phylogenies by maximum likelihood.
Systematic biology, 52(5): 696–704.
- 60 Hahn, M. W. 2009. Distinguishing among evolutionary models for
 61 the maintenance of gene duplicates. *Journal of Heredity*, 100(5):
 62 605–617.
- 63 Hahn, M. W., Han, M. V., and Han, S.-G. 2007. Gene family
 64 evolution across 12 *drosophila* genomes. *PLoS genetics*, 3(11):
 65 e197.
- 66 Hudson, R. R. 1983. Testing the constant-rate neutral allele model
 67 with protein sequence data. *Evolution*, 37(1): 203–217.
- 68 Innan, H. and Kondrashov, F. 2010. The evolution of gene
 69 duplications: classifying and distinguishing between models.
Nature Reviews Genetics, 11(2): 97.
- 70 Knowles, L. L. and Kubatko, L. S. 2011. *Estimating species trees:
 71 practical and theoretical aspects*. John Wiley and Sons.
- 72 Koonin, E. V. 2005. Orthologs, paralogs, and evolutionary
 73 genomics. *Annu. Rev. Genet.*, 39: 309–338.
- 74 Lang, G. I. and Murray, A. W. 2008. Estimating the per-base-pair
 75 mutation rate in the yeast *saccharomyces cerevisiae*. *Genetics*,
 76 178(1): 67–82.
- 77 Liu, L. and Yu, L. 2011. Estimating species trees from unrooted
 78 gene trees. *Systematic Biology*, 60(5): 661–667.
- 79 Liu, L., Yu, L. L., Kubatko, L., Pearl, D. K., and Edwards, S. V.
 80 2009. Coalescent methods for estimating phylogenetic trees.
Mol. Phylogenet. Evol., 53: 320–328.
- 81 Liu, L., Yu, L., and Edwards, S. V. 2010. A maximum pseudo-
 82 likelihood approach for estimating species trees under the
 83 coalescent model. *BMC evolutionary biology*, 10(1): 302.
- 84 Liu, L., Xi, Z., Wu, S., Davis, C. C., and Edwards, S. V. 2015.
 85 Estimating phylogenetic trees from genome-scale data. *Annals
 86 of the New York Academy of Sciences*, 1360(1): 36–53.
- 87 Maddison, W. P. 1997. Gene trees in species trees. *Systematic
 88 Biology*, 46(3): 523–536.
- 89 Mallo, D., de Oliveira Martins, L., and Posada, D. 2015. Simphy:
 90 phylogenomic simulation of gene, locus, and species trees.
Systematic Biology, 65(2): 334–344.
- 91 Nakhleh, L. 2013. Computational approaches to species phylogeny
 92 inference and gene tree reconciliation. *Trends in Ecology and
 93 Evolution*, 28(12): 719–728.
- 94 Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., and Minh, B. Q.
 95 2014. IQ-TREE: a fast and effective stochastic algorithm for
 96 estimating maximum-likelihood phylogenies. *Molecular biology
 97 and evolution*, 32(1): 268–274.
- 98 Pollard, D. A., Iyer, V. N., Moses, A. M., and Eisen, M. B.
 99 2006. Widespread discordance of gene trees with species tree
 100 in *drosophila*: evidence for incomplete lineage sorting. *PLoS
 101 genetics*, 2(10): e173.
- 102 Rambaut, A. and Grassly, N. C. 1997. Seq-gen: an application
 103 for the monte carlo simulation of dna sequence evolution along
 104 phylogenetic trees. *Bioinformatics*, 13(3): 235–238.
- 105 Rannala, B. and Yang, Z. 2003. Bayes estimation of species
 106 divergence times and ancestral population sizes using dna
 107 sequences from multiple loci. *Genetics*, 164(4): 1645–1656.
- 108 Rasmussen, M. D. and Kellis, M. 2011. A Bayesian approach for
 109 fast and accurate gene tree reconstruction. *Molecular Biology
 110 and Evolution*, 28(1): 273–290.
- 111 Rasmussen, M. D. and Kellis, M. 2012. Unified modeling of gene
 112 duplication, loss, and coalescence using a locus tree. *Genome
 113 Research*, 22(4): 755–765.
- 114 Robinson, D. F. and Foulds, L. R. 1981. Comparison of
 115 phylogenetic trees. *Mathematical biosciences*, 53(1-2): 131–147.
- 116 Roch, S. and Snir, S. 2012. Recovering the tree-like trend
 117 of evolution despite extensive lateral genetic transfer: a
 118 probabilistic analysis. *Journal of Computational Biology*, 20(2):
 119 93–112.
- 120 Sawyer, S. A. and Hartl, D. L. 1992. Population genetics of
 121 polymorphism and divergence. *Genetics*, 132(4): 1161–1176.
- 122

- 1 Schrider, D. R., Houle, D., Lynch, M., and Hahn, M. W. 2013.
2 Rates and genomic consequences of spontaneous mutational
3 events in *Drosophila melanogaster*. *Genetics*, pages genetics–
4 113.
- 5 Scornavacca, C. and Galtier, N. 2017. Incomplete lineage sorting in
6 mammalian phylogenomics. *Systematic Biology*, 66(1): 112–120.
- 7 Solís-Lemus, C., Yang, M., and Ané, C. 2016. Inconsistency of
8 species tree methods under gene flow. *Systematic biology*, 65(5):
9 843–851.
- 10 Stamatakis, A. 2006. Raxml-vi-hpc: maximum likelihood-based
11 phylogenetic analyses with thousands of taxa and mixed models.
12 *Bioinformatics*, 22(21): 2688–2690.
- 13 Steel, M., Linz, S., Huson, D., and Sanderson, M. 2013. Identifying
14 a species tree subject to random lateral gene transfer. *Journal
15 of Theoretical Biology*, 322: 81–93.
- 16 Takahata, N. 1989. Gene genealogy in three related populations:
17 consistency probability between gene and population trees.
18 *Genetics*, 122(4): 957–966.
- 19 Than, C. and Nakhleh, L. 2009. Species tree inference by
20 minimizing deep coalescences. *PLoS Computational Biology*,
21 5(9): e1000501.
- 22 Thomas, G. W., Dohmen, E., Hughes, D. S., Murali, S. C.,
23 Poelchau, M., Glastad, K., Anstead, C. A., Ayoub, N. A.,
24 Batterham, P., Bellair, M., *et al.* 2020. Gene content evolution
25 in the arthropods. *Genome Biology*, 21(1): 1–14.
- 26 Wen, D., Yun, Y., Zhu, J., and Nakhleh, L. 2018. Inferring
27 phylogenetic networks using phylogenetic networks. *Systematic Biology*,
28 67(4): 735–740.
- 29 Yu, Y. and Nakhleh, L. 2015. A maximum pseudo-likelihood
30 approach for phylogenetic networks. *BMC Genomics*, 16(10):
31 S10.
- 32 Yu, Y., Warnow, T., and Nakhleh, L. 2011a. Algorithms for
33 MDC-based multi-locus phylogeny inference. In *International
34 Conference on Research in Computational Molecular Biology*,
35 pages 531–545. Springer.
- 36 Yu, Y., Warnow, T., and Nakhleh, L. 2011b. Algorithms for MDC-
37 based multi-locus phylogeny inference: beyond rooted binary
38 gene trees on single alleles. *Journal of Computational Biology*,
39 18(11): 1543–1559.
- 40 Yu, Y., Dong, J., Liu, K. J., and Nakhleh, L. 2014. Maximum
41 likelihood inference of reticulate evolutionary histories. *Proceedings of the National Academy of Sciences*,
42 111(46): 16448–16453.
- 44 Zhang, B. and Wu, Y.-C. 2017. Coestimation of gene trees
45 and reconciliations under a duplication-loss-coalescence model.
46 In *International Symposium on Bioinformatics Research and
47 Applications*, pages 196–210. Springer.
- 48 Zhang, C., Rabiee, M., Sayyari, E., and Mirarab, S. 2018. Astral-
49 iii: polynomial time species tree reconstruction from partially
50 resolved gene trees. *BMC bioinformatics*, 19(6): 153.
- 51 Zhu, J., Yu, Y., and Nakhleh, L. 2016. In the light of
52 deep coalescence: revisiting trees within networks. *BMC
53 Bioinformatics*, 17(14): 415.