

1 **Title:** New approaches for inferring phylogenies in the presence of paralogs

2

3 **Authors:** Megan L. Smith<sup>1\*</sup> and Matthew W. Hahn<sup>1</sup>

4

5 Department of Biology and Department of Computer Science, Indiana University, Bloomington,  
6 IN 47405, USA

7

8 \*Correspondence: [mls16@indiana.edu](mailto:mls16@indiana.edu) (M.L. Smith).

9

10

11

12

13 **Keywords:** phylogenomics, duplication, incomplete lineage sorting, polyploidy

14

15 **Abstract**

16

17 The availability of whole genome sequences was expected to supply essentially unlimited data  
18 for phylogenetics. However, strict reliance on single-copy genes for this purpose has drastically  
19 limited the amount of data that can be used. Here, we review several approaches for increasing  
20 the amount of data used for phylogenetic inference, focusing on methods that allow for the  
21 inclusion of duplicated genes (paralogs). Recently developed methods that are robust to high  
22 levels of incomplete lineage sorting also appear to be robust to the inclusion of paralogs,  
23 suggesting a promising way to take full advantage of genomic data. We discuss the pitfalls of  
24 these approaches, as well as further avenues for research.

## 25 **The search for orthologs**

26

27 The business of phylogeny-building has been transformed by the availability of whole genome  
28 sequences (reviewed in [1]). Indeed, the promise of “phylogenomics” was access to many  
29 thousands of loci [2]. However, the data requirements of most phylogenetic inference methods—  
30 single-copy genes present in almost all species sampled (Figure 1A)—have meant that a growing  
31 number of phylogenomic studies have actually used very small amounts of data. For instance, in  
32 their dataset of 76 arthropod genomes, Thomas et al. [3] found *no* genes that were single-copy  
33 and present in all species. This study is not unique: even with whole-genome data, as the number  
34 of species sampled goes up, the number of single-copy genes found in all taxa goes down [4].

35

36 Phylogeny estimation has long relied on the identification of single-copy orthologous genes,  
37 filtering out paralogous genes found in multiple copies in one or more species (Box 1). Indeed,  
38 when Fitch [5] introduced the terms ortholog and paralog it was in the context of species  
39 phylogeny estimation: “Phylogenies require orthologous, not paralogous, genes.” This sentiment  
40 is echoed repeatedly in the literature [6,7], based on the belief that, since orthologous genes are  
41 related by speciation events alone, their relationships should more accurately reflect the species  
42 phylogeny. Similar claims are made about the privileged use of orthologs in protein-function  
43 prediction [8–10].

44

45 However, accurate methods for inferring species trees using both orthologs and paralogs were  
46 proposed more than 40 years ago [11], and efficient software implementing these approaches has  
47 been around for at least 20 years [12]. Methods using orthologs and paralogs work because gene  
48 trees containing duplication events also include all of the speciation events that follow (Figure  
49 1B). While each duplication event does add a branch not found in the species tree, it also doubles  
50 the amount of information contained about subsequent speciation events. Most significantly,  
51 recent methods developed for phylogeny inference using orthologs [13,14] turn out to be highly  
52 accurate and extremely efficient when applied to datasets including paralogs. Although the  
53 application of these approaches to such datasets is just beginning, their promise for  
54 phylogenomics is clear.

55

56 In this review, we discuss ways to combat the limitation of single-copy orthologs by increasing  
57 the amount of data that can be used in phylogenomics, while still maintaining a high degree of  
58 accuracy. We first discuss the problem of gene tree heterogeneity, and how it affects the  
59 accuracy of species trees. Next, we review two broad approaches for increasing the amount of  
60 data used in phylogenomic inference: one that still includes only orthologs and one that includes  
61 both orthologs and paralogs. We also describe the newly developed phylogenetic methods that  
62 make both of these approaches possible. Finally, we identify some key topics to consider when  
63 inferring phylogenies in the presence of paralogs, including promising future areas of research on  
64 this topic.

65

### 66 **Gene tree heterogeneity and the problem of “hidden paralogy”**

67

68 Gene tree heterogeneity—a mismatch between the topology of a single region and the topology  
69 of a species—is now recognized as common in phylogenetics [15]. This heterogeneity may be  
70 due to a number of biological factors, including incomplete lineage sorting (ILS), introgression,  
71 and gene duplication and loss (GDL) [16], in addition to technical factors such as error in gene  
72 tree reconstruction. This heterogeneity has important consequences for species tree inference, as  
73 if it is not accounted for it can lead to an incorrect phylogeny. Methods developed to deal with  
74 multiple causes of heterogeneity can also help us to infer phylogenies from a broader set of loci.

75

76 In particular, high levels of ILS can mislead many species tree methods, whether they apply  
77 maximum likelihood methods to concatenated alignments of all loci [17] or count gene tree  
78 topologies individually [18]. Partly because of these issues, methods that account for ILS when  
79 estimating species phylogenies have proliferated [19–24]. These gene-tree-based methods  
80 usually construct a separate tree for each locus (excluding the methods in refs. [20] and [23]),  
81 combining these trees together in a principled way to infer a species tree. As with most  
82 phylogenetic approaches, these methods were designed to use datasets consisting of only single-  
83 copy orthologs, as they account only for ILS as a source of gene tree heterogeneity. Importantly,  
84 however, many of these methods also deal naturally with missing data; this will be key for  
85 several of the new approaches described below.

86

87 Gene duplication and loss leads to gene tree heterogeneity by adding duplication events to gene  
88 trees (Box 1). Such events are not expected in histories that follow only the species tree, so trees  
89 that contain more than one copy of a gene are generally removed from phylogenetic datasets.  
90 More insidiously, “hidden paralogs” [25], or “pseudo-orthologs” [26], contain only a single copy  
91 per species due to differential loss of duplicate copies across species (Figure 2) and can be  
92 mistaken for single-copy orthologs. The topologies inferred from pseudo-orthologs can differ  
93 from the species tree via a process that is rarely modeled by phylogenetic methods.

94  
95 Although they are much feared, few studies have actually evaluated the effects of including  
96 pseudo-orthologs on phylogenetic inference, and these found mixed results. Brown and Thomson  
97 [27] suggested that outlier loci supporting a contentious placement of turtles were paralogs, and  
98 that these had an extreme effect on Bayesian inference applied to a concatenated dataset. Many  
99 other studies have shown differences in the species tree inferred from datasets assembled using  
100 different orthology detection tools, differences that are possibly due to the inclusion of pseudo-  
101 orthologs [reviewed in 6]. Some of these studies found substantial differences in the inferred  
102 trees [28], while others found minimal effects [29,30].

103  
104 What is clear from the work briefly summarized here is that there are many causes of gene tree  
105 heterogeneity that have the capacity to mislead phylogenetic inference. With respect to  
106 increasing the types of loci that can be used in phylogenomics, we would like any approaches  
107 using these loci to be robust to the known problems caused by both ILS and hidden paralogy.

108

### 109 **Increasing data availability without including paralogs**

110

111 If only orthologous genes are required, there are multiple ways to increase the total number of  
112 loci used in phylogenetic inference. Below, we discuss two such approaches that can increase the  
113 amount of available data: relaxing filters for missing data (Figure 1C) and sampling lineage-  
114 specific duplicates (Figure 1D).

115

116 *Sampling single-copy orthologs with missing data*

117

118 Often, researchers require that all or most of their taxa are sampled for a locus to be included in  
119 phylogenetic inference. However, the actual effects of including missing data—i.e. loci for  
120 which no sequences exist in one or more species—remain unclear. In concatenated analyses,  
121 simulation studies have demonstrated that there are limited negative effects of missing data [31].  
122 Other studies have argued that the issue is a lack of informative data rather than missing data *per*  
123 *se* [32,33]. Many empirical studies show little effect of including missing data [34,35], and often  
124 the positive effects of including a larger number of loci or sites seem to outweigh the negative  
125 effects of missing data [36,37].

126

127 There has been a lot of recent work on the effects of missing data on gene-tree-based methods  
128 that can account for ILS [14,19,21,24]. Because these methods combine individual gene trees  
129 from each locus, they can naturally accommodate missing taxa in a subset of trees. Studies have  
130 shown that ILS methods can be robust to substantial levels of missing data, whether these are  
131 randomly or non-randomly distributed [38, 39, 40]. Note, however, that these results may break-  
132 down in cases of extreme branch lengths [41,42].

133

134 Based on these considerations, one simple way to drastically increase the amount of data that can  
135 be used for phylogenetic inference is to relax missing data thresholds. For quartet-based methods  
136 such as ASTRAL [13], the minimum number of taxa required from each locus is four (Figure  
137 1C), as a four-taxon unrooted tree is all that is needed to specify phylogenetic relationships.  
138 Empirically, results of relaxing these thresholds can be dramatic. For example, Eaton and Ree  
139 [43] found that requiring a minimum of four taxa increased the number of loci available in a  
140 group of flowering plants nearly 9-fold compared to requiring that all taxa be sampled. The  
141 relative advantage gained by using these methods can only go up as more taxa are included in a  
142 dataset, though researchers should try to ensure that species are represented approximately  
143 evenly across loci to avoid cases where most of the signal for some branches comes from a small  
144 number of genes (e.g. [44]).

145

146 *Sampling orthologs that have lineage-specific duplication events*

147

148 The requirement that only orthologs be sampled for phylogenetic inference does not mean that  
149 we must only include single-copy orthologs. Notably, there is no theoretical reason to exclude  
150 loci that have undergone lineage-specific duplications, as they can have many-to-one  
151 orthologous relationships with single-copy genes (Box 1). For example, in Figure 1D species-  
152 specific duplications have occurred in lineages **a** and **e**. Since the two copies in each species are  
153 both orthologous to the gene copies in all other lineages, if we chose a single gene from each  
154 species the resulting gene tree would include only speciation events. There can be no gene tree  
155 heterogeneity induced by such a sampling scheme, even when there are more than two copies in  
156 each species.

157

158 Surprisingly, this approach has rarely been used in phylogenetics research. The number of loci  
159 that could be included would greatly increase, but the computational burden would increase  
160 slightly, as well (Box 2). These numbers could be increased even further, too: there should be no  
161 negative effect on the inferred topology of including duplications specific to a *pair* of sister  
162 species. In other words, if one or more duplication events occur in the ancestor of a pair of  
163 species, sampling a single copy from each of these species cannot induce gene tree  
164 heterogeneity. This occurs because there is only a single way this pair can be related, and such  
165 gene tree invariance cannot be ensured for duplicates ancestral to three or more species. Though  
166 the inclusion of duplicates specific to a pair of sister species should not affect the inferred  
167 topology, it could affect estimates of terminal branch lengths (see section on “Branch lengths”  
168 below). Broadening sampling to include these genes would lead to a further increase in the  
169 number of loci available for phylogenetic inference.

170

### 171 **Estimating species trees in the presence of paralogs**

172

173 In the methods described thus far we have still limited ourselves to analyses involving only  
174 orthologous loci. If we relax this restriction even more, we can again greatly increase the number  
175 of loci to be used. Below, we review five general approaches for reconstructing species trees in  
176 the presence of paralogs. We largely go through these methods in the chronological order in  
177 which they appeared in the literature, spending the most time at the end on promising new  
178 methods.

179

180 *Gene Tree Parsimony*

181

182 The earliest methods to infer species trees in the presence of gene duplication and loss used gene  
183 tree parsimony (GTP) [11,45,46]. In these approaches the aim is to find the species tree with the  
184 minimum “reconciliation” cost [47] to a collection of input gene trees; i.e. the species tree that  
185 minimizes the distance to all gene trees. Reconciliation costs are calculated based on explicit  
186 biological causes of gene tree heterogeneity, including, but not limited to, GDL. Some software  
187 packages calculate reconciliation costs based on minimizing duplications and losses  
188 [11,46,48,49], while others focus completely on minimizing the number of differences induced  
189 by ILS [50,51], or allow users to choose among these reconciliation costs [52]. Recognizing that  
190 these processes do not act in isolation, recent approaches consider both GDL and ILS [53], with  
191 some additionally incorporating introgression [54,55]. Although these approaches appear to deal  
192 with ILS, they do not completely account for very high levels of ILS when inferring the species  
193 tree [56], and therefore may give misleading results in such cases.

194

195 *Robinson-Foulds-based methods*

196

197 The Robinson-Foulds (RF) distance between two trees measures the number of branches that  
198 must be removed, and the number of subsequent branches that must be added to make them have  
199 the same topology [57]. RF species tree methods try to find the species tree that minimizes the  
200 RF distance to a collection of input gene trees [58]. Although this is a similar approach to gene  
201 tree parsimony, RF-based approaches make no assumptions about the biological processes  
202 leading to heterogeneity between the gene trees and the species tree, and there are therefore no  
203 options to apply different costs to different processes.

204

205 Although RF-based methods as originally described were applicable only to input trees with no  
206 duplicates, interest in applying these methods to multi-copy gene trees (i.e. those with both  
207 orthologs and paralogs) led to several advancements that permitted the calculation of RF  
208 distances between them [59,60]. Chaudhary et al. [61,62] then introduced an approach for  
209 finding a species tree using multi-copy gene trees as input. Their method, MulRF, compares



210 favorably to GTP approaches [63], and has recently been improved by Molloy and Warnow [64].  
211 RF methods appear to perform well under general conditions [63, 64], though, like GTP  
212 methods, they are not accurate under high levels of ILS [65].

213

#### 214 *Probabilistic Methods*

215

216 Several probabilistic approaches have been introduced for inferring species trees in the presence  
217 of gene duplication and loss, but these are often much more computationally intensive than GTP  
218 and RF methods. For example, PHYLOGDOG jointly estimates gene family trees, species trees, and  
219 the number of duplications and losses under a model of GDL by maximizing their likelihood  
220 given a set of alignments [66]. However, PHYLOGDOG does not consider other sources of gene  
221 tree incongruence (e.g. ILS) and the computational costs are high, preventing its application to  
222 large genomic datasets [63].

223

224 De Oliveira Martins et al. [67] introduced *guenomu*, a probabilistic supertree approach to infer  
225 species trees in the presence of both ILS and GDL. *Guenomu* implements a hierarchical  
226 Bayesian model: it takes as input a posterior distribution of gene trees and uses a multivariate  
227 distance metric based on ILS and GDL to infer a posterior distribution of species trees. However,  
228 like PHYLOGDOG, *guenomu* is computationally intensive, and therefore neither approach truly  
229 expands the number of loci one could use in phylogenomics.

230

#### 231 *Methods based on Neighbor Joining (and other clustering approaches)*

232

233 Neighbor Joining (NJ; [68]) and other distance-based approaches are popular methods for  
234 species tree inference using orthologs. Newer application of these approaches can accommodate  
235 ILS by calculating a distance matrix from a collection of gene trees inferred from separate loci,  
236 and then NJ or another clustering algorithm is used to estimate a species tree from this distance  
237 matrix. Distance methods applicable to gene trees can broadly be divided into two classes: those  
238 that construct distance matrices based on sequence distances and those that construct distance  
239 matrices based on internode distances. The former approach includes the methods implemented  
240 in STEAC [69] and METAL [70]. Methods based on internode distances include STAR [69],

241  $NJ_{st}$  [14], and ASTRID [21]. Distance-based approaches have been proven to return the correct  
242 species tree under high levels of ILS [70–72].

243  
244 Extending distance methods to cases including paralogs is straightforward, because distance  
245 matrices can be calculated as averages over multiple samples from a species. Application to  
246 datasets containing orthologs and paralogs has already been done using  $NJ_{st}$  [63,73] and  
247 ASTRID [74]. STAG [4] is another distance method introduced specifically to estimate species  
248 trees from multi-copy gene trees, though it requires that loci have no missing species. Testing the  
249 accuracy of distance methods using orthologs and paralogs, Chaudhary et al. [63] found that  $NJ_{st}$   
250 was outperformed by methods based on GTP, RF distances, and probabilistic models. In  
251 contrast, Yan et al. [73] found that  $NJ_{st}$  performed comparably to quartet-based methods, and  
252 Legried et al. [74] found that ASTRID had similar or higher accuracy than all other methods  
253 evaluated. Overall, distance-based methods appear to be a generally accurate and efficient  
254 method for inferring species trees using paralogs.

#### 255 256 *Quartet-based Methods*

257  
258 Methods to build species trees from quartet sub-trees have been around for some time [75–79],  
259 but have found renewed popularity due to the introduction of more accurate, more efficient  
260 algorithms. These methods scale well to genomic datasets and are robust to both high levels of  
261 ILS [80,81], and, as mentioned earlier, large amounts of missing data. ASTRAL [13,19,80] is  
262 among the most popular of these methods: it infers a species tree from a set of input gene trees,  
263 extracting quartets from them automatically, and finding the phylogeny that maximizes the  
264 number of shared quartet trees. ASTRAL was designed for use with single-copy orthologs, but  
265 can accommodate multiple haplotypes sampled within species (ASTRAL-multi [82]). In these  
266 cases, ASTRAL-multi effectively averages over haplotypes by sampling quartets with at most  
267 one sequence per species.

268  
269 Gene trees with paralogs in them take advantage of the same sampling scheme used by  
270 ASTRAL-multi, and perform very well because the most common quartet in multi-copy gene  
271 trees is still the quartet that matches the species tree (Figure 2; [73,74]). ASTRAL-multi has

272 multiple mathematical guarantees about its accuracy in the presence of both ILS and GDL, at  
273 least under some models [74], and simulation studies have also demonstrated its accuracy  
274 [73,74]. Most recently, a version of the software explicitly built for the inclusion of paralogs,  
275 ASTRAL-Pro, outperformed ASTRAL-multi, MulRF, and GTP methods [65].

276

277 Quartet-based methods are also robust to the hidden paralog problem, as can be illustrated by an  
278 extreme example. Yan et al. [73] suggested that such methods should be accurate even if a single  
279 gene is randomly selected from each species for each gene tree and used as input to ASTRAL (a  
280 sampling scheme that has been referred to as “ASTRAL-ONE” [73,74]). In such a scenario,  
281 there are more combinations of sampled genes that result in pseudo-orthologs than in true  
282 orthologs (Figure 2B). However, one-third of the pseudo-ortholog combinations match the  
283 species tree topology, and the other two-thirds are split evenly between the two alternative  
284 topologies. Together, the orthologs and pseudo-orthologs matching the species tree ensure that  
285 this quartet is always the most common [74,83]; simulations show that with even a few hundred  
286 loci accurate species trees can be recovered using this approach [73]. Although the numbers of  
287 tree topologies given here only involve four species (including the outgroup) and one duplication  
288 event, they should hold for all larger trees since these can be deconstructed into quartets (cf.  
289 [84]). In biological scenarios involving similarly extreme gene loss, both orthologs and pseudo-  
290 orthologs matching the species tree are more likely to be sampled because they require fewer  
291 losses to produce them than the pseudo-ortholog trees that do not match (Figure 2B). This makes  
292 the species tree even more likely to be accurately inferred using quartet methods.

293

294 Because of their relative simplicity, ease-of-use, speed, accuracy, and robustness to multiple  
295 issues that confound other phylogenetic methods, quartet methods have become a mainstay of  
296 standard phylogenetic inference using single-copy orthologs. For all of the same reasons, they  
297 are likely to become widely used when sampling both orthologs and paralogs. We also suspect  
298 that methods related to ASTRAL that have not yet been evaluated under the inclusion of  
299 paralogs (e.g., [20]) will perform equally well under these conditions.

300

301 **Considerations when inferring phylogenies with paralogs**

302

303 Although multiple of the methods discussed here ensure accurate inference of species tree  
304 topologies when paralogs are used, there are important caveats and implications that merit  
305 specific consideration. Below, we discuss several of these.

306

### 307 *Branch lengths*

308

309 Although topology estimates should not be biased by the inclusion of paralogs, the same is not  
310 true for branch lengths. When branch lengths are estimated as substitutions per site [85,86], the  
311 inclusion of pseudo-orthologs will force branches to be longer than they actually are (e.g. Figure  
312 2; [84]). Conversely, when branch lengths are estimated in coalescent units [13,24], the  
313 additional gene tree heterogeneity introduced by paralogs (hidden or not) will result in the  
314 underestimation of branch lengths. No matter what type of branch lengths are to be estimated, we  
315 recommend that the dataset used be restricted to orthologs. Thus, a reasonable approach would  
316 be to estimate a species tree topology using all genes, and then to estimate branch lengths on this  
317 topology with a dataset including only orthologs (allowing for sampling among species-specific  
318 paralogs; Figure 1D).

319

### 320 *Alignment*

321

322 One of the most error-prone, but underappreciated, steps in phylogenomics is alignment.  
323 Automated alignment of thousands of loci means that many errors can creep in, especially when  
324 non-homologous (alternative) exons are sampled from different species. Fortunately, there are  
325 good methods for identifying regions with low alignment quality (e.g. GUIDANCE2; [87]). A  
326 related problem involves deciding how to choose among lineage-specific paralogs (Figure 1D) in  
327 order to maximize alignment length while minimizing alignment error. One promising approach  
328 would be to co-opt methods designed to choose among alternative isoforms at a single locus:  
329 some of these try to pick the set of genes that are most similar in length across species to avoid  
330 the inclusion of non-homologous exons [88]. Combining such methods with tools that identify  
331 and filter unreliable portions of alignments [87,89–92] should minimize error.

332

### 333 *Polyploidy*

334  
335 Polyploidy is a special case of gene duplication and loss in which the whole genome is  
336 duplicated, and offers a particular challenge both to methods for identifying orthologs and to  
337 species tree inference. In autopolyploidy both sets of chromosomes come from the same species,  
338 and gene copies are paralogs that behave in much the same manner as the smaller duplication  
339 events described above. Therefore, the gene tree methods discussed here should not be misled by  
340 autopolyploidy.

341  
342 Allopolyploidy occurs when the chromosome number doubles via hybridization between species;  
343 the resulting gene copies are referred to as homeologs [93]. Since gene copies found in the same  
344 allopolyploid genome are related through speciation between the parental species, homeologs are  
345 not paralogs in the traditional sense. Similarly, there is not a single bifurcating species tree that  
346 describes relationships involving allopolyploids. While this makes it difficult to evaluate the  
347 effect of including homeologs on traditional species tree inference, gene-tree-based methods  
348 should identify one of the two potentially correct species tree topologies as the correct topology  
349 [e.g., 94].

350  
351 *Detecting introgression*

352  
353 Much less consideration has been given to the effect of including paralogs when attempting to  
354 detect introgression. The most commonly used phylogenetic methods for detecting introgression  
355 are based on the expectation that, for any quartet of species, the two minor topologies (i.e. the  
356 topologies that do not match the species tree) should occur at the same frequency; therefore,  
357 asymmetries between topologies can provide evidence for introgression [95–98]. We suggest  
358 here that, for methods that depend on the frequencies of minor topologies to detect introgression,  
359 the inclusion of paralogs should not bias inference. Consider the example shown in Figure 2: as  
360 discussed above, the most common topology matches the species tree. However, four topologies  
361 do not match the species tree. These four potential trees all require three lineage-specific losses  
362 (one in each taxon), and should occur at equal frequency under a model of GDL in the absence  
363 of introgression, similarly to under cases without duplication. Thus, methods for detecting  
364 introgression based on asymmetry in minor topologies should perform well in the presence of

365 paralogs. This proposal merits additional consideration, however, as does the effect of paralogs  
366 on additional methods for detecting introgression not discussed here.

367

### 368 *Concatenation*

369

370 To carry out a concatenated analysis, one gene copy must be sampled per species per locus and  
371 put into a single alignment. If the intention is to include only orthologs (whether single-copy or  
372 not), a small number of pseudo-orthologs can have an extreme, negative influence on  
373 phylogenetic relationships [27,99]. This occurs because pseudo-orthologs have internal branches  
374 that are longer than those of true orthologs (Figure 2B), giving them more phylogenetically  
375 informative changes. To minimize these potential problems, it may in fact help to instead include  
376 all of the data, rather than attempting to include only orthologs. We imagine here a sampling  
377 scheme similar to the approach taken in [73], where a single copy is randomly sampled per  
378 species (i.e. “ASTRAL-ONE”). Not only are more underlying tree topologies guaranteed to  
379 match the species tree topology, but the pseudo-orthologs matching the species tree have longer  
380 internal branches than those matching alternate topologies (Figure 2B). Thus, with enough data,  
381 the topology matching the species tree should be favored by concatenated analyses, even in the  
382 presence of pseudo-orthologs. While certainly not a standard phylogenetic analysis, we suggest  
383 that this may be a fruitful way forward in the future.

384

### 385 **Concluding Remarks**

386

387 Despite the massive amount of genomic data being collected across the tree of life, phylogeny  
388 inference is often restricted to a small portion of this data due to filtering for single-copy  
389 orthologs and minimal missing data. Recent work has demonstrated that several leading methods  
390 for species tree inference perform well in the presence of paralogs, suggesting a source of  
391 additional data for phylogenomic inference. Additionally, recent work has shown that missing  
392 data may not be as much of an issue as feared. Thus, the amount of data available for  
393 phylogenomic inference may be much larger than previously thought. Future work should  
394 consider branch length estimation when paralogs are present, as well as the potential effects of  
395 paralog inclusion on inferences of introgression.

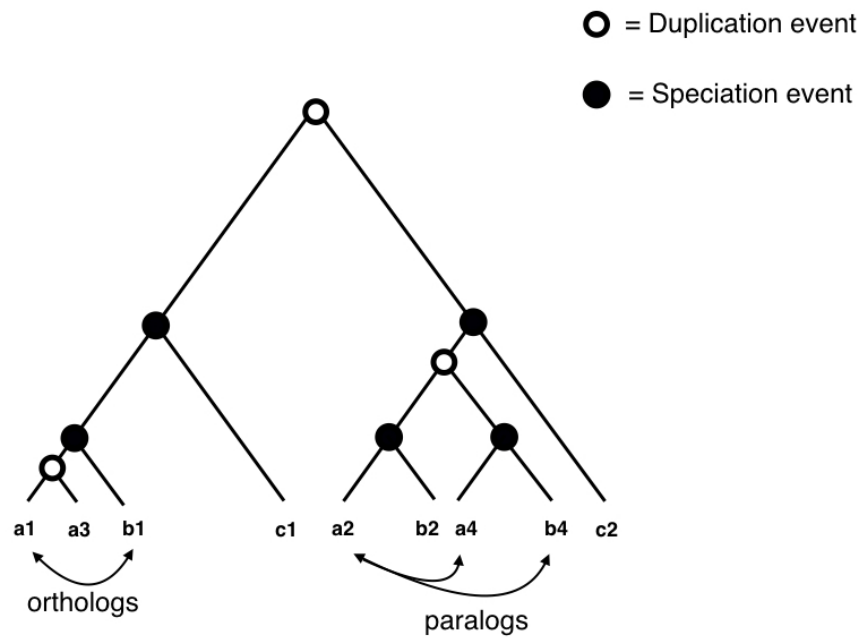
396

397 **Acknowledgements**

398 We thank Rob Lanfear and Erin Molloy for helpful comments. This work was supported by a  
399 National Science Foundation postdoctoral fellowship to MLS (DBI-2009989) and an NSF grant  
400 to MWH (DEB-1936187).

401

402 **Box 1. Types of homologous relationships and implications for phylogenetic inference**  
 403



404  
 405 *Homologous* loci share a common ancestor. *Orthologous* loci share a common ancestor due to  
 406 speciation (e.g. a1 and b1), while *paralogous* loci share a common ancestor due to duplication  
 407 (e.g. a1 and a3; [5]). Orthology relationships can be classified as *one-to-one*, *one-to-many*, and  
 408 *many-to-many* based on whether speciation was followed by duplication in neither, either, or  
 409 both lineages [100]. For example, b1 and c1, are *one-to-one* orthologs. These are the orthologs  
 410 that are typically used in phylogenetic inference. Specifically, researchers target single-copy  
 411 orthologs, which exist in only a single copy in all species considered. However, many-to-one or  
 412 many-to-many orthologs may also be useful. Since the duplication event leading to paralogs a1  
 413 and a3 occurred after the speciation event with b1, they have a *many-to-one* orthologous  
 414 relationship. Such lineage-specific duplications should not affect phylogenetic inference because  
 415 a1 and a3 are *co-orthologous* to b1 and c1, meaning that either copy has an orthologous  
 416 relationship with b1 and c1. Similarly, a2 and a4 have a *many-to-one* orthologous relation to c2.  
 417 The large numbers of complex *many-to-many* relationships that can arise (for instance, the  
 418 relationship between a1, a2, e1, and e2 in Figure 1D) make ortholog group delimitation a  
 419 difficult task, though these loci can still be used in many types of phylogenetic inference.  
 420



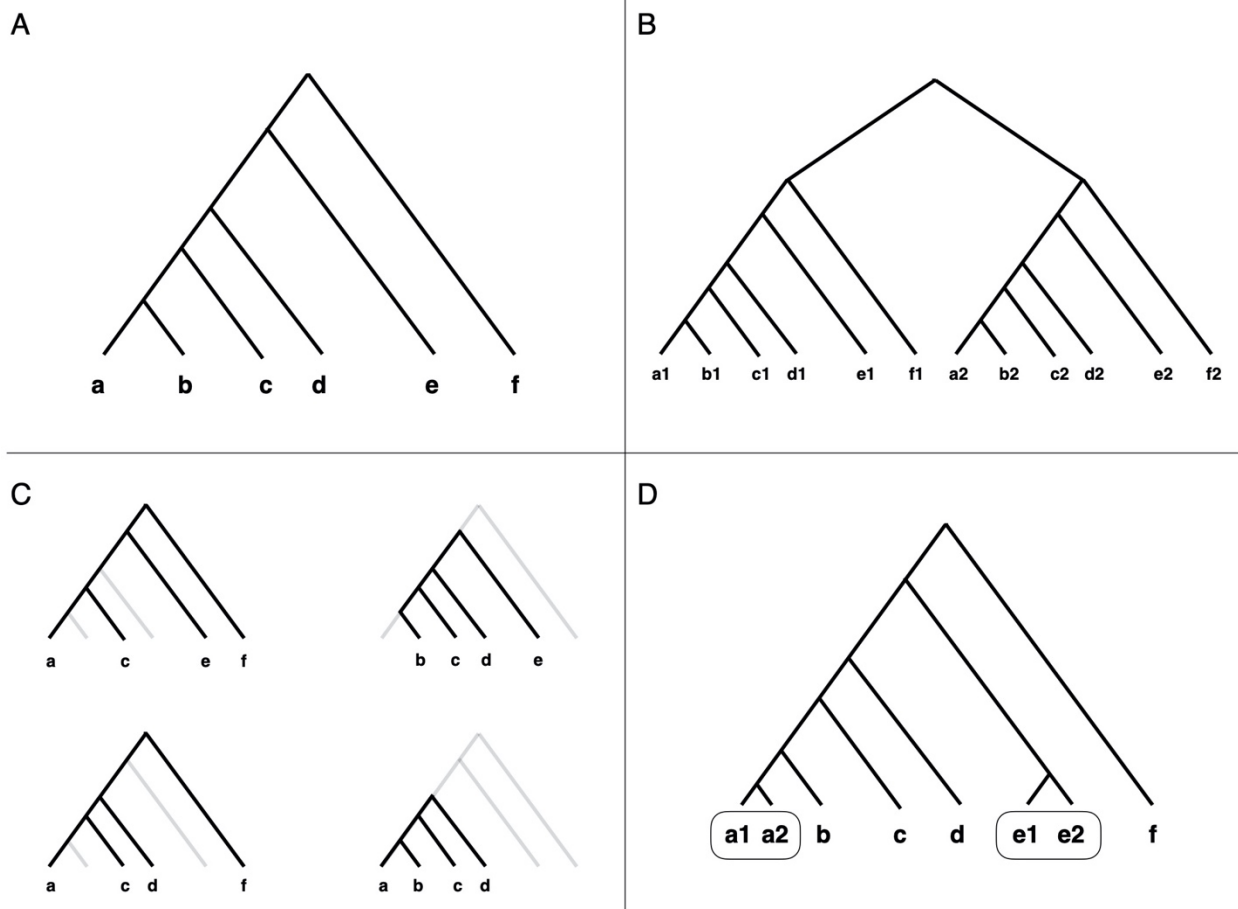
421 **Box 2. Identifying orthologous genes and sampling lineage-specific paralogs.**

422

423 Due to interest in identifying orthologs both for phylogeny reconstruction and for functional  
424 prediction, several methods for ortholog detection have been developed (reviewed in [100]). The  
425 most commonly used approaches for ortholog detection are graph-based approaches [100] which  
426 rely on the identification of reciprocal best hits (RBHs). This is based on the assumption that the  
427 two most closely related homologs between a pair of species should be orthologs. After RBHs  
428 are identified, some approach must be used to construct groups of orthologous sequences; for  
429 example, in OrthoMCL [101] a Markov clustering algorithm is used to identify orthogroups,  
430 which consist of orthologs and recent paralogs. Typically, for downstream phylogenomic  
431 inference, single-copy orthologs present in most species are extracted from these results. While  
432 lineage-specific duplicates need not be excluded from datasets for phylogenetic inference (see  
433 main text), it is not straightforward to extract these from the output of many graph-based  
434 approaches. The most obvious way to identify and include these genes is by reconstructing gene  
435 trees for all orthogroups, identifying lineage-specific duplicates, and selecting one copy per  
436 species for downstream inference. Some recently introduced branch-cutting methods can also  
437 sample such genes from orthogroups containing duplicates. Yang and Smith [102] consider  
438 several different branch-cutting algorithms to extract orthologs appropriate for phylogeny  
439 estimation, and these considerably increase the number of genes available for phylogenetic  
440 inference. For example, in a Hymenoptera dataset analyzed by these authors, the number of  
441 orthologs present in at least eight taxa increased from 4,937 using only single-copy-orthologs to  
442 9,128 under one branch-cutting technique [102]. Thus, even when including paralogs is not  
443 desirable, orthologs can be extracted from many datasets not traditionally considered in  
444 phylogenetic inference.

445

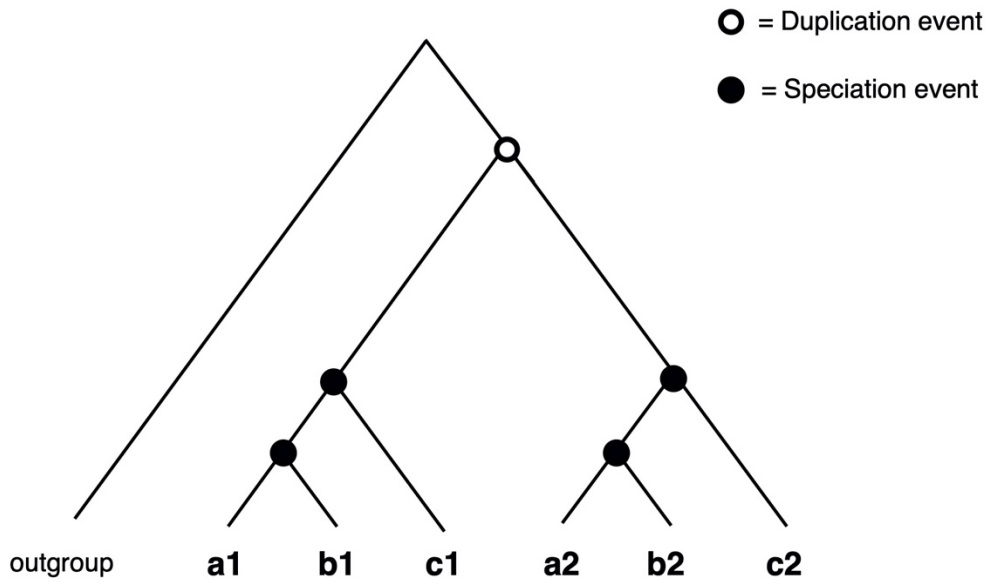
446 **Figure 1. Sampling orthologs and paralogs (Key Figure)**



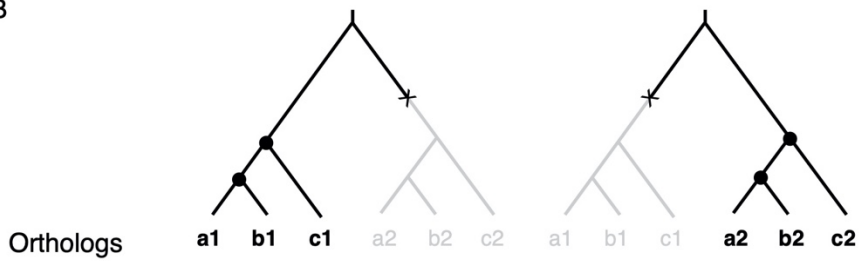
447  
 448 Figure 1. Sampling orthologs and paralogs. There are several potential sampling strategies in  
 449 phylogenetic inference. Here, we illustrate a few of these, although these categories are not  
 450 mutually exclusive. (A) Phylogenies can be constructed from complete sampling of single-copy  
 451 orthologs. (B) Phylogenies can be reconstructed from sets of paralogs. The tree shown has a  
 452 single duplication event in the ancestor of all species. (C) Phylogenies can be constructed from  
 453 genes with missing data, either due to incomplete sampling or to gene loss. (D) Phylogenies can  
 454 be constructed from loci with lineage-specific duplications. Duplications in lineages **a** and **e**  
 455 result in two copies in each of these species in the tree shown. Sampling a single copy from each  
 456 species should not affect phylogenetic inference.  
 457

458 **Figure 2. Orthologs, pseudo-orthologs, and quartet frequencies.**

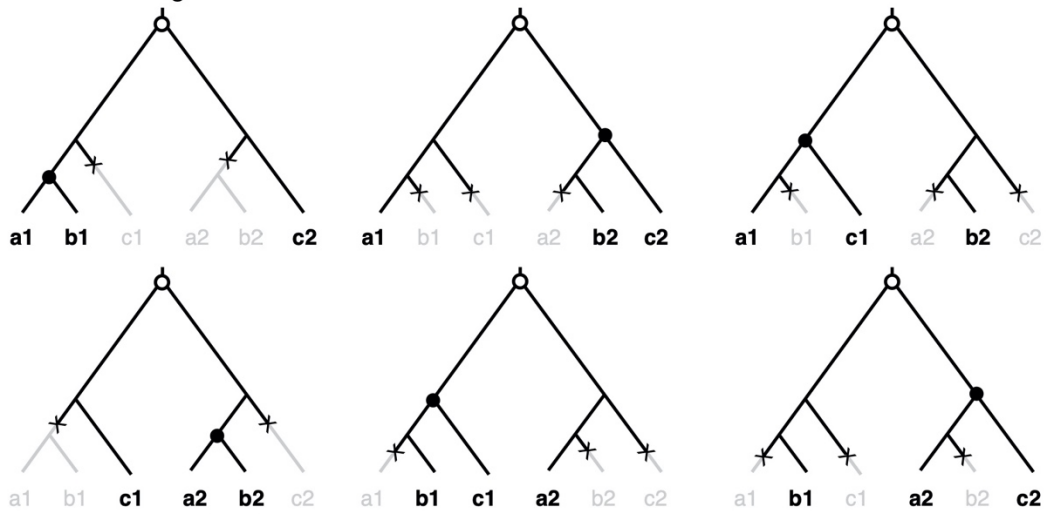
A



B



Pseudo-orthologs



459

460

461 Figure 2. Orthologs, pseudo-orthologs, and quartet frequencies. (A) The full history of a locus in  
462 three species and an outgroup, including one duplication event and two speciation events (which  
463 are shown separately for each set of orthologs). (B) Scenarios where only a single gene copy is  
464 sampled per species; the outgroup is assumed to be sampled in each, but is not shown for clarity.  
465 The single copies may be present because of gene losses (shown here as X's), or simply because  
466 a single copy is randomly chosen per species. The latter case is also what would happen if there  
467 were no missing copies but quartets were sampled from the full gene tree as input to ASTRAL  
468 [73, 74]. There are four quartets that match the species tree: the two orthologs and the two left-  
469 most pseudo-orthologs (“hidden paralogs”). The remaining pseudo-orthologs either place  
470 lineages **b** and **c** sister to one another (center) or **a** and **c** sister to one another (right). Therefore,  
471 quartet methods should perform well even when paralogs are included, because the most  
472 common set of relationships should still match the species tree. Note that if genes are single-copy  
473 because of gene losses, the species tree relationship is likely to become even more common: the  
474 orthologs require only one loss in their history and the matching pseudo-orthologs require two  
475 losses. Pseudo-orthologs not matching the species tree can only be generated when there are  
476 three separate loss events.  
477

478 **References**

- 479 1 Scornavacca, C. et al. (2020) *Phylogenetics in the Genomic Era*, No commercial publisher |  
480 Authors open access book.
- 481 2 Delsuc, F. et al. (2005) Phylogenomics and the reconstruction of the tree of life. *Nat. Rev.*  
482 *Genet.* 6, 361–375
- 483 3 Thomas, G.W.C. et al. (2020) Gene content evolution in the arthropods. *Genome Biol.* 21, 15
- 484 4 Emms, D.M. and Kelly, S. (2018) STAG: Species tree inference from all genes. *bioRxiv*  
485 DOI: 10.1101/267914
- 486 5 Fitch, W.M. (1970) Distinguishing homologous from analogous proteins. *Syst. Zool.* 19, 99–  
487 113
- 488 6 Fernández, R. et al. (2020) Orthology: Definitions, prediction, and impact on species  
489 phylogeny inference. In *Phylogenetics in the Genomic Era* (Scornavacca, C. et al., eds), pp.  
490 2.4:1–2.4:14, No commercial publisher | Authors open access book
- 491 7 Kapli, P. et al. (2020) Phylogenetic tree building in the genomic age. *Nat. Rev. Genet.* 21,  
492 428–444
- 493 8 Nehrt, N.L. et al. (2011) Testing the ortholog conjecture with comparative functional  
494 genomic data from mammals. *PLOS Comput. Biol.* 7, e1002073
- 495 9 Studer, R.A. and Robinson-Rechavi, M. (2009) How confident can we be that orthologs are  
496 similar, but paralogs differ? *Trends Genet.* 25, 210–216
- 497 10 Stamboulian, M. et al. (2019) The ortholog conjecture revisited: The value of orthologs and  
498 paralogs in function prediction. *bioRxiv* DOI: 10.1101/2019.12.27.889691
- 499 11 Goodman, M. et al. (1979) Fitting the gene lineage into its species lineage, a parsimony  
500 strategy illustrated by cladograms constructed from globin sequences. *Syst. Biol.* 28, 132–  
501 163
- 502 12 Chen, K. et al. (2000) NOTUNG: A program for dating gene duplications and optimizing  
503 gene family trees. *J. Comput. Biol.* 7, 429–447
- 504 13 Zhang, C. et al. (2018) ASTRAL-III: Polynomial time species tree reconstruction from  
505 partially resolved gene trees. *BMC Bioinformatics* 19, 153
- 506 14 Liu, L. and Yu, L. (2011) Estimating species trees from unrooted gene trees. *Syst. Biol.* 60,  
507 661–667
- 508 15 Bravo, G.A. et al. (2019) Embracing heterogeneity: Coalescing the tree of life and the future  
509 of phylogenomics. *PeerJ* 7, e6399
- 510 16 Maddison, W.P. (1997) Gene trees in species trees. *Syst. Biol.* 46, 523–536
- 511 17 Kubatko, L.S. and Degnan, J.H. (2007) Inconsistency of phylogenetic estimates from  
512 concatenated data under coalescence. *Syst. Biol.* 56, 17–24
- 513 18 Degnan, J.H. and Rosenberg, N.A. (2006) Discordance of species trees with their most likely  
514 gene trees. *PLOS Genet.* 2, e68
- 515 19 Mirarab, S. and Warnow, T. (2015) ASTRAL-II: Coalescent-based species tree estimation  
516 with many hundreds of taxa and thousands of genes. *Bioinformatics* 31, i44–i52
- 517 20 Chifman, J. and Kubatko, L. (2014) Quartet inference from SNP data under the coalescent  
518 model. *Bioinformatics* 30, 3317–3324
- 519 21 Vachaspati, P. and Warnow, T. (2015) ASTRID: Accurate species trees from internode  
520 distances. *BMC Genomics* 16, S3
- 521 22 Heled, J. and Drummond, A.J. (2010) Bayesian inference of species trees from multilocus  
522 data using \*BEAST. *Mol. Biol. Evol.* 27, 570–580

- 523 23 Bryant, D. et al. (2012) Inferring species trees directly from biallelic genetic markers:  
524 Bypassing gene trees in a full coalescent analysis. *Mol. Biol. Evol.* 29, 1917–1932
- 525 24 Liu, L. et al. (2010) A maximum pseudo-likelihood approach for estimating species trees  
526 under the coalescent model. *BMC Evol. Biol.* 10, 302
- 527 25 Doolittle, W.F. and Brown, J.R. (1994) Tempo, mode, the progenote, and the universal root.  
528 *Proc. Natl. Acad. Sci.* 91, 6721–6728
- 529 26 Koonin, E.V. (2005) Orthologs, paralogs, and evolutionary genomics. *Annu. Rev. Genet.* 39,  
530 309–338
- 531 27 Brown, J.M. and Thomson, R.C. (2017) Bayes factors unmask highly variable information  
532 content, bias, and extreme influence in phylogenomic analyses. *Syst. Biol.* 66, 517–530
- 533 28 Altenhoff, A.M. et al. (2019) OMA standalone: Orthology inference among public and  
534 custom genomes and transcriptomes. *Genome Res.* 29, 1152–1163
- 535 29 Kallal, R.J. et al. (2018) A phylotranscriptomic backbone of the orb-weaving spider family  
536 Araneidae (Arachnida, Araneae) supported by multiple methodological approaches. *Mol.*  
537 *Phylogenet. Evol.* 126, 129–140
- 538 30 Fernández, R. et al. (2018) Phylogenomics, diversification dynamics, and comparative  
539 transcriptomics across the spider tree of life. *Curr. Biol.* 28, 1489–1497
- 540 31 Roure, B. et al. (2013) Impact of missing data on phylogenies inferred from empirical  
541 phylogenomic data sets. *Mol. Biol. Evol.* 30, 197–214
- 542 32 Wiens, J.J. (2003) Missing data, incomplete taxa, and phylogenetic accuracy. *Syst. Biol.* 52,  
543 528–538
- 544 33 Wiens, J.J. (2006) Missing data and the design of phylogenetic analyses. *Syst. Biol.* 39, 34–  
545 42
- 546 34 Philippe, H. et al. (2004) Phylogenomics of eukaryotes: Impact of missing data on large  
547 alignments. *Mol. Biol. Evol.* 21, 1740–1752
- 548 35 Driskell, A.C. et al. (2004) Prospects for building the tree of life from large sequence  
549 databases. *Science* 306, 1172–1174
- 550 36 Hosner, P.A. et al. (2016) Avoiding missing data biases in phylogenomic inference: An  
551 empirical study in the landfowl (Aves: Galliformes). *Mol. Biol. Evol.* 33, 1110–1125
- 552 37 Wiens, J.J. and Morrill, M.C. (2011) Missing data in phylogenetic analysis: Reconciling  
553 results from simulations and empirical data. *Syst. Biol.* 60, 719–731
- 554 38 Nute, M. et al. (2018) The performance of coalescent-based species tree estimation methods  
555 under models of missing data. *BMC Genomics* 19, 286
- 556 39 Xi, Z. et al. (2016) The impact of missing data on species tree estimation. *Mol. Biol. Evol.*  
557 33, 838–860
- 558 40 Molloy, E.K. and Warnow, T. (2018) To include or not to include: The impact of gene  
559 filtering on species tree estimation methods. *Syst. Biol.* 67, 285–303
- 560 41 Rhodes, J.A. et al. (2020) NJst and ASTRID are not statistically consistent under a random  
561 model of missing data. *arXiv preprint arXiv:2001.07844*.
- 562 42 Nute, M. et al. (2020) Correction to: The performance of coalescent-based species tree  
563 estimation methods under models of missing data. *BMC Genomics* 21, 133
- 564 43 Eaton, D.A.R. and Ree, R.H. (2013) Inferring phylogeny and introgression using RADseq  
565 data: An example from flowering plants (*Pedicularis*: Orobanchaceae). *Syst. Biol.* 62, 689–  
566 706

- 567 44 Gatesy, J. et al. (2019) Partitioned coalescence support reveals biases in species-tree methods  
568 and detects gene trees that determine phylogenomic conflicts. *Mol. Phylogenet. Evol.* 139,  
569 106539
- 570 45 Page, R.D.M. (1994) Maps between trees and cladistic analysis of historical associations  
571 among genes, organisms, and areas. *Syst. Biol.* 43, 58–77
- 572 46 Guigo, R. et al. (1996) Reconstruction of ancient molecular phylogeny. *Mol. Phylogenet.*  
573 *Evol.* 6, 189–213
- 574 47 Boussau, B. and Scornavacca, C. (2020) Reconciling gene trees with species trees. In  
575 *Phylogenetics in the Genomic Era* pp. 3.2:1-3.2:23, No commercial publisher | Authors open  
576 access book
- 577 48 Wehe, A. et al. (2008) DupTree: A program for large-scale phylogenetic analyses using gene  
578 tree parsimony. *Bioinformatics* 24, 1540–1541
- 579 49 Bayzid, M.S. and Warnow, T. (2018) Gene tree parsimony for incomplete gene trees:  
580 Addressing true biological loss. *Algorithms Mol. Biol.* 13, 1
- 581 50 Maddison, W.P. and Knowles, L.L. (2006) Inferring phylogeny despite incomplete lineage  
582 sorting. *Syst. Biol.* 55, 21–30
- 583 51 Than, C. and Nakhleh, L. (2009) Species tree inference by minimizing deep coalescences.  
584 *PLOS Comput. Biol.* 5, e1000501
- 585 52 Chaudhary, R. et al. (2010) iGTP: A software package for large-scale gene tree parsimony  
586 analysis. *BMC Bioinformatics* 11, 574
- 587 53 Wu, Y.-C. et al. (2014) Most parsimonious reconciliation in the presence of gene  
588 duplication, loss, and deep coalescence using labeled coalescent trees. *Genome Res.* 24,  
589 475–486
- 590 54 Stolzer, M. et al. (2012) Inferring duplications, losses, transfers and incomplete lineage  
591 sorting with nonbinary species trees. *Bioinformatics* 28, i409–i415
- 592 55 Chan, Y. et al. (2017) Inferring incomplete lineage sorting, duplications, transfers and losses  
593 with reconciliations. *J. Theor. Biol.* 432, 1–13
- 594 56 Than, C.V. and Rosenberg, N.A. (2011) Consistency properties of species tree inference by  
595 minimizing deep coalescences. *J. Comput. Biol.* 18, 1–15
- 596 57 Robinson, D.F. and Foulds, L.R. (1981) Comparison of phylogenetic trees. *Math. Biosci.* 53,  
597 131–147
- 598 58 Bansal, M.S. et al. (2010) Robinson-Foulds supertrees. *Algorithms Mol. Biol.* 5, 18
- 599 59 Puigbo, P. et al. (2007) TOPD/FMTS: A new software to compare phylogenetic trees.  
600 *Bioinformatics* 23, 1556–1558
- 601 60 Marcet-Houben, M. and Gabaldón, T. (2011) TreeKO: A duplication-aware algorithm for the  
602 comparison of phylogenetic trees. *Nucleic Acids Res.* 39, e66
- 603 61 Chaudhary, R. et al. (2013) Inferring species trees from incongruent multi-copy gene trees  
604 using the Robinson-Foulds distance. *Algorithms Mol. Biol.* 8, 28
- 605 62 Chaudhary, R. et al. (2015) MulRF: A software package for phylogenetic analysis using  
606 multi-copy gene trees. *Bioinformatics* 31, 432–433
- 607 63 Chaudhary, R. et al. (2015) Assessing approaches for inferring species trees from multi-copy  
608 genes. *Syst. Biol.* 64, 325–339
- 609 64 Molloy, E.K. and Warnow, T. (2019) FastMulRFS: Fast and accurate species tree estimation  
610 under generic gene duplication and loss models. *bioRxiv* DOI: 10.1101/835553
- 611 65 Zhang, C. et al. (2019) ASTRAL-Pro: Quartet-based species tree inference despite paralogy.  
612 *bioRxiv* DOI: 10.1101/2019.12.12.874727.

- 613 66 Boussau, B. et al. (2013) Genome-scale coestimation of species and gene trees. *Genome Res.*  
614 23, 323–330
- 615 67 De Oliveira Martins, L. et al. (2016) A Bayesian supertree model for genome-wide species  
616 tree reconstruction. *Syst. Biol.* 65, 397–416
- 617 68 Saitou, N. and Nei, M. (1987) The neighbor-joining method: A new method for  
618 reconstructing phylogenetic trees. *Mol. Biol. Evol.* 4, 406–425
- 619 69 Liu, L. *et al.* (2009) Estimating species phylogenies using coalescence times among  
620 sequences. *Syst. Biol.* 58, 468–477
- 621 70 Dasarathy, G. et al. (2015) Data requirement for phylogenetic inference from multiple loci:  
622 A new distance method. *IEEE/ACM Trans. Comput. Biol. Bioinform.* 12, 422–432
- 623 71 Allman, E.S. et al. (2013) Species tree inference by the STAR method and its  
624 generalizations. *J. Comput. Biol.* 20, 50–61
- 625 72 Allman, E.S. et al. (2018) Species tree inference from gene splits by unrooted STAR  
626 methods. *IEEE/ACM Trans. Comput. Biol. Bioinform.* 15, 337–342
- 627 73 Yan, Z. et al. (2018) Species tree inference under the multispecies coalescent on data with  
628 paralogs is accurate. *bioRxiv* DOI: 10.1101/498378
- 629 74 Legried, B. et al. (2020) Polynomial-time statistical estimation of species trees under gene  
630 duplication and loss. *Proceedings of RECOMB 2020: The 24th Annual International  
631 Conference Research in Computational Molecular Biology*, pp. 120-135. Springer, Cham.
- 632 75 Graur, D. et al. (1996) Phylogenetic position of the order Lagomorpha (rabbits, hares, and  
633 allies). *Nature* 379, 333-335
- 634 76 Bryant, D. and Steel, M. (2001) Constructing optimal trees from quartets. *J. Algorithms* 38,  
635 237–259
- 636 77 Strimmer, K. and von Haeseler, A. (1996) Quartet puzzling: A quartet maximum-likelihood  
637 method for reconstructing tree topologies. *Mol. Biol. Evol.* 13, 964–969
- 638 78 Snir, S. and Rao, S. (2012) Quartet MaxCut: A fast algorithm for amalgamating quartet  
639 trees. *Mol. Phylogenetic. Evol.* 62, 1-8
- 640 79 Reaz, R. et al. (2014) Accurate Phylogenetic Tree Reconstruction from Quartets: A Heuristic  
641 Approach. *PLOS One* 9, e104008
- 642 80 Mirarab, S. et al. (2014) ASTRAL: Genome-scale coalescent-based species tree estimation.  
643 *Bioinformatics* 30, i541–i548
- 644 81 Wascher, M. and Kubatko, L. (2019) Consistency of SVDQuartets and maximum likelihood  
645 for coalescent-based species tree estimation. *Syst. Biol.* DOI: 10.1093/sysbio/syaa039
- 646 82 Rabiee, M. et al. (2019) Multi-allele species reconstruction using ASTRAL. *Mol.*  
647 *Phylogenet. Evol.* 130, 286–296
- 648 83 Markin, A. and Eulenstein, O. (2020) Quartet-based inference methods are statistically  
649 consistent under the unified duplication-loss-coalescence model. *arXiv preprint*  
650 *arXiv:2004.04299*.
- 651 84 Siu-Ting, K. et al. (2019) Inadvertent paralog inclusion drives artifactual topologies and  
652 timetree estimates in phylogenomics. *Mol. Biol. Evol.* 36, 1344–1356
- 653 85 Kozlov, A.M. et al. (2019) RAXML-NG: A fast, scalable and user-friendly tool for maximum  
654 likelihood phylogenetic inference. *Bioinformatics* 35, 4453–4455
- 655 86 Huelsenbeck, J.P. and Ronquist, F. (2001) MrBayes: Bayesian inference of phylogenetic  
656 trees. *Bioinformatics* 17, 754–755
- 657 87 Sela, I. et al. (2015) GUIDANCE2: Accurate detection of unreliable alignment regions  
658 accounting for the uncertainty of multiple parameters. *Nucleic Acids Res.* 43, W7–W14



- 659 88 Villanueva-Cañas, J.L. et al. (2013) Improving genome-wide scans of positive selection by  
660 using protein isoforms of similar length. *Genome Biol. Evol.* 5, 457–467
- 661 89 Capella-Gutiérrez, S. et al. (2009) trimAl: A tool for automated alignment trimming in large-  
662 scale phylogenetic analyses. *Bioinformatics* 25, 1972–1973
- 663 90 Castresana, J. (2000) Selection of conserved blocks from multiple alignments for their use in  
664 phylogenetic analysis. *Mol. Biol. Evol.* 17, 540–552
- 665 91 Dress, A.W. et al. (2008) Noisy: Identification of problematic columns in multiple sequence  
666 alignments. *Algorithms Mol. Biol.* 3, 7
- 667 92 Landan, G. and Graur, D. (2007) Heads or tails: A simple reliability check for multiple  
668 sequence alignments. *Mol. Biol. Evol.* 24, 1380–1383
- 669 93 Glover, N.M. et al. (2016) Homoeologs: What are they and how do we infer them? *Trends*  
670 *Plant Sci.* 21, 609–621
- 671 94 Thomas, G.W.C. et al. (2017) Gene-tree reconciliation with MUL-trees to resolve polyploidy  
672 events. *Syst. Biol.* 66, 1007–1018
- 673 95 Huson, D.H. et al. (2005) Reconstruction of reticulate networks from gene trees. *Proceedings*  
674 *of RECOMB 2005: The 9th Annual International Conference Research in Computational*  
675 *Molecular Biology*, pp. 233-249. Berlin: Springer.
- 676 96 Vanderpool, D. et al. (2020) Primate phylogenomics uncovers multiple rapid radiations and  
677 ancient interspecific introgression. *bioRxiv* DOI: 10.1101/2020.04.15.043786
- 678 97 Yu, Y. and Nakhleh, L. (2015) A maximum pseudo-likelihood approach for phylogenetic  
679 networks. *BMC Genomics* 16, S10
- 680 98 Solís-Lemus, C. and Ané, C. (2016) Inferring phylogenetic networks with maximum  
681 pseudolikelihood under incomplete lineage sorting. *PLOS Genet.* 12, e1005896
- 682 99 Shen, X.-X. et al. (2017) Contentious relationships in phylogenomic studies can be driven by  
683 a handful of genes. *Nat. Ecol. Evol.* 1, 126
- 684 100 Altenhoff, A.M. et al. (2019) Inferring orthology and paralogy. In *Evolutionary Genomics:*  
685 *Statistical and Computational Methods* (Anisimova, M., ed), pp. 149–175, Springer
- 686 101 Li, L. (2003) OrthoMCL: Identification of ortholog groups for eukaryotic genomes.  
687 *Genome Research* 13, 2178-2189
- 688 102 Yang, Y. and Smith, S. A. (2014) Orthology inference in nonmodel organisms using  
689 transcriptomes and low-coverage genomes: Improving accuracy and matrix occupancy  
690 for phylogenomics. *Mol. Biol. Evol.* 31, 3081-3092