

1 **Running head**

2 Tree root inference from gene duplication events

3 **Title**

4 STRIDE: Species Tree Root Inference from Gene Duplication Events

5 **Authors**

6 Emms, D.M.<sup>1</sup> and Kelly, S.<sup>1\*</sup>

7 **Affiliations**

8 1) Department of Plant Sciences, University of Oxford, South Parks Road, Oxford, OX1  
9 3RB, UK.

10 \* Corresponding author

11 Phone: +44-1865-275123

12 Email: [steven.kelly@plants.ox.ac.uk](mailto:steven.kelly@plants.ox.ac.uk)

13 **Abstract**

14 The correct interpretation of a phylogenetic tree is dependent on it being correctly rooted.  
15 A gene duplication event at the base of a clade of species is synapomorphic, and thus  
16 excludes the root of the species tree from that clade. We present STRIDE, a fast, effective,  
17 and outgroup-free method for species tree root inference from gene duplication events.  
18 STRIDE identifies sets of well-supported gene duplication events from cohorts of gene  
19 trees, and analyses these events to infer a probability distribution over an unrooted  
20 species tree for the location of the true root. We show that STRIDE infers the correct root  
21 of the species tree for a large range of simulated and real species sets. We demonstrate  
22 that the novel probability model implemented in STRIDE can accurately represent the  
23 ambiguity in species tree root assignment for datasets where information is limited.  
24 Furthermore, application of STRIDE to inference of the origin of the eukaryotic tree

25 resulted in a root probability distribution that was consistent with, but unable to distinguish  
26 between, leading hypotheses for the origin of the eukaryotes. In summary, STRIDE is a  
27 fast, scalable, and effective method for species tree root inference from genome scale  
28 data.

## 29 **Keywords**

30 Phylogenetics; phylogenomics; gene duplication;

## 31 **Introduction**

32 “*Nothing in biology makes sense except in the light of evolution*” (Dobzhansky, T.  
33 2013), “*nothing in evolution makes sense except in the light of phylogeny*” (Sytsma, K.J.  
34 and Pires, J.C. 2001), and nothing in a phylogeny makes sense except in the light of its  
35 root. For example, the phylogeny for four species (Fig. 1A) has five possible roots (Fig. 1B-  
36 F) and each of the different roots corresponds to a different hypothesis as to the  
37 evolutionary history of the species. For the presented tree, identifying a wrong branch as  
38 the root (for example Fig. 1E) would lead us to conclude that elephants are more closely  
39 related to fish and birds than they are to wolves, even though we are using a tree with the  
40 correct topology. A species tree only gives the correct evolutionary relationships when  
41 rooted correctly (Fig. 1B). Thus it is of critical importance to our interpretation of  
42 relationships, and the evolutionary history of life on earth, that we have accurate methods  
43 of inferring the root of species phylogenies.

44 Correct species tree rooting is also of critical importance for the inference of  
45 orthology relationships between genes. Given an unrooted gene tree (Fig. 2A), knowledge  
46 of the correct branching order of the species tree (Fig. 1B) is required to correctly root the  
47 gene tree (Fig. 2B). An incorrect rooting of the species tree (Fig. 1C-F) leads to an  
48 incorrect inference of the root of the gene tree (Fig. 2C-F), and thus incorrect identification  
49 of orthologous genes (Fig. 2G-H). Therefore, our ability to compare the biology of species,

50 through comparisons between orthologous genes, is reliant on accurate methods of  
51 inferring the root of species phylogenies.

52 Although correct root placement is essential for our ability to interpret phylogenies,  
53 almost all models of sequence evolution used for tree inference are time-reversible and  
54 produce unrooted phylogenetic trees. In order to identify the root of a phylogeny extra  
55 information is required, usually knowledge of an extra species that is a suitable (i.e. closely  
56 related) outgroup for the set of species for which the root is unknown. However, outgroup  
57 choice is a common source of error in phylogenetic tree inference, with distantly related  
58 outgroups leading to inaccurate root placement and distortion of the phylogeny due to long  
59 branch attraction (Felsenstein, J. 1981) (Berger, S.A., Krompass, D., et al. 2011). While  
60 time-irreversible models of sequence evolution exist, they do not implicitly provide a  
61 method for accurately inferring the direction of time in a tree (Huelsenbeck, J.P., Bollback,  
62 J.P., et al. 2002, Williams, T.A., Heaps, S.E., et al. 2015). To address this issue, methods  
63 have been developed that can simultaneously infer rooted species and gene trees  
64 (Boussau, B., Szollosi, G.J., et al. 2013). However, these methods are computationally  
65 expensive and do not scale well to moderate or large species sets.

66 “Duplicate gene rooting” has been proposed as an alternative method for rooting  
67 species trees (Donoghue, M.J. and Mathews, S. 1998, Simmons, M.P., Bailey, C.D., et al.  
68 2000). The conceptual basis for this approach is that gene duplication events are time-  
69 irreversible, unlike character substitution, and thus infer the direction of time on the  
70 species tree. Specifically, every node in an unrooted, binary gene tree has three branches  
71 incident upon it. If the node is a speciation node then any of the three incident branches  
72 could be the edge in the direction of the root, with the other two being in the opposite  
73 direction. Thus, speciation nodes are uninformative about the direction of time along the  
74 tree. For a duplication node, however, the symmetry is broken. Two of the edges will  
75 correspond to the two copies of the gene post-duplication, while the third edge will

76 correspond to the gene pre-duplication and thus point towards the root of the tree (Fig. 2A,  
77 node marked ' $\mathcal{D}$ '). In the case of this example tree, it can be inferred that the root of the  
78 species tree must be outside of the subtree containing elephant and wolf. In an idealised  
79 case (with no effects such as incomplete lineage sorting or lateral gene transfer) the two  
80 post-duplication branches can be distinguished from the pre-duplication branch as the  
81 post-duplication branches contain genes from overlapping species sets. Furthermore,  
82 these species sets will be identical if there has been no gene loss or horizontal gene  
83 transfer, and the topology of the duplicate subtrees will recapitulate the species tree  
84 topology. Thus, if gene duplication nodes can be accurately identified in an unrooted gene  
85 tree, then the direction of time can be ascertained for all branches in the post-duplication  
86 subtrees. The direction of time on these branches determines the direction of time on the  
87 corresponding branches of the species tree, and multiple gene duplication events can be  
88 aggregated to determine the direction of time across the whole species tree, thus revealing  
89 the location of the root.

90 Here we present STRIDE, a novel algorithm for Species Tree Root Inference from  
91 gene Duplication Events. STRIDE identifies sets of well-supported gene duplication events  
92 from cohorts of gene trees, and analyses these events to infer a probability distribution  
93 over an unrooted species tree for the location of the true root. We show that STRIDE  
94 correctly identifies the community-accepted root of the majority of species trees.  
95 Additionally, we demonstrate that STRIDE effectively captures uncertainty in root  
96 placement when data is limited or conflicting. Finally, we demonstrate the utility of STRIDE  
97 to challenging phylogenetic problems by providing an outgroup-free root analysis of the  
98 origin of the eukaryotes.

## 99 **Methods**

### 100 ***Problem definition and approach***

101 A branch of an unrooted species tree corresponds to a bipartition that splits the  
102 tree's taxa into two blocks. The presence of a well-supported gene duplication that  
103 respects the topology of the species tree is a synapomorphy that stipulates that the block  
104 in which the duplicates are found is a monophyletic clade. This synapomorphy identifies  
105 the direction of time along the branches within this monophyletic clade. The single  
106 exception to this is the branch in the unrooted tree corresponding to the root in which time  
107 flows in both directions. This is because the branch that spans the root in the unrooted  
108 species tree corresponds to two branches in the rooted species tree and both of its  
109 corresponding blocks are monophyletic clades (Fig. 1A & B). The method presented here  
110 aims to identify this root branch by identifying and analysing a set of well-supported gene-  
111 duplication events. The method identifies the complete set of gene duplication events  
112 contained within a set of user-supplied gene trees and uses these to infer the location of  
113 the root of the species tree. To express uncertainty in the case of limited data or data  
114 conflict, the method uses a probabilistic model of gene-duplication events to calculate a  
115 probability distribution across the branches of the species tree for the location of the root.

### 116 ***Inference of Orthogroups and Gene Trees***

117 For each species set, the protein sequence translations of representative gene  
118 models were downloaded from appropriate online databases. These protein sequences  
119 were then subject to orthogroup inference using OrthoFinder v1.1.4 (Emms, D.M. and  
120 Kelly, S. 2015). The resulting sets of protein sequence orthogroups were aligned using  
121 MAFFT L-INS-I v7.305b (Kato, K. and Standley, D.M. 2013) and subject phylogenetic  
122 inference using IQTREE v1.5.3 (Nguyen, L.T., Schmidt, H.A., et al. 2015). All methods  
123 used their default settings. Parallelisation of MAFFT and IQTREE runs was done using  
124 GNU Parallel (Tange, O. 2011). Alignments were viewed using AliView (Larsson, A. 2014).

125 Trees were viewed using Dendroscope (Huson, D.H. and Scornavacca, C. 2012) and  
126 drawn using the ETE library (Huerta-Cepas, J., Serra, F., et al. 2016).

### 127 ***Identification of Putative Duplications***

128 Gene-duplication events are considered informative if they occur on any branch  
129 other than a terminal branch in the species phylogeny, as a duplicated gene that occurs  
130 only in a single species is not informative to the position of the root of the tree. To identify  
131 informative gene duplication events a novel tree analysis algorithm was developed (Fig. 3).  
132 Prior to analysis of the gene trees, the unrooted species tree was analysed to determine  
133 the species sets in which genes would be expected to occur following a gene duplication  
134 event along any branch of the species tree. For each direction along each branch the sets  
135 of species in the child clades (X & Y) and in the grandchild clades ( $x_1$ ,  $x_2$ ,  $y_1$  &  $y_2$ )  
136 immediately following that branch were identified (Fig. 3A). Then each gene tree was  
137 analysed in turn to identify well-supported gene duplication events within the gene tree as  
138 follows: each node,  $n$ , in the tree was considered in turn, if the node was an unresolved  
139 polytomy it was excluded as such nodes correspond to either a higher-order multiplication  
140 events (e.g. triplication) or an unresolved event in the gene tree (e.g. an amalgamation of  
141 several weakly supported bipartitions). Each analysed node therefore had three branches  
142 incident on it, and any pair of branches could potentially represent duplicate genes (Fig.  
143 3B). For each pair of branches,  $b_1$  and  $b_2$ , the sets of species,  $S_1$  and  $S_2$ , below each  
144 branch were used to identify the locations in the species tree corresponding to these  
145 branches in the gene tree. This was done by identifying the smallest block,  $B_i$  in the  
146 species tree that contains all the species in  $S_i$  ( $i=1,2$ ), thus making the method robust in  
147 the case of subsequent gene loss (Fig. 3B). If there was more than one block satisfying  
148 this criteria, each of these possible blocks were considered. A node,  $n$ , was considered as  
149 a putative gene duplication event if  $B_1=B_2$ .

150 Nodes that were identified as putative gene duplication events were further  
151 examined to reduce the possibility that their existence or location had been misidentified  
152 due to errors in gene tree inference. The criteria were: 1) There must be at least one gene  
153 from each of the expected grandchild clades in both  $S_1$  and  $S_2$  (Fig. 3C). 2) The branching  
154 structure immediately after the gene duplication event on branches  $b_1$  and  $b_2$  must match  
155 the expected branching structure (Fig. 3D), i.e. the first node for each duplicate split the  
156 descendent species into the expected sets X and Y, or subsets thereof. Note that it would  
157 not be valid to check the topology to the level of grandchild clades in step 2 since this  
158 would fail to identify gene duplication events if there were also a subsequent gene  
159 duplication event one branch lower in the species tree. In this case, the observed  
160 grandchild clades would both be subsets of one of the expected child clades rather than  
161 grandchild clades. Steps 1 and 2 check that the observed clades are subsets of the  
162 expected clades (rather than requiring they be equal to) as this is necessary to make the  
163 method robust to subsequent gene loss events.

#### 164 ***Identifying the Maximum-Parsimony Root of the Species Tree***

165 As discussed above, a gene duplication on a bipartition of an unrooted species tree  
166 stipulates the direction of time for all branches of the subtrees derived from that  
167 bipartition. Given a set of gene duplication events, the branch in the species tree that  
168 violates the fewest gene duplication events is identified as the maximum parsimony root. If  
169 multiple such branches exist then they are each identified as equally parsimonious.

#### 170 ***Probability model for the root of the species tree***

171 For any given set of gene-trees, it is possible that errors in gene-tree inference will  
172 lead to false positive inference of gene duplication events that past the filtration criteria. To  
173 account for this, a probability model was developed for the location of the root of the tree  
174 given the set of (potentially conflicting) putative gene duplication events identified. The  
175 model consisted of two parts. The first part, the branch-level model, calculated the

176 probability that a branch was the root given only the duplications observed in either  
177 direction along that branch. The second part, the tree-level model, aggregated all  
178 duplications observed across all branches of tree to give the final probability distribution for  
179 the location of the root taking into account all information obtained from all gene  
180 duplication events observed across the tree.

181 At the branch-level, the set of putative gene duplication events identified on that  
182 branch are modelled by two Poisson processes, one giving rise to true positive gene  
183 duplications and the other to false positive duplications. On a given branch,  $i$ , of a species  
184 tree,  $m_i$  duplications are observed that support time flowing in one direction along the  
185 branch,  $\rightarrow$ , and  $n_i$  duplications supporting time flowing in the opposite direction,  $\leftarrow$ . The  
186 set of duplications on branch  $i$  is then written,  $d_i = \binom{m_i}{n_i}$ , and  $D$  is the set of putative  
187 duplications observed on all branches of the species tree,  $D = \{d_1, d_2, \dots, d_b\}$ .

188 The final tree-level probability distribution  $P(i = root|D)$  takes into account the complete  
189 set of duplications,  $D$ , observed on all branches of the tree rather than just the  
190 duplications,  $d_i$ , observed on a single branch:

$$P(o_i = root|D) = \frac{\prod_j P(o_j^{(i)}|d_j)}{\sum_k \prod_j P(o_j^{(k)}|d_j)} \quad (1)$$

191 where  $o_j^{(i)} \in \{\rightarrow, \leftarrow, root\}$  is the orientation of the branch  $j$  that would be implied by the root  
192 of the tree being branch  $i$ . That is, the probability distribution for the root given all the gene  
193 duplication events on the tree can be expressed in terms of the probabilities for the  
194 orientation of each branch given only the gene duplications on that branch;  $P(\rightarrow |d_i)$ ,  
195  $P(\leftarrow |d_i)$  and  $P(root|d_i)$ .



## 196 **Poisson Model for Gene Duplications**

197 To calculate  $P(o_i|d_i)$  the duplications observed on a branch are modelled as arising  
198 from two Poisson processes. One process describes the number of true positive  
199 duplications (corresponding to the actual direction of time along the branch) and the other  
200 describes the number of false positive duplications. Let  $\alpha$  be a parameter giving the  
201 relative frequency of false positives to true positives across all branches of the tree. Then  
202  $m \sim Po(\lambda)$  and  $n \sim Po(\alpha\lambda)$ , where  $\lambda$  is the expected number of true positives on the branch.  
203 We set the total expected number of duplications on the branch from the two Poisson  
204 processes to match the actual number observed,  $N$ . Thus  $\lambda = N/(1 + \alpha)$ . The relative rate  
205 of false positives to true positives across the whole tree can be estimated from the number  
206 conflicting duplications given the maximum parsimony root of the tree. So as not to over-  
207 penalise false-positive duplications, we take  $\alpha$  to be one tenth of the ratio of conflicting to  
208 non-conflicting duplications of the maximum parsimony root.

209 Bayes' rule gives

$$210 \quad P(o_i) = \frac{P(d_i|o_i)P(o_i)}{P(d_i)}$$

211 where  $P(d_i) = \sum_{o \in \{\rightarrow, \leftarrow, r\}} P(d_i|o)P(o)$ . The priors are given by  $P(r) = 1/b$  and  $P(\rightarrow) =$   
212  $P(\leftarrow) = b - 1/2b$ , where  $b=2t-3$  is the number of branches on an unrooted tree with  $t$  taxa.

213 The probability mass function for the Poisson distribution immediately gives  $P(d|\rightarrow)$  and  
214  $P(d|\leftarrow)$ :

$$215 \quad P(d|\rightarrow) = Po(m; \lambda)Po(n; \alpha\lambda)$$
$$216 \quad = \frac{\lambda^m e^{-\lambda}}{m!} \frac{(\alpha\lambda)^n e^{-\alpha\lambda}}{n!}$$

217 and,

$$218 \quad P(d|\leftarrow) = Po(n; \lambda)Po(m; \alpha\lambda)$$

219 
$$= \frac{\lambda^n e^{-\lambda}}{n!} \frac{(\alpha\lambda)^m e^{-\alpha\lambda}}{m!}$$

220 The branch with the root is more complicated since it actually corresponds to two  
 221 branches on the rooted tree we are attempting to recover. On these two branches time  
 222 flows in opposite directions, away from a central root that separates them. We must allow  
 223 for the  $\binom{m}{n}$  duplications on the branch to actually correspond to  $\binom{m-s}{t}$  duplications on  
 224 one of the two branches and  $\binom{n-t}{s}$  on the other branch (with opposite orientation to the  
 225 first). The number of false positive duplications,  $s$  and  $t$ , are unknown and therefore must  
 226 be summed over. Similarly, the location of root could fall at any point along the length of  
 227 the original branch. If the root were a fraction,  $x$ , along the length of the branch then the  
 228 expected rate of false positive and true positive duplications on that fraction of the branch  
 229 would be  $x\lambda$  and  $x\alpha\lambda$  respectively whereas on the other branch the rates would be  $(1-x)\lambda$   
 230 and  $(1-x)\alpha\lambda$ . Thus, integrating over the position of the root along the branch and summing  
 231 over the distribution of the  $\binom{m}{n}$  putative duplications between true positives and false  
 232 positives on the two resulting branches, we find:

233 
$$P(d|r) = \sum_{s=0}^m \sum_{t=0}^n \int_0^1 P_o^T(m-s; x\lambda) P_o^F(t; x\alpha\lambda) P_o^T(n-t; (1-x)\lambda) P_o^F(s; (1-x)\alpha\lambda) dx$$

234 
$$= \sum_{s=0}^m \sum_{t=0}^n B(m-s+t+1, n-t+s+1) \frac{\lambda^{m-s} e^{-\lambda}}{(m-s)!} \frac{(\alpha\lambda)^{n-t} e^{-\alpha\lambda}}{(n-t)!} \frac{\lambda^{s+t} \alpha^{s+2t-n}}{s! t!}$$

235 Where  $B(, )$  is the beta function.

236 The duplications observed in just one species are uninformative as to the location of  
 237 the root and so should not affect the root probabilities produced by the model. As such, the  
 238 branch model for terminal branches is modified to only model the number of inward  
 239 duplications (those supporting the tree minus the species on the terminal branch as a  
 240 monophyletic clade). The rates  $\lambda_{Term,TP}$  and  $\lambda_{Term,FP}$  are the observed true positive and

241 false positive rates for inward duplications on the terminal branches for the maximum  
242 parsimony root. For the terminal branches, the branch model is:

$$243 \quad P^{Term}(d|\rightarrow) = \text{Po}(m; \lambda_{Term,FP})$$

244 and

$$245 \quad P^{Term}(d|r) = \text{Po}(m; \lambda_{Term,TP}).$$

246

247         The branch-level model takes into account only the duplications observed on a  
248 single branch and these probabilities feed into the tree-level model to give the final  
249 probabilities for the position of the root (Fig. 4). The behaviour of the branch model is in  
250 good agreement with an intuitive understanding of the probabilities that should be  
251 assigned to the three possible orientations for a branch given the number of putative  
252 duplications observed in either direction (Fig. 4A-C). The probability of time flowing to the  
253 left/right increases monotonically with the number of putative duplications supporting it.  
254 The probability of a branch being a root is highest when the number of putative gene  
255 duplications in both directions is the same. Finally, the probability of a branch being a root  
256 remains significantly above zero if there is any number of gene duplications in both  
257 directions (Fig. 4B & C). This reflects the fact that putative gene duplications supporting  
258 the monophyletic nature of both blocks of a bipartition support that bipartition being the  
259 root. The fact that there could be a large difference in the number of gene duplications in  
260 one direction compared to the other due to different branch lengths on the two sides of the  
261 root is accounted for by integrating over the position of the root along the original root  
262 branch. Thus, the probability of a branch being a root is > 30% when there are 20  
263 duplications in one direction compared to 5 in the opposite direction (Fig. 4C). For  
264 comparison, the probability of the orientation of the branch being in the direction of the 5  
265 duplications is vanishingly small ( $\sim 10^{-13}$ ). The branch-level probability model thus gives

266 probabilities for each branch taking into account only the duplications observed on that  
267 branch. The final probabilities for the root of the tree, taking into account all duplications  
268 across the tree are then given by the tree-level model (Equation 1, Fig. 4D & E).

### 269 ***Algorithm implementation and availability***

270 STRIDE is implemented in python. Further information, use instructions, an example  
271 dataset, and a standalone implementation of the algorithm is available under the University  
272 of Oxford Academic Use Licence at <https://github.com/davidemms/STRIDE>. The complete  
273 set of gene trees and species trees required to replicate this analysis are provided for  
274 download from the Zenodo research data archive at  
275 <https://doi.org/10.5281/zenodo.581360>.

### 276 ***Results***

#### 277 ***STRIDE identifies the correct root of species trees given simulated gene tree***

#### 278 ***datasets***

279 The ability of STRIDE to correctly infer the root of a known species tree was tested  
280 using three published, simulated gene tree datasets. The first dataset consisted of 2000  
281 simulated gene trees from 40 species with heterogeneous rates of gene duplication and  
282 loss within trees (Boussau, B., Szollosi, G.J., et al. 2013). The second and third datasets  
283 consisted of 12000 gene trees from 12 species and 7500 gene trees from 17 species,  
284 respectively. These two datasets were similar to the first dataset but also incorporated  
285 incomplete lineage sorting generated using a range of effective population sizes (Wu,  
286 Y.C., Rasmussen, M.D., et al. 2014). Since incomplete lineage sorting can lead to  
287 misidentification of gene duplication and loss events these latter two datasets provided a  
288 good test of STRIDE's robustness in the face of gene-tree/species-tree incongruence. For  
289 all three simulated datasets, STRIDE correctly inferred the root of the species tree and

290 assigned it a probability of 100% (Table 1, Supplementary File 1. Fig. S1-S3). Thus for  
291 these simulated datasets the method performed well.

### 292 ***Application of STRIDE to real species datasets***

293 Simulated datasets generally do not capture all the nuances and difficulties seen in  
294 real biological datasets. These nuances include errors in orthogroup inference, alignment  
295 inference and gene tree inference. Thus to demonstrate the utility of STRIDE, a diverse  
296 range of groups of species were sampled from throughout the eukaryotic domain (Table  
297 1). This included every group of eukaryotes on Ensembl Genomes containing more than 4  
298 genera (Yates, A., Akanni, W., et al. 2016). To expand this group of tests, additional sets  
299 of genomes were obtained for 47 Birds (Jarvis, E.D., Mirarab, S., et al. 2014), 42 Green  
300 Plants (Goodstein, D.M., Shu, S.Q., et al. 2012) and 16 Kinetoplastids (Aslett, M.,  
301 Aurrecochea, C., et al. 2010). In total, this gave 12 species groups with varying levels of  
302 taxon sampling and with estimated divergence times ranging from c. 56 million years for  
303 the Primates (dos Reis, M., Donoghue, P.C.J., et al. 2014) to c. 1500 million years for the  
304 Green Plants (Parfrey, L.W., Lahr, D.J.G., et al. 2011). These species sets thus provided a  
305 diverse group with which to test the utility of STRIDE. Furthermore, for each of these  
306 species sets, there is an accepted consensus on the topology and location of the root of  
307 the species tree (Supplemental File 1). In all cases these topologies and root branches  
308 were assumed to be true when STRIDE's performance was assessed. On average, across  
309 each of the simulated and real dataset in this analysis STRIDE took ~18 seconds to run  
310 using four cores of an Intel Core i7-4770 3.4GHz CPU.

311 Orthogroups for each species set were inferred using OrthoFinder (Emms, D.M.  
312 and Kelly, S. 2015), and gene trees for each orthogroup were inferred using IQTREE  
313 v1.5.3 (Nguyen, L.T., Schmidt, H.A., et al. 2015) from a multiple sequence alignment  
314 generated using MAFFT L-INS-I v7.305b (Kato, K. and Standley, D.M. 2013). For each  
315 species set, STRIDE was run with a published unrooted species tree (without branch

316 lengths) and the complete set of gene trees inferred from all orthogroups identified by  
317 OrthoFinder. The number species, gene trees, informative duplications and other details  
318 are provided in Table 1.

319 In all 12 test cases, there is a single maximum parsimony root. In 9 of the 12 tests  
320 this root agreed with the accepted root of the species set (Table 1). Figures 5 to 7 present  
321 the results of the STRIDE analysis applied to the plant, fungi, and bird datasets. These  
322 datasets correspond to the largest, median and smallest number of informative  
323 duplications per species identified by STRIDE. The results for the remaining datasets can  
324 be found in Supplemental File 1 Figures S4-S12. For the plant dataset, sufficient gene-  
325 duplication events were identified for the probability model to assign a probability of 100%  
326 to the accepted root separating the algae from the land plants (Ruhfel et al. BMC  
327 Evolutionary Biology 2014 14:23) (Fig. 5). A probability of 100% was also assigned for the  
328 correct root in the fungi, even though fewer informative gene duplication events were  
329 identified (Fig. 6, Table 1). In both the plant and fungal datasets, STRIDE also identified  
330 substantial numbers of gene duplication events that support sub-clades within both  
331 species trees (Fig. 5 and Fig. 6).

332 While STRIDE identified the community-accepted root in 75% of the datasets, it  
333 failed to identify this root for the bird (Fig. 7), rodent and Laurasiatheria (Supplementary  
334 File 1 Fig. S11 & S12) datasets. These three datasets had the smallest, second smallest  
335 and fourth smallest number of informative gene duplication events per species respectively  
336 (Table 1). In addition, while there were no conflicting gene duplication events in the bird  
337 dataset, the rodent and Laurasiatheria datasets had the highest and fifth highest ratio of  
338 conflicting to informative duplications (Table 1). Consistent with these observations,  
339 analysis of the factors affecting the accuracy of STRIDE revealed that root probability  
340 assignment was positively correlated with the number of informative duplications per  
341 species ( $R^2=0.17$ , Supplementary File 1 Fig. S13A) and negatively correlated with the

342 proportion of duplications which were in conflict ( $R^2=0.24$ , Supplementary File 1 Fig.  
343 S13B). Furthermore, the proportion of conflicting duplications was negatively correlated  
344 with the number of species ( $R^2=0.36$ , Supplementary File 1 Fig. S13C), suggesting  
345 increased taxon sampling facilitated more accurate identification of gene duplication  
346 events. Thus the ability of STRIDE to detect the true root is affected by taxon sampling  
347 and the number of gene duplication events detected in the dataset.

### 348 ***STRIDE Provides Evidence for Location of the Root of the Eukaryotic Tree***

349         Given the performance of stride on the datasets outlined above it was assessed  
350 whether STRIDE could provide insight into one of the most contentious and difficult tree  
351 rooting problems in biology, the root of the eukaryotic tree (Burki, F. 2014). Here, a set of  
352 45 species that were distributed across the eukaryotic tree were selected. These were  
353 subject to orthogroup and gene tree inference as before and the complete set of 16770  
354 gene trees were submitted for analysis by STRIDE. This identified 2316 putative gene  
355 duplication events excluding the root from (and supporting the monophyly of) major clades  
356 within the eukaryotes including the opisthokonta, fungi, metazoa, and achiplastida (Fig.  
357 8A). Duplication events supporting further subclades within these major groupings were  
358 also abundant (Fig. 8A). In contrast, other major sub-clades including amoebazoa, the  
359 SAR supergroup, and the excavata, did not receive support from gene duplication events  
360 (Fig. 8A). This lack of gene duplication events meant that STRIDE could not exclude the  
361 root of the species tree from the basal branches of these groups and thus could not  
362 provide evidence for or against the five most popular placements for the root of the  
363 eukaryotic tree (Burki, F. 2014). This ambiguity in root assignment is represented  
364 effectively in the probabilities assigned to all putative root-spanning branches (Fig. 8B).

### 365 ***Discussion***

366         STRIDE is an automated method for identifying and analysing gene duplication  
367 events to infer the root of species trees. Through analysis of simulated and real datasets,

368 we show how the performance of STRIDE is affected by data quantity, data conflict, and  
369 taxon sampling. Furthermore, we demonstrate that STRIDE is effective in identifying the  
370 root of species trees for the majority of species datasets and effectively captures the  
371 ambiguity in root assignment given the input data.

372         The aim of STRIDE is to infer a probability distribution over an entire species tree  
373 for the location of its root. This aim is different from algorithms that attempt to reconcile  
374 gene trees with species trees (Szollosi, G.J., Tannier, E., et al. 2015) or model duplication  
375 and loss processes on a tree (Gorecki, P. and Eulenstein, O. 2014). STRIDE identifies and  
376 utilises well-supported gene duplication events and does not evaluate gene loss events for  
377 the following reasons. First, gene trees can distinguish parallel duplication events on  
378 adjacent branches from a single shared duplication event, which is not possible for gene  
379 loss events. Second, the topology of the gene tree post-duplication genes can be  
380 compared with the species tree to confirm the accuracy of the inference, this cannot be  
381 done with gene loss events. Third, most genomes are incomplete and vary considerably in  
382 the quality of their annotation leading to high rates of false positive gene loss (Veeckman,  
383 E., Ruttink, T., et al. 2016, Dunne, M.P. and Kelly, S. 2017).

384         A major advantage of using STRIDE is that sets of species can be analysed without  
385 the inclusion on an outgroup. This is potentially advantageous in situations where inclusion  
386 of an outgroup can effect the topology of gene trees inferred for the ingroup species  
387 (Berger, S.A., Krompass, D., et al. 2011). Moreover, if the outgroup is distantly related to  
388 the ingroup species then additional problems of long branch attraction can lead to incorrect  
389 root placement (Philippe, H., Brinkmann, H., et al. 2011, Kuck, P., Mayer, C., et al. 2012,  
390 Salichos, L. and Rokas, A. 2013). STRIDE is also suitable for large dataset analysis and  
391 for situations where appropriate outgroups are not available. Although STRIDE as  
392 presented is a method for identifying the root of an unrooted species tree, the output from  
393 STRIDE can provide a wealth of useful information. For example, STRIDE maps high



394 confidence gene duplication events to branches in a species tree. These gene duplication  
395 events provide strong evidence for monophyly of the species that share the gene  
396 duplication event. Thus STRIDE can be used to provide additional support to branches in  
397 a species tree that might be weakly supported by molecular sequence data. In this context,  
398 it is worth noting that STRIDE could also be used to evaluate support for alternative  
399 species-tree topologies by providing support for clades from gene duplication events.

400 The application of STRIDE to the eukaryotes was able to exclude the root of the  
401 eukaryotes from the opisthokonts and from a number of other groups, however STRIDE  
402 was unable to uniquely place the root as there were insufficient gene duplication events  
403 identified that could exclude the root from other portions of the tree. It is likely that poor  
404 taxon sampling for some of the groups (e.g. the amoebozoa and excavata), coupled with  
405 genome reduction associated with adaptation to parasitism in many of these species,  
406 impeded the discovery of these gene duplication events. With improved taxon sampling  
407 STRIDE may ultimately be able to provide further insight as to the location of the root of  
408 the eukaryotic tree. Furthermore, as STRIDE produces branch-level probabilities these  
409 could be combined with probabilities obtained from other analyses to perform a multi-data-  
410 type analysis of the origin of the eukaryotes.

411 In summary, STRIDE is a fast and effective method for genome scale phylogenetic  
412 analysis that can be used both to identify high confidence gene duplication events and  
413 identify the root of species trees without the requirement for an outgroup.

## 414 **References**

415 Aslett M, Aurrecochea C, Berriman M, Brestelli J, Brunk BP, Carrington M, Depledge DP,  
416 Fischer S, Gajria B, Gao X, *et al.* 2010. TriTrypDB: a functional genomic resource for the  
417 Trypanosomatidae. *Nucleic Acids Res*, 38:D457-D462.

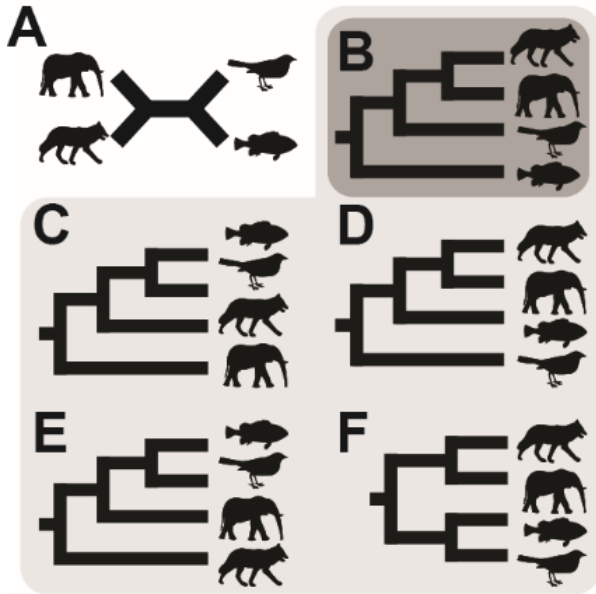
- 418 Berger SA, Krompass D, Stamatakis A. 2011. Performance, Accuracy, and Web Server for  
419 Evolutionary Placement of Short Sequence Reads under Maximum Likelihood. *Syst Biol*,  
420 60:291-302.
- 421 Boussau B, Szollosi GJ, Duret L, Gouy M, Tannier E, Daubin V. 2013. Genome-scale  
422 coestimation of species and gene trees. *Genome Res*, 23:323-330.
- 423 Burki F. 2014. The Eukaryotic Tree of Life from a Global Phylogenomic Perspective. *Csh*  
424 *Perspect Biol*, 6.
- 425 Dobzhansky T. 2013. Nothing in Biology Makes Sense Except in the Light of Evolution.  
426 *Am Biol Teach*, 75:87-91.
- 427 Donoghue MJ, Mathews S. 1998. Duplicate genes and the root of angiosperms, with an  
428 example using phytochrome sequences. *Mol Phylogenet Evol*, 9:489-500.
- 429 dos Reis M, Donoghue PCJ, Yang ZH. 2014. Neither phylogenomic nor palaeontological  
430 data support a Palaeogene origin of placental mammals. *Biol Letters*, 10.
- 431 Dunne MP, Kelly S. 2017. OrthoFiller: utilising data from multiple species to improve the  
432 completeness of genome annotations. *bioRxiv*.
- 433 Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome  
434 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol*, 16.
- 435 Felsenstein J. 1981. Evolutionary Trees from DNA-Sequences - a Maximum-Likelihood  
436 Approach. *J Mol Evol*, 17:368-376.
- 437 Goodstein DM, Shu SQ, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T, Dirks W,  
438 Hellsten U, Putnam N, *et al.* 2012. Phytozome: a comparative platform for green plant  
439 genomics. *Nucleic Acids Res*, 40:D1178-D1186.
- 440 Gorecki P, Eulenstein O. 2014. DrML: Probabilistic Modeling of Gene Duplications. *J*  
441 *Comput Biol*, 21:89-98.

- 442 Huelsenbeck JP, Bollback JP, Levine AM. 2002. Inferring the root of a phylogenetic tree.  
443 Syst Biol, 51:32-43.
- 444 Huerta-Cepas J, Serra F, Bork P. 2016. ETE 3: Reconstruction, Analysis, and  
445 Visualization of Phylogenomic Data. Mol Biol Evol, 33:1635-1638.
- 446 Huson DH, Scornavacca C. 2012. Dendroscope 3: An Interactive Tool for Rooted  
447 Phylogenetic Trees and Networks. Syst Biol, 61:1061-1067.
- 448 Jarvis ED, Mirarab S, Aberer AJ, Li B, Houde P, Li C, Ho SYW, Faircloth BC, Nabholz B,  
449 Howard JT, *et al.* 2014. Whole-genome analyses resolve early branches in the tree of life  
450 of modern birds. Science, 346:1320-1331.
- 451 Katoh K, Standley DM. 2013. MAFFT Multiple Sequence Alignment Software Version 7:  
452 Improvements in Performance and Usability. Mol Biol Evol, 30:772-780.
- 453 Kuck P, Mayer C, Wagele JW, Misof B. 2012. Long Branch Effects Distort Maximum  
454 Likelihood Phylogenies in Simulations Despite Selection of the Correct Model. Plos One,  
455 7.
- 456 Larsson A. 2014. AliView: a fast and lightweight alignment viewer and editor for large  
457 datasets. Bioinformatics, 30:3276-3278.
- 458 Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: A Fast and Effective  
459 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. Mol Biol Evol,  
460 32:268-274.
- 461 Parfrey LW, Lahr DJG, Knoll AH, Katz LA. 2011. Estimating the timing of early eukaryotic  
462 diversification with multigene molecular clocks. Proceedings of the National Academy of  
463 Sciences of the United States of America, 108:13624-13629.

- 464 Philippe H, Brinkmann H, Lavrov DV, Littlewood DTJ, Manuel M, Worheide G, Baurain D.  
465 2011. Resolving Difficult Phylogenetic Questions: Why More Sequences Are Not Enough.  
466 Plos Biol, 9.
- 467 Salichos L, Rokas A. 2013. Inferring ancient divergences requires genes with strong  
468 phylogenetic signals. Nature, 497:327-+.
- 469 Simmons MP, Bailey CD, Nixon KC. 2000. Phylogeny reconstruction using duplicate  
470 genes. Mol Biol Evol, 17:469-473.
- 471 Sytsma KJ, Pires JC. 2001. Plant systematics in the next 50 years - re-mapping the new  
472 frontier. Taxon, 50:713-732.
- 473 Szollosi GJ, Tannier E, Daubin V, Boussau B. 2015. The Inference of Gene Trees with  
474 Species Trees. Syst Biol, 64:E42-E62.
- 475 Tange O. 2011. GNU Parallel - The Command-Line Power Tool. ;login: The USENIX  
476 Magazine, 36:42-47.
- 477 Veeckman E, Ruttink T, Vandepoele K. 2016. Are We There Yet? Reliably Estimating the  
478 Completeness of Plant Genome Sequences. Plant Cell, 28:1759-1768.
- 479 Williams TA, Heaps SE, Cherlin S, Nye TMW, Boys RJ, Embley TM. 2015. New  
480 substitution models for rooting phylogenetic trees. Philos T R Soc B, 370.
- 481 Wu YC, Rasmussen MD, Bansal MS, Kellis M. 2014. Most parsimonious reconciliation in  
482 the presence of gene duplication, loss, and deep coalescence using labeled coalescent  
483 trees. Genome Res, 24:475-486.
- 484 Yates A, Akanni W, Amode MR, Barrell D, Billis K, Carvalho-Silva D, Cummins C,  
485 Clapham P, Fitzgerald S, Gil L, *et al.* 2016. Ensembl 2016. Nucleic Acids Res, 44:D710-  
486 D716.
- 487

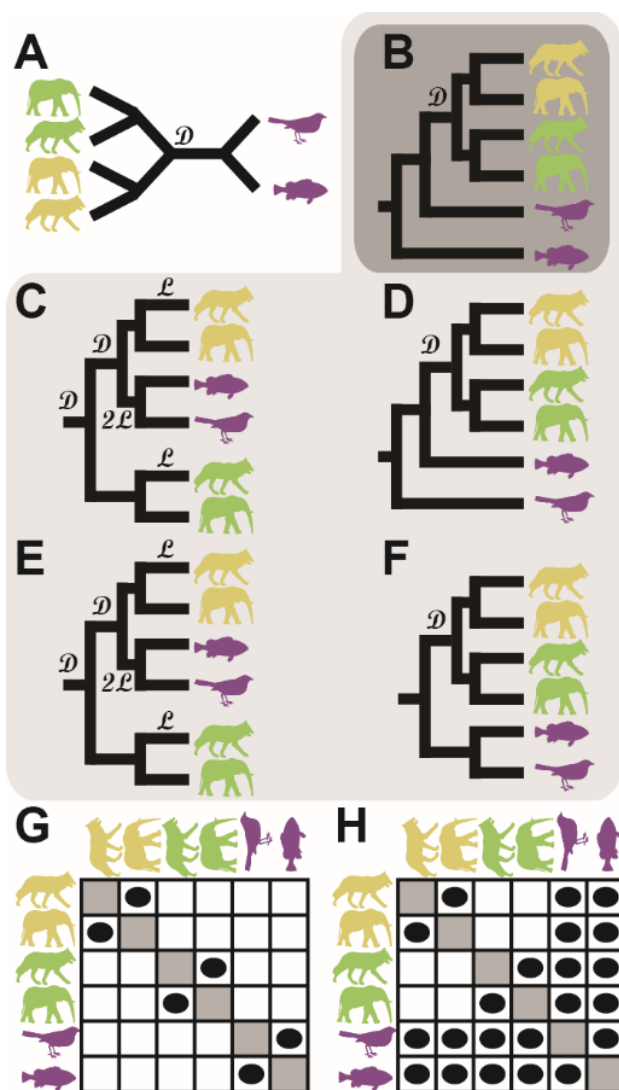
488 **Figures**

489 **Figure 1**



490

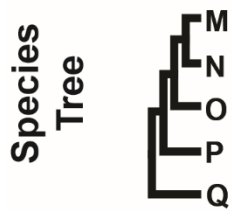
491 **Figure 2**



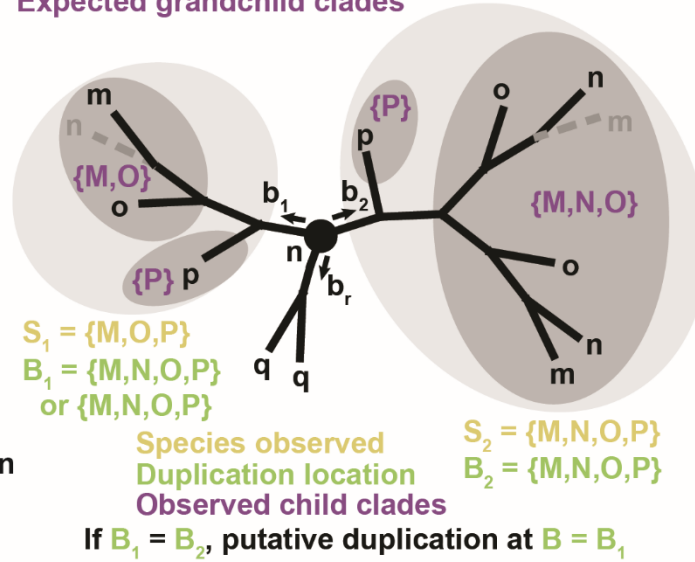
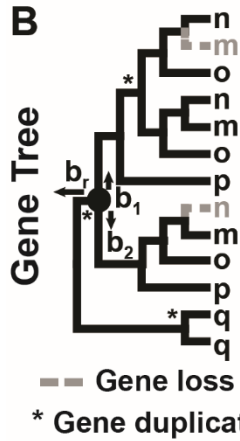
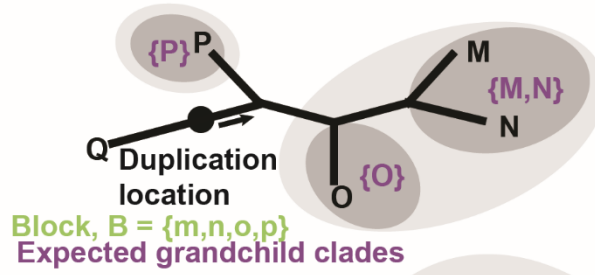
492

493 **Figure 3**

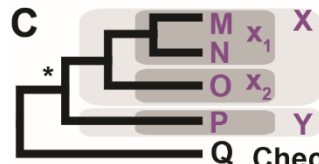
**A Unobserved, rooted trees**



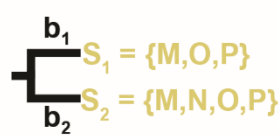
**Observed, unrooted trees**



**Expected**



**Observed**



**Check**

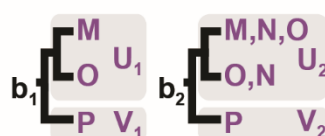
For  $i=1,2$

$X_1 \cap S_i \neq \emptyset$

$X_2 \cap S_i \neq \emptyset$

$Y \cap S_i \neq \emptyset$

Check presence of expected clades



For  $i = 1,2$

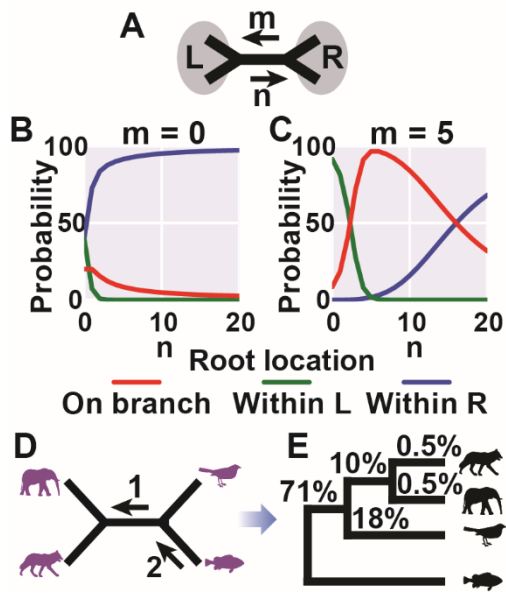
$U_i \subseteq X, V_i \subseteq Y$

Or,  $U_i \subseteq Y, V_i \subseteq X$

Check local topology

494

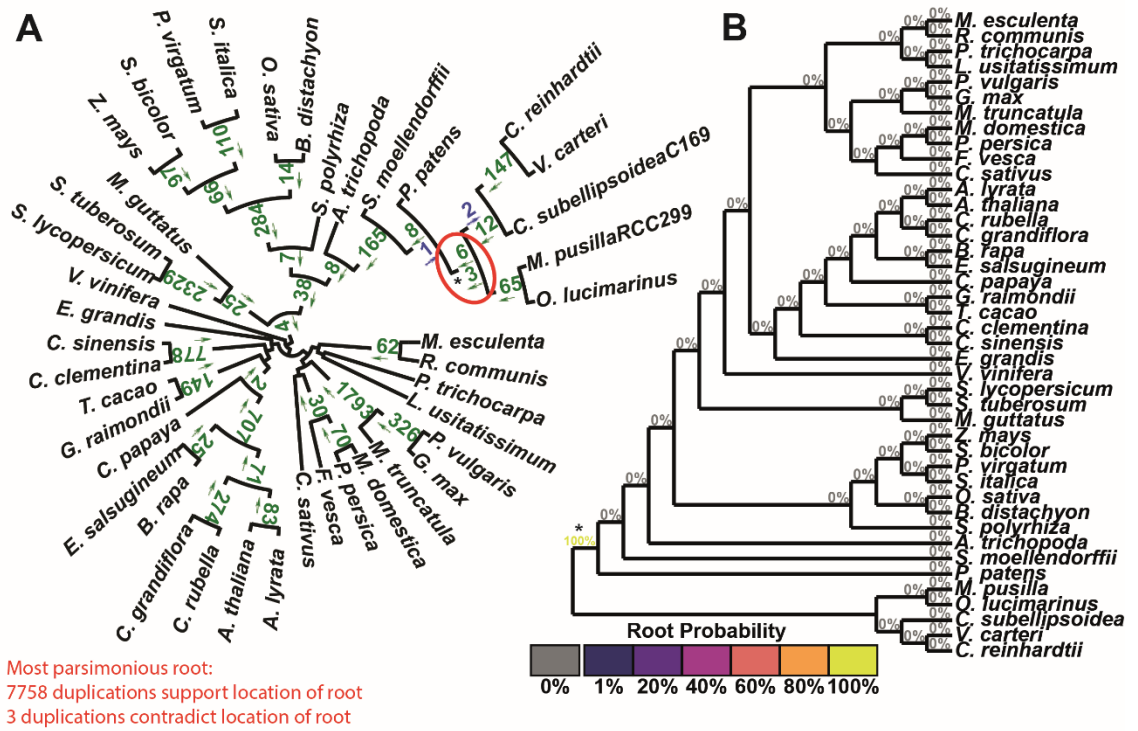
495 **Figure 4**



496

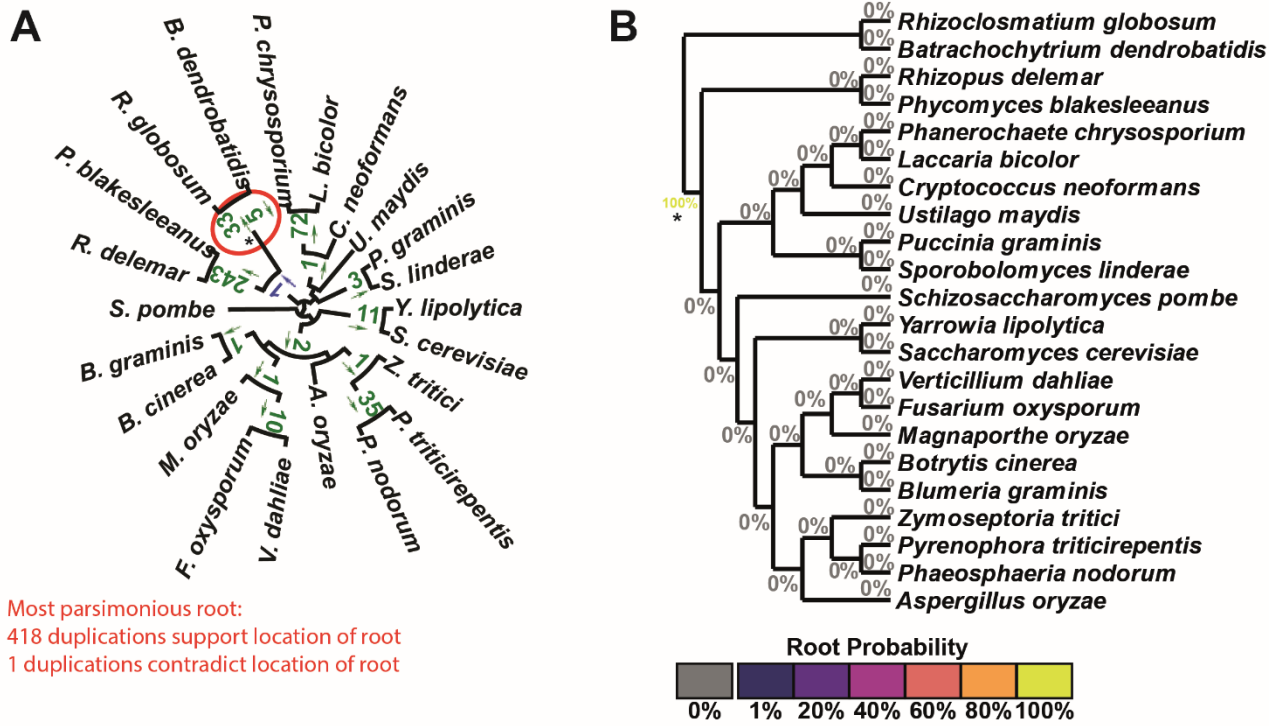


497 **Figure 5**



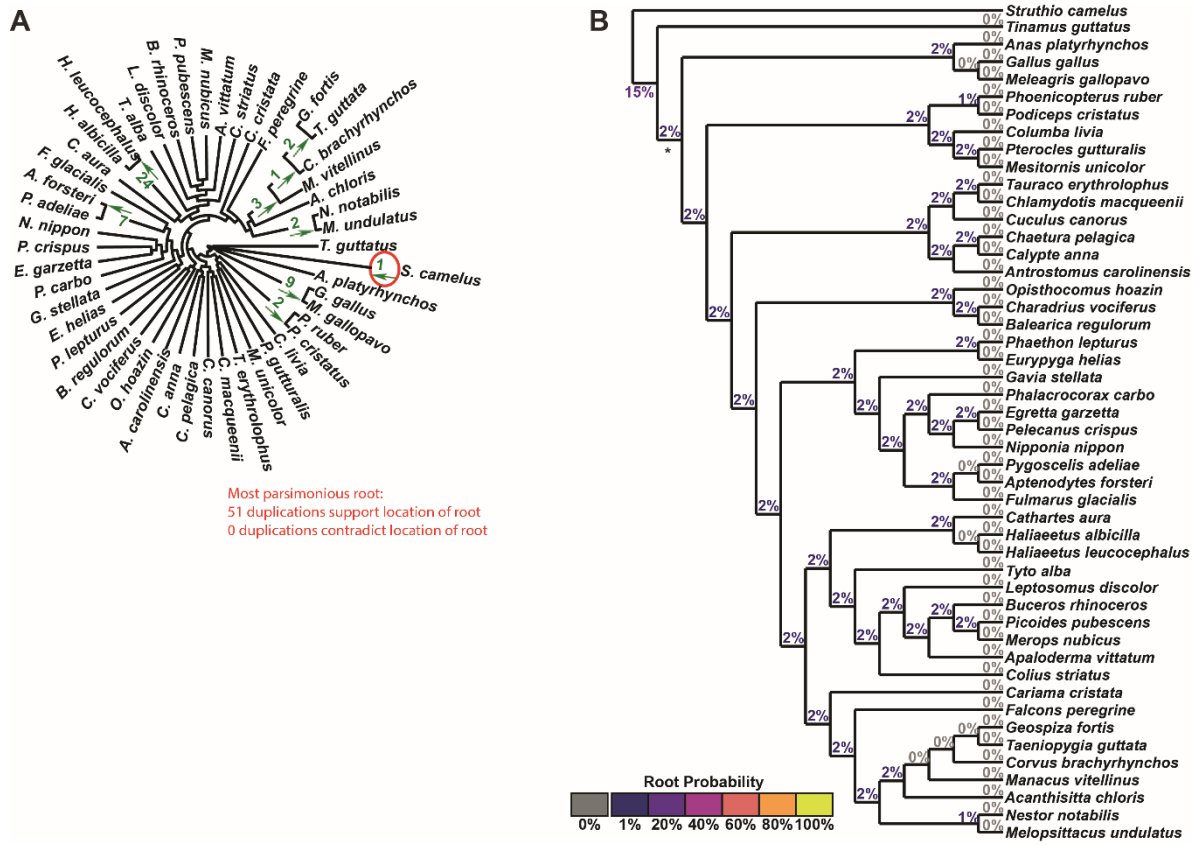
498

499 **Figure 6**



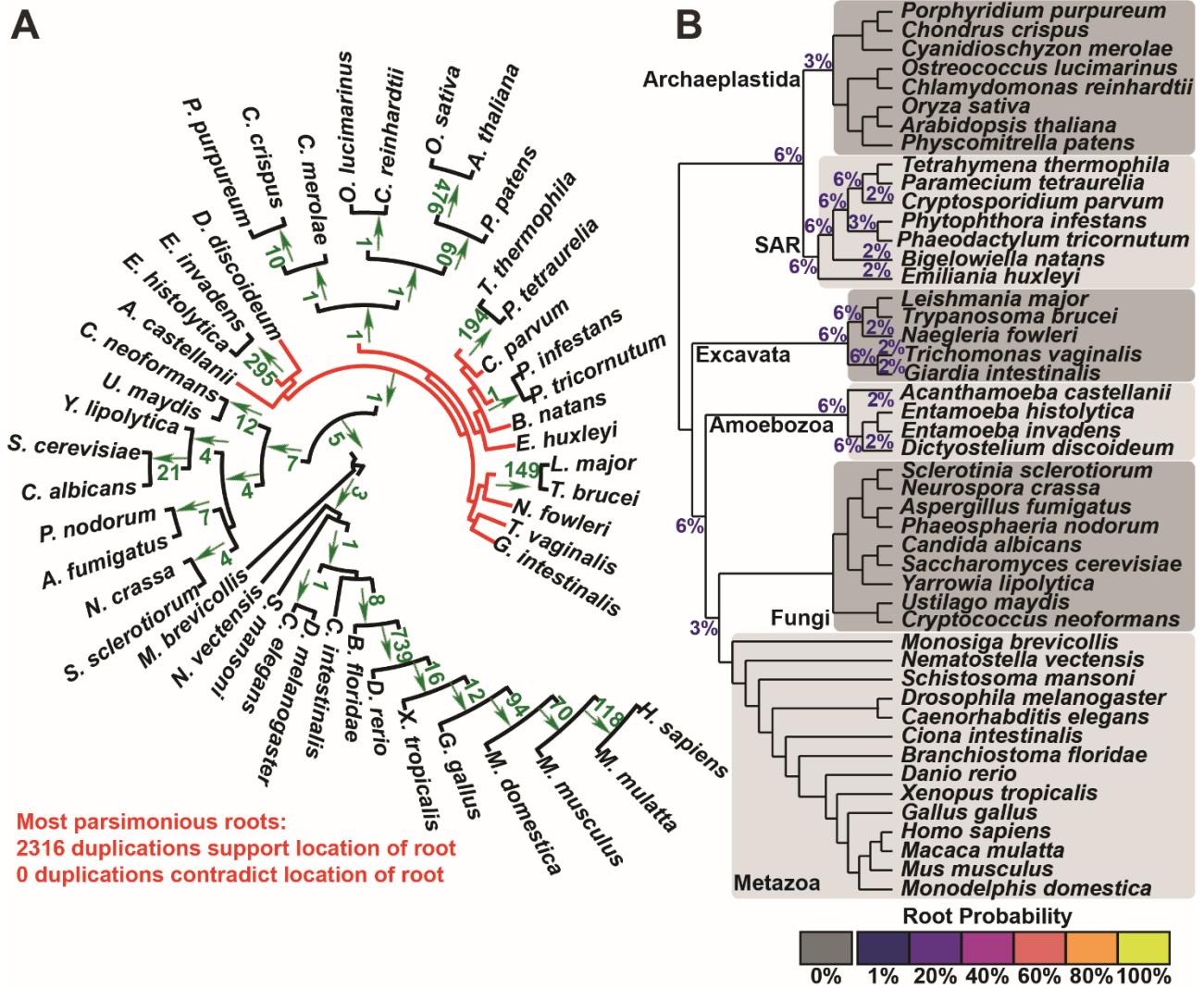
500

501 **Figure 7**



502

503 **Figure 8**



504

505 **Tables**  
506 **Table 1**

Group	Species	Gene Trees	Informative Duplications	Number of Conflicting Duplications	Number of MP Roots	Correct MP Root	Probability for MP Root	Probability for Correct Root
<b>Metazoa + outgroup (simulated)</b>	40	2000	664	0	1	yes	100.0%	100.0%
<b>Drosophila (simulated)</b>	12	12000	1360	1	1	yes	100.0%	100.0%
<b>Primates (simulated)</b>	17	7500	1593	0	1	yes	100.0%	100.0%
<b>Birds (Aves)</b>	47	14454	51	0	1	no	15.0%	2.0%
<b>Flies (Diptera)</b>	7	11688	279	11	1	yes	100.0%	100.0%
<b>Fish</b>	11	16520	650	7	1	yes	100.0%	100.0%
<b>Fungi</b>	21	9325	419	1	1	yes	100.0%	100.0%
<b>Hymenoptera</b>	5	9157	108	7	1	yes	100.0%	100.0%
<b>Kinetoplastids</b>	16	9731	76	4	1	yes	55.0%	55.0%
<b>Laurasiatheria</b>	14	15804	135	7	1	no	100.0%	0.0%
<b>Metazoa</b>	21	13017	2065	0	1	yes	48.0%	48.0%
<b>Nematoda</b>	7	8392	93	2	1	yes	100.0%	100.0%
<b>Primates + outgroup</b>	11	19096	117	11	1	yes	8.0%	8.0%
<b>Rodents</b>	7	15485	22	6	1	no	9.0%	0.5%
<b>Plants</b>	42	28356	7761	3	1	yes	100.0%	100.0%
<b>Eukaryotes</b>	45	16770	2316	0	25	-	-	-
<b>Total simulated</b>	69	21500	3617	1	-	3	-	-
<b>Total real-world</b>	254	187795	14092	59	-	9	-	-
<b>Total</b>	323	209295	17709	60	-	12	-	-

507

508

509

## 510 **Figure Legends**

### 511 **Figure 1**

512 Possible roots for a four-taxa species tree. A) Unrooted species tree for four species:  
513 elephant, wolf, fish & bird. B) The correct rooting of the species tree. B-F) The five possible  
514 rooted species trees for the unrooted species tree in A.

### 515 **Figure 2**

516 Orthologues inferred from gene trees depend on the root. A) An unrooted gene tree  
517 corresponding to an orthogroup with a gene duplication event in the common ancestor of  
518 wolf and elephant. Genes from each species are represented by an image of the species.  
519 B-F) The most parsimonious rootings of the gene trees (fewest duplications and losses) for  
520 each of the five roots of the species tree, as shown in Figure 1B-F.  $\mathcal{D}$  - gene duplication  
521 event,  $\mathcal{L}$  - gene loss event. G) Orthologues inferred from the incorrect trees D & E. H)  
522 Orthologues inferred from the correctly rooted tree B and also the close to correctly rooted  
523 trees D & F.

### 524 **Figure 3**

525 Identification of well-supported gene duplication events. Upper case letters M,N,O,P & Q  
526 are species, lower case m,n,o,p & q are genes from the corresponding species. A) The  
527 unknown, rooted species tree (left) and the observed, unrooted species tree (right). Black  
528 dot and arrow show the location of a single hypothetical gene duplication event on a  
529 branch with time flowing in the direction indicated by the arrow. The branch location and  
530 direction is uniquely identified by the block, B, of species whose common ancestor would  
531 have inherited the duplicate genes. The expected species in the child clades (X & Y) and  
532 grandchild clades ( $x_1$ ,  $x_2$  & Y) after this hypothetical duplication are highlighted with  
533 light/dark grey ellipses respectively. B) The unknown, rooted gene tree (left) and the  
534 observed, unrooted gene tree (right) for a hypothetical gene family with three gene

535 duplication events (marked by \*) and two gene loss events (grey, dotted line). The node  
536 currently being analysed is  $n$  and  $b_r$  is the current, tentative direction towards the root. For  
537 these  $n$  and  $b_r$ ,  $b_1$  and  $b_2$  are analysed to see if they are well-supported gene duplication  
538 branches.  $S_i$  is the set of species below branch  $b_i$ ,  $B_i$  is the smallest block of a bipartition  
539 containing  $S_i$  ( $i=1,2$ ) C) The check that genes from each of the expected grandchild clades  
540 are present on each duplicate branch D) The check that the local topology for each  
541 duplication branch agrees with expected topology.  $U_i$  and  $V_i$  are the observed child clades  
542 on branch  $b_i$ . The observed child clades should not contain genes from any species not in  
543 the expected child clades

#### 544 **Figure 4**

545 The branch-level probability model employed by STRIDE. These branch-level probabilities  
546 are used by the tree probability model to give the overall probabilities for the location of the  
547 root of the species tree. A) A single branch in the tree with  $m/n$  duplications supporting L/R  
548 as monophyletic clades. B) Branch-level model probabilities for position of the root with  
549 respect to the branch when  $m=0$  (the model only takes into account duplications on that  
550 branch). C) As for B with  $m=5$ . D) Hypothetical total number of gene duplication events on  
551 the 4 species phylogeny. One gene duplication event is shared by elephant and dog and 2  
552 are shared by elephant, dog and bird. D) The final tree-level model probabilities for the  
553 location of the root calculated by STRIDE taking into account all the gene duplication  
554 events on all branches in D.

#### 555 **Figure 5**

556 STRIDE analysis applied the set of plant gene trees. A) Numbers of identified gene  
557 duplication events are marked on the branches they are observed on and arrows indicate  
558 the direction in which the duplication occurred. Gene duplication events are in agreement  
559 with the maximum parsimony root of the tree if the arrow points away from the root, and



560 are in green. Those that disagree are in blue. The maximum parsimony root is circled in  
561 red and is in agreement with the correct root, marked with a \*. B) The probabilities for the  
562 location of the root calculated by STRIDE.

563 **Figure 6**

564 STRIDE analysis applied the set of fungi gene trees. A) Numbers of identified gene  
565 duplication events are marked on the branches they are observed on and arrows indicate  
566 the direction in which the duplication occurred. Gene duplication events are in agreement  
567 with the maximum parsimony root of the tree if the arrow points away from the root, and  
568 are in green. Those that disagree are in blue. The maximum parsimony root is circled in  
569 red and is in agreement with the correct root, marked with a \*. B) The probabilities for the  
570 location of the root calculated by STRIDE.

571 **Figure 7**

572 STRIDE analysis applied the set of Bird gene trees. A) Numbers of identified gene  
573 duplication events are marked on the branches they are observed on and arrows indicate  
574 the direction in which the duplication occurred. Gene duplication events are in agreement  
575 with the maximum parsimony root of the tree if the arrow points away from the root, and  
576 are in green. Those that disagree are in blue. The maximum parsimony root is circled in  
577 red and is in agreement with the correct root, marked with a \*. B) The probabilities for the  
578 location of the root calculated by STRIDE, coloured according to the displayed heat map.

579 **Figure 8**

580 STRIDE analysis applied the set of Eukaryotic gene trees. A) Numbers of identified gene  
581 duplication events are marked on the branches they are observed on and arrows indicate  
582 which block of the bipartition the duplicate genes occur in. None of the gene duplication  
583 events contradict each other. The maximum parsimony roots have red branches, the



584 branches from which the root is excluded are black. B) The probabilities for the location of  
585 the root calculated by STRIDE. Major groups of species are marked.