

1
2
3
4
5
6
7
8

Towards Consensus Gene Ages

Benjamin J. Liebeskind^{1,2*}, Claire D. McWhite¹, Edward M. Marcotte¹

¹Center for Systems and Synthetic Biology, Institute for Cellular and Molecular Biology, & Department of Molecular Biosciences, University of Texas at Austin, Austin, TX 78712

²Center for Computational Biology and Bioinformatics, University of Texas at Austin, TX 78712

*To whom correspondence should be addressed: bliebeskind@austin.utexas.edu

9 **Abstract**

10 Correctly estimating the age of a gene or gene family is important for a variety of fields, including
11 molecular evolution, comparative genomics, and phylogenetics, and increasingly for systems biology and
12 disease genetics. However, most studies use only a point estimate of a gene's age, neglecting the
13 substantial uncertainty involved in this estimation. Here, we characterize this uncertainty by investigating
14 the effect of algorithm choice on gene-age inference and calculate consensus gene ages with attendant
15 error distributions for a variety of model eukaryotes. We use thirteen orthology inference algorithms to
16 create gene-age datasets and then characterize the error around each age-call on a per-gene and per-
17 algorithm basis. Systematic error was found to be a large factor in estimating gene age, suggesting that
18 simple consensus algorithms are not enough to give a reliable point estimate. We also found that different
19 sources of error can affect downstream analyses, such as gene ontology enrichment. Our consensus gene-
20 age datasets, with associated error terms, are made fully available at so that researchers can propagate this
21 uncertainty through their analyses (<https://github.com/marcottelab/Gene-Ages>).

22

23 **Introduction**

24 From their inception, whole genome datasets have been used to infer the evolutionary history of gene
25 families [1]. The age of a gene family, its provenance, and its evolutionary history, such as loss and
26 duplication events, can inform us about its function [2]. For instance, gene age has been found to correlate
27 with disease-association [3,4], evolutionary rate [5], and the number of associated protein-interaction
28 partners [6], and a gene's phylogenetic distribution can be used to infer aspects of its function [7]. Gene-
29 ages can also be used to estimate the gene content of ancient organisms, such as the last universal
30 common ancestor (LUCA, [1]), or the last eukaryotic common ancestor (LECA, [8,9]). Accordingly, an
31 analysis of gene family ages on a genomic scale can inform the phylogenetic history of important
32 phenotypes, such as eyes or the nervous system [10,11]. In more recent years, gene age has been used to
33 annotate systems biology datasets [12–14], with the promise of elucidating the evolutionary history of
34 core cellular machinery.

35 Such studies rely first and foremost upon the correct identification of homologs and/or orthologs.
36 These two relationships form the basis of the gene-age determination in nearly all studies, with orthology
37 being the more common criterion [3,15]. Orthology is a pairwise relationship between two genes that
38 occurs when their most recent common ancestor (MRCA) lies at a speciation event in a phylogenetic gene
39 tree. This is in contrast to paralogs, whose MRCA lies at a gene duplication event (nodes on gene trees
40 represent either speciation or gene duplication events, barring horizontal gene transfer)[16,17]. Orthologs
41 tend to display higher functional conservation than paralogs (though perhaps only weakly [18] - see [19]
42 for a review), hence their use as a basis of cross-species comparison. Typically, studies of gene age will
43 consider an orthologous group to be all the descendent lineages of the deepest speciation node, or the
44 divergence between the two most distant homologs, if that is the criterion being used, as in
45 “phylostratigraphy” [4]. Then, the age of the gene group is defined as the MRCA of the species found in
46 that group.

47 Inferring a gene family’s age thus relies on the accuracy of orthology assignment, but inferring
48 correct orthologs is notoriously difficult, with no one of the more than 30 algorithms out-performing all
49 others [20]. In particular, algorithms differ strongly in the tradeoff between recall and precision [20]. Yet
50 most studies on gene age rely on only one kind of algorithm, either using a pre-existing method or
51 establishing an *ad hoc* protocol, most of which resemble one of the pre-existing algorithms [3]. Although
52 methods for probabilistic orthology assignment do exist [21], available methods are not currently scalable
53 to large genomic datasets using protein sequences, and at any rate still rely on a preliminary clustering
54 step to infer gene families. Consensus algorithms also exist, some of which seem to substantially improve
55 performance on established benchmarks [22,23]. However, these methods still give only a point estimate.
56 Another approach is to propagate the uncertainty that necessarily arises in orthology inference through
57 subsequent analyses. However, it is unclear what the relevant sources of uncertainty are in orthology
58 inference, and most consensus algorithms do not keep track of the different sources of error.

59 To remedy this situation, we characterized the error structure of gene-age estimation using 13 popular
60 orthology inference algorithms. In doing so, we identify common types of errors and, after correcting

61 these, present consensus gene-age calls for several model organisms (Table 1). We provide these gene-
62 age estimates along with a detailed analysis and annotation of the uncertainty associated with each age
63 call so that this uncertainty can be propagated through future analyses, as we show for functional term
64 enrichment. The consensus gene ages we calculate can be used for annotating genomic datasets in a
65 variety of fields, and the analysis of error will help prioritize genes for manual annotation and aspects of
66 orthology inference for future study.

67

68 **Results**

69 **Data collection**

70 In order to fairly consider the range of orthology algorithms, we took advantage of the reference
71 datasets managed by the Quest for Orthologs (QFO) consortium. QFO researchers have established
72 community standards and benchmarks for orthology inference and have made their benchmarking results
73 publicly available [20]. Importantly, the algorithms that have contributed to the benchmarking tool are
74 widely used and capture the variety of methods commonly used in the literature to infer orthology and
75 gene age [24–32]. It is therefore expected that nearly every study of gene age, regardless of the method
76 used, will closely resemble the results of at least one of the algorithms we explore here. We downloaded
77 orthology calls for 66 reference proteomes based on 13 orthology inference algorithms from the QFO
78 website and inferred the ages for each human gene by mapping the species in each ortholog group onto a
79 reference species tree from SwissTree, which was derived from a consensus of trees found in the
80 literature [33]. The results below are with reference to the human proteome, but the same methods were
81 applied to a variety of model organism proteomes (Table 1).

82 **The effect of algorithm choice on the distribution of human gene ages**

83 To investigate the effect of algorithm choice on inferred gene age, we broke the reference species tree
84 into eight age categories (Figure 1). These categories form nested clades, with the exception of the
85 category “Euk+Bacteria.” This non-phylogenetic category captures the substantial number of eukaryotic
86 genes that were horizontally transferred from bacteria after eukaryotes diverged from the rest of archaea

87 [34,35], and is defined as genes present in eukaryotes and bacteria but not archaea. For each algorithm,
88 every human gene was assigned to the age category in which the MRCA of the species in its orthogroup
89 falls, and the distributions over the different age categories for the human proteome inferred by that
90 algorithm was calculated.

91 We found that the algorithms fell into two distinct groups with respect to the distribution of age
92 classes. Clustering the algorithms by the average patristic distance between their per-gene age calls
93 recapitulated this grouping (Figure 2), and we define the two groups based on the midpoint root of this
94 tree. One group tended to find that most orthogroups could be traced to the MRCA of vertebrates,
95 whereas the other group found a much older mode age dating back to LECA. We call these two groups
96 the “young” and the “old” group respectively, though of course there are many more subtle and
97 interesting distinctions between the algorithms.

98 Orthology inference algorithms are typically classed into graph-based and tree-based methods [20].
99 However, we found that even though tree-based methods tended to fall in the “old” group, this was not
100 universally the case, nor were all graph-based methods found in the “young” group. The use of species
101 tree information was not a determining factor either (Figure 2). The bimodal nature of the age calls, either
102 “young” or “old”, is therefore not simply a reflection of the graph/tree distinction, although it is clearly
103 correlated. What is the source of this bimodality? One obvious answer is systematic error in the “young”
104 group algorithms, the “old” group, or both. Systematic error in the young group would be equivalent to
105 false negatives, i.e. missing orthology assignments, whereas systematic error in the old group is
106 equivalent to false positives, or spurious orthology assignments. This would have the effect of pushing the
107 age of the group away from or towards the root of the tree, respectively.

108 **Identifying systematic error**

109 We first investigated whether the bimodality of age-calls played out on the single gene level or
110 whether the two groups apparent in Figure 2 were due to the effects of averaging across genes, with error
111 being randomly distributed among proteins. To do so, we calculated a simple statistic that captured how
112 bimodal a protein’s age calls were between the two groups of algorithms (“old” and “young”). This

113 statistic, which we call bimodality, is the difference between age-call variation within the two groups and
114 between them, with more highly bimodal proteins having more variation between groups. Over 80% of
115 proteins had some degree of bimodality corresponding to these two age groups, or none, as is expected
116 given the hierarchical clustering in Figure 1. The remaining genes were anti-correlated with the
117 “old”/“young” groupings. Furthermore, the degree of bimodality between the “young” and “old”
118 algorithm groups correlates well with the amount of error associated with each protein (Spearman’s ρ :
119 .65)(Figure 3). That is, proteins with a large amount of error tend to be more bimodal. The bimodality
120 between algorithms is therefore a systematic phenomenon and a major source of error in these datasets.
121 Unfortunately, in the case of highly polarized genes, we cannot know *a priori* whether the “old” or
122 “young” age is the correct one. It is therefore important to propagate this uncertainty through further
123 analyses, and the bimodality statistic is included with our consensus age estimates.

124 We also investigated whether aspects of the individual proteins contributed to systematic error. For
125 instance, it may be difficult to infer correct evolutionary relationships for small proteins, or those with
126 many domains. At least one orthology inference algorithm uses this idea to “correct” for protein length
127 [36]. However, we found that protein length has a weak positive correlation with age-call error, and that
128 the number of domains also correlates weakly (Spearman’s ρ : $< .2$ in both cases).

129 **Systematic false negatives**

130 What are the causes of systematic false negatives and can we identify them without *a priori*
131 knowledge of the true orthogroup? One clue comes from the different age-category distributions between
132 PANTHER8_all and PANTHER8_LDO [25]. These two sets of orthology calls are based on the same set
133 of gene trees, but differ in their definition of orthology. “LDO” stands for “least diverged ortholog,” and
134 only considers the least diverged among a set of co-orthologs to be the true ortholog of an outgroup. This
135 can be contrasted to the traditional phylogenetic definition of orthology where all co-orthologs are equally
136 orthologous to the outgroup (Figure 4) [17]. Although it may be useful to split co-orthologous groups, as
137 the LDO definition does, in cases where orthology is being used for, e.g. gene function annotation, it is
138 inappropriate for defining the age of a gene or gene family because the age must be in reference to the

139 topology of the phylogenetic tree. The fact that PANTHER8_LDO's age category distribution resembled
140 that of several graph-based methods, and the fact that it clustered with them based on its per-gene age
141 calls (Figure 2), suggests that these methods may be splitting up co-orthologous groups as well.

142 There is no gold standard set of co-orthologs in this dataset, so we used the database PhylomeDB as a
143 reference for identifying co-ortholog over-splitting. To do so, we downloaded PhylomeDB summary files
144 for 10 species in PhylomeDB's model species collection (PhyC2) that overlapped with species in our tree
145 (Table 1), and determined groups of co-orthologs that were then used for the analysis. Briefly, for protein
146 (A), if an algorithm called a younger age (Y) and PhylomeDB an older age (O), and if in the co-orthologs
147 of (A) we could find a protein (B) which that algorithm called at age (O), then (B) was identified as the
148 LDO, age (O) was assumed to be the true age, and that algorithm was determined to be over-splitting the
149 co-ortholog group (Figure 4). This was not carried out for proteins on which PhylomeDB's age call was
150 determined to be a false positive (see below). We note that this method for identifying co-ortholog over-
151 splitting is not ideal, because it relies on a single, imperfect algorithm (PhylomeDB). It is conservative,
152 however, because algorithms will only be trimmed if they give a member of the co-orthologs the *exact*
153 same age on the species tree as that called by PhylomeDB on the focal gene. More thorough analyses of
154 whether graph-based methods are consistently missing co-orthologs will be necessary in the future.

155 **Identifying false positives**

156 If genes of distant organisms are incorrectly inferred to be part of an orthology group, it will drive the
157 age of the orthogroup towards the root of the tree. Recent HGT events are a biological source of such
158 errors, but some algorithmic error is expected to play a role as well. Such problems are perhaps more
159 likely to occur in tree-based algorithms, where slight re-arrangements that don't strongly affect the
160 likelihood of the tree can have an outsized effect on the inference of gene gains and losses [37]. In such
161 cases, the large number of taxa that fall between the true in-group taxa and the false positive out-group
162 taxa will be inferred to have lost the orthogroup. We used this criterion on a per-gene basis to identify
163 algorithms that were likely to have false positives and genes that were likely to be the result of HGT.
164 Algorithms that had an outsized number of taxa missing from an orthogroup (METHODS) were

165 considered false positives and removed from downstream analysis of that orthogroup's age. After
166 trimming these outliers, genes that were in the 95th percentile of inferred losses were flagged as being
167 potential recent HGT events (i.e. horizontally transferred long after LECA). These potential HGT genes
168 are an interesting set in themselves: 66% are from the Euk+Bacteria category, they are hugely enriched
169 for metabolic genes (gProfiler p-value=9.08e⁻¹¹⁶), and several are associated with human diseases.

170 We found that, as expected, algorithms in the “old” group tended to commit more false positive
171 errors, and algorithms in the “young” group committed more false negative errors (Figure 5). Because
172 PhylomeDB was used as a basis for identifying false negatives, its false negative rate could not be
173 quantified.

174 **Consensus**

175 These analyses suggested a way to more robustly estimate consensus gene ages and to calculate a
176 posterior distribution over the estimate. We used the methods described above to identify algorithms that
177 may have committed false positive or false negative errors and then removed these algorithms from
178 consideration on a per-protein basis. After doing so, we generated consensus tables based on the
179 remaining algorithms for the human proteome and for a number of other model eukaryotes (Table 1), and
180 we make these tables available (WEBSITE). Because our tree is best sampled within the opisthokonts
181 (fungi, animals, and closely related protists), we restricted our analyses to this lineage. These tables
182 contain a consensus age category for each protein based on the mode age call of non-trimmed algorithms.
183 Older genes were found to be involved in key components of cell biology. Genes in the Euk+Bac group
184 were found to be highly enriched for mitochondrial function, and genes that date back to the
185 Euk_Archaea node were enriched for translational machinery, as has been shown previously [8,9]. Many
186 of these older genes are also associated with hereditary diseases that represent a deficiency in a cell
187 function associated with that evolutionary epoch. For instance, the cytoskeletal system and cilium date to
188 LECA [9], and genes in this age category are enriched for diseases affecting the cilium, such as primary
189 ciliary dyskinesia and Bardet-Biedl syndrome (Figure 6).

190 These enrichment terms are derived from the point estimates of consensus ages, but we also provide
191 other data that can be used to propagate uncertainty to downstream analyses. For each gene, the
192 distribution over age-calls from the non-trimmed algorithms is given, as well as the number of
193 contributing algorithms and the entropy of the age call distribution. 87% of human proteins had at least 5
194 algorithms contributing after trimming, and 59% had at least 10 out of a total of 13 original algorithms. In
195 addition, the tables contain information on whether the protein was flagged as being a potential horizontal
196 gene transfer event. Finally, we include the node error and bimodality statistics, both of which are
197 measures of uncertainty that reference the reference species tree.

198 We note that in several cases we have made *ad hoc* decisions during the building of the consensus.
199 For instance, algorithms were flagged as false positives if the number of taxa inferred to have lost the
200 orthogroup was two standard deviations above the mean of all algorithms. These decisions were informed
201 by the underlying distributions of values. Nevertheless, we supply the source data files, scripts used for
202 these analyses, as well as interactive iPython notebooks, and we invite researchers to explore and change
203 parameters if they desire (<https://github.com/marcottelab/Gene-Ages>).

204 **Error Propagation**

205 How can our error annotations be used in downstream analyses? Here we give an example of a simple
206 stability analysis for gene ontology enrichment that uses these error terms. It has previously been shown
207 that eukaryotic genes vertically acquired from Archaea are enriched for translation and RNA processing,
208 whereas genes acquired horizontally from bacteria at the root of eukaryotes are enriched for metabolic
209 processes (Figure 6, [8,9]). This conclusion relies on functional term enrichment, but what is the effect of
210 different sources of error on these sorts of enrichment analyses? To investigate the robustness of this
211 conclusion to different sources of error, we used the program g:Profiler [38] to perform functional
212 enrichment analysis on the two age classes “Euk_Archaea” and “Euk+Bacteria” after filtering the datasets
213 at varying levels of stringency (Figure 7A). We found that removing genes that were flagged as a possible
214 late HGT event had a strong effect on the average p-values of functional annotation terms in the
215 Euk+Bacteria age class but not the Euk_Archaea class (Figure 7B). This may be due to these genes being

216 more commonly lost or to many bacterial genes being more recent HGT events (and hence being filtered
217 out). The latter possibility would mean that many genes in this age category could be misidentified as
218 being present in LECA, so these genes are good candidates for manual curation. Notably, filtering on
219 different error terms can increase or decrease the significance of different terms, and, depending on the
220 filtering strategy, the significance ranking of terms can be switched (Figure 7C and D). Analyses that rely
221 on smaller test-sets of genes are likely to be much more strongly affected than these proteome-wide
222 searches.

223 **Discussion**

224 Most studies of gene age use a single point estimate arising from one of a variety of methods. Given
225 our analysis of some of the most popular orthology inference algorithms, we find that point estimates of
226 gene age will be wrong for (at least) thousands of genes in a human-sized proteome (Figure 4). More
227 troubling is the fact that algorithms appear to fall into two classes, each of which presumably has a
228 systematic bias towards either false positives (“old group”) or false negatives (“young group”). This
229 systematic bias happens on a per-gene basis, meaning that simple voting methods will not be able to
230 resolve conflicts. Even with the ideal sampling of algorithms, which we approximate here by exploring a
231 wide diversity of popular algorithms, the effective voting population will still drop to two on highly
232 polarized genes.

233 Many areas of computational biology have faced a similar problem, namely, the need to keep track of
234 error in several components of a workflow, and to correctly propagate this error through the whole
235 analysis [39]. One illustrative example is multiple sequence alignment and phylogenetic inference. The
236 former is a necessary precursor to the latter, and each involves estimation error. Methods have been
237 developed to infer the posterior distributions of both steps simultaneously [40], which is computationally
238 intractable for all but the smallest datasets, or to perform each step iteratively in a maximum likelihood
239 framework [41]. We argue that, eventually, such steps will have to be taken with orthology inference and
240 gene-age estimation. Using a point estimate at each step in the analysis makes the assumption that each

241 inference step has no uncertainty associated with it, which we can clearly reject in the case of gene-age
242 estimation.

243 Some methods for probabilistic orthology inference do exist [21]. These use gene tree models with
244 free parameters for gene duplication, loss, and sometimes HGT, which then contribute to the likelihood
245 along with the multiple sequence alignment. However, these methods are in their infancy, and not usually
246 scalable to large datasets or widely used. In the meantime, it is important to have an understanding of
247 common sources of error in gene-age estimation. We provide that information along with consensus age
248 calls for a variety of model organisms so that researchers can incorporate error propagation into their
249 analyses in a way that is appropriate to their question of interest.

250 Several error terms are likely to be important for a broad range of analyses. The first and most
251 straightforward is the entropy of the age-call estimate after filtering false positives and negatives. This
252 statistic gives a quick idea of how certain an age-call is, with higher entropies being less certain. It is
253 defined with reference to our age categories, so if researchers need to use other age categories, they must
254 use the node age of the gene, which we also provide (Methods). HGT events are also likely to affect some
255 datasets, especially when genes originating in Bacteria are involved (Figure 6). A large number of
256 eukaryotic genes are likely transfers from Bacteria [8], but these may have been transferred at any point
257 on the phylogeny. We define one age category, Euk+Bacteria, to describe all genes transferred before
258 LECA, with later transfers hopefully being caught by our flag. If researchers are primarily interested in
259 HGT, we suggest a much fuller analysis, as our simple method is likely to miss many HGT events.
260 Finally, the bimodality of the age-call between “young” and “old” algorithm types is a key statistic. The
261 systematic biases in the different algorithm types mean that many datasets will be radically different and
262 difficult to compare, and it may account for some of the differences between studies of ancient gene
263 repertoires that used either graph or tree-based methods. Genes that are highly polarized are good
264 candidates for manual curation, because it is unlikely that any *ad hoc* algorithm will differ substantially
265 enough from those we sampled here to be decisive.

266 Although we have characterized only two components of a typical computational biology workflow,
267 orthology inference and gene-age estimation, it would be ideal to characterize error distributions for all
268 the steps in an analysis, which has not been done with gene age data to our knowledge (but see [42] for an
269 interesting example on gene-expression data, and [39] for a general review). The datasets we provide here
270 will hopefully help guide future research efforts aimed at a more formal, probabilistic way to handle error
271 in gene-age estimation, perhaps even in the context of an entire workflow. Until such methods are
272 available, we advocate using our error annotations or a similar analysis in any study incorporating gene-
273 age data.

274

275 **Methods**

276 **Data Collection and Availability**

277 15 algorithms have submitted their estimates on 66 reference proteomes
278 (http://www.ebi.ac.uk/reference_proteomes) to QFO's benchmarking tool
279 (<http://orthology.benchmarkservice.org/cgi-bin/gateway.pl>). We omitted two of these because they either
280 did not have full taxon coverage (RBH), or their results were so different from all the others that it
281 dominated the variance in all downstream analyses (OMA_GETHOGS). Pairwise orthology calls for the
282 13 remaining algorithms were downloaded from the Quest for Orthologs benchmarking website [20].
283 These pairwise calls were converted into tables for each gene, which were then used for subsequent
284 analyses. The reference species tree was downloaded from SwissTree [33] on 06/15/2015
285 (<ftp://ftp.lausanne.isb-sib.ch/pub/databases/SwissTree/speciestree.nhx>) and was pruned to match the taxa
286 in the Quest for Orthologs reference proteomes (http://www.ebi.ac.uk/reference_proteomes). Custom
287 programs were written to perform the analyses below, and these are publicly available, as are iPython
288 notebooks used for plotting. These, and the datasets supporting the conclusions in this article are available
289 on GitHub (<https://github.com/marcottelab/Gene-Ages>) with the following commit id
290 (c1a2862fa894d7da4ccdf3fb8001e1b6b226bd09). Scripts relied heavily on the python packages
291 dendropy [43], BioPython [44], and pandas [45].

292 **Protein Age Calls**

293 All protein ages are referenced to the species tree obtained from SwissTree. The age of a protein is
294 calculated on the species tree by finding the MRCA node of the taxa that have orthologs of that protein.
295 This node is the “node age,” and the age group it falls into is the “binned age.” The binned ages conform
296 to the interior labels given by SwissTree, with the exception of the Euk+Bac age category, which is not
297 phylogenetic, but rather consists of proteins that are present in Bacteria and Eukaryota (and would thus
298 normally be assigned to the oldest age class), but not in Archaea.

299 We calculate several measures of error amongst algorithms. First is an error statistic called “node
300 error” based on the node age calls. Node error is the average number of branches (patristic distance)
301 between the age calls any two algorithms. A similar measure was used to calculate the distance tree in
302 Figure 1. The average patristic distance between age-calls for each pair of algorithms was used as input
303 for a heuristic search in PAUP [46]. Next, because algorithms fall roughly into two groups (“old” and
304 “young”), we calculate the “bimodality” of each protein. This is the difference between the average
305 within group (“old” and “young”) node error and the average between group node error. Note that,
306 although we call this statistic simply “bimodality,” it captures not just the bimodal nature of the age calls,
307 but how different the two peak ages are. Thus the proteins with the highest bimodality score are those for
308 which all the “young” algorithms call one age, all the “old” algorithms call a different age, and these ages
309 are very far apart on the tree (Figure 3).

310 **Filtering False Positives and Negatives**

311 Before calculating a consensus, we flag algorithms that may have committed false positive or false
312 negative errors on a per-gene basis. These algorithms are then removed from consideration of that gene’s
313 age. False positives are orthology calls that are substantially more distant than orthology calls by other
314 algorithms, and have the effect of driving age deeper in the tree. These are found as follows. For each
315 algorithm and each protein: 1.) the node age is calculated 2.) the number of taxa in the species tree
316 descended from this node is found 3.) The number of taxa containing orthologs of the focal protein is
317 subtracted from the number of descendant taxa. This number is the number of taxa without the orthogroup

318 that are descended from an ancestor that putatively had the orthogroup, and is therefore proportional, but
319 not identical, to the number of inferred losses of the orthogroup. For each algorithm and each protein, if
320 this number is two standard deviations above the pooled algorithm mean for the focal protein, that
321 algorithm's age call is considered a false positive and is thrown out.

322 False negatives are cases where an algorithm fails to make an orthology call, driving the inferred age
323 to shallower nodes in the species tree. We identify one possible cause of this, which we call "over-
324 splitting." This is when a group of co-orthologs is not correctly recognized by an algorithm and only one
325 or a few of its members are found as orthologs to a more distant species, while the others are split off into
326 their own orthogroups. The members that are split off would then be called at an incorrectly young age.
327 To identify these errors, we used PhylomeDB's [32] orthogroups as a standard. For each protein and each
328 algorithm (except for PhylomeDB), if the focal algorithm called a younger age than PhylomeDB and a
329 co-ortholog of the focal protein could be found where the focal algorithm called the same node age as
330 PhylomeDB did on the focal protein, then this algorithm was considered to be over-splitting the focal
331 protein, and was not considered in this protein's age call. This error calculation was not performed on
332 proteins where PhylomeDB was flagged as a false positive.

333 **Consensus Ages**

334 We generated consensus binned ages after removing algorithms flagged with false positives and
335 negatives as described above. The number of algorithms favoring each binned age is counted and then
336 normalized by the number of contributing algorithms to give a distribution over age calls. For subsequent
337 analyses, we used the mode of this distribution as the consensus age.

338

339 **Acknowledgements**

340 We would especially like to acknowledge the Quest for Orthologs consortium and those who contributed
341 their algorithms to the benchmarking tool for making their data freely available. B.J.L. was funded by NIH
342 fellowship 1F32GM112504-01A1. E.M.M. acknowledges funding from the NIH, NSF, CPRIT, ARO
343 (61789-MA-MUR), and Welch Foundation (F1515).

344

345

346 **References**

- 347 1. Mushegian AR, Koonin EV. A minimal gene set for cellular life derived by comparison of complete
348 bacterial genomes. *Proc. Natl. Acad. Sci.* 1996;93:10268–73.
- 349 2. Capra JA, Stolzer M, Durand D, Pollard KS. How old is my gene? *Trends Genet.* 2013;29:659–68.
- 350 3. Maxwell EK, Schnitzler CE, Havlak P, Putnam NH, Nguyen A-D, Moreland RT, et al. Evolutionary
351 profiling reveals the heterogeneous origins of classes of human disease genes: implications for modeling
352 disease genetics in animals. *BMC Evol. Biol.* 2014;14:212.
- 353 4. Domazet-Lošo T, Tautz D. An Ancient Evolutionary Origin of Genes Associated with Human Genetic
354 Diseases. *Mol. Biol. Evol.* 2008;25:2699–707.
- 355 5. Wolf YI, Novichkov PS, Karev GP, Koonin EV, Lipman DJ. The universal distribution of evolutionary
356 rates of genes and distinct characteristics of eukaryotic genes of different apparent ages. *Proc. Natl. Acad.*
357 *Sci.* 2009;106:7273–80.
- 358 6. Kim WK, Marcotte EM. Age-Dependent Evolution of the Yeast Protein Interaction Network Suggests
359 a Limited Role of Gene Duplication and Divergence. *PLoS Comput Biol.* 2008;4:e1000232.
- 360 7. Pellegrini M, Marcotte EM, Thompson MJ, Eisenberg D, Yeates TO. Assigning protein functions by
361 comparative genome analysis: Protein phylogenetic profiles. *Proc. Natl. Acad. Sci.* 1999;96:4285–8.
- 362 8. Thiergart T, Landan G, Schenk M, Dagan T, Martin WF. An Evolutionary Network of Genes Present
363 in the Eukaryote Common Ancestor Polls Genomes on Eukaryotic and Mitochondrial Origin. *Genome*
364 *Biol. Evol.* 2012;4:466–85.
- 365 9. Koumandou VL, Wickstead B, Ginger ML, van der Giezen M, Dacks JB, Field MC. Molecular
366 paleontology and complexity in the last eukaryotic common ancestor. *Crit. Rev. Biochem. Mol. Biol.*
367 2013;48:373–96.
- 368 10. Rivera AS, Pankey MS, Plachetzki DC, Villacorta C, Syme AE, Serb JM, et al. Gene duplication and
369 the origins of morphological complexity in pancrustacean eyes, a genomic approach. *BMC Evol. Biol.*
370 2010;10:123.
- 371 11. Liebeskind BJ, Hillis DM, Zakon HH, Hofmann HA. Complex Homology and the Evolution of
372 Nervous Systems. *Trends Ecol. Evol.* 2016;31:127–35.
- 373 12. Conaco C, Bassett DS, Zhou H, Arcila ML, Degnan SM, Degnan BM, et al. Functionalization of a
374 protosynaptic gene expression network. *Proc. Natl. Acad. Sci.* 2012;109:10612–8.
- 375 13. Wan C, Borgeson B, Phanse S, Tu F, Drew K, Clark G, et al. Panorama of ancient metazoan
376 macromolecular complexes. *Nature* [Internet]. 2015 [cited 2015 Sep 9];advance online publication.
377 Available from: <http://www.nature.com/nature/journal/vaop/ncurrent/full/nature14877.html>

- 378 14. Alié A, Hayashi T, Sugimura I, Manuel M, Sugano W, Mano A, et al. The ancestral gene repertoire of
379 animal stem cells. *Proc. Natl. Acad. Sci.* 2015;112:E7093–100.
- 380 15. Gabaldón T. Large-scale assignment of orthology: back to phylogenetics? *Genome Biol.* 2008;9:235.
- 381 16. Fitch WM. Distinguishing Homologous from Analogous Proteins. *Syst. Biol.* 1970;19:99–113.
- 382 17. Fitch WM. Homology a personal view on some of the problems. *Trends Genet. TIG.* 2000;16:227–31.
- 383 18. Chen X, Zhang J. The Ortholog Conjecture Is Untestable by the Current Gene Ontology but Is
384 Supported by RNA Sequencing Data. *PLoS Comput Biol.* 2012;8:e1002784.
- 385 19. Gabaldón T, Koonin EV. Functional and evolutionary implications of gene orthology. *Nat. Rev.*
386 *Genet.* 2013;14:360–6.
- 387 20. Sonnhammer ELL, Gabaldón T, Silva AWS da, Martin M, Robinson-Rechavi M, Boeckmann B, et al.
388 Big data and other challenges in the quest for orthologs. *Bioinformatics.* 2014;30:2993–8.
- 389 21. Ullah I, Sjöstrand J, Andersson P, Sennblad B, Lagergren J. Integrating Sequence Evolution into
390 Probabilistic Orthology Analysis. *Syst. Biol.* 2015;64:969–82.
- 391 22. Pereira C, Denise A, Lespinet O. A meta-approach for improving the prediction and the functional
392 annotation of ortholog groups. *BMC Genomics.* 2014;15:S16.
- 393 23. Maher MC, Hernandez RD. Rock, Paper, Scissors: Harnessing Complementarity in Ortholog
394 Detection Methods Improves Comparative Genomic Inference. *G3 GenesGenomesGenetics.* 2015;5:629–
395 38.
- 396 24. Sonnhammer ELL, Östlund G. InParanoid 8: orthology analysis between 273 proteomes, mostly
397 eukaryotic. *Nucleic Acids Res.* 2015;43:D234–9.
- 398 25. Mi H, Muruganujan A, Thomas PD. PANTHER in 2013: modeling the evolution of gene function,
399 and other gene attributes, in the context of phylogenetic trees. *Nucleic Acids Res.* 2013;41:D377–86.
- 400 26. Altenhoff AM, Škunca N, Glover N, Train C-M, Sueki A, Piližota I, et al. The OMA orthology
401 database in 2015: function predictions, better plant support, synteny view and other improvements.
402 *Nucleic Acids Res.* 2015;43:D240–9.
- 403 27. DeLuca TF, Cui J, Jung J-Y, St. Gabriel KC, Wall DP. Roundup 2.0: enabling comparative genomics
404 for over 1800 genomes. *Bioinformatics.* 2012;28:715–6.
- 405 28. Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, et al. eggNOG 4.5: a
406 hierarchical orthology framework with improved functional annotations for eukaryotic, prokaryotic and
407 viral sequences. *Nucleic Acids Res.* 2015;gkv1248.
- 408 29. Linard B, Allot A, Schneider R, Morel C, Ripp R, Bigler M, et al. OrthoInspector 2.0: Software and
409 database updates. *Bioinforma. Oxf. Engl.* 2015;31:447–8.
- 410 30. Vilella AJ, Severin J, Ureta-Vidal A, Heng L, Durbin R, Birney E. EnsemblCompara GeneTrees:
411 Complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res.* 2009;19:327–35.

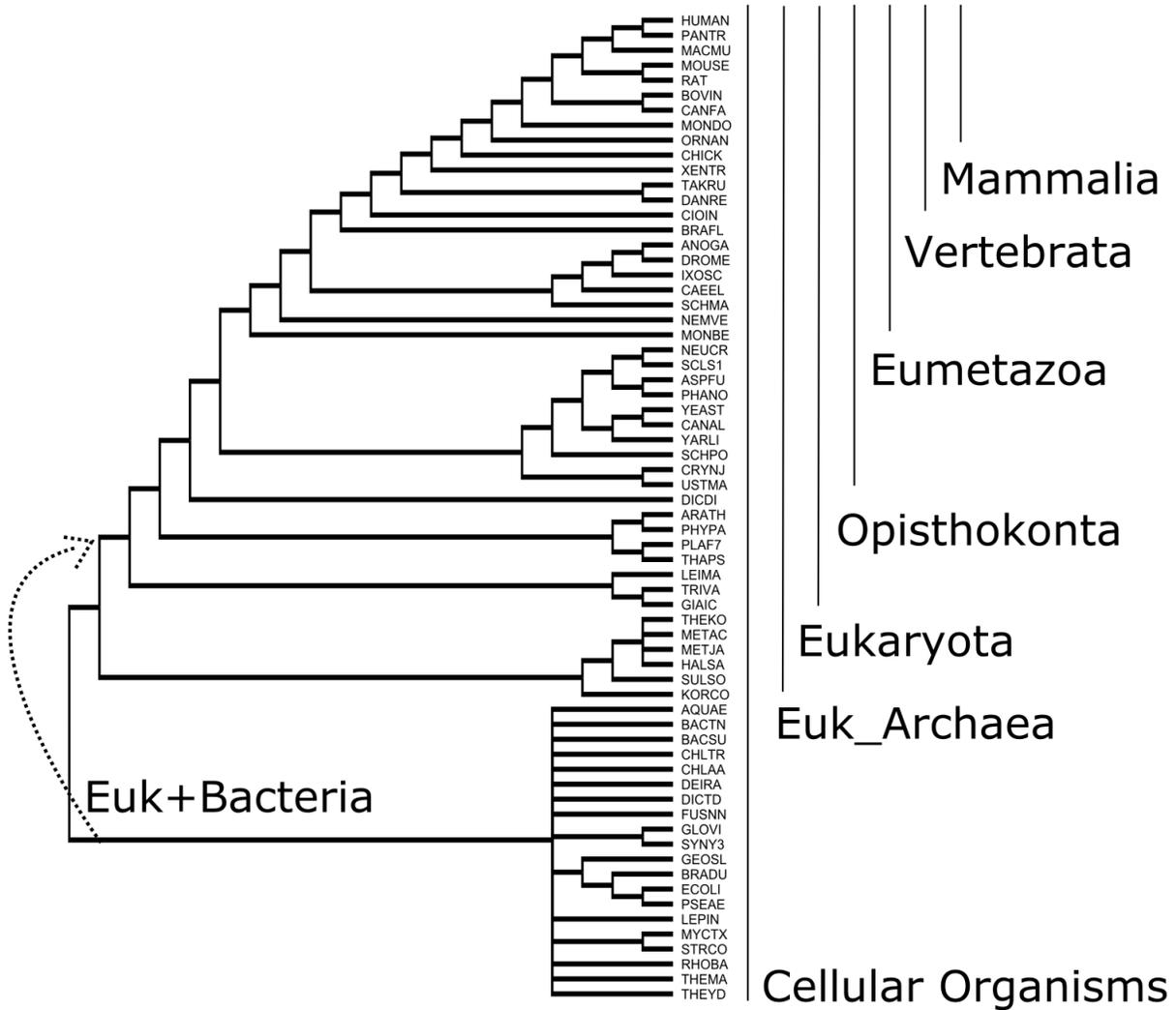
- 412 31. Prysycz LP, Huerta-Cepas J, Gabaldón T. MetaPhOrs: orthology and paralogy predictions from
413 multiple phylogenetic evidence using a consistency-based confidence score. *Nucleic Acids Res.*
414 2011;39:e32–e32.
- 415 32. Huerta-Cepas J, Bueno A, Dopazo J, Gabaldón T. PhylomeDB: a database for genome-wide
416 collections of gene phylogenies. *Nucleic Acids Res.* 2008;36:D491–6.
- 417 33. Boeckmann B, Marcet-Houben M, Rees JA, Forslund K, Huerta-Cepas J, Muffato M, et al. Quest for
418 Orthologs Entails Quest for Tree of Life: In Search of the Gene Stream. *Genome Biol. Evol.*
419 2015;7:1988–99.
- 420 34. Méheust R, Lopez P, Bapteste E. Metabolic bacterial genes and the construction of high-level
421 composite lineages of life. *Trends Ecol. Evol.* 2015;30:127–9.
- 422 35. Pittis AA, Gabaldón T. Late acquisition of mitochondria by a host with chimaeric prokaryotic
423 ancestry. *Nature* [Internet]. 2016 [cited 2016 Feb 25];advance online publication. Available from:
424 <http://www.nature.com/nature/journal/vaop/ncurrent/full/nature16941.html>
- 425 36. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome comparisons
426 dramatically improves orthogroup inference accuracy. *Genome Biol.* 2015;16:1–14.
- 427 37. Hahn MW. Bias in phylogenetic tree reconciliation methods: implications for vertebrate genome
428 evolution. *Genome Biol.* 2007;8:R141.
- 429 38. Reimand J, Kull M, Peterson H, Hansen J, Vilo J. g:Profiler—a web-based toolset for functional
430 profiling of gene lists from large-scale experiments. *Nucleic Acids Res.* 2007;35:W193–200.
- 431 39. Guang A, Zapata F, Howison M, Lawrence CE, Dunn CW. An Integrated Perspective on
432 Phylogenetic Workflows. *Trends Ecol. Evol.* 31:116–26.
- 433 40. Suchard MA, Redelings BD. BAli-Phy: simultaneous Bayesian inference of alignment and
434 phylogeny. *Bioinforma. Oxf. Engl.* 2006;22:2047–8.
- 435 41. Liu K, Raghavan S, Nelesen S, Linder CR, Warnow T. Rapid and Accurate Large-Scale Coestimation
436 of Sequence Alignments and Phylogenetic Trees. *Science.* 2009;324:1561–4.
- 437 42. Thompson A, Vo D, Comfort C, Zakon HH. Expression Evolution Facilitated the Convergent
438 Neofunctionalization of a Sodium Channel Gene. *Mol. Biol. Evol.* 2014;31:1941–55.
- 439 43. Sukumaran J, Holder MT. DendroPy: a Python library for phylogenetic computing. *Bioinformatics.*
440 2010;26:1569–71.
- 441 44. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: Freely Available
442 Python Tools for Computational Molecular Biology and Bioinformatics. *Bioinformatics.* 2009;25:1422–
443 3.
- 444 45. McKinney W. Python for data analysis. Beijing: O’Reilly; 2013.
- 445 46. Swofford, David L. Phylogenetic analysis using parsimony (*and other methods). Sunderland, MA:
446 Sinauer Associates; 2003.
- 447

448

449

450

451 **Figures**

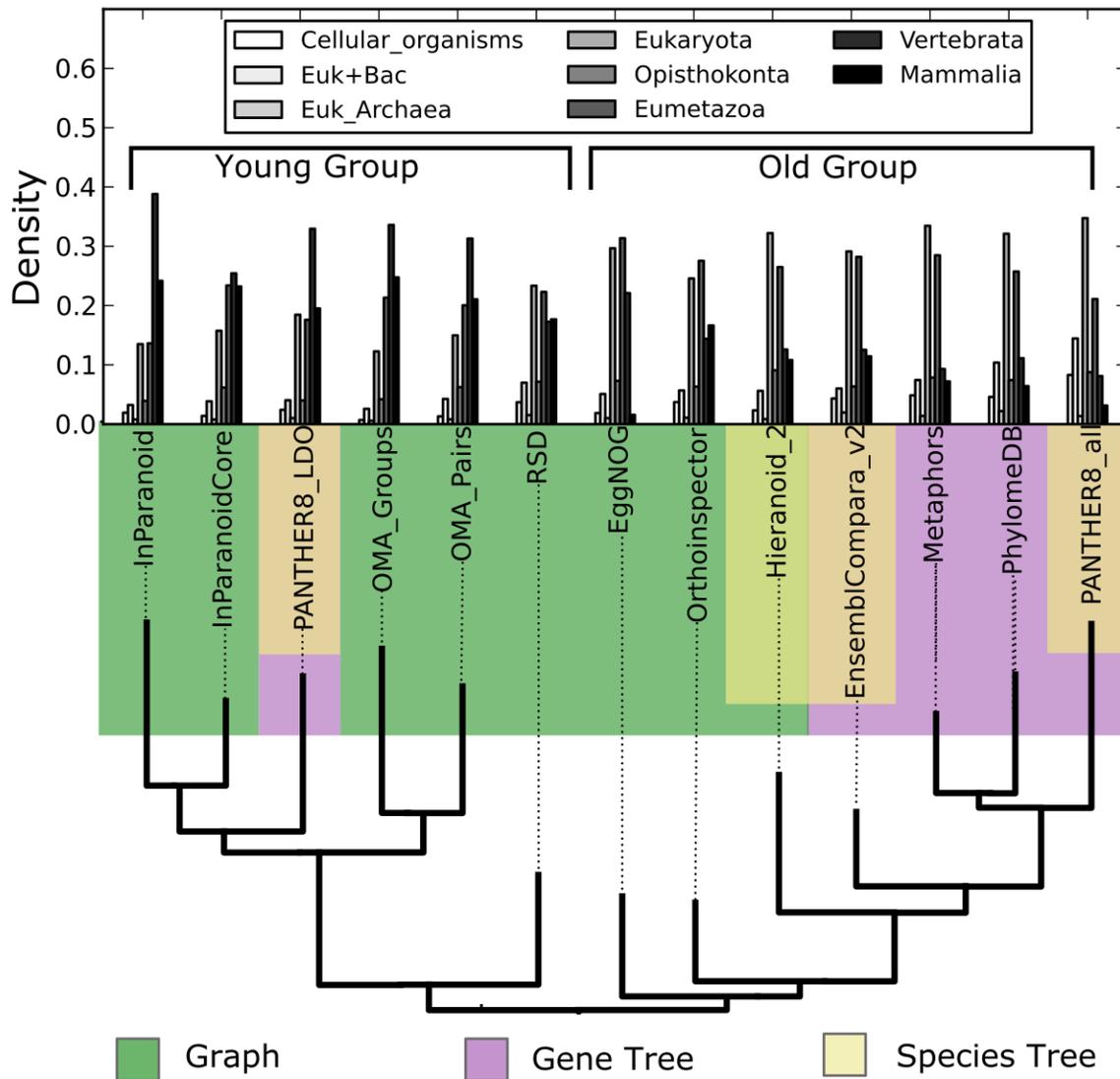


452

453 **Figure 1**

454 The reference species tree and age categories used for gene-age inference. This tree is based on SwissTree

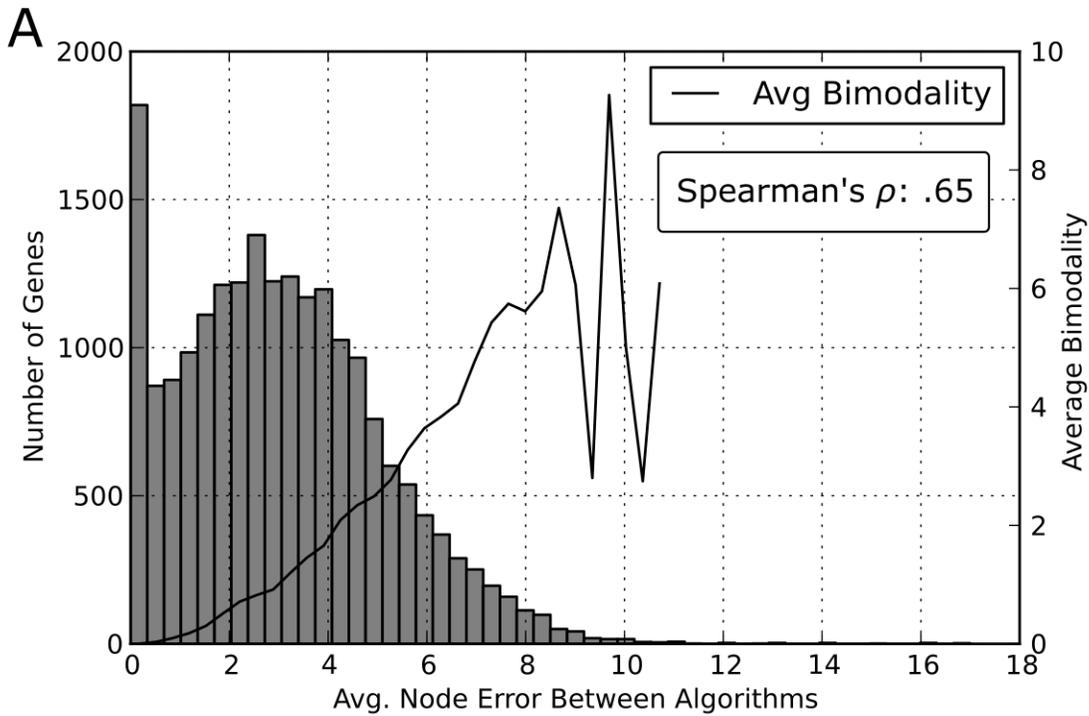
455 [33] and reflects a consensus of recent large-scale phylogenies. Tip names are Uniprot species identifiers.



456

457 **Figure 2**

458 Distribution of age categories in the human proteome inferred by 13 different orthology inference
 459 algorithms. Algorithms were clustered according to the average pairwise distance between their age-calls,
 460 counted in units of braches (patristic distance). The distance tree is rooted at the midpoint. Algorithms are
 461 colored by the methods they use to infer orthology. They either use a graph-based or a gene tree-based
 462 strategy, either with, or without, the use of a species tree.



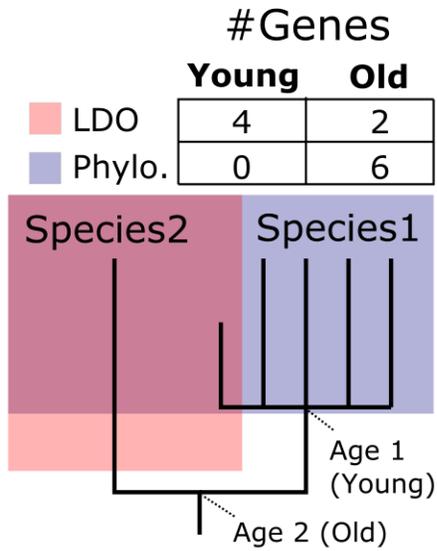
B

	"Young"			"Old"			
	InParanoid	OMA_Groups	PANTHER_LDO	EGGNOG	Hieranoid_2	Metaphors	Bimodality
Q5H910	Mammalia	Mammalia	Mammalia	Eukaryota	Eukaryota	Eukaryota	14
Q7Z5J4	Vertebrata	Vertebrata	Eumetazoa	Vertebrata	Eumetazoa	Eumetazoa	0.3

463

464 **Figure 3**

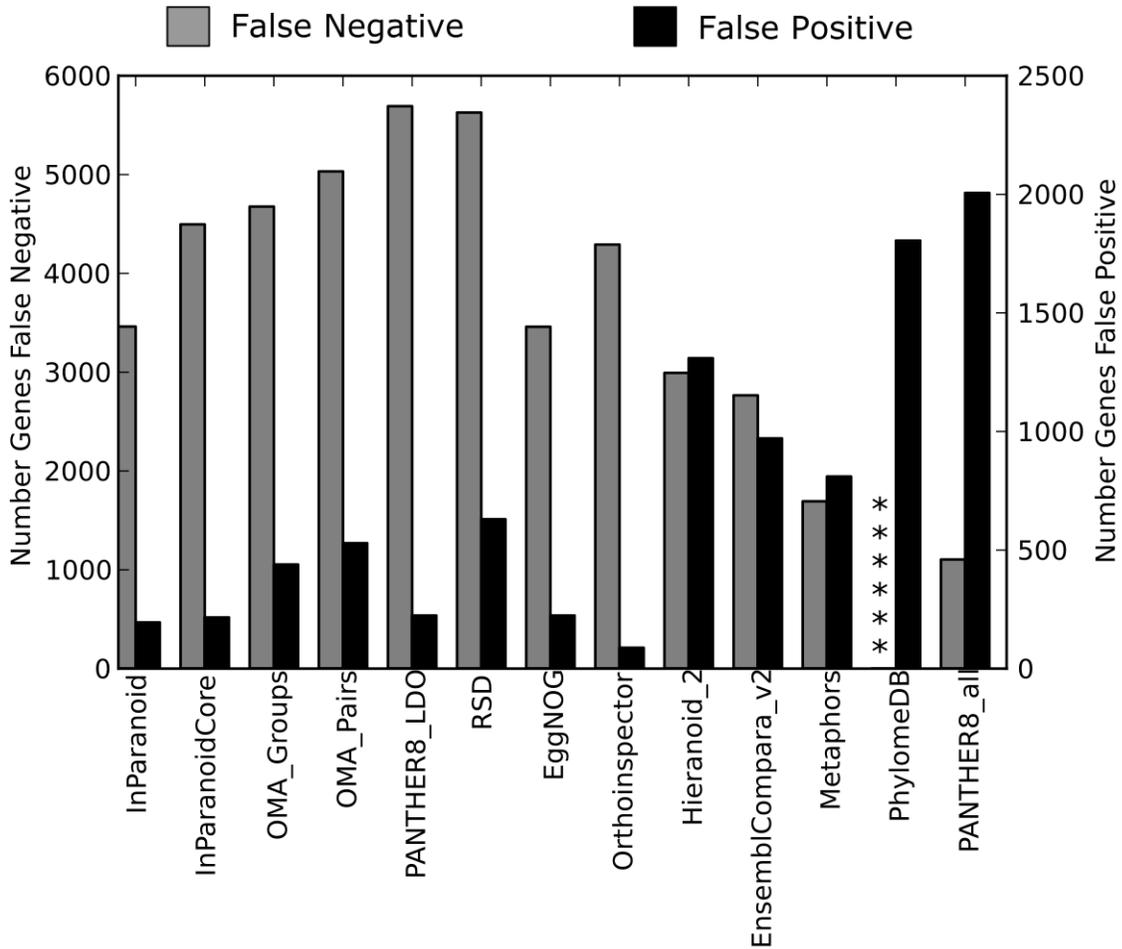
465 Error statistics. (A) The distribution of average node error, a measure of disagreement among the
 466 algorithms for a given gene, is given, along with a plot of the average bimodality in each bin. Genes with
 467 more error tend to be more bimodal between “old” and “young” algorithms. (B) Example of a strongly
 468 bimodal and weakly bimodal gene with a few representative algorithms. The ages are given as categories
 469 for clarity, but the bimodality statistic is calculated according to patristic distance between node age-calls
 470 (Methods).



471

472 **Figure 4**

473 Determination of false negatives due to co-ortholog over-splitting. This tree compares the ages given by
474 least derived orthology (LDO) and traditional, phylogenetic orthology (Phylo.). Given a group of co-
475 orthologs in Species 1, LDO will give only the co-ortholog with the shortest distance to an outgroup
476 (gene in Species 2) the status of ortholog to this outgroup (red box). All others are put in separate
477 orthogroups. Hence, LDO produces more genes that are mapped (incorrectly) to a younger age (**Y**),
478 whereas traditional, phylogenetic orthology (blue box) includes all co-orthologs to the orthogroup,
479 thereby mapping more genes to the older age (**O**).

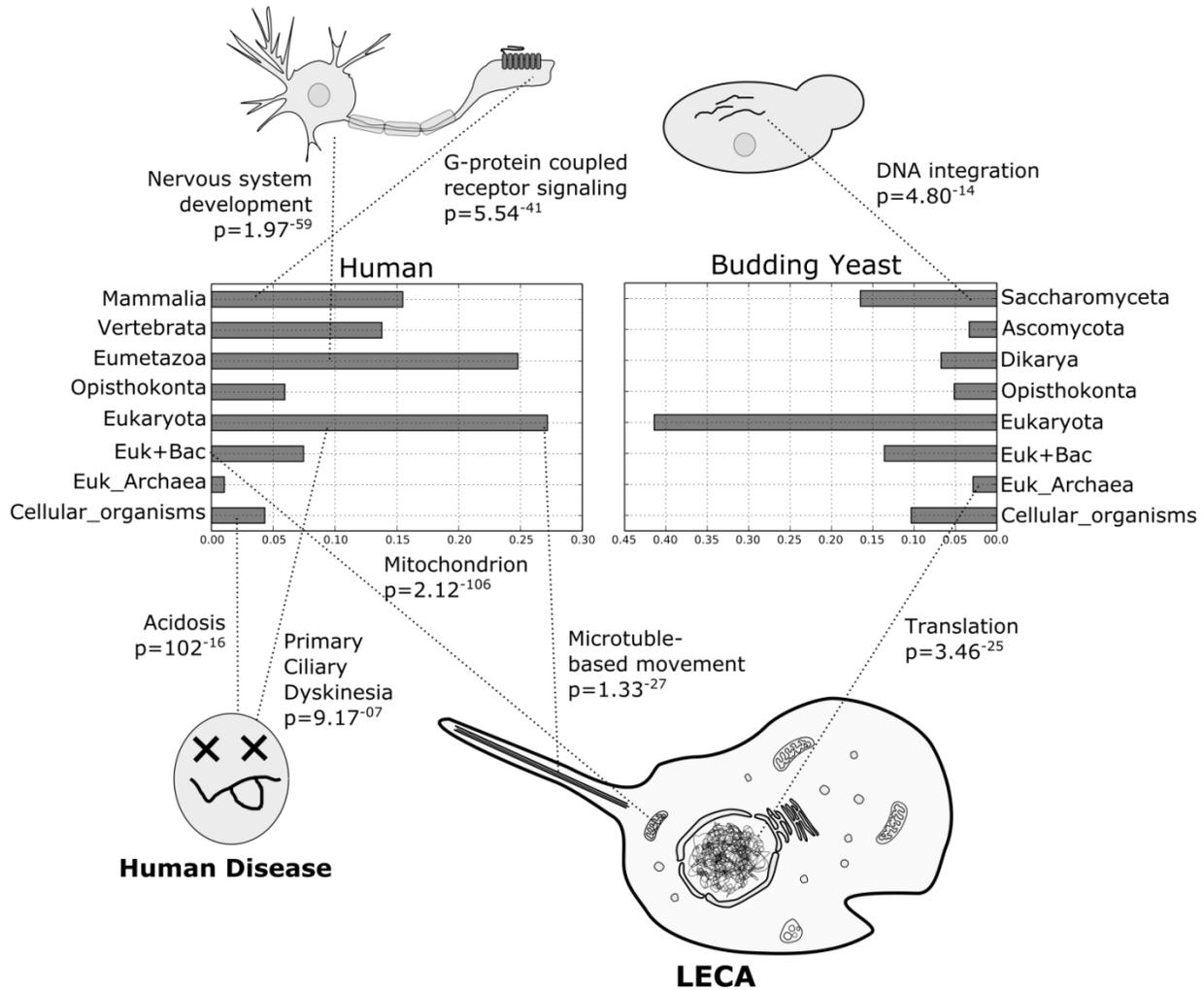


480

481 **Figure 5**

482 Errors committed by different algorithms. False positives and negative are defined in the text.

483 PhylomeDB was used as a standard for false negative, so its false negative count could not be determined.

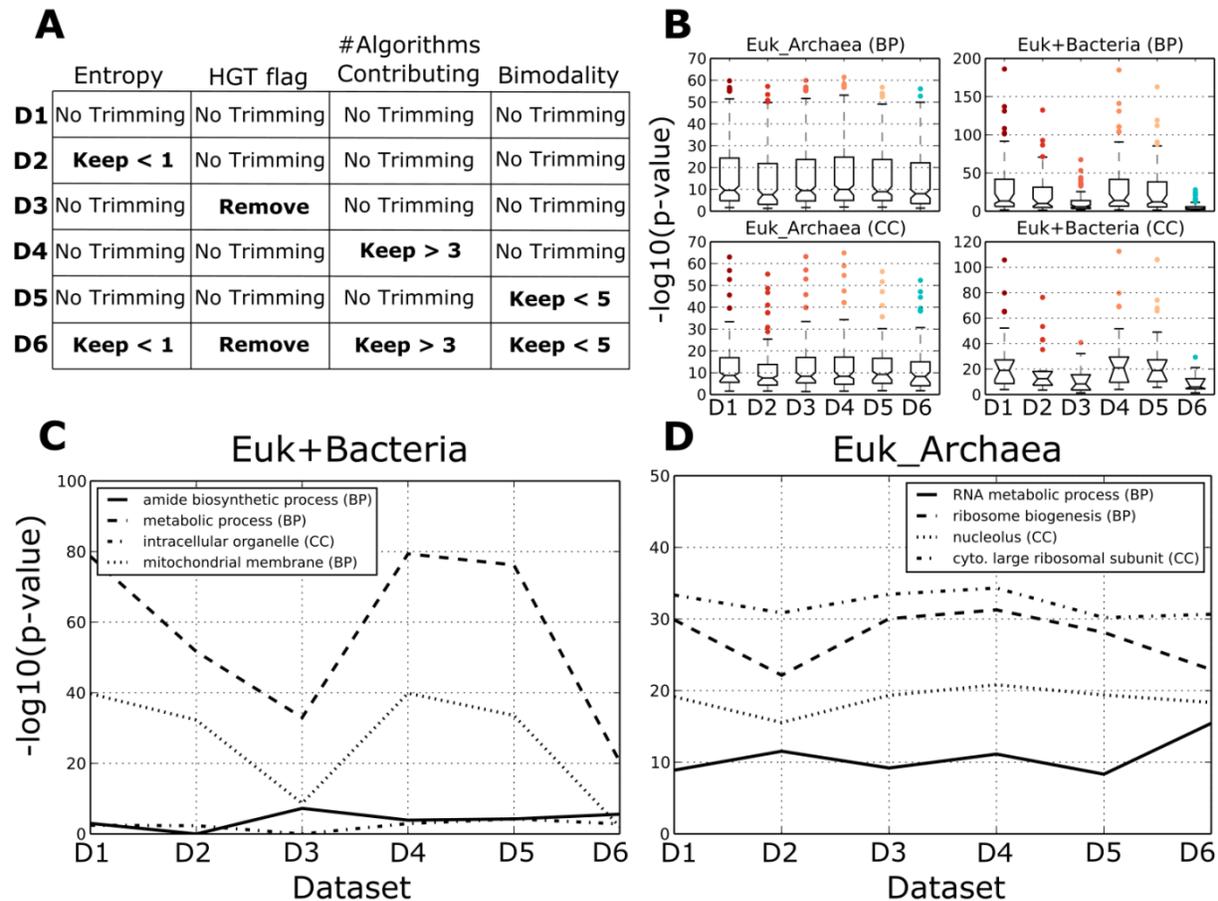


484

485

486 **Figure 6**

487 Enrichment of gene ontology terms and human disease terms (OMIM) in the different consensus age
 488 classes for human and budding yeast (*Saccharomyces cerevisiae*). The distribution of age classes are
 489 shown for each species. Older genes tend to be enriched for core cellular machinery and heritable
 490 diseases. Newer genes are associated with lineage-specific function, such as nervous system development
 491 and olfaction (via G-protein coupled receptors) in mammals, and DNA integration in yeast. P-values are
 492 derived from g:Profiler [38].



493

494 **Figure 7**

495 Effect of filtering on functional term enrichment analysis. (A) Datasets 1-6 were trimmed based on four
 496 sources of error: entropy of age-calls, whether the genes were flagged as potential horizontal gene transfer
 497 (HGT) events, the number of algorithms contributing to the final age call (after filtering algorithms, as
 498 described in Methods and Results), and the polarization of each gene. Dataset 6 was filtered on all four
 499 criteria (B) Negative \log_{10} p-values for the five datasets for two age categories (Euk_Archaea, and
 500 Euk+Bacteria) and two gene ontology terms (Biological Process, and Cellular Compartment). (C-D)
 501 Negative \log_{10} p-values across datasets for eight functional terms in the two age categories. These terms
 502 show the variety of ways that significance can be affected by filtering.

Common Name	Uniprot ID	False-negative analysis
Anopheles gambiae (Mosquito)	ANOGA	No
Bos taurus (Cattle)	BOVIN	No
Branchiostoma floridae (Lancelet)	BRAFL	No
Caenorhabditis elegans (Worm)	CAEEL	Yes
Candida albicans	CANAL	Yes
Canis lupus familiaris (Dog)	CANFA	No
Gallus gallus (Chicken)	CHICK	Yes
Ciona intestinalis (Tunicate)	CIOIN	No
Cryptococcus neoformans	CRYNJ	No
Danio rerio (Zebrafish)	DANRE	Yes
Drosophila melanogaster (Fly)	DROME	Yes
Homo sapiens (Human)	HUMAN	Yes
Ixodes scapularis (Tick)	IXOSC	No
Macaca mulatta (Rhesus macaque)	MACMU	No
Monosiga brevicollis (Choanoflagellate)	MONBE	No
Monodelphis domestica (Opossum)	MONDO	No
Mus musculus (Mouse)	MOUSE	Yes
Nematostella vectensis (Sea anemone)	NEMVE	No
Neurospora crassa	NEUCR	No
Ornithorhynchus anatinus (Platypus)	ORNAN	No
Pan troglodytes (Chimp)	PANTR	No
Phaeosphaeria nodorum	PHANO	No
Rattus rattus (Rat)	RAT	Yes
Saccaromyces cerevisiae (Budding yeast)	YEAST	Yes
Schistosoma mansoni (Blood fluke)	SCHMA	No
Schizosaccharomyces pombe (Fission yeast)	SCHPO	Yes
Sclerotinia sclerotiorum	SCLS1	No
Takifugu rubripes (Pufferfish)	TAKRU	No
Ustilago maydis (Corn smut/Huitlacoche)	USTMA	No
Xenopus tropicalis (Frog)	XENTR	No
Yarrowia lipolytica	YARLI	No

503

504 **Table 1**

505 Species for which final consensus tables were constructed. Tables are available at geneages.org.