

1 **Title**

2 OrthoFinder: scalable phylogenetic orthology inference for comparative genomics

3 **Authors**

4 David M Emms¹, Yi Liu¹, Laurence Belcher¹, Jonathan Holmes¹ and Steven Kelly¹

5 **Abstract**

6 Here, we present a major advance of the OrthoFinder method. This extends OrthoFinder's high
7 accuracy comparative genomic framework to provide substantially enhanced scalability and
8 accuracy. Specifically, we show that enhanced phylogenetic delineation of orthogroups provides a
9 7% relative increase in orthogroup inference accuracy. We further demonstrate that a new gene
10 assignment method substantially reduces overall runtime RAM usage without compromising
11 accuracy. The latest version of OrthoFinder is available at
12 <https://github.com/OrthoFinder/OrthoFinder>.

13 **Affiliations**

14 1 Department of Biology, University of Oxford, South Parks Road, Oxford, OX1 3RB, United
15 Kingdom

16 **Corresponding Author**

17 Name: Steven Kelly

18 Email: steven.kelly@biology.ox.ac.uk

19 Address: Department of Biology, University of Oxford, South Parks Road, Oxford, OX1 3RB, UK

20 **Introduction**

21 Inferring orthology between biological sequences is fundamental to contemporary biological
22 research. It provides the foundation for studying the evolution and diversity of life on Earth and it
23 provides the framework for transfer of biological knowledge between species. Given the central
24 importance of orthology inference to biological research it has been the subject of extensive
25 methodological development for more than 40 years (Fitch, 1970; Gabaldón *et al.*, 2009; Altenhoff
26 *et al.*, 2016; Linard *et al.*, 2021). Diverse approaches to the computational challenge of identifying
27 related genes in different species has resulted in a broad array of methods with varied
28 performance characteristics when applied to diverse datasets (Kristensen *et al.*, 2011; Altenhoff *et al.*,
29 *et al.*, 2016; Nichio, Marchaukoski and Raittz, 2017).

30

31 Although orthology inference has been subject to substantial improvements since the advent of the
32 first automated methods, one of the major challenges facing method development in this field is
33 how to achieve high orthology inference accuracy at scale (Cosentino, Sriswasdi and Iwasaki,
34 2024; Majidian *et al.*, 2025). This challenge has become more pressing in recent years with
35 genome sequencing projects such as The Darwin Tree of Life (The Darwin Tree of Life Project
36 Consortium, 2022) and The Earth BioGenome Project (Lewin *et al.*, 2018), which collectively aim to
37 sequence, assemble and annotate the reference genomes for all two million known species of
38 eukaryote. These pioneering efforts, combined with widespread sequencing efforts that are
39 distributed across many biological communities, are creating an unprecedented scalability
40 challenge for automated orthology inference methods (Linard *et al.*, 2021; Langschied *et al.*, 2024).
41 Thousands of genomes are already available, it is likely that millions more are coming, and there is
42 a real need to be able to accurately and efficiently analyse these resources to maximize their
43 value.

44 There are several distinct computation bottlenecks for inferring orthology for large numbers of
45 species. For example, the entry point to the majority of inference methods is an all-vs-all sequence
46 similarity search from which orthology relationships can be inferred (Emms and Kelly, 2019).
47 Substantial efforts have been focused at increasing the computational efficiency of individual
48 sequence similarity searches with methods such as DIAMOND (Buchfink, Reuter and Drost, 2021),
49 Usearch (Edgar, 2010), and mmseqs (Hauser, Steinegger and Söding, 2016) greatly advancing
50 the capability to efficiently identify similar sequences in large databases. Even though major
51 improvements have been made, all-vs-all sequence searches have an order of n^2 time complexity
52 (where n is the number of species under consideration) and consequently are not well suited to
53 conducting orthology inference as species number increases. Thus, alternative orthology inference
54 approaches that retain high accuracy while improving scalability are required to realize the full
55 value of global sequencing efforts.

56 Previously we developed the OrthoFinder method for phylogenetic orthology inference (Emms and
57 Kelly, 2015, 2019). OrthoFinder first identifies orthogroups, the sets of genes that each descend

58 from a single gene in the last common ancestor of a group of species. OrthoFinder then infers a
59 gene tree for each orthogroup and analyses these gene trees to identify the rooted species tree
60 (Emms and Kelly, 2017, 2018). OrthoFinder also identifies all gene duplication events in the
61 complete set of gene trees and analyses this information in the context of the species tree to
62 provide both gene tree and species tree-level analysis of gene duplication events. Finally,
63 OrthoFinder analyses this cohort of phylogenetic information to identify the complete set of
64 orthologs between all species and provide a suite of comparative genomics statistics. The
65 phylogenetic approach developed in OrthoFinder provided a step change for comparative
66 genomics, enabling the transition from similarity score-based approximations of orthology to tree-
67 based phylogenetic relationships between genes. In recent work, we showed that it was possible to
68 rapidly and accurately place individual sequences into a database of phylogenetic trees and
69 biological sequences produced by an OrthoFinder search (Emms and Kelly, 2022). This approach
70 provided a rapid phylogenetic framework through which individual genes could be added to an
71 existing orthology inference run without incurring a computationally costly all-vs-all re-analysis of
72 existing species. We hypothesized that it should be possible to extend this single-gene approach to
73 support the addition of sets of species and thereby achieve improved scalability without
74 compromising accuracy.

75 Here we present a major update to OrthoFinder that significantly enhances the scalability of the
76 method. We show that it is possible to rapidly search and assign sequence sets from large
77 numbers of species to a phylogenetically partitioned and structured database of biological
78 sequences in near linear time. We demonstrate that using such an accelerated search functionality
79 does not compromise orthology inference accuracy. Moreover, through advancements in
80 phylogenetic interrogation of orthogroups we show that it is possible to achieve scalable inference
81 with higher accuracy than any competitor method or previous versions of OrthoFinder. Finally, we
82 show that the latest version of OrthoFinder is scalable to thousands of genomes on conventional
83 computing resources. As before, the updated version of OrthoFinder is accurate, customizable,
84 and is performed with a simple command using only protein sequences as input.

85 **Results**

86 ***Improved phylogenetic delineation of orthogroups***

87 In previous versions of OrthoFinder, orthologs were defined phylogenetically but orthogroups were
88 defined solely based on MCL clustering of sequence similarity search results (Emms and Kelly,
89 2015, 2019). While this method does correct for inter-species divergence it does not leverage the
90 phylogenetic relationship between the species represented in each orthogroup, and thus does not
91 always identify the true extent of an orthogroup. Errors in orthogroup delineation at this stage
92 propagate through all subsequent steps of the analysis, causing collateral errors in orthology
93 inference. To address this, we developed and implemented a novel phylogenetic re-evaluation of
94 orthogroup membership (Figure 1). First, MCL clustering is performed as in OrthoFinder v2, and a
95 gene tree is inferred for each orthogroup. Next, OrthoFinder applies its high accuracy gene tree
96 species tree reconciliation algorithm (Emms and Kelly, 2019) to identify and map all gene
97 duplication events in each gene tree. By definition, an orthogroup should not contain genes arising
98 from ‘ancient’ duplications that predate the root of the species tree, any such duplications indicate
99 that the genes have descended from multiple independent genes present in the last common
100 ancestor of the set of species under consideration. OrthoFinder splits all gene trees at such ancient
101 gene duplication nodes and in doing so creates a new set of revised phylogenetically defined
102 orthogroups which are then recorded and used for all subsequent downstream analyses.
103 OrthoFinder performs this for all nodes in the species tree creating orthogroups for each node in
104 the species tree. This new step splits erroneously fused orthogroups, and prunes sequences that
105 do not adhere to the phylogenetic definition of the orthogroup (Figure 1).

106 To test the impact of this novel approach on orthogroup inference accuracy, we compared this new
107 version of OrthoFinder to a set of alternative methods on the OrthoBench dataset of expert-curated
108 reference orthogroups (Trachana *et al.*, 2011; Emms and Kelly, 2020). We deployed OrthoFinder
109 in its two most commonly used implementations. One implementation uses multiple sequence
110 alignment and maximum likelihood tree inference to infer orthogroup trees and is provided with a
111 precomputed species tree (OF3_Align_ST), and the other implementation uses the DendroBLAST
112 method to bypass the requirement to construct alignments and alignment-based trees (Kelly and
113 Maini, 2013) (OF3_DB) which is the default implementation in OrthoFinder v2. We also included

114 the new scalable version of OrthoFinder v3 (OF3_Linear) which is discussed in detail in the
115 following section of the manuscript. The list of comparator methods included Broccoli (Derelle,
116 Philippe and Colbourne, 2020), OrthoHMM (Steenwyk *et al.*, 2024), SonicParanoid2 (Cosentino,
117 Sriswasdi and Iwasaki, 2024), OrthoMCL (Li, Stoeckert and Roos, 2003), ProteinOrtho (Klemm,
118 Stadler and Lechner, 2023), Hieranoid (Kaduk and Sonnhammer, 2017), and FastOMA (Majidian
119 *et al.*, 2025).

120 For each method, we calculated seven measures of accuracy: (A) the percentage of missing
121 reference orthogroups, (B) the percentage of missing genes, (C) the percentage of incorrect
122 orthogroup fusion events, (D) the percentage of incorrect orthogroup fission events, and (E, F, G)
123 recall, precision, and entropy. Entropy is the measure of reference orthogroup fragmentation, with
124 higher values indicating greater fragmentation. We then calculated an overall rank-based score by
125 evaluating the relative performance of each method across all seven measures. This revealed that
126 OrthoFinder v3 (in either implementation) outperformed all other orthogroup inference methods
127 including OrthoFinder v2 (Figure 2). OrthoFinder v3 had the highest recall, lowest fraction of
128 missing genes, and the lowest entropy. Moreover, OrthoFinder v3 had 5-7% higher accuracy (f-
129 score harmonic mean of recall and precision. Supplementary S1 Table 1) than the previous version
130 of OrthoFinder run with identical settings (OF2_Align_ST and OF2_DB, Figure 2). While
131 OrthoFinder v3 achieved a lower precision than several other methods, OrthoFinder v3 showed a
132 substantially higher recall, and consequently suffered from a much lower rate of missing data. For
133 example, while SonicParanoid2 (SP_def) showed a 3.7% higher precision than OrthoFinder v3
134 (OF3_Align_ST), the recall of OrthoFinder v3 was 19% higher. Almost all tools were identical in the
135 production of fused reference orthogroups (RefOG), with the exception of FastOMA which
136 erroneously fused 3.5 times more reference orthogroups than other methods (10% compared to
137 2.85%). The spread of reference orthogroup fissions events across the tools varied, with low
138 precision tools such as Broccoli and OrthoHMM making relatively few RefOG fissions. Analysis of
139 the entropy scores, which are similar to a weighted combination of fission and fusion events,
140 revealed that OrthoFinder tends to produce orthogroups whose sizes and membership most
141 accurately recapture the reference dataset (Figure 2). Therefore, the introduction of the

142 phylogenetic delimitation of orthogroups improves overall orthogroup inference accuracy, and
143 OrthoFinder v3 is the most accurate method currently available.

144 ***Achieving enhanced scalability without compromising accuracy***

145 In addition to improvements in the phylogenetic definition of orthogroups, we also developed a
146 novel implementation of OrthoFinder to improve the scalability of the method. This is achieved
147 through adaptation and further development of the SHOOT profile algorithm (Emms and Kelly,
148 2022) which enables rapid addition of new species to existing OrthoFinder runs. This novel
149 scalable version of OrthoFinder is implemented as a two-step process (Figure 3). Prior to starting,
150 the input species set are partitioned into two non-overlapping subsets – a “core” subset and an
151 “assign” subset. It is recommended that the “core” subset contain fewer than 100 species for
152 analysis on conventional computing resources. The “assign” subset can be substantially larger as
153 described below. The first step of this new implementation requires that a conventional
154 OrthoFinder run is performed on the core subset. This creates a phylogenetically partitioned and
155 structured reference database of biological sequences (Emms and Kelly, 2022) that is then used
156 as input for the second stage of the algorithm. The species from the assign subset are then rapidly
157 assigned to the correct orthogroup in the reference database, irrespective of the number of
158 sequences contained within each core orthogroup. Following completion of the assignment step,
159 OrthoFinder’s rapid phylogeny-based orthogroup and ortholog inference steps are performed to
160 produce an expanded rooted species tree, gene trees, phylogenetically determined orthogroups,
161 orthologs, gene duplication events, and comparative genomic statistics (Figure 3).

162 To test the impact on scalability of this new workflow, OrthoFinder v3 was compared against the
163 set of commonly used orthology inference tools. The time each method took to complete analyses
164 of datasets of different sizes was measured, unless a timeout of 7 days was reached before
165 completion. Peak memory usage was also calculated for each analysis. Each tool was run on a
166 Linux server using AMD processors, all tools were allocated 32 threads (16 cores) and up to 200gb
167 of RAM for each run. In order to test performance on a representative set of species, the Ensembl
168 (Harrison *et al.*, 2024) rapid release genomes (Accessed 29th of August 2024) were downloaded.
169 Test datasets ranging in size from 2 - 1024 species we compiled by sampling species from
170 Ensembl dataset using PDA (phylogenetic diversity analyzer) (Chernomor *et al.*, 2015), see

171 'Methods' for more details. All methods were run on the same datasets in their default
172 implementations unless indicated otherwise in the methods or figure legends. OrthoFinder v3 was
173 run using the standard OrthoFinder 2 workflow up to 64 species. Beyond 64 species OrthoFinder
174 v3 was run in the new scalable implementation where the first 64 species were taken to be the
175 core set and all other species were assigned. The final run time for OrthoFinder v3 was calculated
176 using the combined run time for the core 64 species and the assigned dataset.

177 OrthoFinder v3 linear addition outperformed all other tested methods, being the only method able
178 to perform orthology inference on 1024 species within the seven-day cutoff (Figure 4A).
179 OrthoFinder v3 completed this run in 128 hours. SonicParanoid2 (fast mode) and FastOMA were
180 the only two methods able to run 512 proteomes within the cutoff time. OrthoFinder v3 followed a
181 near linear run time trend for 128 – 1024 species. This new implementation of OrthoFinder was
182 faster than each of the other competitor methods tested and significantly (x8) faster than
183 OrthoFinder v2 on datasets larger than 64 species, an important improvement in scalability from
184 the previous version.

185 OrthoFinder v3 also performed well on memory usage (Figure 4B). OrthoFinder v3 outperformed
186 OrthoFinder v2 above 128 species, with increasing memory saving as species number increases.
187 For example, v3 linear shows a 3.4-fold decrease in RAM consumption for 256 species compared
188 to v2 DendroBlast. On the largest datasets that other methods could complete, OrthoFinder v3,
189 used approximately 4-fold lower peak memory consumption (Figure 4B).

190 To demonstrate that the alternative implementation did not compromise orthogroup inference
191 accuracy, the linear species addition method (OF3_Linear) was also tested on the OrthoBench
192 data (Figure 2). OrthoFinder v3 Linear only slightly underperformed compared to the non-scalable
193 version, with slight decreases in some measures observed. However, this scalable version was
194 more accurate than any implementation of OrthoFinder v2 and any other competitor methods
195 (Figure 2), supporting the use of this accelerated implementation of OrthoFinder in comparative
196 genomic analyses.

197 ***OrthoFinder 3 performs accurate ortholog inference on the quest for orthologs***
198 ***benchmarking service.***

199 We also tested the accuracy of the orthologs predicted by OrthoFinder v3 in multiple
200 implementations using the Quest for Orthologs (QfO) benchmarking service (Gabaldón *et al.*,
201 2009; Altenhoff *et al.*, 2016; Linard *et al.*, 2021). QfO is the most widely used orthology
202 benchmarking service, which evaluates a tool's ability to accurately predict orthologs across a
203 diverse set of taxa, including Archaea, Bacteria, and eukaryotes. QfO gives many different
204 measures of accuracy, including comparisons to expert-curated reference orthologs and the
205 accuracy of species trees inferred from a set of orthologs. We used the most recent published set
206 of reference proteomes, which consists of 78 species (23 bacteria, 48 eukaryotes, and 7 archaea)
207 (Nevers *et al.*, 2022). We compared OrthoFinder v3 against the results obtained from 14 orthology
208 inference tools.

209 Analysis of the resulting scores revealed that the new OrthoFinder v3 workflow for high scalability
210 does not compromise accuracy (Figure 5). OrthoFinder v3 is on the Pareto frontier for species tree
211 discordance tests in both eukaryotes (Figure 5A), and bacteria (Figure 5B), demonstrating its
212 ability to perform accurate orthology inference across different domains of life. The only other tool
213 with scalability comparable to OrthoFinder v3 is FastOMA, but this method has substantially lower
214 ortholog inference accuracy (Figure 5). For example, while the Robinson-Foulds distance
215 (measuring disagreement between true and inferred species trees) for eukaryotes is marginally
216 worse in OrthoFinder v3 compared to FastOMA (0.06 vs. 0.05), OrthoFinder v3 has an 80% higher
217 recall (15721 vs. 8686). Similarly, OrthoFinder v3 is marginally better than FastOMA for Robinson-
218 Foulds distance on the bacteria species tree discordance test (0.590 vs. 0.587), but again has a
219 23% higher recall. OrthoFinder v3 is on the Pareto frontier for the enzyme classification test (Figure
220 5C), which measures how well predicted orthologs match to Enzyme Commission numbers (which
221 are given depending on the chemical reaction they catalyze). OrthoFinder v3 is only outperformed
222 in precision by tools that rely on pre-computed databases (e.g. OrthoInspector and PANTHER),
223 and is better than FastOMA in both precision (0.933 vs 0.928) and recall (183368 vs 157049).

224 OrthoFinder v3 is on the Pareto frontier for all three of the human-curated reference sets (Figure 5
225 D-F), scoring particularly well on recall. OrthoFinder v3 has the highest recall of all methods on

226 both the VGNC and SwissTree human-curated reference sets and is only marginally beaten for
227 recall on the TreeFam dataset by OrthoFinder v2 (0.72 vs. 0.74). In summary, similar to the high
228 accuracy of its orthogroup inference, OrthoFinder v3 also exhibits high accuracy ortholog inference
229 across multiple benchmark tests.

230 **Discussion**

231 Planet Earth is home to more than 6,000 species of mammal (Burgin *et al.*, 2018), 300,000 species
232 of plant (Christenhusz and Byng, 2016), 5,000,000 species of insect (Gaston, 1991), and an
233 unknown number of species of unicellular eukaryotes, bacteria and archaea (Mora *et al.*, 2011;
234 Hugenholtz *et al.*, 2021). Inferring the phylogenetic relationships between the biological sequences
235 of these organisms provides a foundation upon which we can study evolution and molecular
236 diversity, and enables us to understand and transfer biological information between organisms. It is
237 likely that a substantial fraction of all known species on Earth will have a representative genome in
238 the coming decades, but methods do not yet exist to enable this scale of data to be analysed.
239 Here, we present a major advance to OrthoFinder that substantially improves the scalability and
240 accuracy of the method. We show it is possible to rapidly assign gene sets from large numbers of
241 species to a phylogenetically partitioned and structured database of biological sequences in near
242 linear time. We also show that using phylogenetic delimitation improves the inference of
243 orthogroups, enhancing the accuracy of OrthoFinder and extending its performance advantage
244 over comparator methods. OrthoFinder remains an easy-to-use, fast, accurate, and fully
245 phylogenetic orthology inference software tool. OrthoFinder also continues to be one of the few
246 tools to provide a wide range of output information including gene duplications, gene trees,
247 sequence alignments and single copy ortholog sequences allowing downstream analysis.

248 The updated OrthoFinder method is now capable of analysing thousands of species using a novel
249 two-step process. Although this process requires an additional step for the user, the only input
250 required is still the set amino acid sequences corresponding to the protein-coding genes for the
251 species of interest. The default parameters for OrthoFinder have been optimized for speed,
252 accuracy, and scalability and enable the combined analysis of thousands of species on commonly
253 available computing resources. OrthoFinder also retains its customizability for expert users, and
254 intermediate steps in the algorithm (such as alignment or tree inference) can be substituted with

255 alternative methods should the user wish. Although there is still some distance to travel, this
256 upgrade to OrthoFinder provides an important step towards the goal of being able to provide high
257 accuracy phylogenetic orthology inference for all species on Earth.

258 **Materials and Methods**

259 ***The OrthoFinder v3 workflow enhanced scalability workflow***

260 The scalable linear species addition framework of OrthoFinder v3 requires the results of a
261 conventional OrthoFinder analysis of an ‘core’ set of species. For each orthogroup, OrthoFinder v3
262 employs the SHOOT profile algorithm (Emms and Kelly, 2022) to select the most representative
263 sequences using k-means clustering applied to an embedding of the sequences from the length L
264 multiple sequence alignment (MSA) of the orthogroup in a 2*L-dimensional space. DIAMOND
265 (Buchfink, Reuter and Drost, 2021) is used to assign the new genes to the core orthogroups using
266 these profiles. Genes not assigned to any core orthogroups are set aside to be analysed
267 separately at a later stage. Once all genes have been assigned, gene trees are inferred. To
268 support the analysis of larger datasets, the MAFFT multiple sequence alignment method has been
269 replaced with the more scalable FAMSA method (Deorowicz, Debudaj-Grabysz and Gudyś, 2016).
270 Users can revert back to previous default alignment method “-A mafft”. Maximum likelihood tree
271 inference of the resulting alignments is then performed using FastTree (Price, Dehal and Arkin,
272 2010) employing the “-fastest” runtime option. The use of other alignment or tree inference
273 programs is supported through use of the config.json file, although users are advised to test if they
274 have sufficient computational resources to use more computationally expensive programs. Once
275 all orthogroup trees are inferred, a revised species tree is then computed using ASTRAL-Pro
276 (Zhang and Mirarab, 2022).

277 As mentioned above, some genes fail to be assigned to any core orthogroups during the initial
278 assignment step. This occurs because the new species (or sets of species) being added will
279 sometimes contain new genes that have arisen *de novo*, and are not found in the core species set.
280 To resolve this, OrthoFinder performs a root to tip tree traversal to identify clades of species
281 (minimum size = 2) that had no representation in the now sparsely sampled initial core species set.
282 Once identified a conventional OrthoFinder analysis is performed on the clade, and the new
283 orthogroups are added to the revised orthogroup set.

284 The original OrthoFinder algorithm applied a per-species-pair gene length normalisation that
285 accounted for both evolutionary-distance and gene-length dependence in pairwise bit-scores
286 between sequences (Emms and Kelly, 2019). However, this normalisation step requires a large
287 number of pairwise hits between the genes in each species in order to fit the required correction
288 functions. To enable the smaller set analysis above, the updated version of OrthoFinder employs
289 an alternative normalisation approach. The approach separates the sequence length bias and
290 species divergences to obtain many independent estimates of the extent of each orthogroup
291 Specifically, bit-scores between two sequences, B_{qh} , are length normalised using geometric mean
292 of the lengths of the query and hit sequences, L_q and L_h : $B_{qh}=B_{qnr}/(L_qL_h)^{1/2}$. As before, normalised
293 reciprocal best hits (RBH) are treated as putative orthologs and are used to estimate the sequence
294 divergence between pairs of species and on a per-orthogroup basis. In the modified version we
295 calculate a set of ratios R_{XY}^* for species X and Y. These estimate the expected ratio of the
296 normalised bit score between the RBH for a pair of genes x and y from species X and Y from some
297 orthogroup and the normalised bit score from the gene x to the most distant gene in that
298 orthogroup. Using these, every RBH observed can be used to estimate a cut-off for the most
299 distant gene in that orthogroup. In this way, if a gene has an RBH in multiple species then it
300 provides multiple independent estimates of the cut-off for the most distant sequence in the
301 orthogroup. Likewise, if multiple genes from a (currently undetermined) orthogroup obtain one or
302 more RBH, then these will each provide multiple independent estimates of the extent of the
303 orthogroup. These estimates are used to construct a graph such that a pair of sequences is
304 connected within the graph if the NBS between those sequences fall within that gene-specific and
305 species-specific cut-off. As such, this graph attempts to connect as densely as possible with edges
306 all sequences within the same orthogroup while having as few edges as possible between
307 sequences in different orthogroups. As before, the MCL algorithm (Van Dongen, 2008) is applied
308 to identify the sets of genes in the graph best satisfying the putative shared orthogroup
309 memberships encoded in the edges of this graph. As with the core orthogroups, the resulting gene
310 trees of the putative orthogroups are analysed to phylogenetically resolve the true extent of each
311 orthogroup (Figure 1). The non-core orthogroups are subjected to phylogenetic analysis to infer

312 rooted gene trees and subsequently these rooted gene trees are analysed to determine the full
313 comparative analysis including orthologs, hierarchical orthogroups, and gene duplication events.

314 **Assessing orthogroup inference accuracy**

315 The full details on the orthogroup benchmarking, including the calculation of the scores, and a
316 complete method by method breakdown of the results are provided in Supplemental File 1 and the
317 complete set of output files generated by each method is provided in Supplemental File 2.

318 **Assessing scalability performance**

319 The Ensembl (Harrison *et al.*, 2024) rapid release dataset was accessed on the 29th August 2024.
320 Amino acid sequences for all protein coding genes in each genome were downloaded. For species
321 which had more than one version available, the most recent proteome file was used. The
322 OrthoFinder v2 script 'primary_transcript.py' was used to extract only the longest variant for each
323 gene giving a dataset with a total of 1789 species. Following collection of the Ensembl proteomes
324 an approximate species tree was derived in order to extract a set of representative proteomes for
325 performance testing across each orthology inference tool. The species tree was generated using
326 orthologs of BUSCO genes (Manni *et al.*, 2021). In brief, the hidden Markov model provided for
327 each BUSCO gene was used to identify high scoring amino acid sequences in each species. The
328 resulting hits from each species to each BUSCO gene were aligned using MAFFT, and the
329 resulting multiple sequence alignments concatenated. An unrooted phylogenetic tree was then
330 inferred from this concatenated alignment using FastTree. This tree was then used as input to
331 Phylogenetic Diversity Analyzer (Chernomor *et al.*, 2015) to select sets of species in the range of
332 2-1024 species such that each subsequent set contained all the species of the previous set, and
333 that each subset maximised the sampled branch length of the inferred tree. The complete result
334 dataset for runtime and scalability of all methods is provided in Supplemental File 3.

335 **Assessing ortholog inference accuracy**

336 OrthoFinder v3 was run on the 2022 edition of the Quest for Orthologs reference proteomes. We
337 ran both OrthoFinder v3 Align (with FAMSA), and OrthoFinder v3 Linear (--core --assign).
338 OrthoFinder v2 along with 13 other orthology inference tools have already been run on this
339 dataset, and is available as part of the publicly available results alongside many methods. The
340 Quest for Orthologs datasets consists of 78 proteomes (48 eukaryotes, 23 prokaryotes, and 7
341 archaea). We ran OrthoFinder on these input datasets and extracted pairs of orthologous genes

342 identified by the method. The file containing orthologous pairs was uploaded to the Quest for
343 Orthologs web server for benchmarking. The Quest for Orthologs papers describe the
344 methodology for benchmarking in full detail (Altenhoff *et al.*, 2016; Nevers *et al.*, 2022).

345 **Acknowledgements**

346 The authors would like to thank the OrthoFinder user community for their feedback on the method
347 and their continued support.

348 **Funding**

349 YL, LB, JH, SK were funded by the Wellcome Trust under grant agreement number
350 226598/Z/22/Z.

351 **Author information**

352 DE developed the method. DE, YL, LB, JH implemented software. DE, YL, LB, JH, SK contributed
353 to the analysis. DE, LB, JH, SK wrote and edited the manuscript.

354 **Ethics declaration**

355 Sk is co-founder of Wild Bioscience Ltd and an employee of Ellison Institute of Technology, Oxford
356 Limited.

357 **References**

- 358 Altenhoff, A.M. *et al.* (2016) 'Standardized benchmarking in the quest for orthologs', *Nature*
359 *Methods*, 13(5), pp. 425–430. Available at: <https://doi.org/10.1038/nmeth.3830>.
- 360 Buchfink, B., Reuter, K. and Drost, H.-G. (2021) 'Sensitive protein alignments at tree-of-life scale
361 using DIAMOND', *Nature Methods*, 18(4), pp. 366–368. Available at:
362 <https://doi.org/10.1038/s41592-021-01101-x>.
- 363 Burgin, C.J. *et al.* (2018) 'How many species of mammals are there?', *Journal of Mammalogy*,
364 99(1), pp. 1–14. Available at: <https://doi.org/10.1093/jmammal/gyx147>.
- 365 Chernomor, O. *et al.* (2015) 'Split diversity in constrained conservation prioritization using integer
366 linear programming', *Methods in Ecology and Evolution*, 6(1), pp. 83–91. Available at:
367 <https://doi.org/10.1111/2041-210X.12299>.
- 368 Christenhusz, M.J.M. and Byng, J.W. (2016) 'The number of known plants species in the world and
369 its annual increase', *Phytotaxa*, 261(3), pp. 201–217. Available at:
370 <https://doi.org/10.11646/phytotaxa.261.3.1>.
- 371 Cosentino, S., Sriswasdi, S. and Iwasaki, W. (2024) 'SonicParanoid2: fast, accurate, and
372 comprehensive orthology inference with machine learning and language models', *Genome Biology*,
373 25(1), p. 195. Available at: <https://doi.org/10.1186/s13059-024-03298-4>.
- 374 Deorowicz, S., Debudaj-Grabysz, A. and Gudyś, A. (2016) 'FAMSA: Fast and accurate multiple
375 sequence alignment of huge protein families', *Scientific Reports*, 6(1), p. 33964. Available at:
376 <https://doi.org/10.1038/srep33964>.

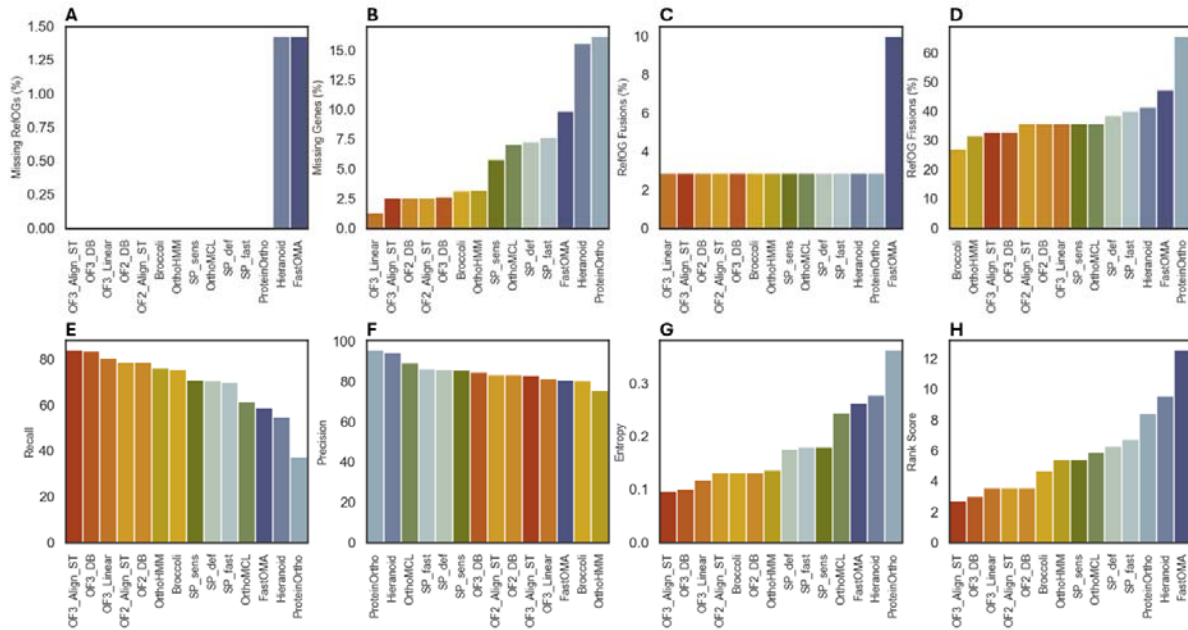
- 377 Derelle, R., Philippe, H. and Colbourne, J.K. (2020) 'Broccoli: Combining Phylogenetic and
378 Network Analyses for Orthology Assignment', *Molecular Biology and Evolution*. Edited by D.
379 Falush, 37(11), pp. 3389–3396. Available at: <https://doi.org/10.1093/molbev/msaa159>.
- 380 Edgar, R.C. (2010) 'Search and clustering orders of magnitude faster than BLAST', *Bioinformatics*,
381 26(19), pp. 2460–2461. Available at: <https://doi.org/10.1093/bioinformatics/btq461>.
- 382 Emms, D.M. and Kelly, S. (2015) 'OrthoFinder: solving fundamental biases in whole genome
383 comparisons dramatically improves orthogroup inference accuracy', *Genome Biology*, 16(1), p.
384 157. Available at: <https://doi.org/10.1186/s13059-015-0721-2>.
- 385 Emms, D.M. and Kelly, S. (2017) 'STRIDE: Species Tree Root Inference from Gene Duplication
386 Events', *Molecular Biology and Evolution*, 34(12), pp. 3267–3278. Available at:
387 <https://doi.org/10.1093/molbev/msx259>.
- 388 Emms, D.M. and Kelly, S. (2018) 'STAG: Species Tree Inference from All Genes'. bioRxiv, p.
389 267914. Available at: <https://doi.org/10.1101/267914>.
- 390 Emms, D.M. and Kelly, S. (2019) 'OrthoFinder: phylogenetic orthology inference for comparative
391 genomics', *Genome Biology*, 20(1), p. 238. Available at: [https://doi.org/10.1186/s13059-019-1832-](https://doi.org/10.1186/s13059-019-1832-y)
392 [y](https://doi.org/10.1186/s13059-019-1832-y).
- 393 Emms, D.M. and Kelly, S. (2020) 'Benchmarking Orthogroup Inference Accuracy: Revisiting
394 Orthobench', *Genome Biology and Evolution*, 12(12), pp. 2258–2266. Available at:
395 <https://doi.org/10.1093/gbe/evaa211>.
- 396 Emms, D.M. and Kelly, S. (2022) 'SHOOT: phylogenetic gene search and ortholog inference',
397 *Genome Biology*, 23(1), p. 85. Available at: <https://doi.org/10.1186/s13059-022-02652-8>.
- 398 Fitch, W.M. (1970) 'Distinguishing Homologous from Analogous Proteins', *Systematic Biology*,
399 19(2), pp. 99–113. Available at: <https://doi.org/10.2307/2412448>.
- 400 Gabaldón, T. *et al.* (2009) 'Joining forces in the quest for orthologs', *Genome Biology*, 10(9), p.
401 403. Available at: <https://doi.org/10.1186/gb-2009-10-9-403>.
- 402 Gaston, K.J. (1991) 'The Magnitude of Global Insect Species Richness', *Conservation Biology*,
403 5(3), pp. 283–296. Available at: <https://doi.org/10.1111/j.1523-1739.1991.tb00140.x>.
- 404 Harrison, P.W. *et al.* (2024) 'Ensembl 2024', *Nucleic Acids Research*, 52(D1), pp. D891–D899.
405 Available at: <https://doi.org/10.1093/nar/gkad1049>.
- 406 Hauser, M., Steinegger, M. and Söding, J. (2016) 'MMseqs software suite for fast and deep
407 clustering and searching of large protein sequence sets', *Bioinformatics*, 32(9), pp. 1323–1330.
408 Available at: <https://doi.org/10.1093/bioinformatics/btw006>.
- 409 Hugenholtz, P. *et al.* (2021) 'Prokaryotic taxonomy and nomenclature in the age of big sequence
410 data', *The ISME Journal*, 15(7), pp. 1879–1892. Available at: [https://doi.org/10.1038/s41396-021-](https://doi.org/10.1038/s41396-021-00941-x)
411 [00941-x](https://doi.org/10.1038/s41396-021-00941-x).
- 412 Kaduk, M. and Sonnhammer, E. (2017) 'Improved orthology inference with Hieranoid 2',
413 *Bioinformatics*, 33(8), pp. 1154–1159. Available at: <https://doi.org/10.1093/bioinformatics/btw774>.
- 414 Kelly, S. and Maini, P.K. (2013) 'DendroBLAST: Approximate Phylogenetic Trees in the Absence
415 of Multiple Sequence Alignments', *PLOS ONE*, 8(3), p. e58537. Available at:
416 <https://doi.org/10.1371/journal.pone.0058537>.

- 417 Klemm, P., Stadler, P.F. and Lechner, M. (2023) 'Proteinortho6: pseudo-reciprocal best alignment
418 heuristic for graph-based detection of (co-)orthologs', *Frontiers in Bioinformatics*, 3. Available at:
419 <https://doi.org/10.3389/fbinf.2023.1322477>.
- 420 Kristensen, D.M. *et al.* (2011) 'Computational methods for Gene Orthology inference', *Briefings in*
421 *Bioinformatics*, 12(5), pp. 379–391. Available at: <https://doi.org/10.1093/bib/bbr030>.
- 422 Langschied, F. *et al.* (2024) 'Quest for Orthologs in the Era of Biodiversity Genomics', *Genome*
423 *Biology and Evolution*, 16(10), p. evae224. Available at: <https://doi.org/10.1093/gbe/evae224>.
- 424 Lewin, H.A. *et al.* (2018) 'Earth BioGenome Project: Sequencing life for the future of life',
425 *Proceedings of the National Academy of Sciences*, 115(17), pp. 4325–4333. Available at:
426 <https://doi.org/10.1073/pnas.1720115115>.
- 427 Li, L., Stoeckert, C.J. and Roos, D.S. (2003) 'OrthoMCL: Identification of Ortholog Groups for
428 Eukaryotic Genomes', *Genome Research*, 13(9), pp. 2178–2189. Available at:
429 <https://doi.org/10.1101/gr.1224503>.
- 430 Linard, B. *et al.* (2021) 'Ten Years of Collaborative Progress in the Quest for Orthologs', *Molecular*
431 *Biology and Evolution*, 38(8), pp. 3033–3045. Available at:
432 <https://doi.org/10.1093/molbev/msab098>.
- 433 Majidian, S. *et al.* (2025) 'Orthology inference at scale with FastOMA', *Nature Methods*, 22(2), pp.
434 269–272. Available at: <https://doi.org/10.1038/s41592-024-02552-8>.
- 435 Manni, M. *et al.* (2021) 'BUSCO Update: Novel and Streamlined Workflows along with Broader and
436 Deeper Phylogenetic Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes',
437 *Molecular Biology and Evolution*, 38(10), pp. 4647–4654. Available at:
438 <https://doi.org/10.1093/molbev/msab199>.
- 439 Mora, C. *et al.* (2011) 'How Many Species Are There on Earth and in the Ocean?', *PLOS Biology*,
440 9(8), p. e1001127. Available at: <https://doi.org/10.1371/journal.pbio.1001127>.
- 441 Nevers, Y. *et al.* (2022) 'The Quest for Orthologs orthology benchmark service in 2022', *Nucleic*
442 *Acids Research*, 50(W1), pp. W623–W632. Available at: <https://doi.org/10.1093/nar/gkac330>.
- 443 Nichio, B.T.L., Marchaukoski, J.N. and Raittz, R.T. (2017) 'New Tools in Orthology Analysis: A
444 Brief Review of Promising Perspectives', *Frontiers in Genetics*, 8, p. 165. Available at:
445 <https://doi.org/10.3389/fgene.2017.00165>.
- 446 Price, M.N., Dehal, P.S. and Arkin, A.P. (2010) 'FastTree 2 – Approximately Maximum-Likelihood
447 Trees for Large Alignments', *PLOS ONE*, 5(3), p. e9490. Available at:
448 <https://doi.org/10.1371/journal.pone.0009490>.
- 449 Steenwyk, J.L. *et al.* (2024) 'OrthoHMM: Improved Inference of Ortholog Groups using Hidden
450 Markov Models'. *bioRxiv*, p. 2024.12.07.627370. Available at:
451 <https://doi.org/10.1101/2024.12.07.627370>.
- 452 The Darwin Tree of Life Project Consortium (2022) 'Sequence locally, think globally: The Darwin
453 Tree of Life Project', *Proceedings of the National Academy of Sciences*, 119(4), p. e2115642118.
454 Available at: <https://doi.org/10.1073/pnas.2115642118>.
- 455 Trachana, K. *et al.* (2011) 'Orthology prediction methods: A quality assessment using curated
456 protein families', *BioEssays*, 33(10), pp. 769–780. Available at:
457 <https://doi.org/10.1002/bies.201100062>.

458 Van Dongen, S. (2008) 'Graph Clustering Via a Discrete Uncoupling Process', *SIAM Journal on*
459 *Matrix Analysis and Applications*, 30(1), pp. 121–141. Available at:
460 <https://doi.org/10.1137/040608635>.

461 Zhang, C. and Mirarab, S. (2022) 'ASTRAL-Pro 2: ultrafast species tree reconstruction from multi-
462 copy gene family trees', *Bioinformatics*, 38(21), pp. 4949–4950. Available at:
463 <https://doi.org/10.1093/bioinformatics/btac620>.

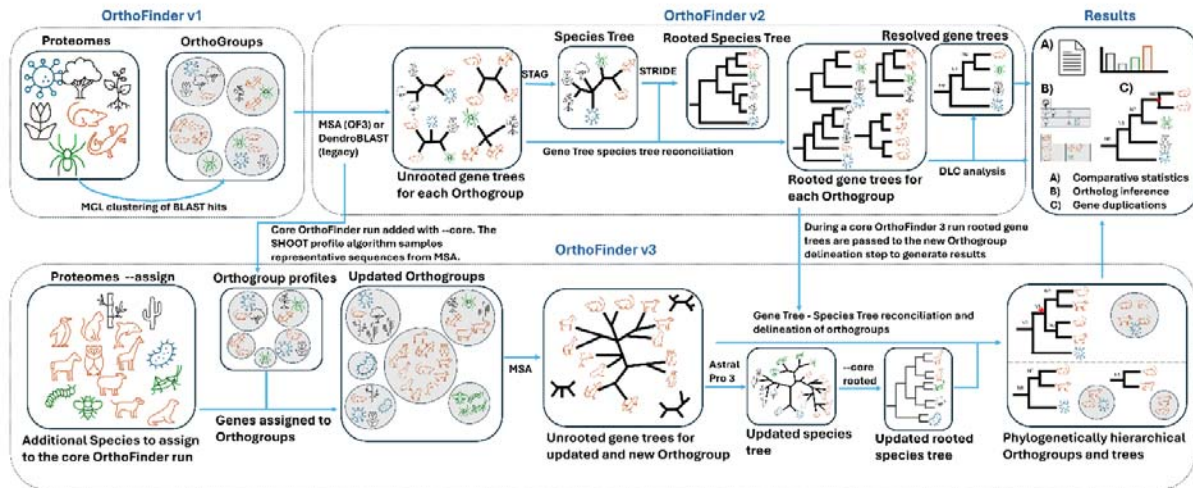
474 **Figure 2**



475

476 **Figure 2:** Comparative analysis of different orthogroup inference tools on the OrthoBench dataset.
 477 Panels A-G each show one of seven performance metrics used to evaluate orthology prediction
 478 accuracy. Lower values indicate better performance for error-based metrics (A-D,G), higher values
 479 are better for recall and precision (E-F). Panel H shows the rank score, combining all seven
 480 metrics.

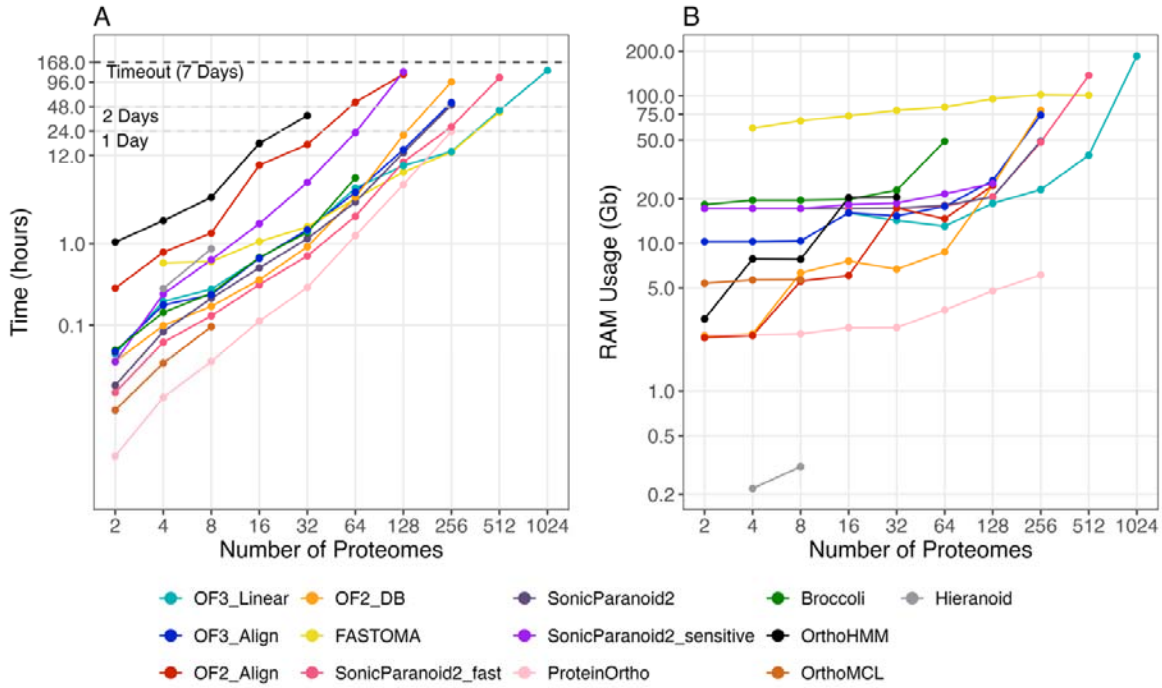
481 **Figure 3**



482

483 **Figure 3:** The OrthoFinder v3 workflow. The first step of OrthoFinder 3 is to perform a conventional
484 OrthoFinder run on a set of core proteomes. The second step then uses an adapted version of the
485 SHOOT profiling algorithm to assign genes from additional species to these core orthogroups and
486 identify new orthogroups not present in the initial input gene set. Gene trees are then inferred and
487 used to infer a revised species tree. All results files are then revised and updated.

488 **Figure 4**



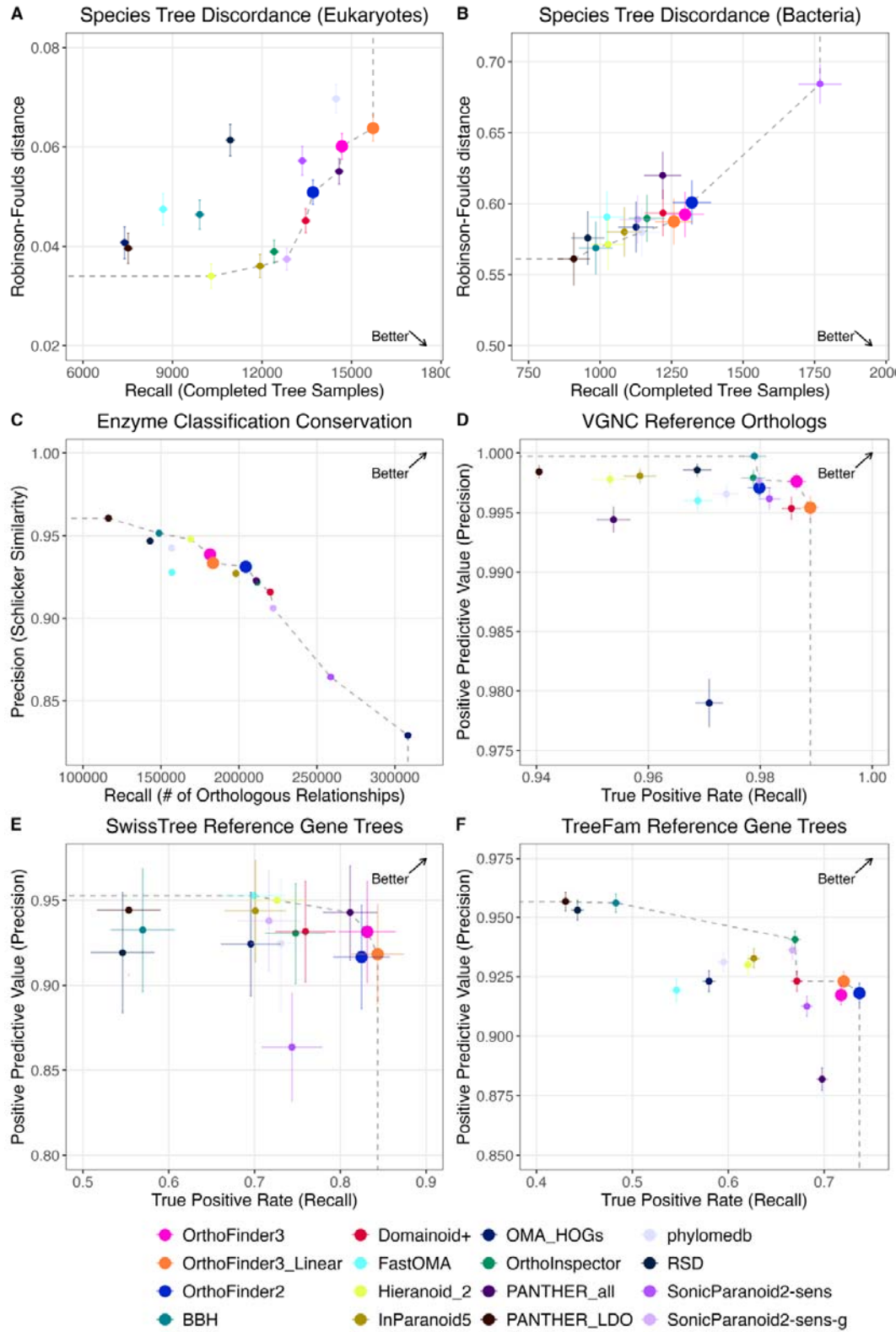
489

490 **Figure 4:** A comparison of runtime and peak memory usage for orthogroup inference tools. (A)
491 Time (log scale) and (B) peak RAM usage required to run different orthology prediction tools
492 across increasing numbers of input proteomes (from 2 to 1024).

493

494

495 **Figure 5**



496

497 **Figure 5:** A comparative analysis of ortholog inference accuracy on the Quest for Orthologs
 498 dataset. Each panel shows a different benchmark from the QfO benchmarking service (2022

499 release). OrthoFinder is compared against other ortholog inference methods which have made
500 their results publicly available. (A–B) Species tree discordance (Robinson-Foulds distance) vs.
501 recall across completed gene trees for eukaryotes (A) and bacteria (B). (C) Conservation of
502 enzyme classification across orthologs (Schlicker similarity) vs. number of orthologous
503 relationships. (D) VGNC reference orthologs: precision vs. recall. (E–F) SwissTree (E) and
504 TreeFam (F) reference gene trees: precision vs. recall. Better-performing methods are toward the
505 top-right for precision measures (C-F), and bottom right for discordance measures (A-B).