

SOFTWARE

Open Access

OrthoFinder: phylogenetic orthology inference for comparative genomics



David M. Emms and Steven Kelly*

Abstract

Here, we present a major advance of the OrthoFinder method. This extends OrthoFinder's high accuracy orthogroup inference to provide phylogenetic inference of orthologs, rooted gene trees, gene duplication events, the rooted species tree, and comparative genomics statistics. Each output is benchmarked on appropriate real or simulated datasets, and where comparable methods exist, OrthoFinder is equivalent to or outperforms these methods. Furthermore, OrthoFinder is the most accurate ortholog inference method on the Quest for Orthologs benchmark test. Finally, OrthoFinder's comprehensive phylogenetic analysis is achieved with equivalent speed and scalability to the fastest, score-based heuristic methods. OrthoFinder is available at <https://github.com/davidemms/OrthoFinder>.

Keywords: Ortholog inference, Gene tree inference, Gene duplication, Comparative genomics

Background

Determining the phylogenetic relationships between gene sequences is fundamental to comparative biological research. It provides the framework for understanding the evolution and diversity of life on Earth and enables the extrapolation of biological knowledge between organisms. Given the central importance of this process to multiple areas of biological research, a diverse array of software tools have been developed that attempt to identify these relationships given sets of user-supplied gene sequences [1–3]. The majority of these software tools try to deduce phylogenetic relationships between gene sequences through heuristic analyses of pairwise sequence similarity scores (or expectation values) obtained from an all-vs-all BLAST [4] search, or accelerated alternatives to BLAST such as DIAMOND [5] or MMseqs2 [6]. Widely used methods include InParanoid [7], OrthoMCL [8], OMA [9], and OrthoFinder [10] all of which take different approaches to interrogating sequence similarity scores, and all of which produce different outputs—some identify orthogroups, some identify orthologs and paralog, and some do both. As they each adopt different approaches to analyzing sequence similarity scores, each of the methods exhibits different performance characteristics on commonly used benchmark databases [1, 11].

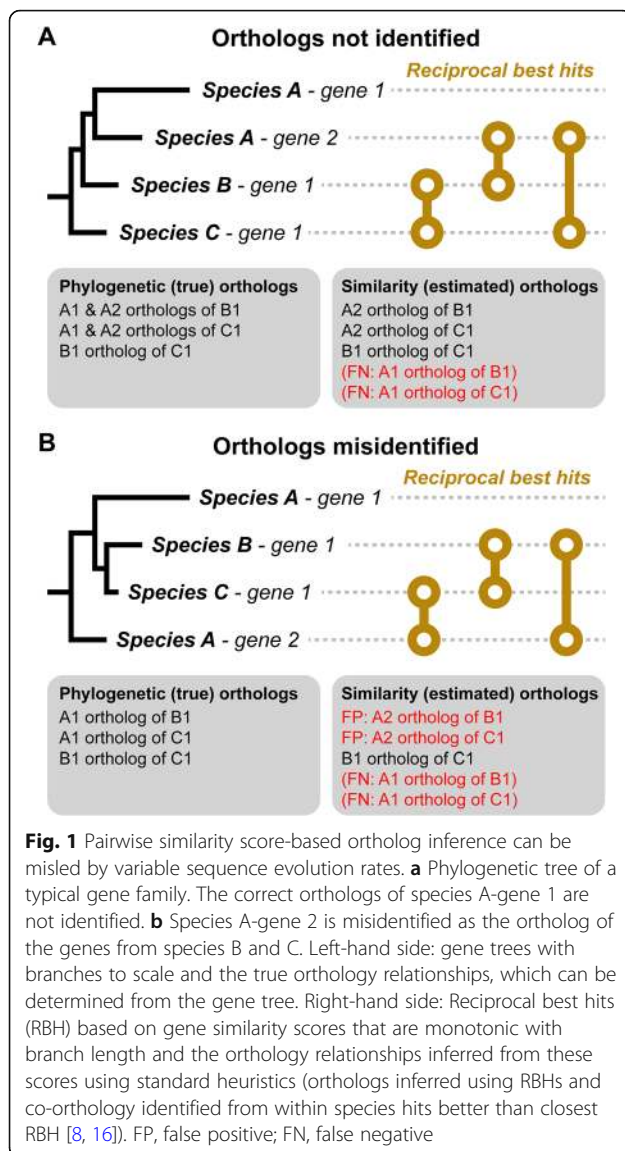
Heuristic analysis of pairwise sequence similarity scores has historically been used to estimate the phylogenetic relationship between genes as it is readily computationally tractable. The central premise underlying their use is that higher scoring sequence pairs are likely to have diverged more recently than lower scoring sequence pairs. Thus, heuristic analysis of sets of pairwise sequence similarity scores can be used to estimate the phylogenetic relationships between sets of genes [7–9, 12, 13]. However, such score-based estimates of the phylogenetic relationship between genes are confounded by multiple factors. For example, variable sequence evolution rates between genes frequently lead to both false-positive and false-negative errors [14, 15] (Fig. 1). Such errors can be mitigated by the analysis of phylogenetic trees of genes [17], as phylogenetic trees are able to distinguish variable sequence evolution rates (branch lengths) from the order in which sequences diverged (tree topology) and hence clarify orthology and paralogy relationships (Fig. 1).

A number of tree-based online databases of orthologs have been developed including PhylomeDB [18], Ensembl-Compara [19], EggNOG [20], and TreeFam [21]. These highly used resources provide the user with the ability to explore the evolutionary history of genes using phylogenetic trees, giving a more complete picture than just pairwise orthology and paralogy relationships alone. Comparative analyses of these methods using standard benchmarking approaches have found no significant difference in ortholog

* Correspondence: steven.kelly@plants.ox.ac.uk

Department of Plant Sciences, University of Oxford, South Parks Road, Oxford OX1 3RB, UK





detection accuracy of these online databases and score-based software tools [1], suggesting that the advantages of a phylogenetic approach have not yet been fully realized. Moreover, the pipelines and methodologies behind these online databases are generally not provided for users to run their own analyses. Thus, there is a need for an automated software tool that effectively exploits the phylogenetic approach to increase accuracy, but with the ease of use, speed, and scalability of a score-based heuristic method.

While an automated software tool for phylogenetic orthology inference from gene sequences is an important goal, the implementation of such a method presents several technical challenges. These comprise the following: (1) inferring a complete set of gene trees for all genes of a given set of species in a time-scale that is competitive with score-based heuristic methods; (2) automatically

rooting these gene trees so that they can be correctly interpreted [22] without requiring the user to know the rooted species tree in advance; and (3) interpreting the gene trees to identify gene duplication events, orthologs, and paralogs while being robust to processes such as gene duplication, loss, incomplete lineage sorting, and gene tree inaccuracies. If these challenges could be addressed in a resource and time-efficient manner, then such a phylogenetic method would provide a step change for orthology inference, enabling the transition from similarity score-based estimates of phylogenetic relationships to phylogenetically delineated phylogenetic relationships between genes.

Some of the challenges listed above have been addressed in isolation by a range of bioinformatic methods. For example, there are a range of methods for identifying orthogroups of genes from user-supplied gene sequences [8–10, 12, 23] and a wide variety of gene tree inference methods that can infer trees from these orthogroups [24–28]. Similarly, there is a range of methods for inferring orthologs from gene trees that also vary in terms of scalability and accuracy [29–32]. However, other critical challenges had no existing solutions. For example, the inference of a complete set of rooted gene trees from a set of species proteomes would be a complex, multi-step process and generally require prior knowledge of the species tree. Equally, methods to infer orthologs from gene trees did not exist that were robust to processes such as incomplete lineage sorting and gene tree inference error while also being scalable to the large-scale analysis required for whole-genome orthology inference across hundreds of species. Thus, substantial technical challenges needed to be addressed to enable fully automated, accurate, and efficient phylogenetic delineation of the phylogenetic relationships between genes.

Here, we present a major update to OrthoFinder that addresses these challenges and significantly extends the scope of the original method. The updated version of OrthoFinder identifies orthogroups as in the original implementation [10] but then uses these orthogroups to infer gene trees for all orthogroups and analyzes these gene trees to identify the rooted species tree. The method subsequently identifies all gene duplication events in the complete set of gene trees and analyzes this information in the context of the species tree to provide both gene tree and species tree-level analysis of gene duplication events. Finally, the method analyzes all of this phylogenetic information to identify the complete set of orthologs between all species and provide extensive comparative genomics statistics. The complete OrthoFinder phylogenetic orthology inference method is accurate, fast, scalable, and customizable and is performed with a single command using only protein sequences as input.

Results

OrthoFinder algorithm overview and summary of results files

The OrthoFinder algorithm is described in detail in the “Methods” section. In brief, it addresses the challenges identified above in five major steps: (a) orthogroup inference, (b) inference of gene trees for each orthogroup, (c and d) analysis of these gene trees to infer the rooted species tree, (e) rooting of the gene trees using the rooted species tree, and (f–h) duplication-loss-coalescence (DLC) analysis of the rooted gene trees to identify orthologs and gene duplication events (mapped to their locations in both the species and gene trees) (Fig. 2). Thus, starting from just gene sequences, OrthoFinder infers orthogroups, orthologs, the complete set of gene trees for all orthogroups, the rooted species tree, and all gene duplication events and computes comparative genomic statistics. To illustrate the standard outputs provided by an OrthoFinder analysis, a graphical example of the complete set of results produced by OrthoFinder for ten metazoan species is shown in Fig. 3a–h.

The default, and fastest, version of OrthoFinder uses DIAMOND [24] for sequence similarity searches. These sequence similarity scores provide both the raw data for orthogroup inference [10] and for gene tree inference of these orthogroups using DendroBLAST [24]. The default implementation of OrthoFinder has been designed to enable a complete analysis with maximum speed and scalability using only gene sequences as input. However, OrthoFinder has also been designed to allow the use of alternative methods for tree inference and sequence search to accommodate user preferences. For example, BLAST [4] can be used for sequence similarity searches in place of DIAMOND. Similarly, gene trees do not need to be inferred using DendroBLAST. Instead, OrthoFinder can automatically infer multiple sequence alignments and phylogenetic trees using most user-preferred multiple sequence alignment and tree inference methods. Moreover, if the species tree is known prior to the analysis, this can also be provided as input, rather than inferred by OrthoFinder. Thus, while OrthoFinder is designed to require minimal inputs and computation, it can be tailored to suit the computational and data resources available to the user.

OrthoFinder has the highest ortholog inference accuracy

The accuracy of key component algorithms of OrthoFinder has been independently assessed in this work and in dedicated publications [5, 10, 22, 24, 33]. To demonstrate the accuracy of the overall method, the orthologs identified by OrthoFinder using its default options, along with multiple different configurations, were submitted to the community-supported *Quest for Orthologs* benchmarking server for the 2011_04 dataset [1] (see the “Methods”

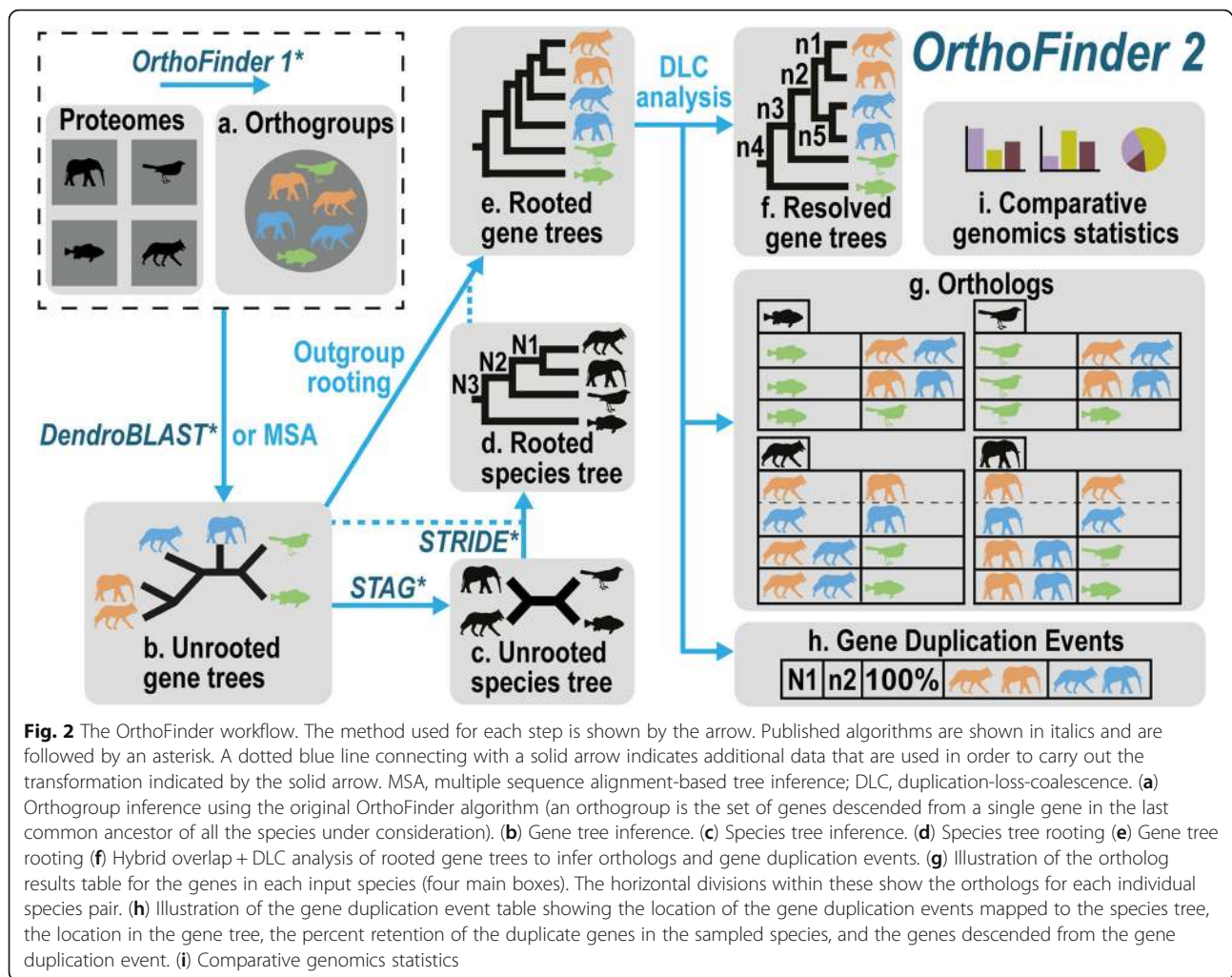
section for details of the tests). This dataset had benchmark results for the largest number of other methods and so allowed the most comprehensive comparison to competitor methods. The results of all of these analyses are shown in Fig. 4a–l and supported by additional analyses in Additional file 1: Figure S1–S3 and Additional file 1: Table S1.

The SwissTree and TreeFam-A tests within *Quest for Orthologs* assess the accuracy of ortholog inference against orthologs from gold-standard trees. For these tests, precision, recall, and *F*-score can be calculated. On these tests, the default, fastest version of OrthoFinder was 3–24% (SwissTree, Fig. 4a) and 2–30% (TreeFam-A, Fig. 4b) more accurate than any other method. The other versions of OrthoFinder were a further 1–3% more accurate than default OrthoFinder. No method was consistently second best to OrthoFinder.

For the *Quest for Orthologs* Standard and Generalized Species Tree Discordance Tests (STDT and GSTDT), no ground truth orthologs are known, and the methods are assessed on the percentage of trials in which a set of orthologs is identified across a set of species and the Robinson-Foulds distance between species tree and the gene tree of the putative orthologs. As such, standard precision, recall, and *F*-score measures cannot be calculated. For these tests, a “pseudo-*F*-score” was calculated using the percentage of the recovered ortholog sets in place of recall and 1 - normalized Robinson-Foulds distance in place of precision (equivalently, the proportion of bipartitions in an agreement between the species tree and the putative orthologs tree). On both STDT and GSTDT, all versions of OrthoFinder had an equal or higher pseudo-*F*-score than all versions of all other methods. The default, fastest version of OrthoFinder was 0–45% (STDT, Fig. 4c) and 10–59% (GSTDT, Fig. 4d) higher scoring than competing methods. The other versions of OrthoFinder were a further 1–6% higher scoring than the default version.

All versions of OrthoFinder, irrespective of algorithmic options, inferred more orthologs (higher recall/recovered ortholog sets) than any other tested method at a similar level of precision (Fig. 4e–l). Across the four tests, the default and fastest version of OrthoFinder (DIAMOND) achieved between 0 (Fig. 4g) and 65% (Fig. 4h) higher recall/recovered ortholog sets than competing methods. It achieved precision/ortholog species tree agreement between 5% lower (Fig. 4h) and 15% higher (Fig. 4g) than competing methods. Similarly, on the latest, 2018, benchmarks, all three versions of OrthoFinder were more accurate than all other methods on all four of the benchmarks: STDT, GSTDT, SwissTree, and TreeFam-A (Additional file 1: Figure S2).

In addition to testing OrthoFinder against competitor methods that can be run on raw sequence data,



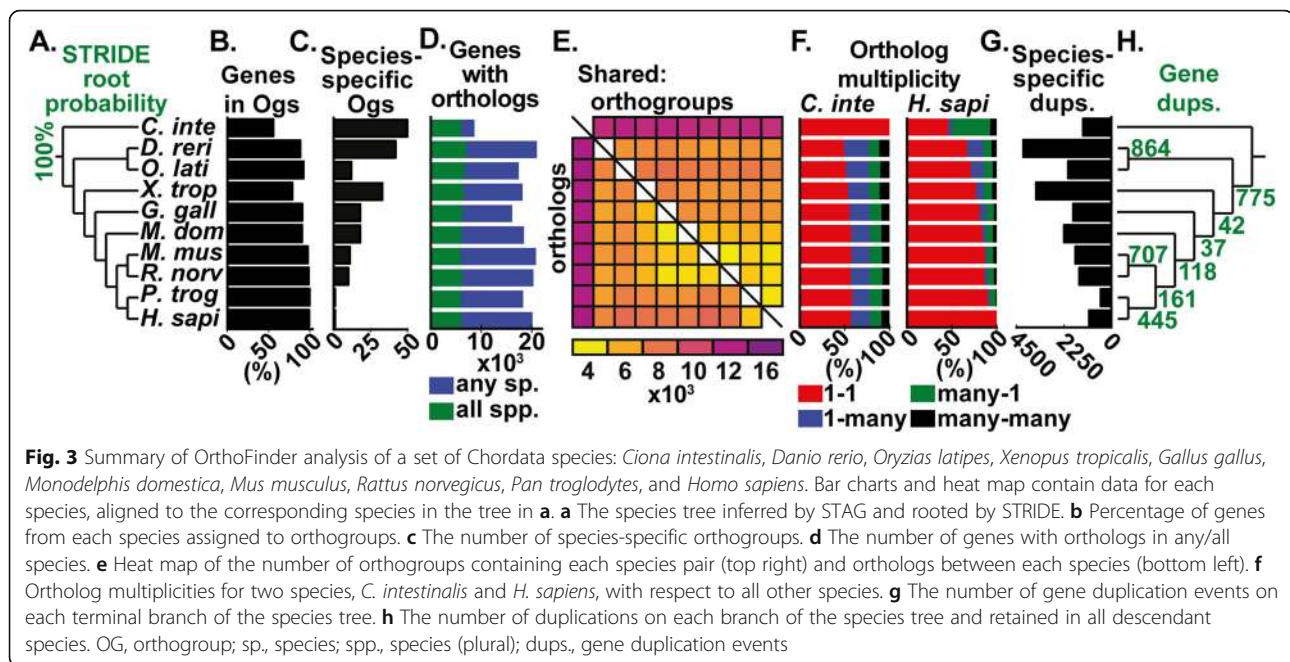
OrthoFinder was also compared with static database methods that involve various levels of human curation. All versions of OrthoFinder, irrespective of algorithmic options, had a higher *F*-score/pseudo-*F*-score across each of the four tests than any of the databases with the one exception of PANTHER on the Species Tree Discordance Test (Additional file 1: Figure S3). Overall, the default version *F*-score/pseudo-*F*-score was between 2 and 14% higher than the database methods. OrthoFinder (BLAST + MSA) scored between 5 and 17% higher than the database methods (Additional file 1: Figure S3). Thus, although OrthoFinder is fully automated and requires no manual curation, it also achieved higher accuracy than curated online database methods.

OrthoFinder is fast and scales well to hundreds of species

To demonstrate the scalability of the OrthoFinder method, it was run on sets of between 4 and 256 fungal species with 16 parallel processes (Fig. 4m). All other publicly available software tools that have been

benchmarked on the *Quest for Orthologs* dataset were similarly tested. The default version of OrthoFinder ran in 192 s on the 4 species and 1.8 days on the 256 species datasets. In this time, it inferred orthogroups, all gene trees, the rooted species tree, orthologs, and gene duplication events (Fig. 4n). Overall, OrthoFinder was the second quickest method, with the fastest method SonicParanoid taking 1.2 days on the same 256 species set. Both OrthoFinder and SonicParanoid scaled well to the largest datasets, both taking less than half the time of the next best method (4.1 days, Fig. 4m).

There was a large range of runtimes across the complete set of methods. Many methods were unsuited to larger species sets, with 64 species being the largest set on which all methods were runnable within the 120 h (5 days) cutoff. At this point of comparison, the slowest method took 200 times longer to run than OrthoFinder. It should also be noted that no competitor method also provides gene trees or identifies gene duplication events (Fig. 4n). Thus, not only is OrthoFinder the most accurate method and the



second fastest method, it also provides the largest quantity of phylogenomic information.

OrthoFinder efficiently and accurately solves the challenge of inferring a rooted species tree from unaligned protein sequence data

Rooted gene trees are required to enable the use of phylogenetic information for ortholog inference, since the correct placement of the root is required for the correct dissection of phylogenetic relationships between genes in the tree [22]. However, the vast majority of tree inference methods infer unrooted trees. Gene trees can be correctly rooted given the knowledge of the underlying rooted species tree, and thus, OrthoFinder first infers and then roots the species tree for the set of species being analyzed. OrthoFinder solves these two challenges (species tree inference and rooting) using two algorithms developed specifically for this purpose.

The species tree is inferred from the set of unrooted orthogroup gene trees using STAG [33], and this species tree is rooted using STRIDE [22]. STAG was developed to leverage the vast amount of phylogenetic information already available in the complete set of orthogroup gene trees inferred by OrthoFinder. It was also developed to be robust to high levels of gene duplication and loss that can hamper methods that rely on sets of single-copy orthologs [33]. It outperformed popular species tree inference methods on benchmark data and scaled well to large datasets [33].

Methods for *ab initio* species tree rooting (i.e., without prior knowledge of a suitable outgroup) have received little attention [22]. STRIDE was similarly

developed to leverage gene duplication events in the complete set of orthogroup gene trees to efficiently determine the root of the species tree and achieved high accuracy on benchmark data [22]. The ability of OrthoFinder to automatically leverage the raw amino acid sequence data to infer the rooted species tree thus enables outgroup rooting of the complete set of orthogroup gene trees for any input set of species and for all gene trees. This is a critical step for enabling phylogenetic orthology inference from gene sequences.

OrthoFinder implements a novel duplication-loss-coalescent algorithm for identifying gene duplication events and orthologs

Given a set of rooted orthogroup gene trees, the final major challenge in accurately dissecting phylogenetic relationships between genes is to account for incomplete lineage sorting and gene tree error. Existing methods for determining if genes within a gene tree are orthologs or paralogs either had poor accuracy or were unable to scale to the number and size of the orthogroup gene trees that must be analyzed. Thus, to address this challenge, a novel, scalable algorithm based on the duplication-loss-coalescent model was developed (see the “Methods” section).

To demonstrate the relative performance characteristics of this method, it was applied to two independent simulated datasets [32, 34] and compared to three popular, comparable methods: GSDI Forester [29], DLCpar (full and search) [32], and species overlap method [31] (Fig. 5). It was also compared to Notung [30], but since branch support values were not available, which Notung

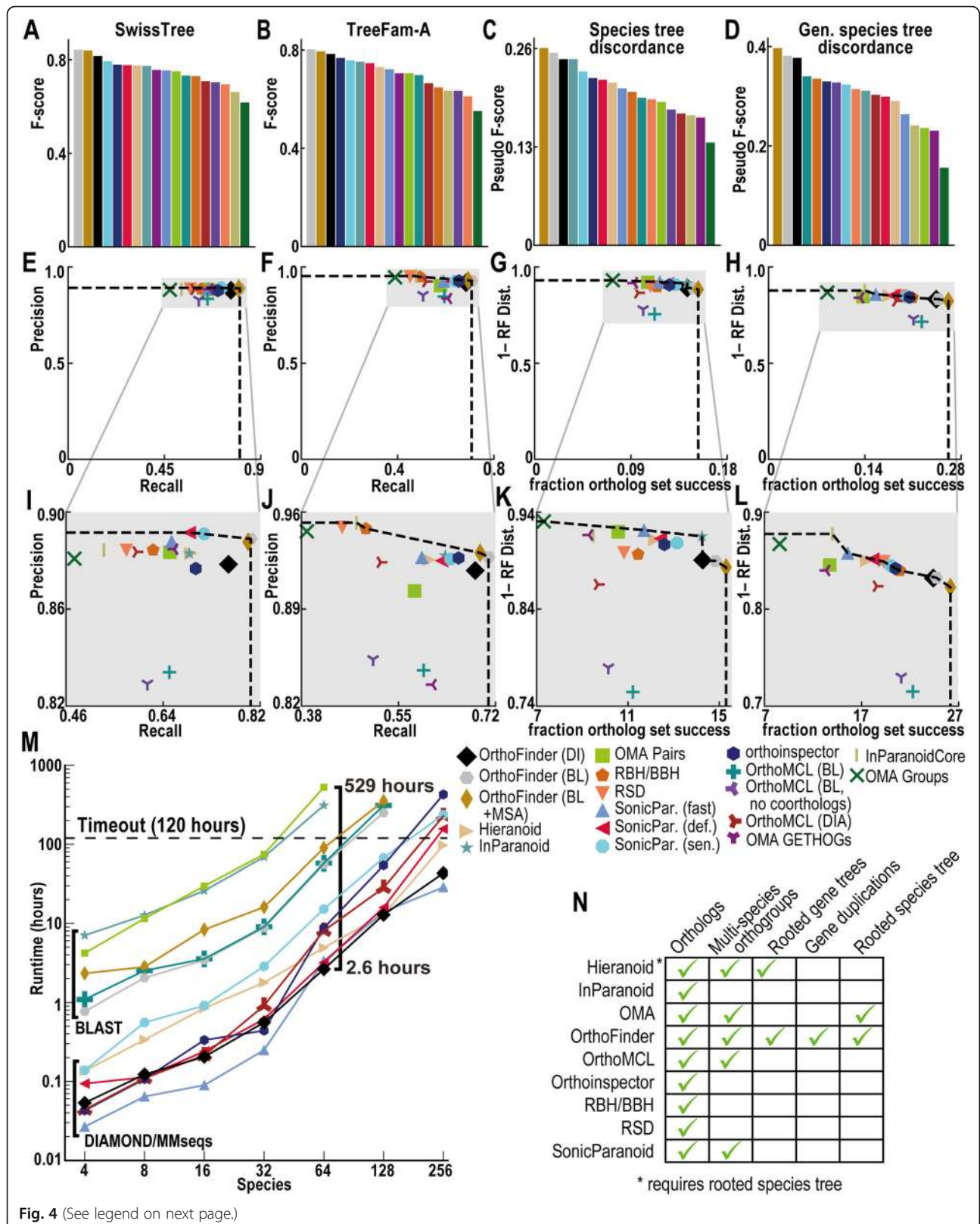


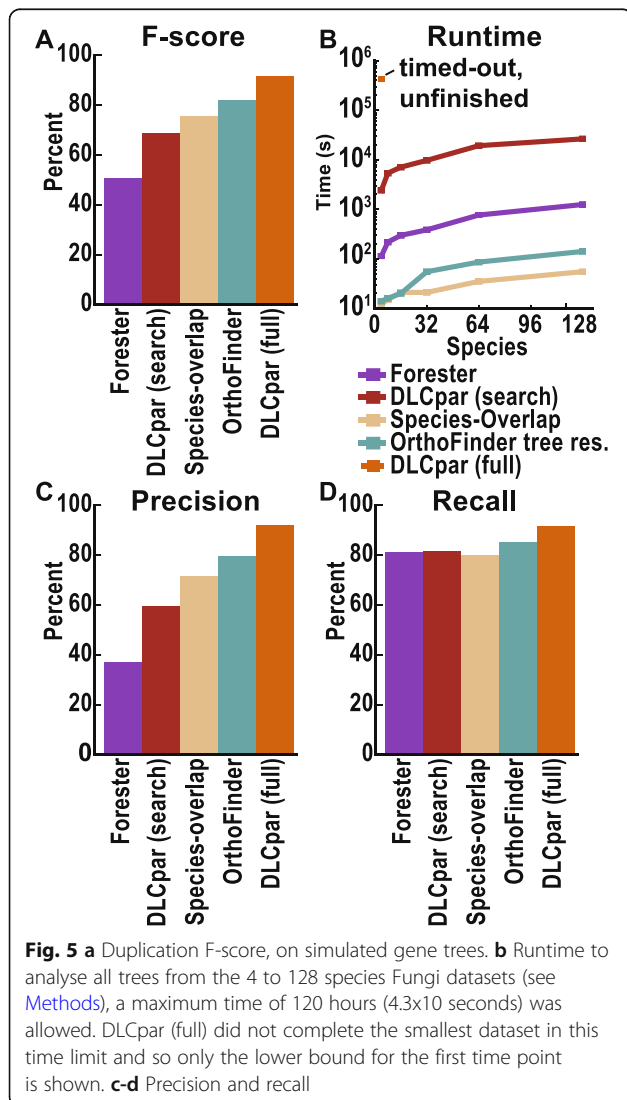
Fig. 4 (See legend on next page.)

(See figure on previous page.)

Fig. 4 a–l Quest for Orthologs 2011_04 benchmarks (see [1]) on 66 species across Eukarya, Bacteria and Archaea for ortholog inference methods. Dotted line shows Pareto frontier. Data for graphs are in Additional file 1: Table S1. **a, b** *F*-score on SwissTree and TreeFam-A tests. **c** “Pseudo-*F*-score” across the two Species Tree Discordance Tests (STDT). **d** “Pseudo-*F*-score” across the four Generalized Species Tree Discordance Tests (GSTDT). **e–f** Agreement of orthologs SwissTree/FreeFam-A trees **g–h** Benchmarks across the STDT & GSTDT. X-axis: Total fraction of randomly selected genes with predicted orthologs in a predefined set of species for the two STDTs & four GSTDTs respectively. Y-axis: Average (1 – normalised Robinson-Foulds distance) between gene tree for putative orthologs and the known species tree across the two STDT & four GSTDT respectively. The four individual GSTDTs and two individual STDTs are shown in Additional file: 1 **i–l** Zoom in of plots **e–h**. See **Methods** section “Ortholog Benchmarking” for details of Quest for Orthologs benchmarks. **m** Runtime for each method with 4-256 input Fungi proteomes. **n** Results returned by methods, a multi-species orthogroup is the set of genes descended from a single gene in the last common ancestor of three or more species

use to improve its accuracy, it achieved identical results to Forester and so is not shown as an additional method here. In terms of accuracy, the novel OrthoFinder method outperformed all methods other than DLCpar (full) (Fig. 5a, Additional file 1: Table S2). However, DLCpar (full) was unable to analyze realistic-sized species datasets. For example, while the OrthoFinder method was able to analyze the complete set of 18,651

orthogroup gene trees (948,449 genes) from 128 fungal species in 141 s, DLCpar (full) was unable to process a considerably smaller, 4-species dataset (2259 trees, 12,958 genes) in 120 h (Fig. 5b). Thus, OrthoFinder is the most accurate method that is scalable to realistic datasets. This algorithm enables accurate interrogation of orthogroup gene trees in a manner that can analyze thousands of gene trees across hundreds of species in minutes on standard computing hardware (Fig. 5b).



Discussion and conclusions

Phylogenetic relationships between gene sequences are defined by their relationship in a gene tree in the context of a species tree. Due to the complexity of conducting phylogenetic orthology inference from raw gene sequences, multiple methods have been developed to bypass phylogeny and approximate phylogenetic relationships from heuristics on pairwise sequence similarity scores. Such approximations are subject to common errors that are avoidable by the analysis of phylogenetic trees of gene sequences. Here, we present a substantial update to OrthoFinder that provides an easy-to-use, fast, accurate, and fully phylogenetic orthology inference software tool.

From testing on community standard benchmarks, we demonstrate that OrthoFinder is the most accurate orthology inference method available. Furthermore, we show that by taking a phylogenetic approach, OrthoFinder provides substantial additional information (including rooted gene trees, rooted species trees, and gene duplication events) that are not provided by heuristic methods. Thus, OrthoFinder is the most accurate and most data-rich orthology inference method for comparative genomics.

The only input required for OrthoFinder is the set of amino acid sequences of the protein-coding genes for the species of interest. OrthoFinder has been designed with ease of use in mind, and the entire analysis is launched with a single command. The default parameters for OrthoFinder are optimized for speed and scalability and enable the combined analysis of hundreds of species on commonly available computer resources. However, OrthoFinder is also designed with the expert user in mind, and intermediate steps in the algorithm can be substituted with other methods for multiple

sequence alignment and tree inference should the user wish. We illustrate the time-accuracy trade-off associated with changes in the internal steps of the algorithm and show that the fastest and least accurate implementation OrthoFinder is still more accurate than any other orthology inference method.

Methods

OrthoFinder workflow

A gene tree is the canonical representation of the evolutionary relationships between the genes in a gene family. Thus, ortholog inference from gene trees is an important goal. However, no automated software tools are available that provide genome-wide ortholog inference from gene trees. A number of challenges had to be addressed to enable this. These included the efficient partitioning of genes into small, non-overlapping sets such that all orthologs of a gene are contained in the same set as the original gene; scalable and accurate inference of gene trees from these gene sets; automatic rooting of these gene trees without a user-provided species tree; and robust ortholog inference in the presence of imperfect gene tree inference. The OrthoFinder workflow was designed to address each of these challenges and is described in detail below.

By default, OrthoFinder infers orthologs from the orthogroup trees (a gene tree for the orthogroup) using the steps shown in Fig. 2. Input proteomes are provided by the user using one FASTA file per species. Each file contains the amino acid sequences for the proteins in that species. Orthogroups are inferred using the original OrthoFinder algorithm [10]; an unrooted gene tree is inferred for each orthogroup using DendroBLAST [24]; the unrooted species tree is inferred from this set of unrooted orthogroup trees using the STAG algorithm [33]; this STAG species tree is then rooted using the STRIDE algorithm by identifying high-confidence gene duplication events in the complete set of unrooted orthogroup trees [22]; the rooted species tree is used to root the orthogroup trees; orthologs and gene duplication events are inferred from the rooted orthogroup trees by a novel hybrid algorithm that combines the “species-overlap” method [31] and the duplication-loss-coalescent model [32] (described below); and comparative statistics are calculated. All major steps of the algorithm are parallelized to allow optimal use of computational resources. Only the orthogroup inference was provided in the original implementation of OrthoFinder [10]; all other subsequent steps are new and described below.

Use of orthogroups for gene tree inference

Orthologs are the set of genes in a species pair descended from a single gene in the last common ancestor of those two species. An orthogroup is the set of genes

from multiple species descended from a single gene in the last common ancestor (LCA) of that set of species. Thus, an orthogroup is the natural extension of orthology to multiple species.

For ortholog inference, orthogroups are the optimum partitioning of genes for gene tree inference: An orthogroup is the smallest set of genes such that, for all genes it contains, the orthologs of these genes are also in the same set. Since gene tree inference scales super-linearly with the number of genes, partitioning genes into the smallest possible sets is the most efficient way of constructing a set of gene trees that encompass all orthology relationships. Although partitioning genes into larger sets (e.g., gene families containing gene duplication events prior to the LCA) would decrease the number of gene trees to be inferred, the super-linear scaling of gene tree inference would result in a longer overall runtime for the complete set of trees. The original OrthoFinder orthogroup inference method is still the most accurate method on the independent Orthobench test set [10] and thus is used for this step.

Customizable steps in the OrthoFinder method

There are two customizable steps in the OrthoFinder method: (1) the sequence search method and (2) the orthogroup tree inference method. The default option for step 1 is DIAMOND [5]. The default option for step 2 is DendroBLAST [24]. The default options are recommended by the authors as they are fast and achieve high accuracy on the Quest for Orthologs benchmarks [1] (Fig. 4a–d). However, the user is free to substitute any alternative methods for these steps. Currently, supported methods for step 1 include BLAST [4] and MMseqs2 [6]. Similarly, any combination of multiple sequence alignment and tree inference method can be substituted in for step 2. For illustrative purposes, the default multiple sequence alignment method is MAFFT [35] and the default tree inference method is FastTree [25]; this combination is benchmarked above. It is impossible for the authors to test all possible combinations of multiple sequence alignment and tree inference methods, and the selected methods were chosen because of their speed and scalability characteristics [25, 35]. OrthoFinder provides flexibility for the user to select their preferred method. More accurate multiple sequence alignment and tree inference methods should give more accurate ortholog inference, and many studies exist comparing the accuracy and runtime characteristics of the available methods [36, 37]. A user-editable configuration file is provided in JSON format that allows new sequence search, multiple sequence alignment, and tree inference methods to be added to OrthoFinder. To facilitate the trialing of alternative multiple sequence alignment and tree inference methods, OrthoFinder provides the option to restart an existing

analysis after the orthogroup inference stage. This skips the requirement to compute the all-vs-all sequence search and orthogroup inference and thus accelerates testing of different internal steps.

Species tree inference and rooting

The rooted species tree is required in order to identify the correct out-group in each orthogroup tree, as correct gene tree rooting is critical for the orthology assessment from that tree [22]. Since orthogroups can potentially contain any subset of the species in the analysis, it is not sufficient to simply know the out-group for the complete species set. Instead, the complete rooted species tree is required. If the user knows the rooted species tree for the set of species being analyzed, then it is recommended to specify this tree manually at the command line to remove the possibility of species tree inference error. Such a tree can be provided as a Newick format text file. In the event that a species tree is not provided (or not known), then OrthoFinder automatically infers it.

Sets of one-to-one orthologs that are present in all species are often used for species tree inference; however, in real-world large-scale analyses, these can be rare [33]. A new algorithm, Species Tree from All Genes (STAG), was developed to allow species tree inference even for species sets with few or no complete sets of one-to-one orthologs present in all species [33]. Without this algorithm, species tree inference could fail if there were no sets of one-to-one orthologs present in all species. STAG infers the species tree using the most closely related genes within single-copy or multi-copy orthogroups. In benchmark tests, STAG [24] had higher accuracy than other leading methods for species tree inference, including maximum likelihood species tree inference from concatenated alignments of protein sequences, ASTRAL [38] and NJst [39].

The Species Tree Root Inference from Duplication Events (STRIDE) algorithm [22] is used to root the species tree in OrthoFinder. STRIDE was developed to enable the rooting of the species tree using only information available in the set of gene trees. STRIDE does this by identifying the set of well-supported in-group gene duplication events in the complete set of unrooted orthogroup trees, and using these events to infer a probability distribution over an unrooted STAG species tree for the location of its root. Similarly to STAG, STRIDE has been shown to identify the correct root of the species tree in multiple large-scale molecular phylogenetic data sets spanning a wide range of time scales and taxonomic groups [22]. In some cases, it is possible that there could be few duplications within the gene trees, and so STRIDE will not be able to identify the root of the species tree, or will only be able to exclude the root from clades in which gene duplication events are

observed. In this case, ortholog inference should still not be significantly impacted since the rooting of the gene tree only affects ortholog inference in cases where gene duplication events are present [22]. This makes the STRIDE approach particularly suited to gene tree rooting for ortholog inference.

Gene tree rooting

Tree inference methods infer unrooted gene trees. A gene tree must be correctly rooted in order for it to show the correct evolutionary history of the gene family and thus to allow correct ortholog inference. The orthogroup trees could contain any subset of the input species. In general, the rooted species tree, inferred as described above, can be used to root the orthogroup trees by identifying the out-group clade in each orthogroup tree and placing the root on the branch separating this out-group from the remaining genes.

However, species tree and gene tree topologies can arise in which this simple approach will not work, and so, a robust generalization of this outgroup rooting method is required in order to be able to root any potential gene tree. Firstly, in the species tree, the out-group could consist of a single species or multiple species. Secondly, in the gene tree, the genes from the out-group could be in a monophyletic clade or there may be no bipartition in the tree that separates all the genes from the out-group from all remaining genes. Thirdly, a gene duplication event could have occurred in the gene tree prior to the divergence of the out-group from the remaining species. Thus, the most ancient bipartition of the gene tree would be a gene duplication event separating the genes into two clades rather than a bipartition separating the out-group from the in-group. Such a gene tree should be rooted on this bipartition. Both of these two descendant clades could then potentially contain genes from both the out-group and in-group species. Thus, there will be no bipartition in such a tree that separates the genes of the out-group species from the genes of the in-group species.

The algorithm used by OrthoFinder searches for the correct bipartition on which to place the root. For each bipartition in the gene tree, it calculates two scores. The first, S_{AD} , quantifies how well the bipartition corresponds to an ancient duplication prior to the divergence of the species. The second, S_{IO} , quantifies how well the bipartition corresponds to the divergence of the out-group species from the in-group species. Both S_{IO} and S_{AD} range between 0 and 1. Let O be the set of species in the out-group and I be the set of species in the in-group. For a bipartition in the unrooted gene tree, let A be the set of species with genes on one side of the bipartition and let B be the set of species with genes on the other side of the bipartition. Then:

$$S_{IO} = \frac{|OnA| |InB|}{|O| |I|} \left(1 - \frac{|OnB|}{|O|}\right) \left(1 - \frac{|InA|}{|I|}\right),$$

$$S_{AD} = \frac{|OnA| |InB| |OnB| |InA|}{|O| |I| |O| |I|}.$$

Each of the four terms in these equations quantifies the proportion of in-/out-group species the bipartition correctly includes/excludes from clade A/B of the gene tree (giving the $2^3 = 8$ terms in total across the two equations). The bipartition with the highest score for either S_{IO} or S_{AD} is the optimal root for the gene tree using this measure.

The effectiveness of these scores at identifying the correct root can be seen by considering the following. A bipartition with a value of 1 for S_{IO} implies that it perfectly divides the tree into an in-group and out-group and implies a value of 0 for S_{AD} for all bipartitions in the tree (thus, there are no potential bipartitions corresponding to an ancient duplication). This is the correct bipartition on which to root the tree since it separates the in-group from the out-group genes. Conversely, a bipartition with a value of 1 for S_{AD} implies that the bipartition is a duplication event before the divergence of any of the species, with all species present for both duplicates. It implies a value of 0 for S_{IO} for all bipartitions in the tree (thus, there is no bipartition that corresponds to a first speciation event that splits the genes into an out-group clade and an in-group clade). The highest value for either S_{IO} or S_{AD} across the tree shows that the corresponding bipartition is close to one of these perfect cases and is the best root for the gene tree.

Ortholog inference and identification of gene duplication events from gene trees

A number of methods were considered for distinguishing orthologs from paralogs in gene trees. Duplication and loss reconciliation, e.g., Forester, uses a rooted species tree and rooted gene tree to determine if each node in the gene tree is a speciation or a duplication event. Genes that diverged at a speciation event are orthologs whereas those that diverged at a duplication event are paralogs. DLCpar [32] uses a model for duplication-loss-(deep) coalescent (DLC) that addresses incongruence between the gene and species trees to increase accuracy. It exists in two versions which we label DLCpar (full) and DLCpar (search). DLCpar (full) considers the complete space of possible reconciliations to find the maximum parsimony solution under the DLC model but can have large runtimes even for relatively small gene trees. DLCpar (search) instead employs an iterative search for a locally optimal solution, which can differ from the globally optimal solution. A third approach, here referred to as the species-overlap method, is employed in a

number of ortholog databases [20, 31] and was originally described in a method for determining orthologs of human genes [31]. In this method, nodes in the gene tree are identified as duplication nodes if the sets of species below its child nodes overlap; otherwise, the node is a speciation node. Genes that diverged at a speciation node are orthologs, and those that diverged at a duplication node are paralogs.

These methods were tested on the fungal orthogroups (in parallel, using 16 cores) to determine their runtime on sets of typical orthogroup trees derived from sets of between 4 and 128 species. Our implementation of the species-overlap method was the fastest, taking 55 s to analyze the largest dataset (Fig. 5). This dataset consisted of the 18,651 orthogroup trees containing 948,449 genes and corresponded to the complete set of orthogroup trees for the 128 fungal species. Forester was 21 times slower, and DLCpar (search) was over 500 times slower. DLCpar (full) was unable to complete the analysis of the smallest input dataset in 120 h and so was not tested on any of the larger datasets. To put this time in context, all steps in the OrthoFinder algorithm for this dataset collectively take less than 4 min in total (i.e., orthogroup inference, gene tree inference, species tree inference, species tree rooting, gene tree rooting).

To compare the accuracy of the above methods, they were each tested for their precision and recall in identifying gene duplication events on simulated “flies” and “primates” datasets [32] and a simulated “metazoa” dataset [34]. Since for all methods tested a node in a gene tree is either a duplication or speciation event, the identification of all gene duplication events is equivalent (by complementation) to the identification of all speciation events. Thus, the overall accuracy at identifying gene duplication events is equivalent to the overall accuracy at identifying orthologs. The most accurate method on the simulated data was DLCpar (full) with an F -score of 91.8% followed by the species-overlap method with an F -score of 75.5%.

Since DLCpar (full) was the most accurate method on the simulated datasets but was unsuitable for analyzing gene trees with more than four species a novel hybrid algorithm was developed. This aimed to combine the strengths of the highest accuracy DLCpar (full) method with simplifications from the species-overlap method to achieve high accuracy in a reasonable runtime.

In the DLC model, clades of genes containing no duplicates are analyzed to find the most parsimonious reconciliation with the species tree. This is required since the goal for DLCpar is a complete reconciliation of the gene tree with the species tree. However, in the species-overlap method, clades of single-copy genes are identified as orthologs without further analysis of the topology of their relationship. This assumption is reasonable,

since trees of single-copy orthologs are frequently topologically distinct from the species tree. For example, in an analysis of 1030 gene trees of one-to-one orthologs from 23 fungi species, all 1030 gene trees were topologically distinct from each other and from the species tree [40]. The analysis of such clades under the DLC model is likely to be computationally costly with no benefit in terms of accuracy of ortholog inference.

On the other hand, when a gene duplication event has occurred, it is important to accurately identify the genes affected by this event since the location of the event determines which genes are orthologs and which are paralogs. In the hybrid algorithm developed for OrthoFinder, these nodes, for which there is evidence of a gene duplication event through overlapping species sets, are analyzed under the DLC model. The DLC model is used to attempt to find the most parsimonious interpretation of this node in terms of which genes diverged at the gene duplication event and which diverged at a speciation event.

As described, this method would still require exploring a large search space for the nodes under consideration, and the reduction in runtime would not be significant. Thus, to accelerate the process, duplication and loss events are inferred directly using the species-overlap method. A duplication event is inferred from an overlap in the species sets below a node and a loss event is inferred by the presence of a gene from a species in one of the descendant clades but not in the other. The analysis can then be accelerated by classifying a node according to the species overlaps of its subclades up to a maximum total topological depth of two below the node being analyzed (clades O, Additional file 1: Figure S4A). The possible sub-cases for the overlaps between these clades have been enumerated (Additional file 1: Figure S4B). For each sub-case, the most parsimonious interpretation under the DLC model has been pre-calculated (Additional file 1: Figure S4C) and can thus be corrected without the need for a topology search.

The algorithm implemented in OrthoFinder is as follows. A post-order traversal of the orthogroup tree is performed (a node is not visited until all its descendant nodes have been visited), analyzing each node of the orthogroup tree in turn. A given node is analyzed to identify if the species sets below its child nodes overlap. If there is an overlap, the smallest sub-clade below each child node that contains the complete set of overlapping species is identified up to a maximum total topological depth of two below the node (clades O, Additional file 1: Figure S4A). The node is assigned to the corresponding sub-case (Additional file 1: Figure S4B). If a more parsimonious interpretation of the sub-case is available under the DLC model, then the sub-tree below the node is rearranged to match this interpretation (Additional file 1:

Figure S4C). After the node has been analyzed, the next node in the post-order traversal is analyzed. Note, the choice of a post-order traversal allows the traversal to be continued unimpeded despite any such rearrangements below the node being analyzed. The resulting gene trees are referred to as “resolved” gene trees and correspond to the “locus tree” under the DLCpar model [32]. Orthologs and gene duplication events are determined from the resolved gene tree according to the species overlap method.

Although only a single traversal of the tree is employed, rather than the iterative search and rearrangement employed by DLCpar, the post-order traversal enables more parsimonious interpretations of child clades below a node to be identified prior to the analysis of the parent node. Thus, the analysis of sub-trees below a node informs the subsequent analysis of the node itself. In theory, nodes could be categorized to sub-cases based on the overlaps of clades at a greater topological depth than that employed here. This conservative approach was taken since the number of subcases increases exponentially, and a total topological depth of two proved sufficient to achieve a higher accuracy for the method compared to the simple species overlap. The analysis of clades to this depth proved sufficient to increase the *F*-score from 72% with just the species-overlap method to 80% with the hybrid algorithm (Fig. 5a). The pre-calculated solutions for each sub-case removed the need for costly, iterative search using random (i.e., unguided) tree rearrangement operations thus accelerating the analysis considerably. The hybrid algorithm was able to analyze the complete set of orthogroup trees for the 128 fungi species in 141 s; this was 9 times faster than Forster and 187 times faster than DLCpar (search) (Fig. 5d). The hybrid method also outperformed both methods in terms of accuracy (Fig. 5a). Note that the species tree is not required for the hybrid model used by OrthoFinder. The only use of the species tree is in determining the root for each orthogroup tree. All gene tree processing is performed using the python ETE toolkit [41].

Simulation tests of OrthoFinder gene duplication event inference accuracy

The tests for gene duplication event inference accuracy were performed on the simulated “flies” and “primates” dataset from [32] and a simulated “metazoa” dataset from [34]. To model real data, the flies and primate datasets used known species trees, parameters for divergence times, duplication rates, loss rates, population sizes, and generation times. Trees were simulated with varying effective population sizes and duplication rates so as to model incomplete lineage sorting [32, 34]. The flies dataset consisted of 12,000 trees with 12 species and 12,032 gene duplication events. The primates dataset

consisted of 7500 trees with 17 species and 16,066 gene duplication events. The metazoa dataset intended to emulate the complexity of real data by using heterogeneity in rates of duplication and loss, a complex model of sequence evolution, and then inferring trees with a homogenous, simple model [34]. It consisted of 2000 gene trees with 40 species and 4967 gene duplication events. For comparison, Forester [29], DLCpar (full), DLCpar (search) [32], and the overlap algorithm (i.e., without OrthoFinder's tree resolution) were also tested.

All methods were provided with the input rooted gene tree and, where appropriate, the rooted species tree (Forester and DLCpar). No other parameters required specification for any of the other methods. The rooted gene trees were provided as part of the simulated data for the flies and primates datasets. Multiple sequence alignment (MSA) files were provided for the metazoa dataset. For this dataset, gene trees were inferred from the MSAs using FastTree so as to also include a potential level of tree inference error and were rooted with reconroot [32]. The OrthoFinder rooting algorithm was not used so as to avoid inadvertently biasing the results in favor of OrthoFinder. All methods were provided with the same input rooted gene trees. The complete set of gene duplication events identified by each of the methods was compared against the ground truth gene duplication events. An inferred gene duplication was identified as correct if the two sets of genes observed post-duplication exactly matched the two sets of genes post-duplication from the ground truth data.

The performance testing of the methods for identifying gene duplication events was performed on the orthogroup trees from the 4- to 128-species Fungi datasets as inferred by OrthoFinder with default parameters. The commands for Forester and DLCpar were run in parallel using GNU Parallel [42] using 16 threads on these gene trees. The OrthoFinder method was run via the "scripts/resolve.py" program included as part of the OrthoFinder distribution. To allow testing, the species-overlap method was also implemented in OrthoFinder and was run using the same program with the option "--no_resolve."

Ortholog benchmarking

Orthogroup inference accuracy of OrthoFinder has already been tested using the independent Orthobench dataset [11]. This showed to be the most accurate method tested in terms of overall *F*-score (although other methods scored higher in terms of either precision or recall while scoring proportionally worse in the other) [10]. The community developed "Quest for Orthologs" benchmarks [1] were used to assess the accuracy of the newly developed OrthoFinder ortholog inference using the 2011_04 dataset. This dataset had benchmarks for the largest set of methods and so provided the widest comparison with other methods.

OrthoFinder was tested using the default method (DIAMOND sequence search and DendroBLAST trees, no additional options). It was also tested with the BLAST replacing DIAMOND (options: "-S blast") and with both BLAST search and multiple sequence alignment and maximum likelihood tree inference (options: "-S blast -M msa"). In the latter, MAFFT [35] and FastTree [25] were used for multiple sequence alignment and tree inference as described above. For each of these three cases, OrthoFinder was run on the 66 reference proteomes of the Quest for Orthologs test set with a single command ("-f Proteomes/" + options), and the inferred orthologs were submitted to the Quest for Orthologs web server for benchmarking.

The Quest for Orthologs benchmarks are described in detail in [1]. The Species Tree Discordance Test and the generalized version of this test both consider a set of species partitioned into clades with a known species tree topology connecting the clades. The benchmarking consists of a repeated test. For one of the clades of species, a gene is selected at random for each instance of the test. If the orthology inference method under scrutiny predicts an ortholog for that gene for at least one species from each of the remaining clades, then the test is recorded as a "successful ortholog set." For each successful ortholog set, an MSA is constructed and a gene tree inferred using RAxML [28]. The normalized Robinson-Foulds (RF) distance is calculated between this tree and the known species tree. The result of the benchmark is the fraction of successful ortholog sets and the average RF distance for these successful sets. A higher fraction of success and a lower average RF distance indicates a better ortholog inference method under this test. The benchmarks include two different Species Tree Discordance Tests (STDT) across two different species sets and four Generalized Species Tree Discordance Tests (GSTDT) across four different species sets. In Fig. 4g, h, the total fraction of successful ortholog sets and the average normalized RF distance across these successful ortholog sets across the two/four species sets are reported for the STDT and GSTDT. The individual GSTDT and STDT results for the four individual species sets are given in Additional file 1: Figure S1.

Minor changes have been made to the labeling and orientation of the axes compared to the presentation in the Quest for Orthologs paper [1] to improve the consistency with the SwissTree and TreeFam-A benchmarks. The altered *y*-axis for the GSTDT and STDT presented here is (1 - normalized RF distance) so that higher *y* values always correspond to the better agreement with the species tree for all benchmark figures. The number of completed ortholog set successes for the STDT and GSTDT is reported as a fraction rather than the total number. For the SwissTree and TreeFam-A tests, the axes are labeled as "precision" instead of "pos. predictive value rate" and

“recall” instead of “true positive rate” as this is more standard terminology for the quantities reported by the tests.

The full set of benchmarks, the input files, and the ortholog inference results can be seen online at <http://orthology.benchmarkservice.org/>. A comprehensive summary of the benchmarks, as described above, is shown in Fig. 4a–l for ortholog prediction software tools. The corresponding comparisons against online databases are shown in Additional file 1: Figure S2 and Additional file 1: S3. The complete datasets are available to download from Zenodo research archive at <https://doi.org/10.5281/zenodo.1481147> [43].

Performance testing

We constructed sets of fungal proteomes of increasing size for performance testing. Ensembl Genomes was interrogated on 6 November 2017 using its REST API [44] to identify all available fungal genomes. To achieve an even sampling of species, we selected 1 species per genera and excluded genomes from candidate phyla or phyla with fewer than 3 sequenced genomes. This gave a set of 272 species which were downloaded from the Ensembl FTP site [45]. We created datasets of increasing size by randomly selecting 4, 8, 16, 32, 64, 128, and 256 species such that the last common ancestor was the same for each dataset. Each dataset was analyzed using a single Intel E5-2640v3 Haswell node (16 cores) on the Oxford University ARCUS-B server using 16 parallel threads for OrthoFinder with DIAMOND (arguments: “-S diamond -t 16 -a 16”). The complete datasets for all analyzed species subsets are available for download from Zenodo at <https://doi.org/10.5281/zenodo.1481147>. All methods submitted to Quest for Orthologs that provided a user-runnable implementation of the method were tested on the same fungi datasets and the same ARCUS-B server nodes and run in parallel using 16 threads (when supported by the method).

Chordata dataset

The data for the OrthoFinder analysis of the ten Chordata species for the illustration of the results of an OrthoFinder analysis (Fig. 2a–h) are provided in the Zenodo archive <https://doi.org/10.5281/zenodo.1481147>. This includes the input proteomes, the OrthoFinder results, and the script used to generate the figures from the results. OrthoFinder was run with default settings (DIAMOND sequence search and DendroBLAST gene trees).

Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s13059-019-1832-y>.

Additional file 1. Supplementary figures and tables.

Additional file 2. Review history.

Acknowledgements

The authors would like to acknowledge the use of the University of Oxford Advanced Research Computing (ARC) facility in carrying out this work.

Review history

The review history is available as Additional file 2.

Additional information

Peer review information: Barbara Cheifet and Tim Sands were the primary editors on this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

Authors' contributions

DE and SK conceived and designed the project. DE developed the algorithms. DE and SK discussed the results and wrote the manuscript. Both authors read and approved the final manuscript.

Funding

This work was supported by the European Union's Horizon 2020 research and innovation program under grant agreement number 637765. SK is a Royal Society University Research Fellow.

Availability of data and materials

The OrthoFinder source code and executables are available at <https://github.com/davidemms/OrthoFinder> [46] and are released under the GNU General Public License, GPL-3.0 license. The 256 Fungi proteomes and the Chordata proteomes datasets were prepared for this study and are available in the Zenodo research data archive at <https://doi.org/10.5281/zenodo.1481147> [42]. The reference proteomes for the Quest for Orthologs benchmarks are available from ftp://ftp.ebi.ac.uk/pub/databases/reference_proteomes/previous_releases/. The orthologs inferred by all methods and all benchmark results are available from the Quest for Orthologs Benchmark Server [1]: <https://orthology.benchmarkservice.org>. The metazoa dataset [34] and flies and primates [32] datasets used for the genome duplication accuracy analysis are available from the respective authors.

Ethics approval and consent to participate

Not applicable

Consent for publication

Not applicable

Competing interests

The authors declare that they have no competing interests.

Received: 24 April 2019 Accepted: 23 September 2019

Published online: 14 November 2019

References

- Altenhoff AM, Boeckmann B, Capella-Gutierrez S, Dalquen DA, DeLuca T, Forslund K, Huerta-Cepas J, Linard B, Pereira C, Przytycki LP, et al. Standardized benchmarking in the quest for orthologs. *Nature Methods*. 2016;13:425.
- Nichio BTL, Marchaukoski JN, Raittz RT. New tools in orthology analysis: a brief review of promising perspectives. *Front Genet*. 2017;8:165.
- Kristensen DM, Wolf YI, Mushegian AR, Koonin EV. Computational methods for gene orthology inference. *Brief Bioinform*. 2011;12:379–91.
- Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. BLAST+: architecture and applications. *BMC Bioinformatics*. 2009;10:421.
- Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat Methods*. 2015;12:59–60.
- Steinberger M, Soding J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol*. 2017;35:1026–8.
- Ostlund G, Schmitt T, Forslund K, Kostler T, Messina DN, Roopra S, Frings O, Sonnhammer ELL. InParanoid 7: new algorithms and tools for eukaryotic orthology analysis. *Nucleic Acids Res*. 2010;38:D196–203.
- Li L, Stoeckert CJ, Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic genomes. *Genome Res*. 2003;13:2178–89.

9. Altenhoff AM, Schneider A, Gonnet GH, Dessimoz C. OMA 2011: orthology inference among 1000 complete genomes. *Nucleic Acids Res.* 2011;39:D289–94.
10. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biol.* 2015;16:157.
11. Trachana K, Larsson TA, Powell S, Chen WH, Doerks T, Muller J, Bork P. Orthology prediction methods: a quality assessment using curated protein families. *Bioessays.* 2011;33:769–80.
12. Cosentino S, Iwasaki W. SonicParanoid: fast, accurate and easy orthology inference. *Bioinformatics.* 2019;35:149–51.
13. Linard B, Thompson JD, Poch O, Lecompte O. OrthoInspector: comprehensive orthology analysis and visual exploration. *Bmc Bioinformatics.* 2011;12:11.
14. Lafond M, Miardan MM, Sankoff D. Accurate prediction of orthologs in the presence of divergence after duplication. *Bioinformatics.* 2018;34:366–75.
15. Fitch WM. Distinguishing homologous from analogous proteins. *Sys Zool.* 1970;19:99.
16. Remm M, Storm CEV, Sonnhammer ELL. Automatic clustering of orthologs and in-paralogs from pairwise species comparisons. *J Mol Biol.* 2001;314:1041–52.
17. Dalquen DA, Dessimoz C. Bidirectional best hits miss many orthologs in duplication-rich clades such as plants and animals. *Genome Biol Evol.* 2013; 5:1800–6.
18. Huerta-Cepas J, Capella-Gutierrez S, Pryszcz LP, Marcet-Houben M, Gabaldon T. PhylomeDB v4: zooming into the plurality of evolutionary histories of a genome. *Nucleic Acids Res.* 2014;42:D897–902.
19. Herrero J, Muffato M, Beal K, Fitzgerald S, Gordon L, Pignatelli M, Vilella AJ, Searle SMJ, Amode R, Brent S, et al. Ensembl comparative genomics resources. *Database.* 2016;2016:baw053. <https://academic.oup.com/database/article/doi/10.1093/database/baw053/2630361>.
20. Powell S, Forslund K, Szklarczyk D, Trachana K, Roth A, Huerta-Cepas J, Gabaldon T, Rattei T, Creevey C, Kuhn M, et al. eggNOG v4.0: nested orthology inference across 3686 organisms. *Nucleic Acids Res.* 2014;42: D231–9.
21. Li H, Coghlan A, Ruan J, Coin LJ, Heriche JK, Osmotherly L, Li RQ, Liu T, Zhang Z, Bolund L, et al. TreeFam: a curated database of phylogenetic trees of animal gene families. *Nucleic Acids Res.* 2006;34:D572–80.
22. Emms DM, Kelly S. STRIDE: species tree root inference from gene duplication events. *Mol Biol Evol.* 2017;34(12):3267–78.
23. Schreiber F, Sonnhammer ELL. Hieranoid: hierarchical orthology inference. *J Mol Biol.* 2013;425:2072–81.
24. Kelly S, Maini PK. DendroBLAST: approximate phylogenetic trees in the absence of multiple sequence alignments. *PLoS One.* 2013;8:e58537.
25. Price MN, Dehal PS, Arkin AP. FastTree 2-approximately maximum-likelihood trees for large alignments. *PLoS One.* 2010;5:e9490.
26. Lefort V, Desper R, Gascuel O. FastME 2.0: a comprehensive, accurate, and fast distance-based phylogeny inference program. *Mol Biol Evol.* 2015;32: 2798–800.
27. Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol.* 2015;32:268–74.
28. Stamatakis A. RAXML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics.* 2014;30:1312–3.
29. Zmasek CM, Eddy SR. A simple algorithm to infer gene duplication and speciation events on a gene tree. *Bioinformatics.* 2001;17:821–8.
30. Chen K, Durand D, Farach-Colton M. NOTUNG: a program for dating gene duplications and optimizing gene family trees. *J Comput Biol.* 2000;7:429–47.
31. Huerta-Cepas J, Dopazo H, Dopazo J, Gabaldon T. The human phylome. *Genome Biol.* 2007;8(6):R109.
32. Wu YC, Rasmussen MD, Bansal MS, Kellis M. Most parsimonious reconciliation in the presence of gene duplication, loss, and deep coalescence using labeled coalescent trees. *Genome Res.* 2014;24:475–86.
33. Emms D, Kelly S. STAG: species tree inference from all genes. *bioRxiv.* 2018. <https://www.biorxiv.org/content/10.1101/267914v1>.
34. Boussau B, Szollosi GJ, Duret L, Gouy M, Tannier E, Daubin V. Genome-scale coestimation of species and gene trees. *Genome Res.* 2013;23:323–30.
35. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol.* 2013;30:772–80.
36. Thompson JD, Linard B, Lecompte O, Poch O. A comprehensive benchmark study of multiple sequence alignment methods: current challenges and future perspectives. *PLoS One.* 2011;6:e18093.
37. Zhou XF, Shen XX, Hittinger CT, Rokas A. Evaluating fast maximum likelihood-based phylogenetic programs using empirical phylogenomic data sets. *Mol Biol Evol.* 2018;35:486–503.
38. Mirarab S, Warnow T. ASTRAL-II: coalescent-based species tree estimation with many hundreds of taxa and thousands of genes. *Bioinformatics.* 2015; 31:44–52.
39. Liu L, Yu LL. Estimating species trees from unrooted gene trees. *Syst Biol.* 2011;60:661–7.
40. Salichos L, Rokas A. Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature.* 2013;497:327.
41. Huerta-Cepas J, Serra F, Bork P. ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data. *Molecular Biology and Evolution* 2016, 33:1635–38.
42. Tange O. GNU Parallel - the command-line power tool. *login.* 2011;36:42–7.
43. Emms D, Kelly S. Supplemental dataset for: OrthoFinder2: fast and accurate phylogenomic orthology analysis from gene sequences. <https://doi.org/10.5281/zenodo.1481147> 2019.
44. Yates A, Beal K, Keenan S, McLaren W, Pignatelli M, Ritchie GRS, Ruffier M, Taylor K, Vullo A, Flicek P. The Ensembl REST API: Ensembl data for any language. *Bioinformatics.* 2015;31:143–5.
45. Cunningham F, Amode MR, Barrell D, Beal K, Billis K, Brent S, Carvalho-Silva D, Clapham P, Coates G, Fitzgerald S, et al. Ensembl 2015. *Nucleic Acids Res.* 2015;43:D662–9.
46. Emms D, Kelly S. OrthoFinder. GitHub. <https://github.com/davidemms/OrthoFinder>. 2019. Accessed 21 Oct 2019.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

