1 **orthofisher: a broadly applicable tool for automated gene identification and retrieval**

2

3 Jacob L. Steenwyk[1,*] & Antonis Rokas[1,*]

4

5 [1] Vanderbilt University, Department of Biological Sciences, VU Station B#35-1634, Nashville,

6 TN 37235, United States of America

7

8 **ORCiDs**

9 J. L. Steenwyk: 0000-0002-8436-595X

10 A. Rokas: 0000-0002-7248-6551

11

12 *Correspondence should be addressed to: jacob.steenwyk@vanderbilt.edu or

13 antonis.rokas@vanderbilt.edu

14

15 **Running title:**        orthofisher: automated gene retrieval

16

19

20   **Abstract**

21   Identification and retrieval of genes of interest from genomic data is an essential step for many

22   bioinformatic applications. We present orthofisher, a command-line tool for automated

23   identification and retrieval of genes with high sequence similarity to a query profile-Hidden

24   Markov Model sequence alignment across a set of proteomes. Performance assessment of

25   orthofisher revealed high accuracy and precision during single-copy orthologous gene

26   identification. orthofisher may be useful for assessing gene annotation quality, identifying single-

27   copy orthologous genes for phylogenomic analyses, estimating gene copy number, and other

28   evolutionary analyses that rely on identification and retrieval of homologous genes from

29   genomic data. orthofisher comes complete with comprehensive documentation

30   (https://jlsteenwyk.com/orthofisher/), is freely available under the MIT license, and is available

31   for download from GitHub (https://github.com/JLSteenwyk/orthofisher), PyPi

32   (https://pypi.org/project/orthofisher/), and the Anaconda Cloud

33   (https://anaconda.org/jlsteenwyk/orthofisher).

34

35

36   **Introduction**

37   Sequence similarity searches of genomic data are commonly employed in diverse fields of

38   biology. Several pieces of software have been designed to infer statistically homologous

39   sequences from databases of sequence data, such as BLAST, DIAMOND, and HMMER

40   (Camacho *et al.* 2009; Eddy 2011; Madden 2013; Buchfink *et al.* 2015). One frequent use of

41   sequence similarity search methods is for the identification of orthologs, sequences present in the

42   common ancestor of two species, and homologs, sequences that stem from the same common

43   ancestral sequence (Gabaldón and Koonin 2013). For example, the OrthoFinder software

44   conducts BLAST all-vs-all searches across proteomes to infer groups of putatively orthologous

45   genes (Emms and Kelly 2019). Similarly, the BUSCO software aims to identify putatively

46   orthologous genes using a predetermined set of profile Hidden Markov Model sequence

47   alignments (pHMMs) derived from single-copy orthologous proteins from the OrthoDB database

48   (Waterhouse *et al.* 2013, 2018).

49

50    The results of these or similar pieces of software can facilitate diverse downstream analyses

51    (Remm *et al.* 2001; Li *et al.* 2003; Train *et al.* 2017; Waterhouse *et al.* 2018; Emms and Kelly

52    2019). However, global analyses, such as those conducted by OrthoFinder, are computationally

53    expensive and may be beyond the scope of a research project (e.g., studies focused on a few

54    genes). Similarly, software that rely on databases, such as BUSCO, are constrained to the

55    orthologs therein. As a result, there is a need for bioinformatic software that can conduct

56    automated identification and retrieval of putative homologs and orthologs across sequence

57    databases using user-specified query sequences and output files that facilitate downstream

58    analyses.

59

60    We introduce orthofisher, a command-line toolkit for automated identification of highly similar

61    sequences across proteomes using custom pHMMs. orthofisher facilitates downstream analyses

62    by creating multi-FASTA files populated with highly similar sequences identified during pHMM

63    searches. Default parameters are designed to identify sequences with the highest sequence

64    similarity (i.e., putative orthologous genes), but users can customize its use to best fit their

65    research question (e.g., relaxed thresholds can be used to obtain all putatively homologous genes;

66    similarly, searches in databases that contain gene isoforms can be used to retrieve all isoforms of

67    a particular gene). We demonstrate the efficacy of orthofisher by evaluating the precision and

68    recall for identification of sequences with high similarity to query pHMMs in a multiple

69    sequence FASTA (multi-FASTA) files from animals, plants, and fungi. Comparison of

70    orthofisher, BUSCO, and OrthoFinder revealed similar performance in identification of

71    sequences with high sequence similarity. Thus, orthofisher aims to streamline gene identification

72    and retrieval from genomic data, which is the first step of many bioinformatic analyses and

73    projects. We anticipate orthofisher will be of interest to diverse fields of computational biology

74    and to biologists and bioinformaticians.

75

76    **Methods**

77    orthofisher requires two files as input (Figure 1). One file—specified with the -m, --hmm

78    argument—provides the paths to query pHMMs that will be used during sequence similarity

79    search; the other file—specified with the -f, --fasta argument—provides the paths to FASTA files

80    that will be used as the sequence search database. orthofisher then loops through each FASTA

81 file and uses each pHMM to search for similar sequences using HMMER3 (Eddy 2011) with an

82 expectation-value threshold of 0.001 (which can be modified with the -e, --evalue argument).

83 orthofisher then parses the resulting HMMER3 output using biopython (Cock *et al.* 2009) and

84 identifies top hits. Top hits are defined following criteria used in the BUSCO pipeline

85 (Waterhouse *et al.* 2018) wherein all sequences with scores that are greater than or equal to 85%

86 of the score of the best hit are maintained. Users can modify this threshold using the -b, --

87 bitscore argument. Top hits are considered homologous genes.

88

89 orthofisher outputs three directories and two text files that enable researchers to easily evaluate

90 results from sequence similarity search and facilitate downstream analyses. The three directories

91 are

92     • *hmmsearch_output*: HMMER3 output files,

93     • *all_sequences*: one multi-FASTA file per pHMM, which are populated with

94     homologous sequences identified during the sequence similarity search step, and

95     • *scog*: one multi-FASTA file per pHMM, which are populated with only those

96     homologous sequences that are present at most only once in each genome.

97 The two text files are

98     • *short_summary.txt*: the number and percentage of sequences present in single-copy,

99     multi-copy, or absent sequences per pHMM search, and

100     • *long_summary.txt*: the homologous sequences identified during pHMM search for every

101     query and sequence database.

102 Contents of output files will be heavily dependent on user parameters, the pHMMs used, and the

103 input files. For example, transcriptomic data may require additional processing steps such as

104 collapsing isoforms into a single representative sequence per gene. The intent of orthofisher—

105 which is to identify single-copy orthologous genes—is flexible enough to capture paralogous

106 sequences as well. A tutorial for how to use orthofisher is publicly available as part of the online

107 documentation https://jlsteenwyk.com/orthofisher/tutorial.

108

109 Nearly 30% of bioinformatic tools fail to install (Mangul *et al.* 2019), which poses a nontrivial

110 problem for the reproducibility of computational experiments. To remedy this issue, we

111 implemented state-of-the-art standards of software development practices and design principles

112    (Darriba *et al.* 2018) following previously established protocol (Steenwyk *et al.* 2020, 2021). For

113    example, whenever changes to code are made, faithful function of orthofisher is tested using a

114    continuous integration pipeline, a process that automatically builds, packages, and tests

115    installation and function using Python versions 3.6, 3.7, and 3.8. We also wrote several unit and

116    integration tests that span 95% of the orthofisher code.

117

118    **<u>Results and Discussion</u>**

119    To determine the similarities and differences between orthofisher and other algorithms that

120    identify putative orthologs, we compared results obtained from orthofisher with that of BUSCO

121    and OrthoFinder. BUSCO and OrthoFinder are both widely adopted methods of identifying

122    orthologous genes across multiple proteomes. As noted in the introduction, each software differs

123    – more specifically, BUSCO conducts homology searches using a predefined set of pHMMs and

124    OrthoFinder conducts proteome-wide analysis to identify groups of orthologous genes. Thus, we

125    expect that if orthofisher can identify putative orthologs across proteomes, it will identify the

126    same genes BUSCO identifies during its sequence similarity search. Given that both algorithms

127    conduct pHMM-based searches, we anticipate that both will exhibit near identical performances.

128    When comparing orthofisher and BUSCO to OrthoFinder, we anticipate the sequences identified

129    during sequence similarity search by orthofisher and BUSCO will be in the same orthologous

130    group of genes inferred by OrthoFinder.

131

132    **orthofisher and BUSCO obtain similar results**

133    To evaluate the efficacy of orthofisher, we compared results obtained from orthofisher to those

134    obtained from BUSCO, v4.0.4 (Waterhouse *et al.* 2018). To do so, both algorithms were used to

135    identify 255 near-universally single-copy orthologous genes obtained from the Eukaryota

136    OrthoDB, v10 (Waterhouse *et al.* 2013), database across the proteomes of animals (*Homo*

137    *sapiens*: GCF_000001405.39; *Mus musculus*: GCF_000001635.27), plants (*Arabidopsis*

138    *thaliana*, NCBI accession: GCA_000001735.2; *Solanum lycopersicum*: GCF_000188115.4), and

139    fungi (*Saccharomyces cerevisiae*, NCBI accession: GCA_000146045.2; *Candida albicans*:

140    GCA_000182965.3). Measures of precision and recall were calculated as follows:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

141     where *TP* represents true positives, *FP* represents false positives, and *FN* represents false

142     negatives of single-copy orthologous genes. Precision and recall values range from 0 to 1 and

143     higher values reflect better performance.

144

145     Near perfect values of precision and recall (0.98 or [231 / [231 + 4]] and 1.0 or [231 / [231 + 0]],

146     respectively) reveal orthofisher is able to automate the identification and retrieval of sequences

147     with high similarity to the query pHMM. A low false positive rate of 0.02 was observed. The

148     difference in the performance of BUSCO and orthofisher stems from an additional set of gene-

149     specific score and length thresholds used by the BUSCO software, which are not implemented in

150     orthofisher. These results demonstrate that orthofisher can accurately identify homologous genes.

151

152     To demonstrate the importance of using a score threshold of 85% of the score observed in the

153     best hit following the BUSCO pipeline (Waterhouse *et al.* 2018), we highlight an example where

154     absence of a score threshold would have led to identification of additional putatively orthologous

155     genes. A HMMER search using the query BUSCO pHMM 1001705at2759 and a e-value

156     threshold of 1e-10 in the proteome of *A. thaliana* reports the gene as multi-copy whereas both

157     orthofisher and BUSCO report this gene to be single-copy. More specifically, when using only

158     an e-value threshold of 1e-10, the following nine genes are reported: AEE76455.1, AEE78573.1,

159     AEC10322.1, ANM68500.1, AED93406.1, AEE76521.1, AEE82221.1, AED98328.1, and

160     AEE29324.1; however, AEE76455.1 has a score of 242.5 and the next best hit, AEE78573.1, has

161     a score of 64.5. Thus, a score threshold of 85% of the best hit (in this case 242.5*0.85) is helpful

162     during sequence similarity searches.

163

164     **orthofisher and BUSCO perform similarly to OrthoFinder**

165     Comparison of the results of BUSCO and orthofisher to OrthoFinder, a global (or whole

166     proteome) ortholog calling algorithm revealed BUSCO, orthofisher, and OrthoFinder produce

167     similar results. To perform these comparisons, we first used OrthoFinder, v2.3.8 (Emms and

168     Kelly 2019), to identify putative orthologous groups of genes in the same animal, plant, and

169     fungal proteomes described above using an inflation parameter of 1.5 and DIAMOND,

170   v0.9.24.125 (Buchfink *et al.* 2015). Then, we determined if genes identified as multi-copy are

171   part of the same or different orthologous group(s) of genes and also assessed if genes identified

172   as single-copy in BUSCO or orthofisher were also single-copy in OrthoFinder.

173

174   Among multi-copy genes, we found BUSCO and OrthoFinder had nearly identical performance

175   in the proteomes of *A. thaliana*, *S. lycopersicum*, and *C. albicans*. For *S. cerevisiae*, one gene,

176   1545004at2759, out of 255 differed between BUSCO and OrthoFinder wherein BUSCO

177   identified two homologs and OrthoFinder split these two genes into different orthologous groups

178   of genes. A similar scenario was observed among 12 / 255 and 3 / 255 genes in the human and

179   mouse proteomes, respectively. For orthofisher, a similar scenario was observed for 1 / 255

180   genes in *S. lycopersicum*; 1 / 255 genes in *A. thaliana*; 8 / 255 genes in *S. cerevisiae*; 4 / 255

181   genes in *C. albicans*; 13 / 255 genes in the human proteome; and 4 / 255 genes in the mouse

182   proteome. We note that isoforms of the same gene sequence were present in the analysed

183   proteomes and were accounted for in these analyses.

184

185   Among single-copy genes, we observed a few instances where single-copy genes in BUSCO

186   were multi-copy in OrthoFinder. More specifically, this was observed for 8 genes in *S.*

187   *lycopersicum*; 16 genes in *A. thaliana*; 2 genes in *S. cerevisiae*; 2 genes in *C. albicans*; 36 genes

188   in the human proteome; and 26 genes in the mouse proteome. Similar results were observed for

189   orthofisher. More specifically, 16 / 255 genes in *A. thaliana* were identified as single-copy by

190   orthofisher but were in multi-copy orthologous groups of genes in OrthoFinder. The same

191   observation was made for 7 / 255 genes in *S. lycopersicum*; 1 / 255 gene in *S. cerevisiae*; 2 / 255

192   genes in *C. albicans*; 35 / 255 genes in the human proteome; and 24 / 255 genes in the mouse

193   proteome.

194

195   In summary, sequence similarity searches of 255 genes in 6 proteomes identified differences

196   among 105 genes (6.86%; 105 / 1,530) between BUSCO and OrthoFinder; similarly, we

197   identified differences among 116 genes (7.58%; 116 / 1,530) between orthofisher and

198   OrthoFinder. These differences likely stem from differences in the approach of each algorithm to

199   identify putative orthologs. Specifically, OrthoFinder uses DIAMOND and Markov clustering to

200   identify orthologous groups, BUSCO uses pHMM-based search and gene-specific score and

201 length thresholds using OrthoDB, and orthofisher uses pHMM-based similarity search

202 thresholds. Also, these differences are in part driven by each algorithm reporting different results

203 (i.e., OrthoFinder reports groups of putatively orthologous genes and BUSCO and orthofisher

204 report putative orthologous genes).

205

206 **orthofisher is helpful for estimating the number of members in a gene family**

207 To demonstrate how to use orthofisher to estimate the number of gene family members, we

208 estimate the number of DNA photolyase (PFam: PF00875) and zinc finger, C2H2 type (PFam:

209 PF00096) homologs in *S. cerevisiae*, *C. albicans*, two species from the *Hanseniaspora* genus (*H.*

210 *uvarum* NRRL Y1614 and *H. vineae* NRRL Y17529, both of which are known to lack DNA

211 photolyases (Steenwyk *et al.* 2019)), and three *Aspergillus* species (*A. niger* CBS 513.88, *A.*

212 *fumigatus* Af293, and *A. flavus* NRRL 3357). When estimating gene family number, we

213 recommend lowering the score threshold to, for example, 25% of the best hit, which we have

214 done here. In line with previous reports, we found that *Hanseniaspora* species lacked DNA

215 photolyases whereas *S. cerevisiae*, *C. albicans*, and all *Aspergillus* species had one or two DNA

216 photolyases. In contrast, proteins with Zinc finger domains are more abundant across all species

217 with copies ranging from 16 (*H. vineae*) to 39 (*A. flavus*).

218

219 <u>**Practical considerations**</u>

220 <u>The intended use of orthofisher is to help identify orthologous genes across species using</u>

221 <u>accurate and sensitive pHMM-based searches. We encourage users to evaluate results produced</u>

222 <u>by orthofisher using additional approaches (e.g., phylogenetic inference) to infer precise</u>

223 <u>relationships of orthology and paralogy among sequences. We note that orthofisher is not</u>

224 <u>explicitly designed to identify a single-representative sequence if multiple isoforms encoded by</u>

225 <u>one gene sequence are present in a proteome. Thus, we also suggest users collapse isoforms prior</u>

226 <u>to or after orthofisher analysis following standard protocol in many transcriptomics studies.</u>

227

228 In summary, orthofisher is a command-line tool for automated identification and retrieval of

229 genes of interest from genomic data. We anticipate orthofisher will be useful for evaluating

230 genome completeness, performing phylogenomic inferences, estimating gene family size, and

231 other analyses that rely on identification and retrieval of homologous genes from genomic data.

232

## **Web resources**

orthofisher comes complete with comprehensive documentation
(https://jlsteenwyk.com/orthofisher/), is freely available under the MIT license, and is available
for download from GitHub (https://github.com/JLSteenwyk/orthofisher), PyPi
(https://pypi.org/project/orthofisher/), and the Anaconda Cloud
(https://anaconda.org/jlsteenwyk/orthofisher). The proteomes, pHMMs, and outputs of
orthofisher, BUSCO, and OrthoFinder are available through figshare (doi:
10.6084/m9.figshare.14399150).

## Literature Cited
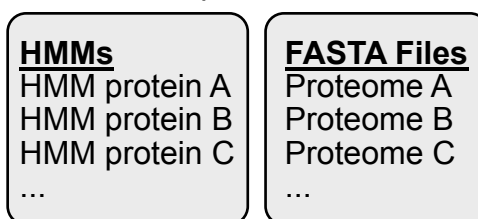
252   **Literature Cited**

253   Buchfink, B., C. Xie, and D. H. Huson, 2015 Fast and sensitive protein alignment using
254      DIAMOND. Nat. Methods 12: 59–60.

255   Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos *et al.*, 2009 BLAST+:
256      architecture and applications. BMC Bioinformatics 10: 421.

257   Cock, P. J. A., T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox *et al.*, 2009 Biopython: freely
258      available Python tools for computational molecular biology and bioinformatics.
259      Bioinformatics 25: 1422–1423.

260   Darriba, D., T. Flouri, and A. Stamatakis, 2018 The State of Software for Evolutionary Biology
261      (K. Crandall, Ed.). Mol. Biol. Evol. 35: 1037–1046.

262   Eddy, S. R., 2011 Accelerated Profile HMM Searches (W. R. Pearson, Ed.). PLoS Comput. Biol.
263      7: e1002195.

264   Emms, D. M., and S. Kelly, 2019 OrthoFinder: phylogenetic orthology inference for comparative
265      genomics. Genome Biol. 20: 238.

266   Gabaldón, T., and E. V. Koonin, 2013 Functional and evolutionary implications of gene
267      orthology. Nat. Rev. Genet. 14: 360–366.

268   Li, L., C. J. Stoeckert, and D. S. Roos, 2003 OrthoMCL: Identification of ortholog groups for
269      eukaryotic genomes. Genome Res. 13: 2178–2189.

270   Madden, T., 2013 The BLAST sequence analysis tool. BLAST Seq. Anal. Tool 1–17.

271   Mangul, S., T. Mosqueiro, R. J. Abdill, D. Duong, K. Mitchell *et al.*, 2019 Challenges and
272      recommendations to improve the installability and archival stability of omics computational
273      tools. PLOS Biol. 17: e3000333.

274   Remm, M., C. E. V. Storm, and E. L. L. Sonnhammer, 2001 Automatic clustering of orthologs
275      and in-paralogs from pairwise species comparisons. J. Mol. Biol. 314: 1041–1052.

276   Steenwyk, J. L., T. J. Buida, A. L. Labella, Y. Li, X.-X. Shen *et al.*, 2021 PhyKIT: a broadly
277      applicable UNIX shell toolkit for processing and analyzing phylogenomic data (R.
278      Schwartz, Ed.). Bioinformatics.

279   Steenwyk, J. L., T. J. Buida, Y. Li, X.-X. Shen, and A. Rokas, 2020 ClipKIT: A multiple
280      sequence alignment trimming software for accurate phylogenomic inference (A. Hejnol,
281      Ed.). PLOS Biol. 18: e3001007.

282   Steenwyk, J. L., D. A. Opulente, J. Kominek, X.-X. Shen, X. Zhou *et al.*, 2019 Extensive loss of
283      cell-cycle and DNA repair genes in an ancient lineage of bipolar budding yeasts (S.
284      Kamoun, Ed.). PLOS Biol. 17: e3000255.

285   Train, C.-M., N. M. Glover, G. H. Gonnet, A. M. Altenhoff, and C. Dessimoz, 2017 Orthologous
286      Matrix (OMA) algorithm 2.0: more robust to asymmetric evolutionary rates and more
287      scalable hierarchical orthologous group inference. Bioinformatics 33: i75–i82.

288   Waterhouse, R. M., M. Seppey, F. A. Simão, M. Manni, P. Ioannidis *et al.*, 2018 BUSCO
289      Applications from Quality Assessments to Gene Prediction and Phylogenomics. Mol. Biol.
290      Evol. 35: 543–548.

291   Waterhouse, R. M., F. Tegenfeldt, J. Li, E. M. Zdobnov, and E. V. Kriventseva, 2013 OrthoDB:
292      a hierarchical catalog of animal, fungal and bacterial orthologs. Nucleic Acids Res. 41:
293      D358–D365.

294

295 <u>**Figure Legend**</u>

296

297 **Figure 1. Workflow overview for orthofisher.** orthofisher takes two files as input, which

298 specify the location of query pHMMs and the FASTA files wherein sequence similarity searches

299 will be performed. orthofisher then outputs three directories and two text files that summarize

300 results and facilitate downstream analyses.