

1   **Title:** Universal orthologs infer deep phylogenies and improve genome quality

2   assessments

3   **Authors:** Md Nafis Ul Alam<sup>1,2</sup>, Cristian Román-Palacios<sup>3</sup>, Dario Copetti<sup>1</sup>, Rod A. Wing<sup>1,4</sup>

4   <sup>1</sup> Arizona Genomics Institute, School of Plant Sciences, University of Arizona, Tucson, AZ,

5   USA

6   <sup>2</sup> Plant Biotechnology Laboratory, Department of Biochemistry and Molecular Biology,

7   University of Dhaka, Dhaka, Bangladesh

8   <sup>3</sup> College of Information Science, University of Arizona, Tucson, AZ, USA

9   <sup>4</sup> Biological and Environmental Sciences and Engineering Division (BESE), King Abdullah

10   University of Science and Technology (KAUST), Thuwal, 23955-6900, Saudi Arabia

11

12

13

14

15

16

17

18

## 19 Abstract

20 Universal single-copy orthologs are the most conserved components of genomes. Although they  
21 are routinely used for studying evolutionary histories and assessing new assemblies, current  
22 methods do not incorporate information from available genomic data. Here, we first determine  
23 the influence of evolutionary history on universal gene content in plants, fungi and animals. We  
24 find that across 11,098 genomes comprising 2,606 taxonomic groups, 215 groups significantly  
25 vary from their respective lineages in terms of their BUSCO (Benchmarking Universal Single  
26 Copy Orthologs) completeness. Additionally, 169 groups display an elevated complement of  
27 duplicated orthologs, likely as an artifact of whole genome duplication events. Secondly, we  
28 investigate the extent of taxonomic congruence in BUSCO-derived whole-genome phylogenies.  
29 For 275 suitable families out of 543 tested, sites evolving at higher rates produce at most 23.84%  
30 more taxonomically concordant, and at least 46.15% less terminally variable phylogenies  
31 compared to lower-rate sites. We find topological differences between BUSCO concatenated and  
32 coalescent trees to be marginal and conclude that higher rate sites from concatenated alignments  
33 produce the most congruent and least variable phylogenies. Finally, we show that BUSCO  
34 misannotations can lead to misrepresentations of assembly quality. To overcome this issue, we  
35 filter a Curated set of BUSCOs (CUSCOs) that provide up to 6.99% fewer false positives  
36 compared to the standard BUSCO search and introduce novel methods for comparing assemblies  
37 using BUSCO synteny. Overall, we highlight the importance of considering evolutionary  
38 histories during assembly evaluations and release the phyca software toolkit that reconstructs  
39 consistent phylogenies and reports phylogenetically informed assembly assessments.

40 **Keywords:** Phylogenomics; Genomics; BUSCO; CUSCO; Orthologs; Gene Annotation

## 41 Introduction

42 High-quality reference genomes are becoming available for earth's flora and fauna at an  
43 accelerating rate. For example, between 12 August 2022 and 21 August 2023, 7,845 new  
44 organism genomes were released by NCBI alone (Sayers et al., 2024; Sayers et al., 2023). With  
45 advancements in long-read sequencing, nuclear conformation capture and optical mapping, the  
46 reconstruction of high-quality telomere-to-telomere assemblies (Li & Durbin, 2023; Rautiainen  
47 et al., 2023) is now becoming routine across all extant clades in the tree of life (Garg et al., 2024).

48 Conserved single-copy orthologs are used to create phylogenies (Van Damme et al., 2022) and  
49 evaluate the completeness of new assemblies (Manni et al., 2021), yet current tools and  
50 databases remain mostly oblivious to their varying evolutionary histories and taxonomic biases.  
51 For instance, OrthoDB (Kriventseva et al., 2019) is an established database of universal  
52 orthologs, but does not specifically explore the genome-wide variations in gene presence within  
53 major taxonomic groups. Similarly, OrthoFinder (Emms & Kelly, 2019) is used to reconstruct  
54 gene trees and species phylogenies, but does not analyze phylogenetic conflicts within and  
55 between gene features in alignment sites. Moreover, the detrimental effects of disregarding  
56 information about evolutionary history when using universal orthologs for assembly  
57 completeness tests (Cunha et al., 2023) has been overlooked in popular methods (Manni et al.,  
58 2021). Hence, a systematic exploration of public genomic data has the potential to improve  
59 existing methods for the utilization of universal orthologs in phylogenomics and assembly  
60 quality assessments.

61 Universal single-copy orthologs are the most stable components of genomes as they remain  
62 identifiably conserved in higher eukaryotes that diverged over millions of years ago (Gundappa

63 et al., 2022). A query set of universal single-copy orthologs (BUSCOs) (Manni et al., 2021)  
64 serves as a standard method for benchmarking gene content in newly assembled genomes.  
65 Fluctuations in BUSCO gene incidence is seen in some taxonomic groups (Cunha et al., 2023)  
66 but the full extent of BUSCO gene absence across genomically well-represented lineages has not  
67 been the subject of a focused or recent study. Although these genes remain under an evolutionary  
68 constraint of being maintained as single copies to balance dosage, polyploids (Fornasiero et al.,  
69 2024) and descendants of recently genome duplicated ancestors (Liu et al., 2020; Mansfeld et al.,  
70 2021; Wighard et al., 2022) carry fractionally elevated copy numbers. As such, BUSCO copy  
71 number variations have not been cataloged in detail across taxonomies or by gene identity.  
72 BUSCO gene sets have been the basis for some deep molecular phylogenies (Timilsena et al.,  
73 2022; Van Damme et al., 2022). BUSCOphylo (Sahbou et al., 2022) allows users to create  
74 BUSCO phylogenies, but it is not computationally feasible for gigabase-scale genomes or a large  
75 number of taxa. It also does not explore the accuracies or inconsistencies of BUSCO-derived  
76 phylogenies. Moreover, from the perspective of molecular phylogenetics, while substitution  
77 models have been trained on empirical sequences (Jarvis et al., 2015; Misof et al., 2014; Ran et  
78 al., 2018) to improve likelihood estimates, there have been limited efforts in incorporating  
79 divergent reference genome data (Armstrong et al., 2020) to derive improved inferences. Among  
80 many unknowns, there are known sources of model inadequacies that violate basic phylogenetic  
81 assumptions. For instance, gene histories are often obscured by incomplete lineage sorting (Yan  
82 et al., 2021), horizontal gene transfer (Schrempf & Szöllősi, 2020) or hybridization (Komarova  
83 & Lavrenchenko, 2022) and sites in gene alignments may support conflicting histories due to  
84 alignment errors (Edgar, 2021), recombination, long-branch attractions (Susko & Roger, 2021)  
85 or node-density artifacts (Venditti et al., 2006). Furthermore, alignment concatenation has been

86 shown to be statistically inconsistent for tree reconstructions (Kubatko & Degnan, 2007). This  
87 has led to many researchers assaying both concatenated and coalescent trees (Jarvis et al., 2015;  
88 Luo et al., 2022). Therefore, an explorative study that decouples sites in large, concatenated  
89 alignments from gene structures based on the column's rate of evolution has the potential to  
90 improve current methods of phylogenomic reconstructions.

91 In this study, we compiled BUSCO statistics for all plant, fungal and animal genomes cataloged  
92 in NCBI Genome (Sayers et al., 2022) up to January of 2024. Our objective was to improve  
93 methods for the utilization of BUSCO genes in phylogenomics and genome completeness  
94 evaluations. Under a wide range of rate and site configurations, we assessed the capacity of  
95 BUSCO genes in reconstructing taxonomically congruent phylogenies. We tested individual  
96 trees for taxonomic concordance, and tree distributions under the same conditions for variations  
97 in terminal leaf bifurcations. Through the constructed BUSCO database, we provided evidence  
98 for 2.25% to 13.33% mean lineage-wise gene misidentifications using the most widely used  
99 default BUSCO search parameters. Categorically, we procured a Curated set of BUSCO  
100 orthologs (CUSCOs) that attains a higher specificity for 10 major BUSCO eukaryotic lineages,  
101 namely Viridiplantae, Liliopsida, Eudicots, Chlorophyta, Fungi, Ascomycota, Basidiomycota,  
102 Metazoa, Arthropoda and Vertebrata. For robust comparisons and evaluations of closely related  
103 assemblies, a syntenic BUSCO metric was derived that offers higher contrast and better  
104 resolution than standard BUSCO gene searches. Our results, data and source code have been  
105 made available through a public database and software module named phyca.

106

107 **Results**

## 108 **BUSCO gene content is influenced by evolutionary history**

109 We compiled 11,098 eukaryotic genome assemblies from NCBI and observed that genomes for  
110 new animal genera were being released at a greater rate than plants and fungi (Figure 1A). The  
111 majority of NCBI genome assemblies contained a complete or near-complete complement of  
112 single and duplicated BUSCO genes (Figure 1B). Plant lineages had a much higher mean  
113 BUSCO duplication rate at 16.57% compared to fungi and animals at 2.79% and 2.21%  
114 respectively (Figure 1B and 1C). It is known that genomes of higher ploidy are often assembled  
115 into variable sets of pseudomolecules (Fornasiero et al., 2024; Healey et al., 2024) and this is  
116 reflected in our database (Supplementary Figure 1). The mean number of observed copies for the  
117 complete BUSCO gene set had 99.05% linear correlation with the number of copies of  
118 pseudomolecules in phased and partly phased assemblies (Supplementary Figure 1). There were  
119 169, 165 and 258 taxonomic groups out of 2,606 total that had significantly elevated means for  
120 duplicated BUSCO genes, mean BUSCO copy numbers and log assembly size respectively  
121 (Supplementary Table 1). For example, among the well-represented fungal classes, all 13  
122 assemblies of the family Backusellaceae had duplicated BUSCOs significantly greater than other  
123 fungal groups with a minimum of 11.42% and mean of 12.18%. For the 25 assemblies in the  
124 Mucoraceae family, the minimum and mean for duplicated BUSCOs were 5.1% and 6.54%  
125 respectively. The assembly counts, mean, minimum and maximum number of BUSCO metrics  
126 for every taxonomic group including Mann-Whitney U test p-values for deviation from group  
127 means are provided in Supplementary Table 1.

128 Extended drops in BUSCO completeness in Figure 1C are a result of bulk genome sequencing  
129 projects that resulted in large numbers of draft genome assemblies, e.g., Ellis et al., 2021 (Ellis et

130 al., 2021) who submitted 822 *de novo* butterfly genomes, Ronco et al., 2021 (Ronco et al., 2021)  
131 who submitted 539 cichlid fish genomes. Some taxonomic groups do show a predisposition to  
132 comparatively lower BUSCO completeness, as outlined in Supplementary Table 1. For instance,  
133 a number of *Incertae sedis* fungi-like organisms (mostly microsporidia) were found to contain  
134 <25% BUSCO genes and are seen as a dip at the trail of the fungal bars in Figure 1C  
135 (Supplementary table 1). In terms of taxonomy, it was found that across all BUSCO lineages and  
136 taxonomic levels, 215 groups had significantly different mean BUSCO completeness. The  
137 complete database, along with taxonomic classifications, assembly and BUSCO statistics are  
138 available to download and view at [www.phyca.org](http://www.phyca.org).

139 **Sites evolving at higher rates produce more taxonomically congruent  
140 phylogenies**

141 From our compiled data, we sought to determine the best way to utilize BUSCO genes to create  
142 broad whole-genome phylogenies spanning large evolutionary distances. Individual phylogenies  
143 were tested for agreement with NCBI taxonomic classifications. To assess taxonomic congruence,  
144 we created 3,566 phylogenetic trees for the 5 largest BUSCO lineages in terms of assembly and  
145 gene count. Our tests were focused on the Eudicots, Ascomycota, Basidiomycota, Arthropoda  
146 and Vertebrata lineages. Gene alignments for divergent taxa varied significantly based on  
147 parameters passed to the alignment algorithm (Supplementary figure 2). Different lineages had  
148 different rate profiles for aligned sites (Supplementary figure 3). Algae, fungi and early  
149 diverging metazoans displayed greater site heterogeneity in their alignments (Supplementary  
150 figures 2 and 3).

151 Phylogenetic trees under different evolutionary rates and alignment lengths were compared for  
152 taxonomic congruity. Variations of the LG and JTT substitution models (Le & Gascuel, 2008)  
153 with different rate categories were consistently found to have the highest likelihood under all  
154 conditions (Supplementary Table 2). The top 5 best substitution models based on Bayesian  
155 Information Criterion (BIC) for each condition with model comparison metrics are included in  
156 Supplementary Table 3. The number of unique amino acid residues in an alignment column was  
157 used as a proxy for site evolutionary rate. Sites evolving at higher rates together with longer  
158 alignments generally produced more taxonomically concordant trees (Figure 2A and  
159 Supplementary figure 4). Taxonomic concordance was predominant in eudicots with either 68 or  
160 69 out of 69 total families (98.55-100%) being reconstructed as monophyletic above 4,000  
161 alignment length and 5 or more unique amino acids. In arthropods and vertebrates, up to 113 out  
162 of 125 (90.40%) and 187 out of 225 (83.11%) respectively were reconstructed as monophyletic.  
163 In ascomycetes and basidiomycetes, only up to 60 out of 97 (61.86%) and 63 out of 88 (71.59%)  
164 respectively were found monophyletic in any single condition. The lineage and condition-wise  
165 monophyly counts are presented in Supplementary Figure 4. For each lineage, a consistent  
166 number of families were resolved as monophyletic in most of the trees, whilst some families  
167 precariously only appeared monophyletic at certain conditions (Figure 2B and Supplementary  
168 table 4). Alignments with greater numbers of sites and unique residues almost always resolved  
169 greater numbers of families (Figure 2C). Rate effects were more potent than alignment length  
170 (Figure 2D). For instance, 32 families out of 543 total were monophyletic under all tested  
171 conditions. Of the remaining 511, 59.47%, 84.61% and 86.53% were monophyletic when  
172 reconstructed with 2, 8 and 14 (low, moderate and high) unique amino acids per column  
173 respectively and 67.18%, 80.12% and 83.32% were reconstructed as monophyletic with 1,000,

174 5,000 and 10,000 alignment lengths respectively. Under conditions where the alignments did not  
175 provide sufficient information to accurately resolve tree topology (Supplementary figure 5),  
176 likelihoods were higher at greater rates and alignment lengths (Figure 2E). Variations in  
177 taxonomic concordance receded with increasing site counts and evolutionary rates in eudicots,  
178 arthropods and vertebrates, but the pattern was less prominent in ascomycetes and  
179 basidiomycetes (Figure 2F-G and Supplementary Figure 6).

180 We observed that all five tested lineages showed a similar trend where 462 out of 543 families  
181 were found monophyletic at the most informative condition with 14-character columns and an  
182 alignment length of 10,000 (Figure 2C and Supplementary table 4). Of the remaining 81, 42  
183 families could not be resolved as monophyletic (0 out of 50 trees) and the monophyly status of  
184 the remaining 39 families remained inconsistent. Rate preferences for monophyly in the queried  
185 families were not observed. The Petroicidae family of birds was the only family that yielded  
186 monophyletic trees across all 50 trees at rate condition 8, but was not consistently monophyletic  
187 in the higher rate condition of 14 with monophyly in 49 out of 50 trees in alignments of length  
188 10,000 (Supplementary table 4).

189 To interpret the relationship between tree likelihoods and taxonomic concordance, we  
190 recomputed likelihoods for all trees under a fixed set of alignments. Correlations between mean  
191 tree likelihood and taxonomic concordance diminished with longer alignments and faster  
192 evolving sites (Supplementary figure 5). At the same time, tree topologies were more stable at  
193 the terminal taxa for all lineages at higher evolutionary rates and greater site counts  
194 (Supplementary figure 7). The sets of trees created from sites with 8 unique characters and an  
195 alignment length of 10,000 for eudicots, arthropods and vertebrates were compared to BUSCO  
196 coalescent trees to contrast tree concordance between the two popular methods. No significant

197 variations in taxonomic agreement between concatenated trees and trees created under the  
198 multispecies coalescent model were observed (Figure 2H).

199 **A filtered BUSCO set provides improved assembly assessments**

200 Across all 10 lineages, on average 2.25% to 13.33% of BUSCO genes were misidentified in  
201 genomes where all BUSCO genes had been removed (Figure 3A and Table 1). Misidentification  
202 implies that a default BUSCO search would not identify divergent copies of these genes and the  
203 absence of the identified BUSCO gene in a query assembly would result in the inadvertent  
204 identification of the divergent copy. The magnitude of misidentification rates varies by lineage  
205 and was observed to be lowest across the fungal assemblies and highest across vertebrate and  
206 plant assemblies (Figure 3A and Table 1). Roughly 10% of BUSCO genes in all 10 lineages  
207 were misidentified at a far greater number of assemblies than others (Supplementary figure 8).  
208 Assessment of BUSCO completeness with these genes removed resulted in reduced numbers  
209 (Table 1) of BUSCO gene misidentifications in all lineages (Figure 3B). The reduction in false  
210 hits was more pronounced in the Vertebrata, Liliopsida, Eudicots and Chlorophyta lineages  
211 (Figure 3B and Table 1). For clarity, the Curated set of BUSCO genes has been named CUSCOs  
212 and the remaining Misannotation-prone BUSCO genes are hereon abbreviated as MUSCOs.

213 We analyzed the incidence of BUSCO misannotations by assembly and gene identity to  
214 extrapolate the source of this phenomenon. Gene misannotations were found to be more  
215 weighted towards the query gene rather than the query genome assembly (Figure 4A). Removal  
216 of MUSCOs resulted in better assembly assessment metrics and shifted the assembly quantiles of  
217 BUSCO misidentification towards the gene quantities (Figure 4B). Correlation analysis of  
218 lineage-wise misannotation rates with assembly metrics revealed that BUSCO gene

219 misidentifications correlated most with the mean number of BUSCO copies in the assembly, a  
220 metric we termed inflation (Figure 4C). Other variables showing the highest correlations were  
221 the number of miniProt hits (MPH) and the log of assembly size, being more pronounced in  
222 chlorophytes and vertebrates, respectively (Figure 4C).

223 Given the observed preponderance of misannotation rates in complex genomes in terms of  
224 assembly size, gene hits and BUSCO inflation (Figure 4C), we analyzed the syntenic patterns of  
225 identified and misidentified BUSCO genes to query potential evolutionary origins. For  
226 computational feasibility, all possible permutations of identified and misidentified BUSCO genes  
227 in 10 sets of gene blocks harboring up to 10 genes were tested. Gene block analysis revealed that  
228 beyond the species level, misidentified BUSCO genes are preserved in syntenic order at the  
229 highest rates in the Liliopsida, Viridiplantae and Eudicots lineages at 4.07%, 3.97% and 3.78%  
230 respectively. The fourth highest rate of syntenic misidentifications was in the Basidiomycota at  
231 just 0.88% and the lowest was in Arthropoda at 0.14%. Two such representative gene blocks  
232 from the Eudicots and Vertebrata lineages are shown in Figure 4D top and bottom respectively.  
233 This suggests that some misidentified BUSCO genes are remnants of gene duplication events  
234 where the syntenic copy became more divergent. Details for all computed gene blocks are  
235 available to download at [www.phyca.org/data.html](http://www.phyca.org/data.html). The syntenic analysis was extended to our  
236 complete data set with syntenic gene pairs to determine whether CUSCO and MUSCO genes  
237 contained pairs with one and two remnant genes in similar proportions. CUSCO syntenic  
238 doublets were progressively found in lower proportions with one and two remnant genes (Figure  
239 4E). However, MUSCO syntenic doublets appeared in similar proportions with pairs of  
240 identified and pairs of remnant genes (Figure 4E). MUSCO genes are therefore more syntenic in  
241 the remnant-remnant configuration compared to CUSCO genes.

242

243

244 **BUSCO collinearity is an indicator of pseudomolecule quality**

245 To demonstrate the utility of BUSCO synteny in assembly comparisons, we compiled and  
246 compared 1035 pairs of genomes of the same species with contrasting quality metrics. We  
247 employed an adjusted Intersection Over Union (IoU) metric with BUSCO gene doublets found in  
248 the same order and orientation to compare two assemblies. The denominator is adjusted by the  
249 difference in the number of contigs such that highly fragmented assemblies with the same gene  
250 order and orientation would be syntenically equivalent to highly contiguous assemblies. Hence,  
251 the syntenic doublet metric is designed to only capture differences in gene synteny and to not be  
252 influenced by varying numbers of contigs in query assemblies (Supplementary Figure 9).

253 BUSCO syntenic connections were able to capture far greater contrast in the assembly pairs  
254 compared to simply the difference in BUSCO completeness (Figure 5A). Syntenic BUSCO  
255 connections decayed exponentially with phylogenetic distance in our six non-overlapping  
256 BUSCO lineages (Figure 5B and 5C). We further compiled the 40 least contiguous NCBI  
257 assemblies of *Oryza sativa*, *Mus musculus*, *Drosophila melanogaster*, *Ovis aries* and  
258 *Arabidopsis thaliana* to represent the BUSCO syntenic distance between the assemblies as a  
259 dendrogram. Metrics for the full set of assemblies are provided in Supplementary Figures 10, 11,  
260 12, 13 and 14 respectively. An example of a dendrogram with 8 fragmented *Mus musculus*  
261 assemblies and a highly contiguous reference assembly is shown in Figure 5D. Less contiguous  
262 assemblies were found to be at greater syntenic distances to the higher quality assembly,  
263 implying greater numbers of BUSCO misidentification events or more extensive misassemblies.

264 To further assess how BUSCO synteny can indicate assembly quality, we visualized  
265 chromosome-wise BUSCO collinearity in a set of *Oryza* assemblies as a case study. The *Oryza*  
266 genus is genetically well characterized with several state-of-the-art chromosome level  
267 assemblies (Fornasiero et al., 2024). We demonstrate with a draft assembly (GenBank ID:  
268 GCA\_009805545.1) and a high-quality assembly of *Oryza longistaminata* (Reuscher et al., 2018)  
269 that BUSCO synteny can provide greater contrast between assemblies of varying quality  
270 compared to BUSCO metrics alone (Figure 6). Between the two *O. longistaminata* assemblies,  
271 although the number of curated BUSCO genes identified was comparable (98.82% and 93.17%),  
272 BUSCO collinearity was not preserved across the closely related sister taxa within the genus  
273 (Figure 6). These observed syntenic deviations are quantified by our adjusted IoU metric based  
274 on BUSCO gene connections (Supplementary Figure 9) and the syntenic distance between the  
275 two *O. longistaminata* assemblies was 82.25%. The full set of chromosomes for this test case is  
276 available on the phyca website at [www.phyca.org/data.html](http://www.phyca.org/data.html). The phyca software package allows  
277 users to similarly compare and visualize syntenic distances between assemblies and query  
278 genomes.

279

## 280 **Discussion**

281 Here, we presented our studies across three facets. First, we determined the prevalence of  
282 BUSCO gene variations by taxonomy through the compilation of available plant, fungi and  
283 animal genomes in the public domain. Second, we optimized site conditions for consistent  
284 phylogenomic reconstructions by maximizing taxonomic congruity and minimizing tree set  
285 variability. We then created large whole-genome phylogenies under the best determined

286 conditions for 10 major BUSCO lineages. Third, we provided evidence for BUSCO  
287 misannotations with the current software defaults and filtered a curated set of BUSCO genes for  
288 better genome quality assessments. To mitigate the effects of BUSCO misannotations during  
289 assembly evaluations, we described a novel method of comparing assemblies with BUSCO  
290 synteny that provides much better contrast for closely related assemblies of varying quality.

291

292 **BUSCO completeness and copy number variations**

293 Universal genes have been instrumental for querying gene space completeness and assembly  
294 quality (Manni et al., 2021). Our results show that the evolutionary history of a genome  
295 influences its BUSCO score and that this influence is prevalent in many taxonomic groups rather  
296 than just a few (Cunha et al., 2023). It was also observed that some groups vary more  
297 dramatically than others in BUSCO metrics (Supplementary table 1). Therefore, for assemblies  
298 from early diverging groups with few extant taxa or available genomes, BUSCO genes may  
299 provide an inadequate representation of gene space completeness. Given these observations, we  
300 propose that it is necessary to consider the evolutionary history of related taxa when evaluating  
301 the gene content of new genome assemblies.

302 Assembly gene content is influenced drastically by evolutionary history. Polyploid organisms are  
303 known for being able to maintain multiple sets of single-copy orthologs (Fornasiero et al., 2024)  
304 and genomes fractionate at varying rates post-duplication (Garsmeur et al., 2014). It is likely that  
305 groups that were found to harbor large sets of duplicated BUSCO genes in haploid assemblies  
306 have either experienced recent whole-genome duplication events or have adjusted their gene  
307 regulation to accommodate an inflated complement of some single-copy orthologs. The set of

308 genes that are more likely to be misidentified (Supplementary figure 8) are likely tolerated more  
309 in genomes at greater copy numbers. This is supported by the high correlation of gene  
310 misannotations to the BUSCO inflation metric shown in Figure 4C and the preservation of some  
311 syntenic remnant genes across large phylogenetic distances (Figures 4D and 4E). It is probable  
312 that misannotation-prone genes duplicated and subsequently functionalized in ancient ancestral  
313 genomes multiple times. Some of the duplicated copies may have taken up important functions  
314 that prevented the sequences from diverging drastically and the shared homology is now  
315 responsible for the observed false hits. The availability of a consolidated database of BUSCO  
316 results from public genomes allows researchers to derive meaningful copy number expectations  
317 for BUSCO genes in new assemblies based on evolutionary history.

318

319 **Decoupling aligned sites from gene features and a case for fast evolving columns**

320 Likelihood estimation in phylogenetics assumes that all sites evolve independently (Yang, 2006).  
321 Since this is not biologically meaningful (Nasrallah et al., 2010), advanced tree search algorithms  
322 split columns into invariant sites (Yang, 1996) and several rate categories (Yang, 1994) to  
323 address rate heterogeneity. We assumed that unique amino acid counts in aligned columns could  
324 serve as a proxy for evolutionary rate at that site and filtering sites by evolutionary rate would  
325 decouple sites from intragenic evolutionary influences. In practice, researchers often select fast  
326 evolving sites for dense phylogenies (Matschiner et al., 2020) and conversely, for deep  
327 phylogenies, they tend to use slowly evolving sites to optimize information content in the  
328 alignment (Misof et al., 2014). Our study broadly highlights the practical effects of rate variation  
329 and alignment information content on tree reconstruction. Rosenberg and Kumar, 2001

330 (Rosenberg & Kumar, 2001) showed that the number of sites have greater effect on tree accuracy  
331 compared to substitution rates. On the contrary, our results show that when an adequate number  
332 of sites are sampled (Figure 2D), site evolutionary rate has a greater effect on tree accuracy in  
333 terms of taxonomic congruity. In our studies, higher rate sites were generally found to produce  
334 better trees and there was minimal hindrance caused by long-branch attraction biases and  
335 heterotachy (Figure 2C and 2D).

336 Slow evolving sites have been favored throughout the history of molecular phylogenetics (Pisani,  
337 2004). Slow-fast analysis was popularized for phylogenetic reconstructions in the context of  
338 substitution saturation and long-branch biases (Cummins & McInerney, 2011; Kostka et al.,  
339 2008). Similarly, chi-squared tests are employed to detect compositional heterogeneity in  
340 alignments (Boudinot et al., 2023; Foster, 2004). The primary goal of these analyses has been to  
341 identify and prune fast evolving sites to improve phylogenies (Pisani, 2004). Such practices have  
342 recently been perceived with scrutiny (Superson & Battistuzzi, 2022) and Rangel and Fournier,  
343 2023 (Rangel & Fournier, 2023) has shown that fast evolving alignment sites can be highly  
344 informative. We show in Figure 2 (and Supplementary table 4) that higher rate sites improve  
345 taxonomic concordance across almost all 543 families tested, and always increase tree set  
346 consistencies (Supplementary figure 7) compared to lower rate sites. Therefore, contrary to  
347 popular practices, our results suggest that with adequate taxon sampling, faster rates for protein  
348 characters may produce more accurate phylogenies regardless of node depth.

349

350 **Phylogenies within the kingdom Fungi and recalcitrant evolutionary histories**

351 Some taxonomic classifications in the fungal domain are based on molecular ITS data (Carbone  
352 et al., 2017). Although ITS-based primers are commonly used for phylogenetic placement, the  
353 drawbacks of ITS sequences are apparent. RNA code has fewer letters than protein code and the  
354 ITS sequences are much shorter than most protein coding genes. Further, rRNA genes appear in  
355 large copy numbers (Lavrinienko et al., 2021; Lofgren et al., 2019) making them amenable to  
356 multiple evolutionary histories at greater divergence times. In contrast, single-copy orthologs  
357 exist under dosage restraints and this generally prevents copy number variations from persisting  
358 throughout evolutionary timescales (Garsmeur et al., 2014). Additionally, sampling greater  
359 numbers of taxa generally has a strong positive effect on phylogenetic accuracy (Heath et al.,  
360 2008) and BUSCO genes offer the means to include highly divergent clades. For these reasons, it  
361 is reasonable that BUSCO genes would be able to resolve deeper phylogenies with greater  
362 precision than ITS sequences.

363 We found taxonomic classifications to be more obscure for the kingdom fungi. Although tree  
364 entropy at the termini reduced by about 50% (Supplementary figure 7), we did not observe the  
365 same level of gradual reductions in the variance of monophyletic counts as seen from plants and  
366 higher animals (Supplementary figure 6). One likely explanation for these complications is their  
367 significantly higher rate of evolution and shorter generation times compared to other clades  
368 (Naranjo-Ortiz & Gabaldón, 2020). This can be seen in the greater fraction of high-rate sites  
369 shown in the state frequency spectra in Supplementary figure 3. This effect in conjunction with  
370 their compact genome sizes, relatively higher rates of gene flow (Gonçalves & Gonçalves, 2022)  
371 and very short generation times compared to higher eukaryotes makes the accurate  
372 reconstruction of fungal evolutionary histories challenging. Despite these challenges, the fungal  
373 families did follow the same trend as the higher eukaryotes in response to increasing

374 evolutionary rates in Figures 2C and 2D, albeit a greater fraction of families seemed to have  
375 members descended from more than one most recent common ancestor. The greater fraction of  
376 non-monophyletic groups could be an artifact of the limitations of the standard ITS-based  
377 classification scheme. These views are supported by a 9.72% observed higher fraction of  
378 monophyly in the higher fungi, basidiomycetes compared to the lower fungal phyla, ascomycetes  
379 (Supplementary table 4). Furthermore, compared to the other three lineages tested, ascomycetes  
380 and basidiomycetes show noticeably greater numbers of monophyletic groups with alignments of  
381 slowly evolving sites (Supplementary tables 4 and 5). One cause behind this could be higher  
382 rates of alignment errors in more distantly related taxa. In this regard, we did not consider the  
383 consistently reproduced alignments in Supplementary Figure 2 to be infallible since they are  
384 biased by the heuristics of multiple-sequence alignment algorithms (Edgar, 2021). Additionally,  
385 the higher range of monophyly in shorter alignments (Supplementary figure 4) could be  
386 explained by ITS-derived taxonomic classifications since those alignments resemble ITS  
387 alignments better in terms of length, slower rates of evolution and overall information content.  
388 Because of these ambiguities, and a higher fraction of families found not to be monophyletic, the  
389 ascomycetes and basidiomycetes weren't included in the coalescent study in Figure 2H.

390

### 391 **BUSCO provides a standard for whole-genome phylogenies**

392 At present, both concatenated and coalescent phylogenies are used in practice (Jarvis et al., 2015;  
393 Luo et al., 2022). The multispecies coalescent model corrects for incomplete lineage sorting to  
394 resolve ancestral relationships in higher taxa speciating from large populations. Jian et al., 2019  
395 (Jian et al., 2019) showed that the multispecies coalescent outperforms concatenation across a

396 range of metazoan groups. Our results suggest that such differences are usually marginal in terms  
397 of taxonomic congruity (Figure 2H). In the arthropods for instance, the coalescent tree set  
398 demonstrates a greater range of variation. It is also important to note that the total number of  
399 sites in the coalescent trees were far greater than the concatenated trees since up to 75 whole  
400 genes were included. For comparison, the vertebrate tree likelihoods were still improving at the  
401 10,000 site count mark (Supplementary figure 5). We propose that when there is adequate  
402 information content in the alignments, the high dimensional likelihood surface flattens out  
403 harboring several vicinal and localized peaks and valleys. This results in the distribution of  
404 alternate topologies with varying model likelihoods spread out within a range of monophyly  
405 counts in the correlation plots shown in Supplementary figure 5. We thus conclude that the  
406 multispecies coalescent offers a powerful framework, but results should still be interpreted with  
407 caution, and our BUSCO concatenation method offers a robust alternative when suitable.

408 The search space for phylogenetic trees grows faster than exponentials with increasing numbers  
409 of terminal nodes (Yang, 2006). Our smallest tested tree had 592 terminal nodes which equates  
410 to a search space of  $\frac{((2 \times 592) - 5)!}{2^{(592-3)} \times (592-3)!} = 2.12 \times 10^{1556}$ . This high number of taxa makes the tree  
411 space numerically intractable even with the best available heuristics. The exact same tree  
412 topology was never reproduced in our results under any condition. From our evaluations of the  
413 tree distributions, we suggest that: 1) consistent reconstruction of a greater number of groups as  
414 monophyletic offers support for internal nodes, and 2) reduced terminal variability in tree  
415 distributions provides confidence for accuracy of overall tree topology. Combined, ancestral  
416 histories reconstructed from our method of sampling high-rate sites from whole-genome BUSCO  
417 data should be deemed more reliable than ITS or gene trees, and on par with coalescent-based  
418 trees. In the phyca website and software package, the 39 clades with undetermined monophyly

419 status have been shared in Supplementary Table 4 to alert users to be cautious about drawing  
420 interpretations. It is important to be aware that with large datasets, model inadequacies (Delsuc et  
421 al., 2005) could result in erroneous topologies having high support values. It is therefore possible  
422 that for any individual taxa or clade, the reduced terminal variability in our tree sets may have  
423 reinforced erroneous placements. We recommend that researchers with more nuanced  
424 evolutionary questions should consider rebuilding subtrees within their clade of interest. For this  
425 purpose, phyca provides a user-friendly implementation of our proposed methods to construct  
426 phylogenies from user defined sets of query taxa.

427

## 428 **Shortcomings of homology-based and probabilistic gene predictions**

429 BUSCO has been the unrivaled standard for gene space completeness tests since 2019 (Seppey et  
430 al., 2019). BUSCO relies on sequence homology searches through sequence alignments and  
431 subsequent refinement of search results by trained hidden Markov models (Manni et al., 2021).  
432 In general, alignment-based methods for gene identification are employed using arbitrary cutoffs  
433 (Levy Karin et al., 2020) and probabilistic models are used with empirically trained probabilities  
434 (Edgar, 2021; Wheeler & Eddy, 2013). BUSCO gene prediction by Compleasm (Huang & Li,  
435 2023), a better implementation of BUSCO, starts with a miniProt (Li, 2023) search that is  
436 restricted to report duplicate genes only if the alignment score is at least 95% of the best  
437 alignment. Compleasm has four additional threshold parameters for secondary hits, gene identity,  
438 fraction and completeness respectively. These thresholds have been empirically optimized by the  
439 developers to maximize precision and recall (Huang & Li, 2023). Almost all user-reported  
440 BUSCO results are reported based on default parameters (Ellis et al., 2021; Fornasiero et al.,

441 2024; Healey et al., 2024; Liu et al., 2020; Mansfeld et al., 2021). Readjustment of these  
442 parameters would adversely alter the preoptimized tunings, and for experimental explorations,  
443 there would be an inordinate number of permutations to consider. Our method of removing genes  
444 and rerunning under default settings mimics the effect of natural gene loss events. Our analysis  
445 of false positive hits revealed a set of less reliable BUSCO genes with a significantly higher  
446 propensity of being misannotated (Supplementary figure 8). We surmise that for gene predictions  
447 there may be no “one glove fits all” method that will work for all genes across all possible  
448 lineages. With this view in mind, integrative approaches have been suggested in the past to  
449 improve gene prediction accuracies (Alam & Chowdhury, 2020). We conclude that putative gene  
450 prediction is a tricky endeavor and demonstrate in Figure 3B and Table 1 that removing the less  
451 reliable genes from the BUSCO gene set improves precision without compromising recall.

452

## 453 Conclusion

454 Universal orthologs are critical inferential tools for evolutionary genomic research. To improve  
455 the utilization of BUSCO genes in this field, we first compiled and comprehensively analyzed  
456 their presence and copy number variations within the expansive higher eukaryotic domain. Based  
457 on our findings, we suggest that evolutionary histories must be considered for proper  
458 interpretation of BUSCO completeness metrics. Second, we determined the extent to which the  
459 ancestral histories of major eukaryotic lineages could be resolved through universal single-copy  
460 orthologs. Our results imply that columns evolving at higher rates in alignments of protein  
461 characters are more robust for deep phylogenomic reconstructions. We described a novel way to  
462 consider phylogenetic accuracy using taxonomy and a simplified way to express tree set

463 variability by enumerating terminal leaf bifurcations. In light of our findings, we produced the  
464 largest unified nuclear genome-based phylogenies for 10 major taxonomic groups in the plant,  
465 fungi and animal kingdoms to date. Within these phylogenies, we highlighted clades that were  
466 consistently reconstructed as monophyletic with respect to their taxonomic labels and  
467 distinguished clades that demonstrated more recalcitrant ancestral histories. Finally, our database  
468 yielded a filtered set of BUSCO orthologs that provide a better representation of assembly gene  
469 content compared to the standard BUSCO search. We showed that more robust evaluation of  
470 genome quality can be attained through the incorporation of BUSCO synteny information from  
471 related assemblies. Our processed data and tools have been made easily accessible for robust  
472 phylogenomic reconstructions, rapid placement of query assemblies by appending BUSCOs to  
473 large, precomputed alignments and for deriving phylogenetically informed assembly quality  
474 evaluations.

475

## 476 **Materials and methods**

### 477 **Database compilation and classification**

478 Metadata for plant, fungi and animal genome assemblies were sourced from the NCBI genome  
479 database (Sayers et al., 2022) accessed on January 14, 2024. Assemblies flagged by NCBI as  
480 partial and contaminated were not used. Special characters (\(\)-/#:=+[]) were removed from  
481 organism names to avoid software errors during automation. The assembly metadata were sorted  
482 by level of assembly set by NCBI (complete, chromosome, scaffold, contig), date of release  
483 (newest to oldest) and assembly size (largest to smallest) respectively. Only the top entry for

484 identical organism names was kept. Batch downloads were executed using the cURL application  
485 ([www.curl.se](http://www.curl.se)). The NCBItax2lin software (<https://github.com/zxyue/ncbitax2lin>) was used to  
486 assign taxonomic classifications at the phylum, class, order, family and genus levels to the  
487 assemblies. The Mann-Whitney test was used to test the hypotheses of whether assemblies  
488 within a taxonomic group had a significantly different mean for a metric compared to all  
489 assemblies in the BUSCO lineage. A Bonferroni correction of  $\frac{0.05}{2 \times (\text{total assembly count})}$  was carried  
490 out to determine the p-value cutoff thresholds.

## 491 **Finding and aligning universal orthologs**

492 Searches for universal orthologs was executed using Compleasm version 0.2.5 (Huang & Li,  
493 2023) with OrthoDB version 10 reference sequences (Kriventseva et al., 2019) for the  
494 Viridiplantae, Chlorophyta, Liliopsida, Eudicots, Fungi, Ascomycota, Basidiomycota, Metazoa,  
495 Arthropoda and Vertebrata lineages using the default settings. For duplicated universal single-  
496 copy orthologs, the ortholog that was more syntenic with the database was selected. Gene copies  
497 sharing adjacent BUSCO orthologs at greater frequency within the database were defined as  
498 more syntenic. For equally syntenic duplicates, the gene with greater sequence identity was  
499 retained. Assemblies that did not contain 90% of the BUSCO orthologs in each lineage were  
500 included in the database but dropped from the subsequent phylogenetic analysis for suboptimal  
501 quality. All identified orthologs for each gene in each lineage were aligned using MUSCLE  
502 version 5.1 (Edgar, 2021). Alignments for the Viridiplantae, Fungi, Metazoa and Arthropoda  
503 lineages were done with 16 total combinations of four parameter perturbations and four guide  
504 tree permutations to create a stratified ensemble of multiple sequence alignments (Edgar, 2021).

505 Confidence for each column in the alignment was computed using the addconfseq flag in  
506 MUSCLE v5.

507

## 508 **Phylogenetic assessment**

509 For Eudicots, Ascomycota, Basidiomycota, Arthropoda and Vertebrata lineages, aligned sites  
510 were filtered by the number of unique amino acids in the column as a proxy for rate of evolution  
511 at that site. For rate categories 2 to 15, we selected between 1,000 and 20,000 sites at 1,000 site  
512 increments. This resulted in a total of  $14 \times 20 = 280$  alignments per lineage. For the Arthropoda  
513 lineage, we could only select up to 14,000 sites per category because of the relatively lower  
514 number of aligned sites. We had  $14 \times 14 = 196$  total alignments for arthropods. Assemblies that  
515 had fewer than 90% BUSCO genes and aligned sites that comprised more than 10% gaps were  
516 removed. IQ-TREE version 2.1.2 (Minh et al., 2020) was used with default settings including  
517 built-in ModelFinder2 (Kalyaanamoorthy et al., 2017) to create maximum likelihood trees for  
518 every alignment. There were 198 total nuclear substitution models to test including alternate site  
519 specifications and rate category variations. A total of  $280 \times 4 + 196 = 1316$  individual trees were  
520 created and tested in this step. Trees were assessed for taxonomic congruity by counting the  
521 number of families that descended monophyletically from a common ancestor. For the terminal  
522 and central rates 2, 8 and 14, five sets of alignments were sampled for site counts 1,000, 5,000  
523 and 10,000. We carried out 10 independent searches on the tree space for each alignment with a  
524 different random seed, resulting in a total of  $10 \times 5 = 50$  trees for  $3 \times 3 = 9$  conditions in 5 lineages.  
525 The total number of trees at this stage was  $50 \times 9 \times 5 = 2250$ . Each individual tree was assessed for  
526 congruity by counting the number of monophyletic families. The set of trees in each condition

527 was assessed for entropy or degree of variation at the terminal leaves by counting the total  
528 number of unique terminal bifurcations in the set. The 5 alignment sets at site rate 8 and site  
529 count 5,000 were used to compute likelihoods for all 2,250 trees using IQ-TREE. The mean  
530 likelihood score of 5 alignments was used as the likelihood for each individual tree. Gene trees  
531 were created for the 200 longest genes in the Eudicots, Arthropoda and Vertebrata lineages.  
532 From 5 to 75 genes were selected with increments of 5 genes at random to create 15 coalescent  
533 trees under the multi-species coalescent model in Astral-pro3 version 1.19.3.5 (Zhang & Mirarab,  
534 2022).

535

### 536 **Assessing misidentified BUSCOs**

537 For BUSCO misidentification studies, all single and duplicate BUSCO genes identified by  
538 Compleasm were first removed using scripts available on the phyca GitHub page and  
539 Compleasm was rerun on the genome set. Genes found in fragments were not considered. The  
540 curated BUSCO gene set was selected manually by looking at the frequency at which each  
541 BUSCO gene was misidentified. For each assembly, genome inflation was defined as the  
542 average number of times the BUSCO gene set was found in the assembly. Polyploid genomes  
543 shown in Supplementary Figure 1 were labeled manually according to literature through searches  
544 done by the species names. Assembly level for chromosome scale assemblies was determined by  
545 the labels assigned to the pseudomolecules.  
546 Gene blocks were traced with all possible permutations of identified and remnant BUSCO genes  
547 up to 11 genes in length using phyca scripts. To compute CUSCO and MUSCO proportions,

548 Remnant-Identified gene doublets were considered syntenic when they were matched in gene  
549 identity and orientation by a Identified-Identified doublet within the same lineage. Remnant-  
550 Remnant gene doublets were considered syntenic when they were matched by either a Remnant-  
551 Identified doublet or an Identified-Identified doublet. For each set of BUSCO doublets, fraction  
552 of doublets where both genes were CUSCO genes was defined as the CUSCO proportion and the  
553 fraction of doublets where both genes were MUSCO genes were defined as MUSCO proportions.  
  
554 For comparisons of BUSCO gene content and syntenic distance, two assemblies of the highest  
555 and lowest N50 were selected for organisms with more than one available genome assembly  
556 from NCBI Genome. Only pairs where the difference in N50 was greater than 200Kb were  
557 considered. Assemblies with an N50 of less than 1Mb or less than 80% BUSCO content were  
558 filtered out. Syntenic distance and distance matrices were computed by phyca. Exponential  
559 curves were fit using the curve\_fit function from SciPy (Virtanen et al., 2020) version 1.14.1.  
560 Distance matrices were converted to newick trees using scikit-bio version 0.6.2  
561 (<https://scikit.bio>).  
  
562 The *Oryza alta* assembly was from Yu et al., 2021 (Yu et al., 2021) and *Oryza coarctata* was  
563 from Fornasiero et al., 2024 (Fornasiero et al., 2024). Pseudomolecules of two subgenomes of  
564 the polyploid *Oryza* species were separated through their sequence headers. All dendograms and  
565 cladograms were created using BioNick version 0.0.3 (<https://pypi.org/project/BioNick/0.0.3/>).  
566 The phyca website uses phylotree.js (<https://phylotree.hphy.org/>) for dynamic tree  
567 visualizations.  
  
568

569 **Acknowledgements**

570 We would like to acknowledge Robert C. Edgar and Derrick Zwickl for their valuable insights  
571 and suggestions on alignment and phylogenetic methods. We acknowledge all members of the  
572 Arizona Genomics Institute and Data Diversity Lab teams for their continued support and  
573 encouragement. We acknowledge Abid Mahmood and his team for their help with the project  
574 website development. We also thank Chandler Sobel-Sorenson and the University of Arizona  
575 High Performance Computing team for maintaining and assisting us with necessary  
576 computational resources and software.

577

## 578 **Funding**

579 This work was supported by the Bud Antle Endowed Chair of Excellence in Agriculture & Life  
580 Sciences awarded to Rod A. Wing at the University of Arizona.

581

582

## 583 **Author information**

### 584 **Authors and Affiliations**

585 Arizona Genomics Institute, School of Plant Sciences, University of Arizona, Tucson, AZ, USA

586 Md Nafis Ul Alam, Dario Copetti, Rod A. Wing

587 Plant Biotechnology Laboratory, Department of Biochemistry and Molecular Biology,

588 University of Dhaka, Dhaka, Bangladesh

589 Md Nafis Ul Alam

590 College of Information Science, University of Arizona, Tucson, AZ, USA  
591 Cristian Román-Palacios  
  
592 Center for Desert Agriculture, Biological and Environmental Sciences and Engineering Division  
593 (BESE), King Abdullah University of Science and Technology (KAUST), Thuwal, 23955-6900,  
594 Saudi Arabia  
595 Rod A. Wing

596

597

598 **Contributions**

599 RAW and MNUA conceived and planned the project. RAW, CRP and DC supervised the work.  
600 CRP reviewed the phylogenetic methods and helped design further experiments to validate the  
601 results. MNUA implemented the methods, compiled the data set, developed the algorithms and  
602 wrote the scripts and manuscript. All authors reviewed and edited the manuscript.

603

604 **Corresponding author**

605 Please direct correspondence to Rod A. Wing and Cristian Román-Palacios.

606 **Ethics declarations**

607 **Ethics approval and consent to participate**

608 Not applicable.

609    **Consent for publication**

610    Not applicable.

611    **Competing interests**

612    The authors declare that they have no competing interests.

613

614

615

616

617

618

619

620

621 **References**

- 622 Alam, M. N. U., & Chowdhury, U. F. (2020). Short k-mer abundance profiles yield robust  
623 machine learning features and accurate classifiers for RNA viruses. *PLoS One*, 15(9),  
624 e0239381. <https://doi.org/10.1371/journal.pone.0239381>
- 625 Armstrong, J., Hickey, G., Diekhans, M., Fiddes, I. T., Novak, A. M., Deran, A., Fang, Q., Xie,  
626 D., Feng, S., Stiller, J., Genereux, D., Johnson, J., Marinescu, V. D., Alföldi, J., Harris, R.  
627 S., Lindblad-Toh, K., Haussler, D., Karlsson, E., Jarvis, E. D., . . . Paten, B. (2020).  
628 Progressive Cactus is a multiple-genome aligner for the thousand-genome era. *Nature*,  
629 587(7833), 246-251. <https://doi.org/10.1038/s41586-020-2871-y>
- 630 Boudinot, B. E., Fikáček, M., Lieberman, Z. E., Kusy, D., Bocak, L., Mckenna, D. D., & Beutel,  
631 R. G. (2023). Systematic bias and the phylogeny of Coleoptera—A response to Cai et al.  
632 (2022) following the responses to Cai et al. (2020). *Systematic Entomology*, 48(2), 223-  
633 232. <https://doi.org/https://doi.org/10.1111/syen.12570>
- 634 Carbone, I., White, J. B., Miadlikowska, J., Arnold, A. E., Miller, M. A., Kauff, F., U'Ren, J. M.,  
635 May, G., & Lutzoni, F. (2017). T-BAS: Tree-Based Alignment Selector toolkit for  
636 phylogenetic-based placement, alignment downloads and metadata visualization: an  
637 example with the Pezizomycotina tree of life. *Bioinformatics*, 33(8), 1160-1168.  
638 <https://doi.org/10.1093/bioinformatics/btw808>
- 639 Cummins, C. A., & McInerney, J. O. (2011). A method for inferring the rate of evolution of  
640 homologous characters that can potentially improve phylogenetic inference, resolve deep  
641 divergence and correct systematic biases. *Systematic Biology*, 60(6), 833-844.
- 642 Cunha, T. J., de Medeiros, B. A. S., Lord, A., Sørensen, M. V., & Giribet, G. (2023). Rampant  
643 loss of universal metazoan genes revealed by a chromosome-level genome assembly of

- 644 the parasitic Nematomorpha. *Curr. Biol.*, 33(16), 3514-3521.e3514.
- 645 <https://doi.org/10.1016/j.cub.2023.07.003>
- 646 Delsuc, F., Brinkmann, H., & Philippe, H. (2005). Phylogenomics and the reconstruction of the  
647 tree of life. *Nature Reviews Genetics*, 6(5), 361-375. <https://doi.org/10.1038/nrg1603>
- 648 Edgar, R. C. (2021). Muscle5: High-accuracy alignment ensembles enable unbiased assessments  
649 of sequence homology and phylogeny. *Nature Communications*, 13(1), 6968.  
650 <https://doi.org/10.1038/s41467-022-34630-w>
- 651 Ellis, E. A., Storer, C. G., & Kawahara, A. Y. (2021). De novo genome assemblies of butterflies.  
652 *Gigascience*, 10(6). <https://doi.org/10.1093/gigascience/giab041>
- 653 Emms, D. M., & Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for  
654 comparative genomics. *Genome Biol.*, 20(1), 238. [https://doi.org/10.1186/s13059-019-1832-y](https://doi.org/10.1186/s13059-019-<br/>655 1832-y)
- 656 Fornasiero, A., Feng, T., Al-Bader, N., Alsantely, A., Mussurova, S., Hoang, N. V., Misra, G.,  
657 Zhou, Y., Fabbian, L., Mohammed, N., Rivera Serna, L., Thimma, M., Llaca, V.,  
658 Parakkal, P., Kudrna, D., Copetti, D., Rajasekar, S., Lee, S., Talag, J., . . . Wing, R. A.  
659 (2024). Oryza genome evolution through a tetraploid lens. *bioRxiv*.  
660 <https://doi.org/10.1101/2024.05.29.596369>
- 661 Foster, P. G. (2004). Modeling Compositional Heterogeneity. *Systematic Biology*, 53(3), 485-  
662 495. <https://doi.org/10.1080/10635150490445779>
- 663 Garg, V., Bohra, A., Mascher, M., Spannagl, M., Xu, X., Bevan, M. W., Bennetzen, J. L., &  
664 Varshney, R. K. (2024). Unlocking plant genetics with telomere-to-telomere genome  
665 assemblies. *Nat. Genet.*, 1-12. <https://doi.org/10.1038/s41588-024-01830-7>

- 666 Garsmeur, O., Schnable, J. C., Almeida, A., Jourda, C., D'Hont, A., & Freeling, M. (2014). Two  
667 evolutionarily distinct classes of paleopolyploidy. *Mol. Biol. Evol.*, 31(2), 448-454.  
668 <https://doi.org/10.1093/molbev/mst230>
- 669 Gonçalves, P., & Gonçalves, C. (2022). Horizontal gene transfer in yeasts. *Current Opinion in  
670 Genetics & Development*, 76, 101950.  
671 <https://doi.org/https://doi.org/10.1016/j.gde.2022.101950>
- 672 Gundappa, M. K., To, T.-H., Grønvold, L., Martin, S. A. M., Lien, S., Geist, J., Hazlerigg, D.,  
673 Sandve, S. R., & Macqueen, D. J. (2022). Genome-wide reconstruction of  
674 rediploidization following autopolyploidization across one hundred million years of  
675 Salmonid evolution. *Mol. Biol. Evol.*, 39(1). <https://doi.org/10.1093/molbev/msab310>
- 676 Healey, A. L., Garsmeur, O., Lovell, J. T., Shengquiang, S., Sreedasyam, A., Jenkins, J., Plott, C.  
677 B., Piperidis, N., Pompidor, N., Llaca, V., Metcalfe, C. J., Doležel, J., Cápal, P., Carlson,  
678 J. W., Hoarau, J. Y., Hervouet, C., Zini, C., Dievart, A., Lipzen, A., . . . D'Hont, A.  
679 (2024). The complex polyploid genome architecture of sugarcane. *Nature*, 628(8009),  
680 804-810. <https://doi.org/10.1038/s41586-024-07231-4>
- 681 Heath, T. A., Hedtke, S. M., & Hillis, D. M. (2008). Taxon sampling and the accuracy of  
682 phylogenetic analyses. *Journal of systematics and evolution*, 46(3), 239.
- 683 Huang, N., & Li, H. (2023). compleasm: a faster and more accurate reimplementation of BUSCO.  
684 *Bioinformatics*, 39(10). <https://doi.org/10.1093/bioinformatics/btad595>
- 685 Jarvis, E. D., Mirarab, S., Aberer, A. J., Li, B., Houde, P., Li, C., Ho, S. Y. W., Faircloth, B. C.,  
686 Nabholz, B., Howard, J. T., Suh, A., Weber, C. C., da Fonseca, R. R., Alfaro-Núñez, A.,  
687 Narula, N., Liu, L., Burt, D., Ellegren, H., Edwards, S. V., . . . Avian Phylogenomics, C.

- 688 (2015). Phylogenomic analyses data of the avian phylogenomics project. *Gigascience*, 4,  
689 4. <https://doi.org/10.1186/s13742-014-0038-1>
- 690 Jian, X., Edwards, S., & Liu, L. (2019). The multispecies coalescent model outperforms  
691 concatenation across diverse phylogenomic data sets. *Systematic Biology*, 69, 795-812.  
692 <https://doi.org/10.1093/sysbio/syaa008>
- 693 Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A., & Jermiin, L. S. (2017).  
694 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods*,  
695 14(6), 587-589. <https://doi.org/10.1038/nmeth.4285>
- 696 Komarova, V. A., & Lavrenchenko, L. A. (2022). Approaches to the detection of hybridization  
697 events and genetic introgression upon phylogenetic incongruence. *Biol. Bull. Rev.*, 12(3),  
698 240-253. <https://doi.org/10.1134/s2079086422030045>
- 699 Kostka, M., Uzlikova, M., Cepicka, I., & Flegr, J. (2008). SlowFaster, a user-friendly program  
700 for slow-fast analysis and its application on phylogeny of Blastocystis. *BMC  
701 Bioinformatics*, 9(1), 341. <https://doi.org/10.1186/1471-2105-9-341>
- 702 Kriventseva, E. V., Kuznetsov, D., Tegenfeldt, F., Manni, M., Dias, R., Simão, F. A., &  
703 Zdobnov, E. M. (2019). OrthoDB v10: sampling the diversity of animal, plant, fungal,  
704 protist, bacterial and viral genomes for evolutionary and functional annotations of  
705 orthologs. *Nucleic Acids Res.*, 47(D1), D807-D811. <https://doi.org/10.1093/nar/gky1053>
- 706 Kubatko, L. S., & Degnan, J. H. (2007). Inconsistency of phylogenetic estimates from  
707 concatenated data under coalescence. *Syst. Biol.*, 56(1), 17-24.  
708 <https://doi.org/10.1080/10635150601146041>
- 709 Lavrinienko, A., Jernfors, T., Koskimäki, J. J., Pirttilä, A. M., & Watts, P. C. (2021). Does  
710 Intraspecific Variation in rDNA Copy Number Affect Analysis of Microbial

- 711 Communities? *Trends in Microbiology*, 29(1), 19-27.
- 712 <https://doi.org/10.1016/j.tim.2020.05.019>
- 713 Le, S. Q., & Gascuel, O. (2008). An Improved General Amino Acid Replacement Matrix.
- 714 *Molecular Biology and Evolution*, 25(7), 1307-1320.
- 715 <https://doi.org/10.1093/molbev/msn067>
- 716 Levy Karin, E., Mirdita, M., & Söding, J. (2020). MetaEuk-sensitive, high-throughput gene
- 717 discovery, and annotation for large-scale eukaryotic metagenomics. *Microbiome*, 8(1), 48.
- 718 <https://doi.org/10.1186/s40168-020-00808-x>
- 719 Li, H. (2023). Protein-to-genome alignment with miniprot. *Bioinformatics*, 39(1), btad014.
- 720 Li, H., & Durbin, R. (2023). Genome assembly in the telomere-to-telomere era. *ArXiv*.
- 721 <http://arxiv.org/abs/2308.07877>
- 722 Liu, J., Shi, C., Shi, C.-C., Li, W., Zhang, Q.-J., Zhang, Y., Li, K., Lu, H.-F., Shi, C., Zhu, S.-T.,
- 723 Xiao, Z.-Y., Nan, H., Yue, Y., Zhu, X.-G., Wu, Y., Hong, X.-N., Fan, G.-Y., Tong, Y.,
- 724 Zhang, D., . . . Gao, L.-Z. (2020). The chromosome-based rubber tree genome provides
- 725 new insights into spurge genome evolution and rubber biosynthesis. *Mol. Plant*, 13(2),
- 726 336-350. <https://doi.org/10.1016/j.molp.2019.10.017>
- 727 Lofgren, L. A., Uehling, J. K., Branco, S., Bruns, T. D., Martin, F., & Kennedy, P. G. (2019).
- 728 Genome-based estimates of fungal rDNA copy number variation across phylogenetic
- 729 scales and ecological lifestyles. *Molecular Ecology*, 28(4), 721-730.
- 730 <https://doi.org/https://doi.org/10.1111/mec.14995>
- 731 Luo, J., Chen, J., Guo, W., Yang, Z., Lim, K.-J., & Wang, Z. (2022). Correction: Luo et al.
- 732 Reassessment of Annamocarya sinesis (Carya sinensis) Taxonomy through

- 733 Concatenation and Coalescence Phylogenetic Analysis. *Plants* 2022, 11, 52. *Plants*,  
734 11(23), 3282. <https://doi.org/10.3390/plants11233282>
- 735 Manni, M., Berkeley, M. R., Seppey, M., & Zdobnov, E. M. (2021). BUSCO: Assessing  
736 Genomic Data Quality and Beyond. *Current Protocols*, 1(12).
- 737 <https://doi.org/10.1002/cpz1.323>
- 738 Mansfeld, B. N., Boyher, A., Berry, J. C., Wilson, M., Ou, S., Polydore, S., Michael, T. P.,  
739 Fahlgren, N., & Bart, R. S. (2021). Large structural variations in the haplotype-resolved  
740 African cassava genome. *Plant J.*, 108(6), 1830-1848. <https://doi.org/10.1111/tpj.15543>
- 741 Matschiner, M., Böhne, A., Ronco, F., & Salzburger, W. (2020). The genomic timeline of cichlid  
742 fish diversification across continents. *Nature Communications*, 11(1), 5895.
- 743 <https://doi.org/10.1038/s41467-020-17827-9>
- 744 Minh, B. Q., Schmidt, H. A., Chernomor, O., Schrempf, D., Woodhams, M. D., von Haeseler, A.,  
745 & Lanfear, R. (2020). IQ-TREE 2: New Models and Efficient Methods for Phylogenetic  
746 Inference in the Genomic Era. *Mol. Biol. Evol.*, 37(5), 1530-1534.
- 747 <https://doi.org/10.1093/molbev/msaa015>
- 748 Misof, B., Liu, S., Meusemann, K., Peters, R. S., Donath, A., Mayer, C., Frandsen, P. B., Ware,  
749 J., Flouri, T., Beutel, R. G., Niehuis, O., Petersen, M., Izquierdo-Carrasco, F., Wappler,  
750 T., Rust, J., Aberer, A. J., Aspöck, U., Aspöck, H., Bartel, D., . . . Zhou, X. (2014).  
751 Phylogenomics resolves the timing and pattern of insect evolution. *Science*, 346(6210),  
752 763-767. <https://doi.org/10.1126/science.1257570>
- 753 Naranjo-Ortiz, M. A., & Gabaldón, T. (2020). Fungal evolution: cellular, genomic and metabolic  
754 complexity. *Biological Reviews*, 95(5), 1198-1232.
- 755 <https://doi.org/https://doi.org/10.1111;brv.12605>

- 756 Nasrallah, C. A., Mathews, D. H., & Huelsenbeck, J. P. (2010). Quantifying the Impact of  
757 Dependent Evolution among Sites in Phylogenetic Inference. *Systematic Biology*, 60(1),  
758 60-73. <https://doi.org/10.1093/sysbio/syq074>
- 759 Pisani, D. (2004). Identifying and removing fast-evolving sites using compatibility analysis: an  
760 example from the Arthropoda. *Systematic Biology*, 53(6), 978-989.
- 761 Ran, J.-H., Shen, T.-T., Wang, M.-M., & Wang, X.-Q. (2018). Phylogenomics resolves the deep  
762 phylogeny of seed plants and indicates partial convergent or homoplastic evolution  
763 between Gnetales and angiosperms. *Proc. Biol. Sci.*, 285(1881).  
764 <https://doi.org/10.1098/rspb.2018.1012>
- 765 Rangel, L. T., & Fournier, G. P. (2023). Fast-Evolving Alignment Sites Are Highly Informative  
766 for Reconstructions of Deep Tree of Life Phylogenies. *Microorganisms*, 11(10), 2499.  
767 <https://www.mdpi.com/2076-2607/11/10/2499>
- 768 Rautiainen, M., Nurk, S., Walenz, B. P., Logsdon, G. A., Porubsky, D., Rhie, A., Eichler, E. E.,  
769 Phillippe, A. M., & Koren, S. (2023). Telomere-to-telomere assembly of diploid  
770 chromosomes with Verkko. *Nat. Biotechnol.*, 41(10), 1474-1482.  
771 <https://doi.org/10.1038/s41587-023-01662-6>
- 772 Reuscher, S., Furuta, T., Bessho-Uehara, K., Cosi, M., Jena, K. K., Toyoda, A., Fujiyama, A.,  
773 Kurata, N., & Ashikari, M. (2018). Assembling the genome of the African wild rice  
774 Oryza longistaminata by exploiting synteny in closely related Oryza species. *Commun  
775 Biol.*, 1, 162. <https://doi.org/10.1038/s42003-018-0171-y>
- 776 Ronco, F., Matschiner, M., Böhne, A., Boila, A., Büscher, H. H., El Taher, A., Indermaur, A.,  
777 Malinsky, M., Ricci, V., Kahmen, A., Jentoft, S., & Salzburger, W. (2021). Drivers and

- 778 dynamics of a massive adaptive radiation in cichlid fishes. *Nature*, 589(7840), 76-81.
- 779 <https://doi.org/10.1038/s41586-020-2930-4>
- 780 Rosenberg, M. S., & Kumar, S. (2001). Incomplete taxon sampling is not a problem for  
781 phylogenetic inference. *Proceedings of the National Academy of Sciences*, 98(19),  
782 10751-10756. <https://doi.org/doi:10.1073/pnas.191248498>
- 783 Sahbou, A.-E., Iraqi, D., Mentag, R., & Khayi, S. (2022). BuscoPhylo: a webserver for Busco-  
784 based phylogenomic analysis for non-specialists. *Sci. Rep.*, 12(1), 17352.  
785 <https://doi.org/10.1038/s41598-022-22461-0>
- 786 Sayers, E. W., Beck, J., Bolton, E. E., Brister, J. R., Chan, J., Comeau, D. C., Connor, R.,  
787 DiCuccio, M., Farrell, C. M., Feldgarden, M., Fine, A. M., Funk, K., Hatcher, E.,  
788 Hoeppner, M., Kane, M., Kannan, S., Katz, K. S., Kelly, C., Klimke, W., . . . Sherry, S. T.  
789 (2024). Database resources of the National Center for Biotechnology Information.  
790 *Nucleic Acids Res*, 52(D1), D33-d43. <https://doi.org/10.1093/nar/gkad1044>
- 791 Sayers, E. W., Bolton, E. E., Brister, J. R., Canese, K., Chan, J., Comeau, D. C., Connor, R.,  
792 Funk, K., Kelly, C., Kim, S., Madej, T., Marchler-Bauer, A., Lanczycki, C., Lathrop, S.,  
793 Lu, Z., Thibaud-Nissen, F., Murphy, T., Phan, L., Skripchenko, Y., . . . Sherry, S. T.  
794 (2022). Database resources of the national center for biotechnology information. *Nucleic  
795 Acids Res.*, 50(D1), D20-D26. <https://doi.org/10.1093/nar/gkab1112>
- 796 Sayers, E. W., Bolton, E. E., Brister, J. R., Canese, K., Chan, J., Comeau, D. C., Farrell, C. M.,  
797 Feldgarden, M., Fine, A. M., Funk, K., Hatcher, E., Kannan, S., Kelly, C., Kim, S.,  
798 Klimke, W., Landrum, M. J., Lathrop, S., Lu, Z., Madden, T. L., . . . Sherry, S. T. (2023).  
799 Database resources of the National Center for Biotechnology Information in 2023.  
800 *Nucleic Acids Res*, 51(D1), D29-d38. <https://doi.org/10.1093/nar/gkac1032>

- 801 Schrempf, D., & Szöllősi, G. (2020). The sources of phylogenetic conflicts. *Phylogenetics in the*  
802 *genomic era*, 3-1. <https://hal.science/hal-02535482/>  
803 [https://hal.science/hal-02535482/file/chapter\\_3.1\\_Schrempf\\_Szollosi.pdf](https://hal.science/hal-02535482/file/chapter_3.1_Schrempf_Szollosi.pdf)  
804 Seppey, M., Manni, M., & Zdobnov, E. M. (2019). BUSCO: assessing genome assembly and  
805 annotation completeness. *Gene prediction: methods and protocols*, 227-245.  
806 Superson, A., & Battistuzzi, F. (2022). Exclusion of fast evolving genes or fast evolving sites  
807 produces different archaean phylogenies. *Molecular Phylogenetics and Evolution*, 170,  
808 107438.  
809 Susko, E., & Roger, A. (2021). Long Branch Attraction Biases in Phylogenetics. *Syst. Biol.*  
810 <https://doi.org/10.1093/sysbio/syab001>  
811 Timilsena, P. R., Wafula, E. K., Barrett, C. F., Ayyampalayam, S., McNeal, J. R., Rentsch, J. D.,  
812 McKain, M. R., Heyduk, K., Harkess, A., Villegente, M., Conran, J. G., Illing, N.,  
813 Fogliani, B., Ané, C., Pires, J. C., Davis, J. I., Zomlefer, W. B., Stevenson, D. W.,  
814 Graham, S. W., . . . dePamphilis, C. W. (2022). Phylogenomic resolution of order- and  
815 family-level monocot relationships using 602 single-copy nuclear genes and 1375  
816 BUSCO genes [Original Research]. *Frontiers in Plant Science*, 13.  
817 <https://doi.org/10.3389/fpls.2022.876779>  
818 Van Damme, K., Cornetti, L., Fields, P. D., & Ebert, D. (2022). Whole-genome phylogenetic  
819 reconstruction as a powerful tool to reveal homoplasy and ancient rapid radiation in  
820 waterflea evolution. *Syst. Biol.*, 71(4), 777-787. <https://doi.org/10.1093/sysbio/syab094>  
821 Venditti, C., Meade, A., & Pagel, M. (2006). Detecting the node-density artifact in phylogeny  
822 reconstruction. *Syst. Biol.*, 55(4), 637-643. <https://doi.org/10.1080/10635150600865567>

- 823 Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski,  
824 E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J.,  
825 Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., . . . SciPy,  
826 C. (2020). SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature  
methods*, 17(3), 261-272. <https://doi.org/10.1038/s41592-019-0686-2>
- 827 Wheeler, T. J., & Eddy, S. R. (2013). nhmmer: DNA homology search with profile HMMs.  
828 *Bioinformatics*, 29(19), 2487-2489. <https://doi.org/10.1093/bioinformatics/btt403>
- 830 Wighard, S. S., Athanasouli, M., Witte, H., Rödelsperger, C., & Sommer, R. J. (2022). A New  
831 Hope: A hermaphroditic nematode enables analysis of a recent whole genome duplication  
832 event. *Genome Biol. Evol.*, 14(12). <https://doi.org/10.1093/gbe/evac169>
- 833 Yan, Z., Smith, M. L., Du, P., Hahn, M. W., & Nakhleh, L. (2021). Species tree inference  
834 methods intended to deal with incomplete lineage sorting are robust to the presence of  
835 paralogs. *Syst. Biol.*, 71, 367-381. <https://doi.org/10.1093/sysbio/syab056>
- 836 Yang, Z. (1994). Maximum likelihood phylogenetic estimation from DNA sequences with  
837 variable rates over sites: approximate methods. *Journal of Molecular evolution*, 39, 306-  
838 314.
- 839 Yang, Z. (1996). Among-site rate variation and its impact on phylogenetic analyses. *Trends in  
840 ecology & evolution*, 11(9), 367-372.
- 841 Yang, Z. (2006). *Computational molecular evolution*. OUP Oxford.
- 842 Yu, H., Lin, T., Meng, X., Du, H., Zhang, J., Liu, G., Chen, M., Jing, Y., Kou, L., Li, X., Gao, Q.,  
843 Liang, Y., Liu, X., Fan, Z., Liang, Y., Cheng, Z., Chen, M., Tian, Z., Wang, Y., . . . Li, J.  
844 (2021). A route to de novo domestication of wild allotetraploid rice. *Cell*, 184(5), 1156-  
845 1170.e1114. <https://doi.org/10.1016/j.cell.2021.01.013>

846 Zhang, C., & Mirarab, S. (2022). ASTRAL-Pro 2: ultrafast species tree reconstruction from  
847 multi-copy gene family trees. *Bioinformatics*, 38(21), 4949-4950.  
848 <https://doi.org/10.1093/bioinformatics/btac620>

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863 **Figure Legends**

864 **Figure 1. BUSCO database statistics.** **A.** Genome assemblies for new genera and species are  
865 growing linearly for plants and fungi and rapidly for animals, especially in recent years. **B.**  
866 BUSCO statistics vary for plants, fungi and animals. The fraction of single-copy and duplicated  
867 genes are complementary. More duplications are observed in plants and less variation is notable  
868 for the fungi. **C.** Some taxonomic groups, such as ascomycetes and insects are better represented  
869 in NCBI genome. Assemblies from bulk genome sequencing projects with relatively low cost per  
870 genome appear as a stretch with lower BUSCO completeness. Duplicated fractions are more  
871 prominent in plants owing primarily to higher duplication rates and greater incidence of  
872 polyploidy.

873

874 **Figure 2. Higher rates are more informative and produce better phylogenies overall.** **A.**  
875 Taxonomic concordance across 13 rate profiles and 20 alignment lengths. Sites evolving at  
876 higher rates and longer alignments share more agreement with taxonomic groupings. **B.** Most  
877 families in 4 groups are resolved as monophyletic in most trees whereas a smaller number of  
878 families are more sporadic and appear to be monophyletic either randomly or under specific rates.  
879 **C.** Ve, Eu, Ba, As and Ar represent Vertebrate, Eudicots, Basidiomycota, Ascomycota and  
880 Arthropoda lineages respectively. Each vertical bar is a unique family. With few exceptions,  
881 families are more likely to be found monophyletic at greater rates and sites. **D.** Increasing rates  
882 have a greater effect on tree concordance relative to increasing sites. **E.** Under optimum tree  
883 search conditions, tree likelihoods correlate with taxonomic agreement. **F.** Ascomycota  
884 represents the fungal clade which shows increased variance in tree sets and is less responsive to

885 rate and site adjustments. **G.** For the Arthropoda lineage, increasing rates and sites increases  
886 concordance and reduces tree set variance. **H.** Differences between concatenated and coalescent  
887 species trees are marginal.

888 **Figure 3. Removal of erratic BUSCO genes reduces BUSCO misidentification rates. A.**  
889 BUSCO genes are misidentified at different rates in different lineages. Median fraction of false  
890 identification is around 15% for most plants and vertebrates, but much lower in fungi. **B.** Only  
891 considering our Curated set of BUSCO genes (CUSCOs) markedly reduces false hits in some  
892 lineages.

893

894 **Figure 4. Misidentification events are weighted more towards the identity of the gene**  
895 **rather than assembly and false hits correlate most with assembly complexity and gene**  
896 **content. A.** A graph of gene quantiles against assembly quantiles for false hit counts shows that  
897 the majority of assemblies show some false gene hits but the gene quantiles rise more shapely. **B.**  
898 Considering only the curated BUSCO set shifts the assembly quantiles at the lower range  
899 towards the genes. CUSCO genes are misidentified in far fewer assemblies and do not show  
900 assembly preference. **C.** False identification rates correlate most with the number of miniProt  
901 hits (MPH) and mean BUSCO copy counts (Inflation). Moderate correlation to the log of  
902 assembly size is also observed. **D.** Two example blocks of 8 genes conserved beyond the species  
903 level for eudicots (top) and vertebrates (bottom) showing misidentified/remnant BUSCO genes  
904 in syntenic order. **E.** CUSCO and MUSCO proportions for syntenic doublets with 0, 1 and 2  
905 remnant genes. Remnant proportions gradually recede for CUSCOs, but rise back up in remnant  
906 doublets for MUSCOs.

907 **Figure 5. BUSCO syntenic distance offers greater contrast than BUSCO content, decays**  
908 **exponentially with phylogenetic distance and serves as a robust metric to compare closely**  
909 **related assemblies. A.** Boxplot showing differences in BUSCO completeness and BUSCO  
910 syntenic distance between 1035 sets of assemblies that vary in quality. **B.** General trend and  
911 histogram of BUSCO syntenic distance and BUSCO completeness differences. BUSCO syntenic  
912 differences can offer far greater contrast. **C.** Exponential decay and curve function of BUSCO  
913 syntenic similarity for Arthropoda, Vertebrata, Ascomycota, and Basidiomycota lineages **D.**  
914 Exponential decay and curve function for Liliopsida and Eudicots lineages **E.** Eight highly  
915 fragmented *Mus musculus* assemblies compared against a highly contiguous assembly through  
916 BUSCO syntenic distance and quality assembly metrics.  
917

918 **Figure 6. Phylogenetic and syntenic information improves assembly assessment. A.**  
919 BUSCOs in chromosome 1 of *Oryza longistaminata* and *O. meyeriana* assemblies are less  
920 syntenic to sister taxa. A chromosomal translocation event from chromosome 3 to 1 in *O. alta*  
921 subgenome C is also visualized. **B.** Assessment of an improved *O. longistaminata* assembly  
922 reveals that BUSCO genes were either misidentified or contigs were scaffolded poorly in the  
923 inferior assembly. Chromosomes 1 and 7 are visualized at the top and bottom respectively.

924

925

926

927

928

929 **Table 1. BUSCO and CUSCO misidentification rates**

Lineage	BUSCO completeness (mean)	BUSCO completeness (SD)	CUSCO completeness (mean)	CUSCO completeness (SD)	BUSCO false hits (mean)	BUSCO false hits (SD)	CUSCO false hits (mean)	CUSCO false hits (SD)
Viridiplantae	91.88	15.90	90.67	18.37	11.35	8.69	7.52	8.88
Liliopsida	87.30	24.47	90.06	19.54	12.60	8.01	5.90	7.70
Eudicots	92.37	15.00	91.81	16.52	13.34	7.30	6.35	7.40
Fungi	94.66	12.99	94.93	12.91	2.86	4.01	1.64	4.02
Ascomycota	96.74	6.92	96.84	7.13	2.25	2.97	0.69	3.01
Basidiomycot a	95.59	8.08	95.29	9.27	3.00	3.98	2.02	3.86
Metazoa	83.41	23.86	82.32	25.47	6.02	7.44	4.00	6.69

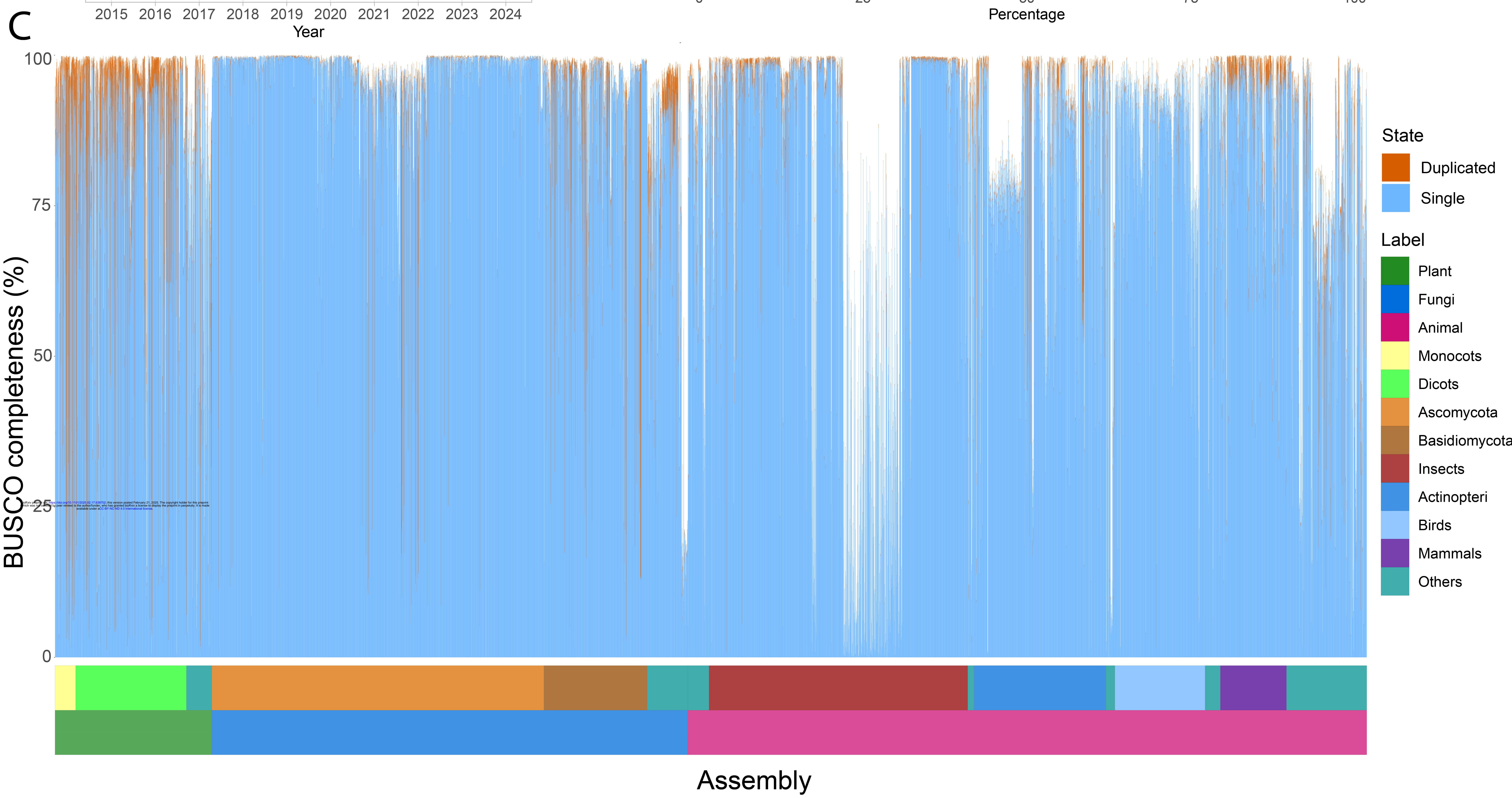
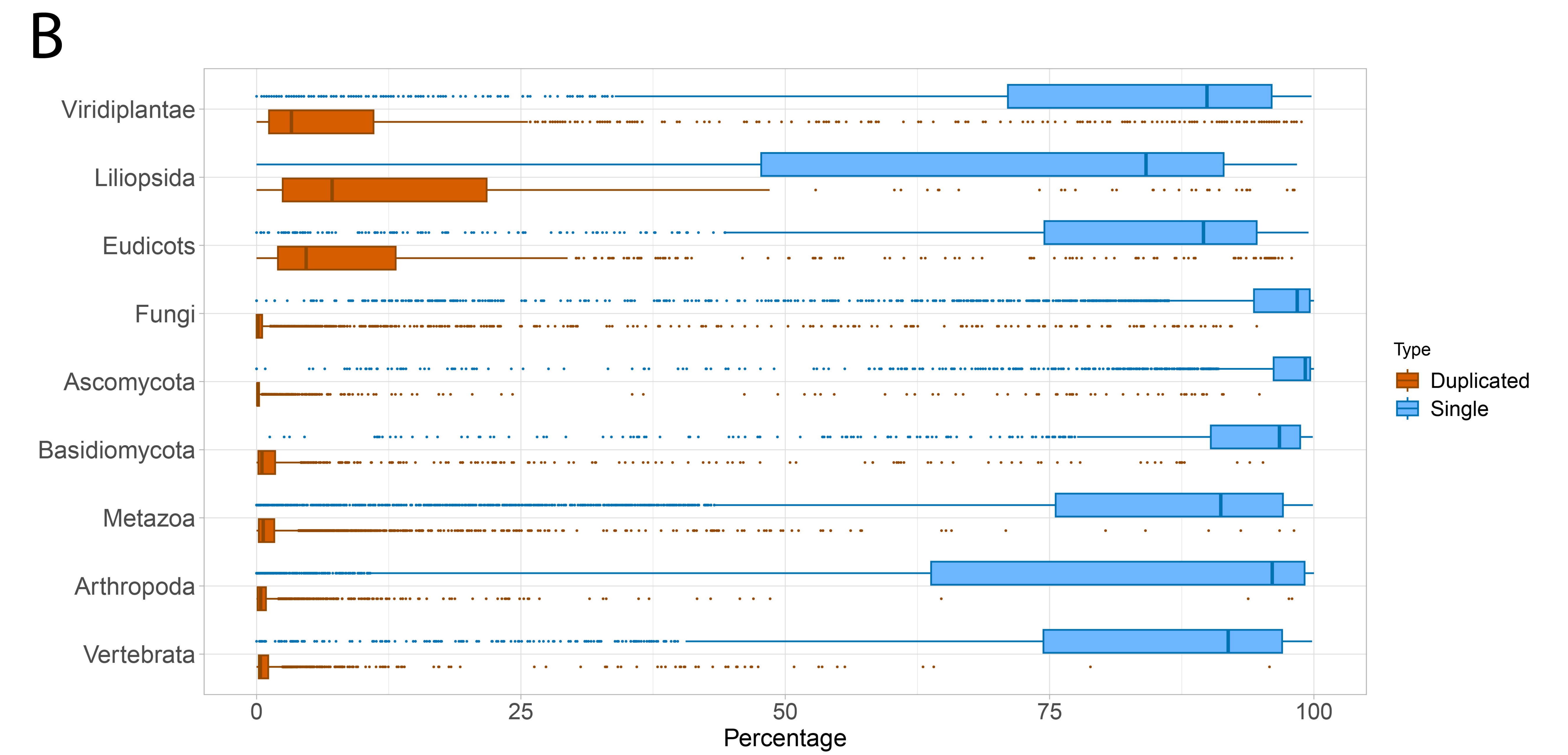
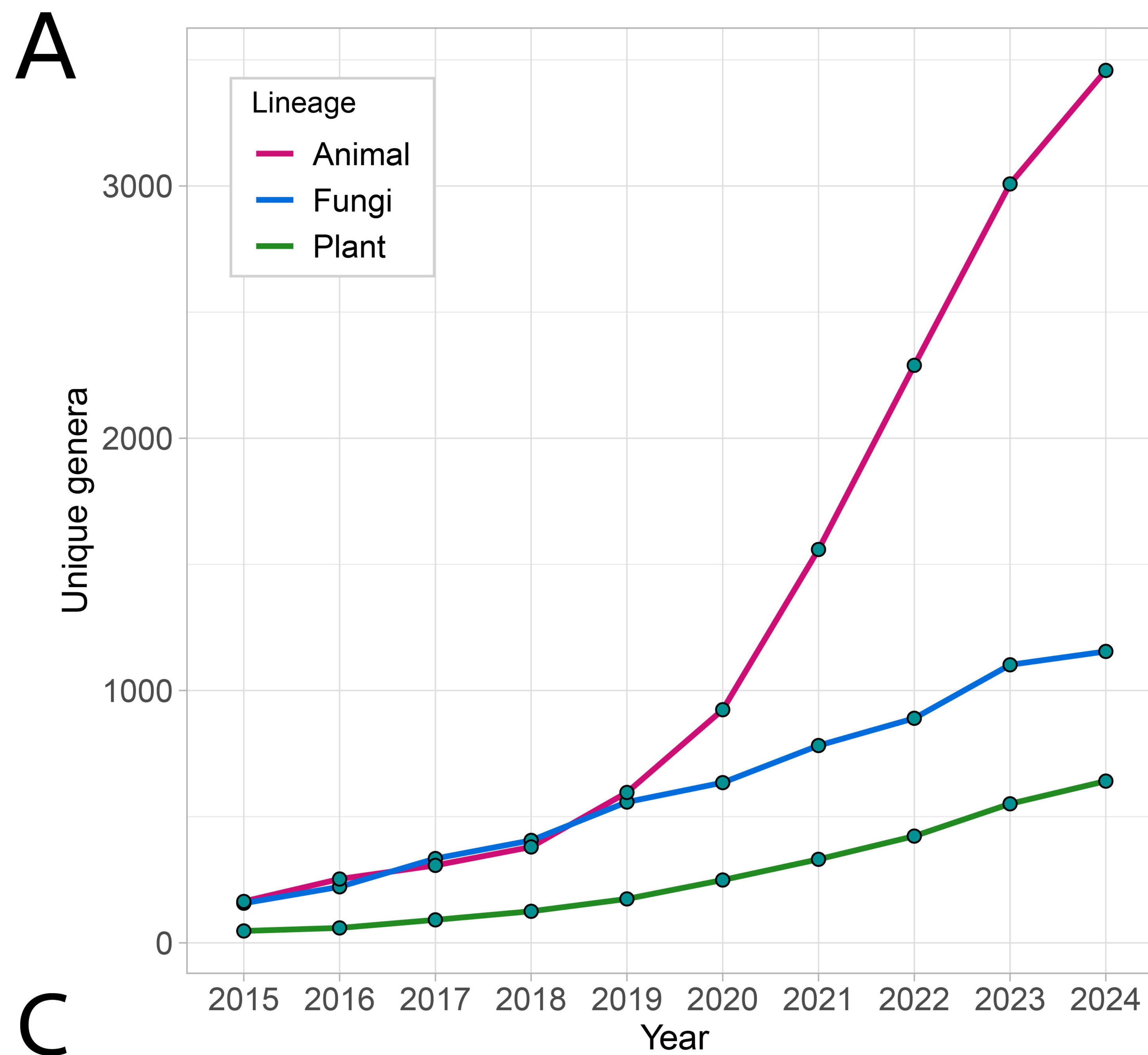
Arthropoda	81.86	29.35	78.86	32.58	4.55	3.92	1.92	3.75
Vertebrata	85.09	19.63	84.72	20.07	9.57	5.12	2.17	3.74
Chlorophyta	86.09	14.29	85.19	16.37	8.12	6.80	3.57	6.09

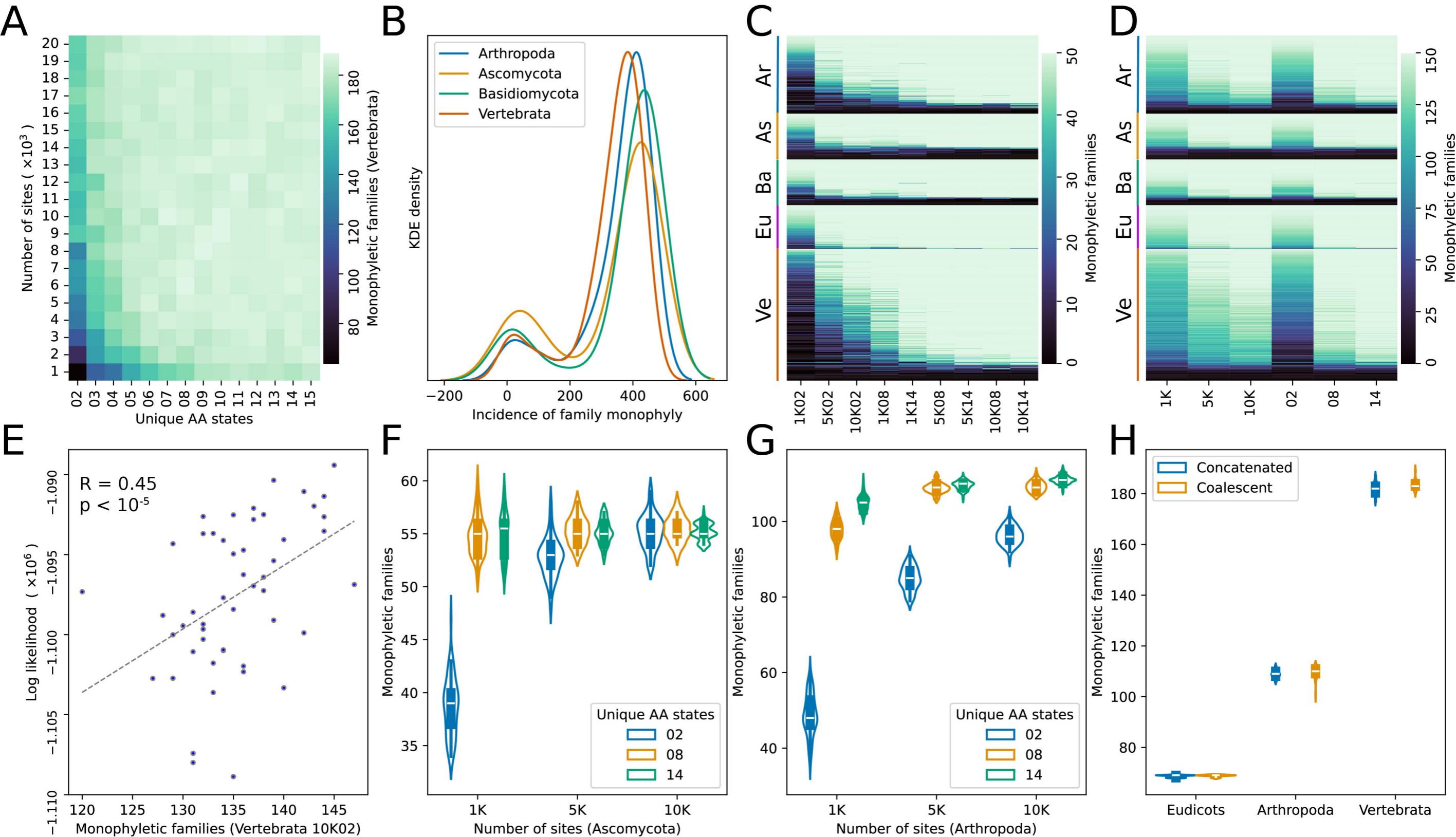
930 \*all numbers are in percentiles

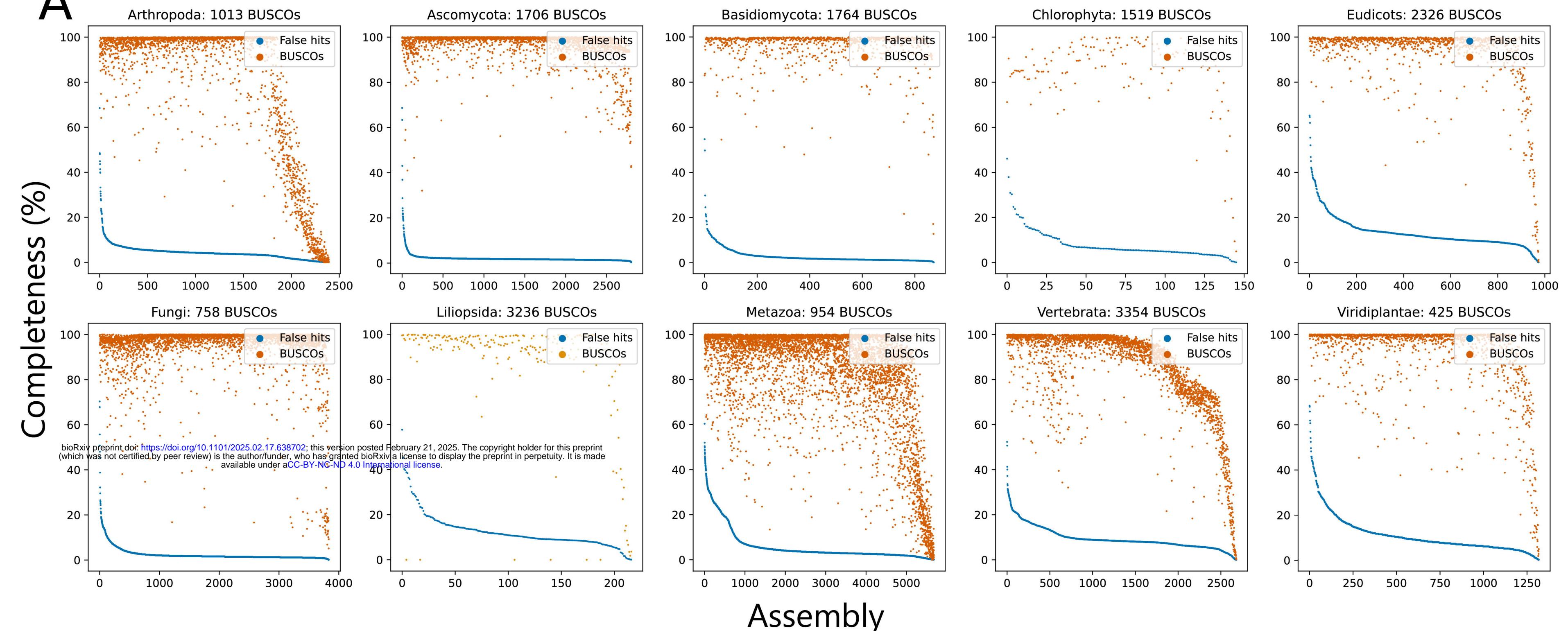
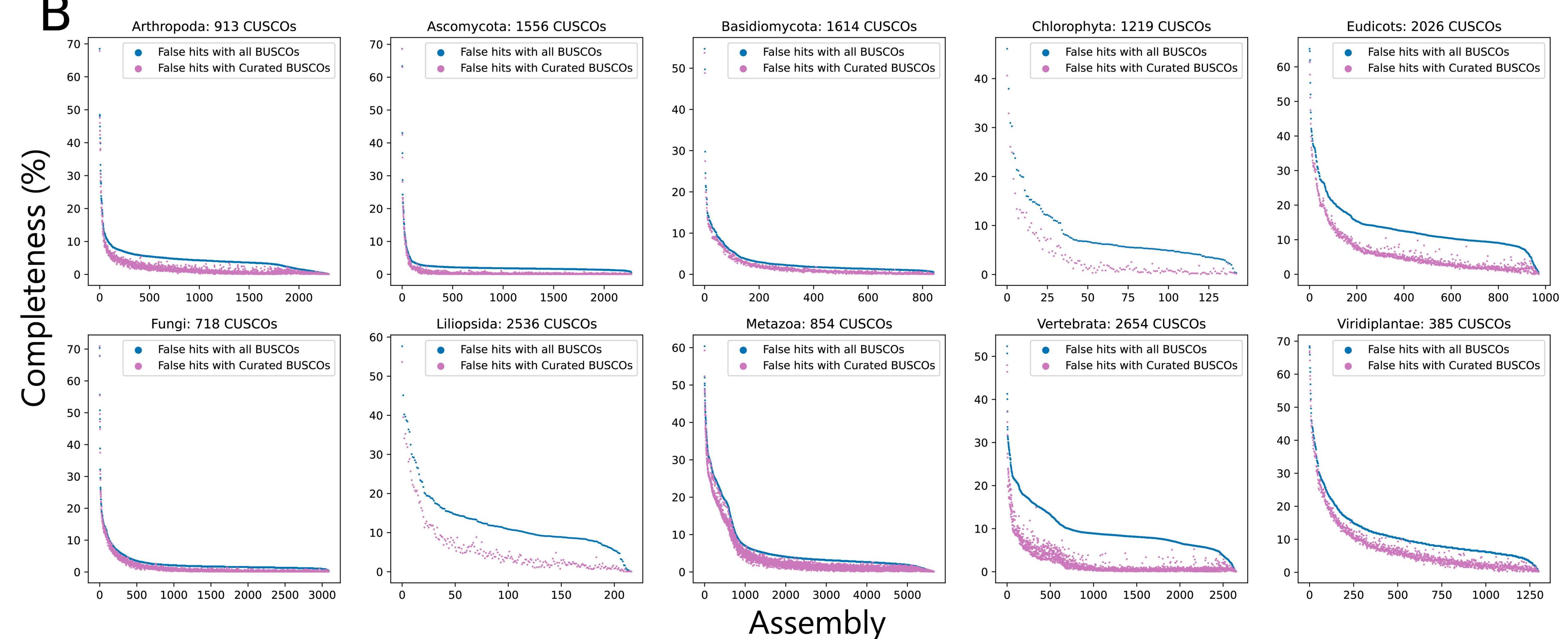
931

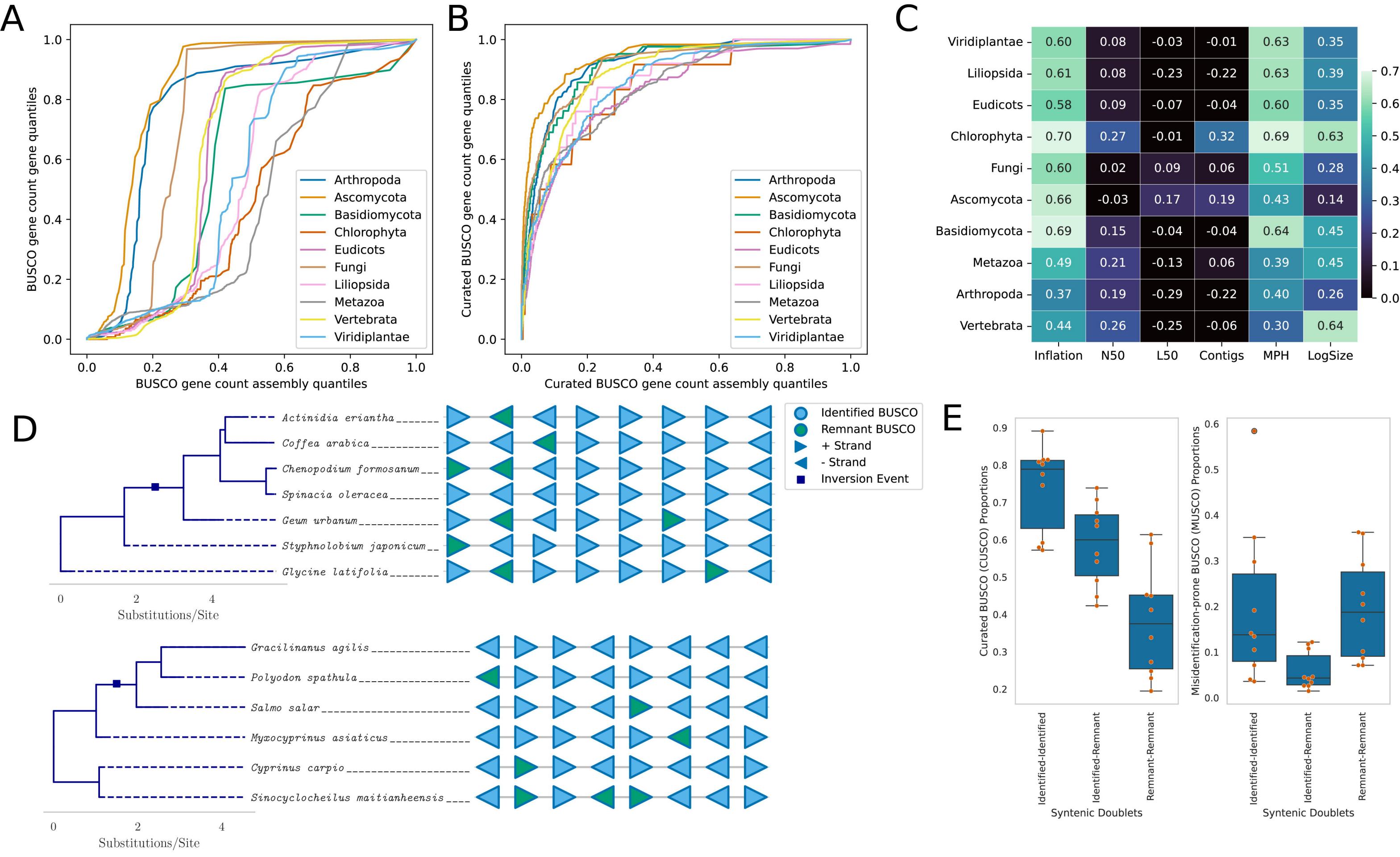
932

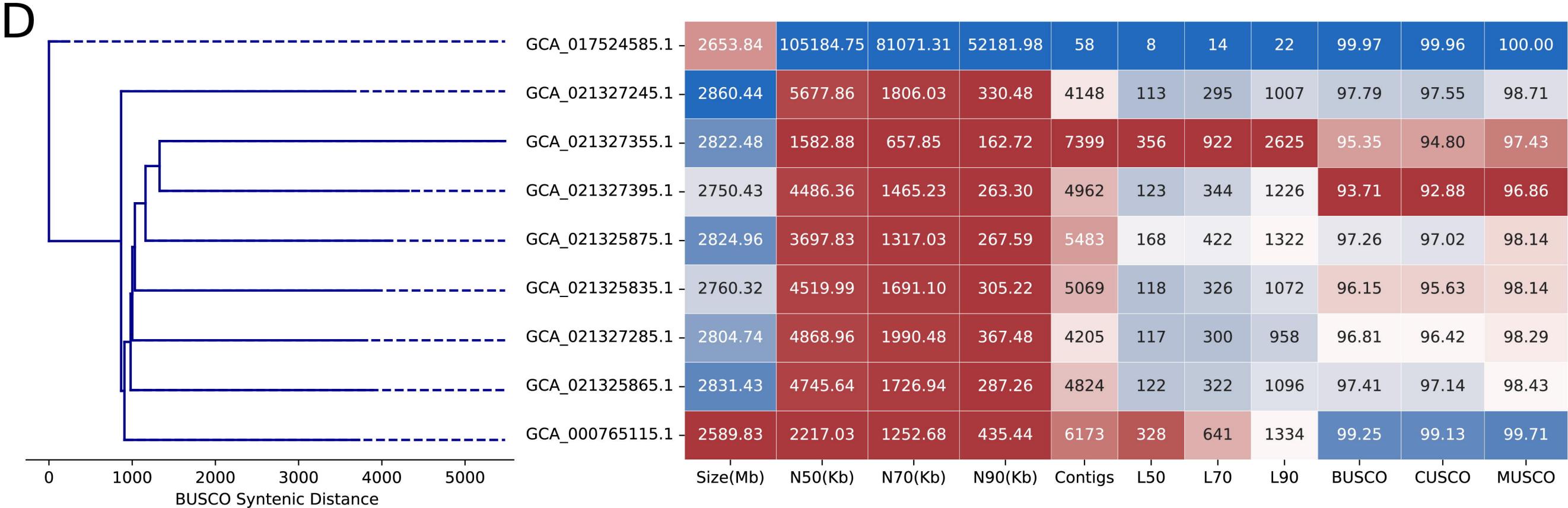
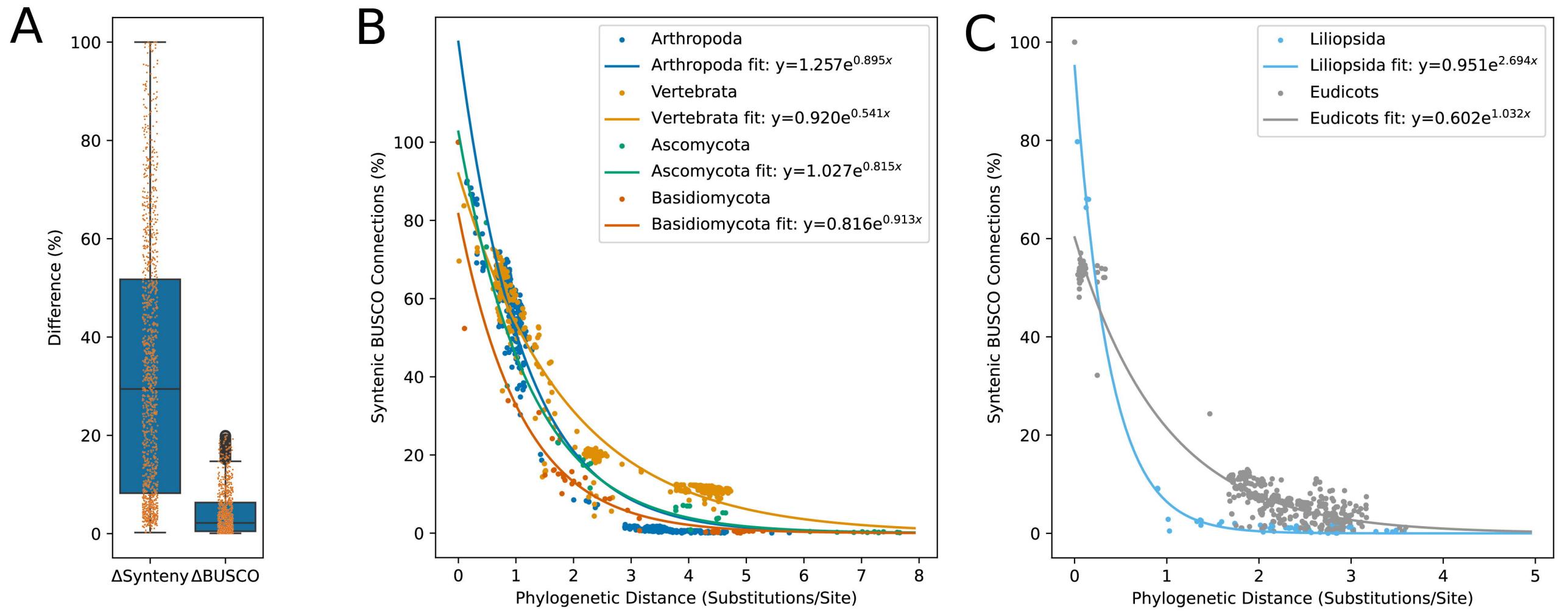
933

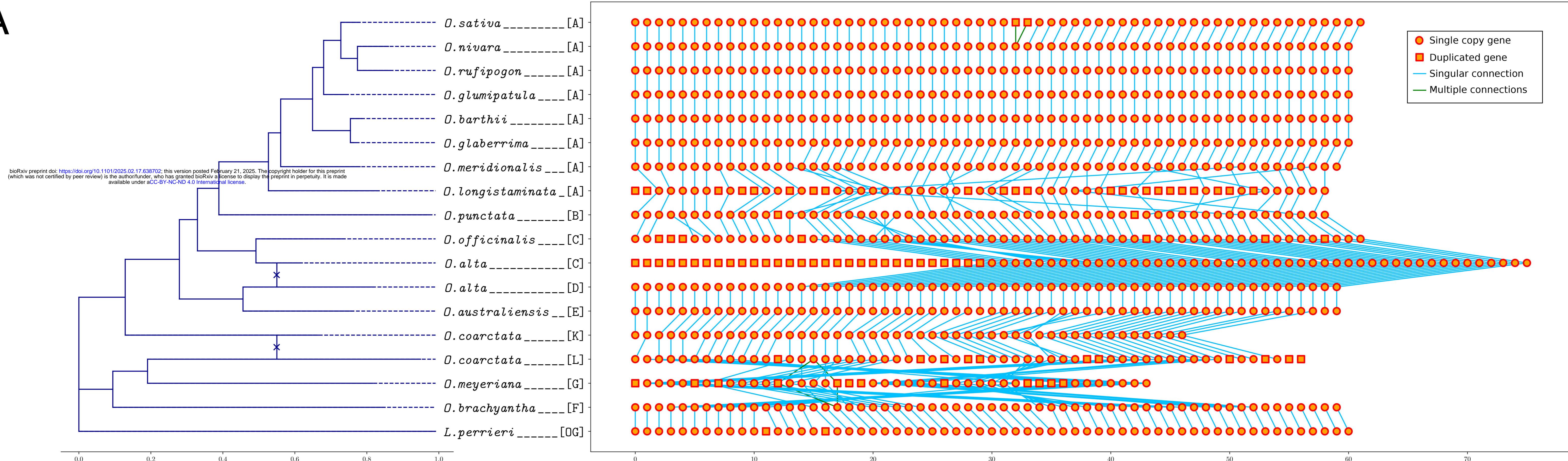




**A****B**





**A****B**