

1

2

3

4

5

6

7

8

Phylogenetic profiling in eukaryotes: The effect of  
species, orthologous group, and interactome selection on  
protein interaction prediction

9

10

11

Eva S. Deutekom<sup>1</sup>, Teunis J.P. van Dam<sup>1</sup>, Berend Snel<sup>1\*</sup>

12

<sup>1</sup> Theoretical Biology and Bioinformatics, Department of Biology, Science faculty, Utrecht  
University, Utrecht, The Netherlands

13

14

\* Corresponding author

15

Email: [B.Snel@uu.nl](mailto:B.Snel@uu.nl)

16

## 17 **Abstract**

18           Phylogenetic profiling in eukaryotes is of continued interest to study and predict  
19 the functional relationships between proteins. This interest is likely driven by the  
20 increased number of available diverse genomes and computational methods to infer  
21 orthologies. The evaluation of phylogenetic profiles has mainly focussed on reference  
22 genome selection in prokaryotes. However, it has been proven to be challenging to obtain  
23 high prediction accuracies in eukaryotes. As part of our recent comparison of orthology  
24 inference methods for eukaryotic genomes, we observed a surprisingly high performance  
25 for predicting interacting orthologous groups. This high performance, in turn, prompted  
26 the question of what factors influence the success of phylogenetic profiling when applied  
27 to eukaryotic genomes.

28           Here we analyse the effect of species, orthologous group and interactome  
29 selection on protein interaction prediction using phylogenetic profiles. We select species  
30 based on the diversity and quality of the genomes and compare this supervised selection  
31 with randomly generated genome subsets. We also analyse the effect on the performance  
32 of orthologous groups defined to be in the last eukaryotic common ancestor of eukaryotes  
33 to that of orthologous groups that are not. Finally, we consider the effects of reference  
34 interactome set filtering and reference interactome species.

35           In agreement with other studies, we find an effect of genome selection based on  
36 quality, less of an effect based on genome diversity, but a more notable effect based on  
37 the amount of information contained within the genomes. Most importantly, we find it is  
38 not merely selecting the correct genomes that is important for high prediction  
39 performance. Other choices in meta parameters such as orthologous group selection, the

40 reference species of the interaction set, and the quality of the interaction set have a much  
41 larger impact on the performance when predicting protein interactions using phylogenetic  
42 profiles. These findings shed light on the differences in reported performance amongst  
43 phylogenetic profiles approaches, and reveal on a more fundamental level for which types  
44 of protein interactions this method has most promise when applied to eukaryotes.

45

## 46 **Introduction**

47 The post-genomic era has provided us with a wealth of eukaryotic genomes of  
48 diverse and underrepresented phyla [1]. Most of the sequences in these new genomes  
49 are without precise function assignment and a challenge remains, protein function and  
50 interaction discovery [2,3]. Computational approaches that are available for large scale  
51 analyses of protein function and interactions include phylogenetic profiling. Phylogenetic  
52 profiling uses correlations of the presences and absences of groups of orthologous  
53 proteins (orthologous groups) across a set of species [4]. Phylogenetic profiling is a  
54 seemingly straightforward method proven to be a valuable alternative resource for  
55 studying functional relationships between proteins. Recently the method has played an  
56 integral part in identifying the cellular functional role of CENATAC that is a key player in  
57 a rare aneuploidy condition in humans [5], identifying eukaryotic reproduction genes [6],  
58 and identifying eukaryotic novel recombination repair genes [7].

59 The information in the phylogenetic profiles given by presence and absence  
60 patterns, are shaped by a diverse range of evolutionary forces. These forces include  
61 horizontal gene transfer, secondary endosymbiosis and gene loss. The method relies on

62 the principle that proteins with a similar profile indicate that the proteins co-evolved due  
63 to them belonging to the same functional pathway or complex. There are countless  
64 observations that co-occurring proteins tend to interact [8–10]. Phylogenetic profiling can  
65 be a powerful tool for function prediction. By comparing, or even clustering, profiles of  
66 proteins with unknown function to those with known function enables us to infer to which  
67 complexes or functional pathways the uncharacterized proteins likely belong and, in turn,  
68 infer their function.

69 Multiple studies have shown the effectiveness of phylogenetic profiling in large  
70 scale analyses of eukaryotes [9,11], which has become possible with the large increase  
71 in genomic data and computational methods to (automatically) infer orthologies [12–14]  
72 or cluster genes [11]. However, benchmarking and analysing the performance of large-  
73 scale phylogenetic profiling has been limited to prokaryotes, for which good performance  
74 can be obtained when predicting protein interactions [15–18]. The performance  
75 decreases when benchmarking is done solely with eukaryotes or when eukaryotes are  
76 combined with prokaryotes [15,16]. Likely, the performance reduction is caused by the  
77 different forces driving eukaryotic genome evolution, compared to the dynamic pan  
78 genomes of prokaryotes where the interplay of rampant horizontal gene transfer of  
79 operons and loss of genes that create highly informative patterns.

80 We recently obtained a high protein interaction prediction performance in a large  
81 set of eukaryotes in the context of evaluating a diverse set of orthologous group inference  
82 methods [19]. The surprisingly high prediction performance only marginally depended on  
83 the orthologous group inference methods (which was the focus of the study), suggesting  
84 that its cause could be any of the other underlying choices. Therefore, a more elaborate

85 analysis of the choices made for phylogenetic profiling is warranted. Here we evaluate in-  
86 depth the meta parameters influencing the performance of phylogenetic profiles in  
87 eukaryotes.

88 Multiple studies have understandably focused their analysis on reference genome  
89 selection or the amount of genomes/data needed to increase prediction performance [16–  
90 18,20,21]. Besides genome diversity and quality, we analyse orthologous groups and  
91 reference interactome selection. Our results demonstrate that an interplay of biological  
92 and technical aspects influence phylogenetic profiling. Most importantly, our results show  
93 that prediction performance is influenced not only by genome selection but mostly by  
94 orthologous and interactome selection.

## 95 **Results**

96 Each results section describes the analyses of meta parameters encompassing  
97 five main concepts: genome quality, genome diversity, performance directed genome  
98 selection, orthologous group selection, and reference interactome selection. To rule out  
99 any orthology specific issues, we performed the analyses using two orthology inference  
100 methods, Sonicparanoid [13] and Broccoli [14]. Sonicparanoid performed the best in our  
101 previous study using phylogenetic profiles for protein interaction prediction [19]. We chose  
102 Sonicparanoid as the primary method, while broccoli serves to determine to what extent  
103 the results are contingent on a specific orthology method. The results for Broccoli can be  
104 found in Supplementary figures and are overall in agreement with the results of  
105 Sonicparanoid.

## 106 **1. Lesser quality genomes have more effect on the prediction** 107 **performance than higher-quality genomes**

108 Phylogenetic profiles can be noisy due to multiple technical reasons, such as gene  
109 annotation and genome assembly errors. Consequently, the quality of genomes can be  
110 an essential factor, as profiles with a lot of noise would be akin to noisy gene expression  
111 or protein interaction measurements. We expect noisy genomes to give much weaker  
112 prediction performance. Given this expectation, the first meta parameter assessed was  
113 genome quality. We calculated genomes quality using two independent metrics, BUSCO  
114 [22] and one of our design (Supplementary figures and Methods and Materials). For  
115 clarity, we use only the BUSCO metric in the main text since both metrics generally agree  
116 with each other.

117 The BUSCO metric assesses genome completeness based on the (in our case)  
118 absence of single-copy orthologs that are highly conserved among eukaryotic species.  
119 The absences of these orthologs can result from incomplete draft genomes or false  
120 negatives in gene prediction, which in both cases leads to false absences of orthologs.  
121 We selected 50 high-quality genomes with the lowest BUSCO values, i.e., genomes with  
122 the least number of unexpected absences. We also selected 50 lower quality genomes  
123 with the highest BUSCO values, i.e., genomes with the most number unexpected  
124 absences (Fig 1.A.). We compared the quality filtered genome sets with 1000 randomly  
125 generated genome sets of 50 genomes each to see if quality-based selection differs from  
126 any random sampling of genomes.

127

128           **Fig 1. Lesser quality genomes have more impact on protein interaction**  
129 **prediction performance. A.** BUSCO absences as a function of retained Last Eukaryotic  
130 Common Ancestor (LECA) orthologous groups in different species. Filled data points are  
131 the selected genomes for the prediction accuracy calculations. **B.** Receiver-operator  
132 Curve of two species sets (n = 50) with the most and least BUSCO absences. The inset  
133 gives the Area Under the Curve (AUC) values compared with the random backdrop of  
134 1000 random species sets (violin plot) and the initial species set (teal diamond).

135  
136           The results show that the performance using the highest quality genomes with the  
137 least suspect absences falls within the distribution of random genome prediction  
138 performance (AUC: 0.765). In contrast, the lower quality genomes fall below the  
139 distribution of random prediction performance (AUC: 0.748) (inset Fig 1.B.). This suggests  
140 that it is more beneficial to filter out lesser-quality genomes than it is to select for high-  
141 quality genomes. This result is consistent between two independent scores of genome  
142 quality (S1 Fig).

143           With these results, it seems prudent to select genomes only based on quality when  
144 applying phylogenetic profiles. However, there is an inherent bias between genome  
145 quality and phylogenetic distribution (Fig 1.A). For instance, eukaryotes belonging to the  
146 Opisthokonta supergroup have overall lower BUSCO absences, biasing the selection of  
147 good genomes towards one eukaryotic supergroup. *A priori*, species diversity seems  
148 another meta parameter in genome selection with potential impact. In the next section,  
149 we will look at the diversity of species and how that influences phylogenetic profiling.

## 150 **2. Genome diversity has little effect on prediction** 151 **performance in eukaryotes**

152 The diversity of species plays a role in the performance of phylogenetic profiles in  
153 prokaryotes [16]. We also expect high species diversity to improve how informative  
154 profiles are by giving high-resolution information on how genes co-evolve in different  
155 organisms. More species diversity allows to maximally discern the effect of evolutionary  
156 forces shaping co-evolving proteins, which might not be apparent in, e.g., an animal only  
157 data set. There will be no discernible and informative phylogenetic pattern in a  
158 homogeneous species set where most ancestral protein complexes are not frequently  
159 lost. A previous study showed that the maximum phylogenetic diversity in Bacteria gives  
160 the best predictive performance [18]. Here we want to test how maximal and minimal  
161 diversity affects prediction performance in eukaryotes.

162 We analysed the impact of eukaryotic diversity by selecting two sets of 50  
163 genomes, one containing the most similar species (Fig 2.A.) and the other the most  
164 diverse species (Fig 2.B.) from our initial species set. The (dis)similarity was measured  
165 using an iterative all-vs-all comparison using the cosine distance between genomes and  
166 their orthologous group content. We started with the most diverse or similar species pairs  
167 and iteratively added to this set the species with the highest (dis)similarity until we  
168 obtained 50 genomes (Materials and Methods). We recalculated the protein-interaction  
169 prediction performance for both these sets. The prediction performance is lower than the  
170 initial set for both sets, but not worse than any randomly selected genome sets (AUC:  
171 0.760 for the dissimilar set and AUC: 0.764 for the similar set) (Fig 2. C. inset).

172



173           **Fig 2. Both high and low diversity sets have little impact on protein**  
174 **interaction prediction performance. A.** The most similar species form more clusters  
175 and are overall more similar to each other. **B.** The most diverse species show no  
176 clustering and are overall less similar to each other. **C.** Receiver-operator Curve of two  
177 species sets (n = 50) with the most diverse and most similar species. The inset gives the  
178 Area Under the Curve (AUC) values compared with the random backdrop of 1000 random  
179 species sets (violin plot) and the initial species set (green diamond).

180

181           Similar species will naturally show more cohesion in profiles, with little separation  
182 of protein co-evolution. Highly diverse species will naturally show more discordance, with  
183 little information left to see protein co-evolution. In both cases, there will not be a gene-  
184 specific signal. A combination of the two should give good separation of actual co-  
185 evolving genes. Together with the effects of genome quality, genome diversity can be an  
186 important factor for the performance of phylogenetic profiles. However, the interplay  
187 between these two factors is complex, and as we previously determined, high genome  
188 quality corresponds with lower phylogenetic diversity. To look into this further, we  
189 investigate the influence of single genomes on predictive performance.

190 **3. Single influential genomes and their combined effect on**  
191 **prediction performance reveal the importance of the type of**  
192 **information in the profiles**

193           Diversity and quality both impact performance and we expect it to have a combined  
194 influence on phylogenetic profiling. Instead of a priori selecting genomes based on a

195 measure for each of these two criteria, we can also objectively evaluate the prediction  
196 performance by removing genomes from the initial species set one-by-one (Fig 3. A.).  
197 Genomes that decrease prediction performance when removed from the initial set we can  
198 consider as advantageous to phylogenetic profiling, while genomes that increase  
199 prediction performance when removed from the initial set we can consider as  
200 disadvantageous to phylogenetic profiling. We selected the top 50 advantageous and top  
201 50 disadvantageous genomes to see whether these genomes together in their respective  
202 sets also influence the prediction performance.

203

204 **Fig 3. Influential genomes and their combined effect. A.** Recalculated Area  
205 Under the Curve (AUC) values when a single species is removed from the initial species  
206 set. Genomes that increase the AUC value when removed can be considered  
207 disadvantages compared to the initial set when predicting protein interactions with  
208 phylogenetic profiles. Genomes that decrease the AUC value when removed can be  
209 considered advantageous for predicting protein interactions. Top 50 advantageous and  
210 top 50 disadvantageous genomes shown with the black fill in the scatter plot. **B.** Receiver-  
211 operator Curve of two species sets (n=50) with the most advantageous and  
212 disadvantageous genomes. The inset gives the Area Under the Curve (AUC) values  
213 compared with the random backdrop of 1000 random species sets (violin plot) and the  
214 initial species set (green diamond). **C.** Comparison of the counts (histogram) and kernel  
215 density estimates (line plot) of (I) illogical absence ratios (illogical absences divided by  
216 total interaction absences (co-absences + illogical absences)), (II) present interactions,  
217 (III) the cosine distance to human, and (IV) total shared orthologous groups with human.

218

219           For both the advantageous and disadvantageous set we can see a large difference  
220 in prediction performance (Fig 3. B.) and a larger difference than the selection based on  
221 the measures for either quality or diversity. With the advantageous genome set, the  
222 performance increases (AUC: 0.801). In contrast, for the disadvantageous set the  
223 performance drops (AUC: 0.730). Both values fall well outside the distribution of 1000  
224 randomly generated genome subsets.

225           Although a large cumulative effect on performance because we used the genomes'  
226 performance to select the genomes, it is still very interesting to see what these genomes  
227 share if it is not quality or diversity. We therefore examined the role of different genomes  
228 in these genome sets. Comparison of a large number of factors (S5 Fig) revealed that  
229 that the difference in prediction performance of the advantageous and disadvantageous  
230 genome sets is related to the (human) interactions retained in the genomes (Fig 3. C. and  
231 S5. A. Fig). The illogical absence ratios and the complete interactions present (or co-  
232 presences) (Fig 3. C. I & II) show intermediate values for the disadvantageous genome  
233 set. At the same time, these values are either high or low for the advantageous genome  
234 set.

235           We can also directly relate the differences between these genome sets to how  
236 close the genomes in the sets are to the human genome. The cosine distance and the  
237 shared orthologous groups of the genomes with the human genome (Fig 3.C. III & IV)  
238 show intermediate values for the disadvantageous set, while the values are either high or  
239 low for the advantageous genome set. For the orthologous groups inferred by Broccoli

240 this signal is even more pronounced (S5. B. Fig). A surprising finding is that the  
241 advantageous set contains numerous parasitic organisms (S1 Table).

242 In other words, genomes boosting the performance share either a lot or a little  
243 similarity with the reference interactome across a range of dimensions. Thus,  
244 phylogenetic profiling in eukaryotes benefits from genomes with a little or a lot of  
245 interactions present with regards to the reference interactome. These results reveal the  
246 importance of selecting genomes based on the evolutionary information contained within  
247 them relative to the query species, and is critical for high performance when predicting  
248 interacting proteins.

#### 249 **4. Orthologous group (pre-)selection improves prediction** 250 **performance by (inadvertently) enriching co-evolving** 251 **proteins in profiles**

252 Phylogenetic profiling benefits from clear modular co-evolution of proteins and  
253 subsets of proteins showing similar evolutionary behaviour [16,23]. A myriad of factors  
254 limit the modular co-evolution of interacting protein [24–26]. In previous research [19],  
255 which provides the starting meta parameters of this study, we evaluated orthology  
256 methods by their ability to recapitulate gene family dynamics in the Last Eukaryotic  
257 Common Ancestor (LECA). Consequently, the results so far are based on orthologous  
258 groups estimated to be in LECA. To see if this selection criterion was a factor in the strong  
259 performance, we performed phylogenetic profiling with other orthologous groups  
260 selections: groups estimated to be post-LECA, or groups not filtered on any criteria (post-  
261 LECA + LECA), i.e., the raw output of the orthology inference methods. We compared

262 these orthologous group sets with 1000 subsets of randomly selected LECA orthologous  
263 groups. The prediction performance was indeed reduced (AUC: 0.691 post-LECA and  
264 AUC: 0.734 all orthologous groups) compared to that of LECA orthologous groups or any  
265 randomly selected set of orthologous groups (Fig 4.). After some reflection, a myriad of  
266 explanations likely factor into this effect. Profiles of LECA proteins have many losses, and  
267 thus a lot of information (entropy) (S6 and S8 Figs). Profiles of post-LECA proteins have  
268 less loss and, by definition, are restricted to specific lineages, and thus contain less  
269 information. Combining LECA and post-LECA orthologous groups produce a set of  
270 phylogenetic profiles with an overall much lower similarity.

271

272 **Fig 4. Orthologous group selection has a large impact on prediction**  
273 **performance.** Receiver-operator Curve of post-LECA orthologous groups and unfiltered  
274 orthologous groups. The inset gives the Area Under the Curve (AUC) values compared  
275 with the random backdrop of randomly selected LECA OGs (violin plot) and the initial  
276 species set (green diamond).

277

278 We now have identified a key meta parameter choice explaining why our previous  
279 research found such high performance. However, it is unclear what the reason for this  
280 effect is and why for specific pairs of proteins, one protein was in LECA and the other not.  
281 This separation could be biological reality, i.e., innovations in the evolution from LECA to  
282 human, or issues in orthology assignment, i.e., one protein is evolving much more rapidly  
283 that causes the protein's predicted orthologous group to give an artifactually lineage-  
284 specific distribution in the profile. Consequently, the protein is falsely inferred as a more

285 recent addition or innovation. Manual inspection of this set (S9 Fig) does not obviously  
286 point towards one of the explanations. It is likely a combination of factors, including  
287 orthology prediction errors (e.g., oversplitting) and actual lineage-specific  
288 additions/inventions. In any case, the meta parameter of orthologous group selection is  
289 perhaps easily overlooked or made implicitly in the OG creation itself. Still, it is highly  
290 impactful, and our results show that OG selection improves prediction performance by  
291 enriching co-evolving proteins in profiles.

## 292 **5. Choice of reference interactome and interaction filtering**

### 293 **improves prediction performance by increasing the amount**

### 294 **of co-evolving proteins and quality of interactions**

295 Phylogenetic profiling attempts to predict which pairs of proteins are part of the  
296 same function, pathway, or complex. The performance of phylogenetic profiles can be  
297 measured using a data set of proteins that interact or are otherwise functionally linked.  
298 For example, we can take KEGG pathways as measuring units, as done in the STRING  
299 database [27]. However, these pathways often have an excess of 30 proteins and not all  
300 of them are expected to have the mutual functional dependence that results in co-  
301 evolution. This unwantedly biases the predictor by having supposedly interacting proteins  
302 with little correlation. Similarly, if we would take a very small well-curated set of  
303 compact/short linear metabolic pathways as was used to seed the CLIME searches [11],  
304 then the choice of what protein pairs to count as false negatives becomes difficult. Hence,  
305 our decision in previous work was to parse human interactions from BioGRID to contain  
306 only interactions found in at least five independent studies (Methods and Materials). This

307 filtering of interactions has been repeatedly demonstrated to effectively increase the  
308 quality and reduce the noise in the interaction data [27,28]. Moreover, the same data set  
309 contains a very good indication of which proteins are not functionally related. Proteins  
310 that are well studied and repeatedly surface in high throughput assays and are subject to  
311 repeated investigations are indeed likely to have no functional relation since these  
312 proteins are evidently never identified to interact.

313         The results in section 3 (Fig 3.C.) reinforce the notion that the reference interaction  
314 set plays a role in the performance of predicting interacting proteins. For these reasons,  
315 we analysed how the filtering and choice of reference interactome influences protein  
316 interaction prediction performance in eukaryotes. Using an unfiltered human protein  
317 interaction dataset reduces the prediction performance from an AUC of 0.779 to an AUC  
318 of 0.638 (Fig 5. A.). This performance is also lower than any set of randomly selected  
319 LECA orthologous groups (inset). The quality of the interaction data used clearly plays a  
320 role in prediction performance, i.e., if we take a noisy “ground truth” it turns out to be  
321 difficult to predict this truth. It is difficult to predict interactions with a set littered with false,  
322 virtually random, pairs.

323

324         **Fig 5. Interactome selection is important for prediction performance. A.**  
325 Receiver-operator Curve of post-LECA orthologous groups and unfiltered orthologous  
326 groups. The inset gives the Area Under the Curve (AUC) values compared with the  
327 random backdrop of randomly selected LECA orthologous groups (violin plot) and the  
328 initial species set (green diamond). **B.** GO-enrichment analysis for genes enriched in  
329 interactions present in only human compared to interactions present in human and yeast.

330 Orthologous groups can contain multiple genes. We randomly selected genes from an  
331 orthologous group to generate a new sample and population sets ten times and  
332 recalculated the enrichment (shown by multiple points in the figure rows).

333

334 We further analysed the choice of reference organism for protein interactions.  
335 Specifically for eukaryotes, the prediction performance was sensitive to the reference  
336 species for protein interactions [21]. Yeast has been the organism of choice as the  
337 reference interaction set for eukaryotes. Yeast is a popular model organism that has been  
338 extensively researched, and it is with yeast that many protein-protein interaction high  
339 throughput methods were pioneered. As a result, we also expect the interaction data of  
340 yeast to be of higher quality than that of human and, consequently, interaction predictions  
341 to be better.

342 We used *Saccharomyces cerevisiae* interactions from BioGRID (Materials and  
343 Methods) filtered with the same number of publications strictness criterion. Surprisingly,  
344 and contrary to for instance [6], the human interaction set performed better with an AUC  
345 of 0.779 compared to the yeast interaction set with an AUC of 0.713 (Fig 5. A.). One reason  
346 could be that ascomycete fungi and yeast in particular, has lost many co-evolving LECA  
347 complexes found in most eukaryotes [29,30]. These losses include Complex I, essential  
348 functions in chromatin modification [31], spliceosomal introns and RNAi machinery giving  
349 patchy patterns of canonical Dicer and Argonaute [32], ciliary genes [8,33], and the  
350 WASH complex [9,11,34]. These observations prompted us to look at the GO term  
351 enrichment of interacting LECA orthologous groups that contain only human genes  
352 versus interacting LECA orthologous groups that have both human and yeast genes.



353 We indeed find evidence of multiple genes belonging to ancestral complexes  
354 enriched in the human interaction set (Fig 5. B.), including enrichment in more  
355 straightforward GO terms related to mitochondria and respiration (e.g., GO:0005747,  
356 GO:0006120, GO:0032981 and GO:0070469), cilium (e.g., GO:0005929) and  
357 spliceosomal components (e.g., SMN complex GO:0032797). We also find evidence in  
358 higher-level GO terms that at lower levels reflect complexes known to be present in  
359 human and absent in yeast (S2 Table), such as chromatin modification (e.g.,  
360 GO:0042127). For the broccoli inferred orthologous groups, more enriched GO terms  
361 reflect at lower-level complexes known to be present in human and not in yeast:  
362 Argonaute and Dicer (GO:0010629, GO:0048471 and GO:0030426), WASH  
363 (GO:0005814, GO:0005856 and GO:0043005) and chromatin modification (e.g.,  
364 GO:0007399) (S10 Fig and S3 Table).

365 Even though there are more yeast than human interactions present in multiple  
366 species, the entropies of the profiles participating in yeast interactions are lower (S11  
367 Fig). This observation and the GO analysis reveal a clear reason why the performance  
368 we reported is high relative to others. Namely, we use the human reference interaction  
369 set with ancestral complexes that have been frequently lost throughout eukaryotic  
370 evolution and are absent in yeast.

## 371 **Discussion**

372 Phylogenetic profiling is complicated due to many biological and technical issues.  
373 These issues include the complex histories of proteins and the choices in the meta  
374 parameters for phylogenetic profiling, such as the quantity, quality and diversity of

375 reference genomes and annotations. Meta parameter choices in phylogenetic profiling  
376 has been extensively studied in mostly prokaryotes, where generally the focus is on the  
377 choice of reference genomes, phylogenetic profile methods, and/or the amount of data.  
378 We focus on eukaryotes and investigate qualitative different meta parameters for  
379 phylogenetic profiling. We showed that phylogenetic profile performance when predicting  
380 protein interactions is influenced by a complex interplay of multiple technical and  
381 biological parameters.

382         Genome diversity plays an important role in prediction performance for prokaryotes  
383 [16,18]. In contrast, our measures of eukaryotic diversity did not significantly influence  
384 prediction performance. Selecting lesser-quality genomes has a larger effect on  
385 prediction performance, while selecting higher-quality genomes does not. Genome  
386 selection and the interplay between quality and diversity does matter. However, other  
387 meta parameters have a much larger impact on prediction performance, such as the  
388 amount of information in the phylogenetic profiles in relation to the reference interaction  
389 dataset. This discrepancy suggests that more complex feature selection procedures  
390 should be explored for reference genome set selection, especially since (non-linear)  
391 interactions between subsets of genomes and combinations of subsets could drastically  
392 boost performance.

393         Other meta parameters, such as orthologous group and reference interactions,  
394 have a much larger effect than genome selection. Some results make a lot of sense from  
395 a technical point of view. For example, low quality/noisy functional data (unfiltered  
396 BIOGRID) or mixing phylogenetic profiles that are at least 50% inconsistent (post-LECA  
397 + LECA orthologous groups) have poor performance. A drawback to filtering out post-

398 LECA orthologous groups is that we remove lineage-specific interactions that are still a  
399 part of a protein complex and show clear co-evolution. Our analysis shows that we should  
400 consider these often hidden choices when encountering large differences in performance  
401 between reported studies.

402         One very counterintuitive finding is that the yeast interaction set showed lower  
403 predictive performance. Compared to the human interaction set, yeast should be of equal  
404 quality by all accounts, arguably even better. Together with the observation that LECA  
405 orthologous groups performed better than post-LECA orthologous groups, this suggests  
406 that the performance of phylogenetic profiles in eukaryotes is optimal for modules that  
407 fulfil a very particular set of conditions. These modules (i) were present in LECA, (ii) were  
408 repeatedly lost in eukaryotic lineages, and (iii) the genes in the module conserved most  
409 of their function. This observation fits with notable examples from the WASH complex and  
410 cilium [8,9,11], or proteins with great success in predicting its components like the minor  
411 spliceosome [5] and RNAi machinery genes including Dicer and Argonaute [32]. These  
412 biological patterns should explain the very strong signal found by studies such as [9,11].  
413 Note, both studies show very strong signals for complexes as well as pathways, which  
414 we excluded due to the problem of defining a quality negative interaction set.

415         In conclusion, we find that for eukaryotes more genomes and better-quality  
416 genomes are not necessarily better. It is instead the type of information in the genomes.  
417 The information in these genomes is not directly related to larger genomes, for instance  
418 parasites increase prediction performance. Instead, the information is related to the  
419 interactions of the reference species present in a given genome. Genome selection has  
420 a minor influence compared to orthologous groups selection and interactome selection,

421 which both greatly improve the performance when predicting protein interactions.  
422 Interactome and orthologous group selection is likely the major source for the large  
423 variance in reported performances. Ancestral complexes that are repeatedly lost are  
424 responsible for the strong performance of phylogenetic profiles in eukaryotes and it is  
425 these hidden choices in orthologous group selection that we should consider when we  
426 find large differences in performance between studies.

## 427 **Material and Methods**

### 428 **1. Initial datasets and methods**

429 We started our investigation from the analysis done in our previous work [19], to  
430 investigate the influence of different parameters on the performance of predicting protein-  
431 protein interactions using phylogenetic profiles. We showed a relatively high prediction  
432 performance using a large set of diverse eukaryotes and orthologous groups inferred to  
433 be in the Last Eukaryotic Common Ancestor (LECA). This reference set is called the initial  
434 set. Any changes that we made are changes in this initial set. In the sections below, we  
435 will briefly describe the composition of this initial set and the methods we used to obtain  
436 it.

#### 437 **1. 1. Large scale eukaryotic dataset and LECA orthologous groups**

438 We inferred orthologous groups on a diverse genome set of 167 eukaryotes using  
439 different orthologous group inference methods in our previous work. For this analysis, we  
440 chose the best performing method regarding protein interaction prediction, Sonicparanoid

441 (version 1.3.0) [13]. To rule out any large orthology specific issues during our current  
442 analyses, we chose at least one other method: Broccoli (version 1.0) [14].

443 Ancestral eukaryotic complexes have been lost together multiple times [35].  
444 Phylogenetic profiles should benefit from this clear modular evolution of proteins.  
445 Therefore, we selected orthologous groups estimated to be in LECA. Briefly, we inferred  
446 LECA orthologous groups using the Dollo parsimony approach [36] with additional strict  
447 inclusion criteria [19]. The Dollo parsimony method assumes that genes can be gained  
448 only once while minimizing gene loss. Before we assigned an orthologous group to LECA,  
449 it must be in at least three supergroups (See Supp. Table X) distributed over the  
450 Amorphae and Diaphoretickes (previously known as opimoda and diphoda) [37].

## 451 **1. 2. Phylogenetic profiling and measuring co-occurrence of proteins**

452 We constructed phylogenetic profiles by determining the presence (1) and  
453 absence (0) of orthologous groups in 167 species. To evaluate prediction accuracy, we  
454 obtained a higher quality reference interaction set by filtering the human BioGRID  
455 interaction database (version 3.5.172 May 2019) [38,39]. BioGRID contains physical  
456 interactions between proteins. We filtered this interaction set to keep non redundant  
457 interaction pairs found in at least five independent studies (PubMed ID's). The number of  
458 independent studies is a measure of how thoroughly these proteins were investigated  
459 and how receptive the proteins are to high-throughput measurements. We mapped the  
460 interacting genes to their corresponding orthologous groups.

461 We used the best performing negative protein interaction set from our previous  
462 analyses [19]. We inferred this negative set by taking pairs of interacting proteins that  
463 were found to be interacting at least five times, but not with each other. This excludes the

464 possibility that the negative set contains interacting proteins that were not found due to  
465 manifold technical reasons.

466 To calculate the (dis)similarities between phylogenetic profiles we used the from  
467 our previous analysis best performing distance measure, the cosine distance.

## 468 **2. Genome selection procedures**

469 We compared the results of all the genome selection procedures to 1000 sets of  
470 genomes randomly selected to exemplify that the differences in prediction accuracies are  
471 not due to random variations in genome composition. We calculated the protein  
472 interaction prediction performance for each of these random genome sets.

### 473 **2. 1. Selecting better and worse quality genomes**

474 To measure the quality of the genomes, we used two quality metrics. The first  
475 metric is the out of the box BUSCO metric that works by calculating the absences of highly  
476 conserved single-copy orthologs [22]. The BUSCO Eukaryota database (odb9) was  
477 aligned to the genomes using the hmmsearch alignment tool from the HMMER package  
478 3.1b2 (dated February 2015) [40]. We took the HMM specific quality score given by  
479 BUSCO to validate the hits in the alignments.

480 The second metric is of our own devising. The second quality metric we used was  
481 the Illogical absences (IA) metric of our design. We added this second independent metric  
482 to remove the dependence of quality on a single measure to establish the completeness  
483 of the genomes and gene prediction. The IA metric calculates the number of absences of  
484 protein interaction partners, which we termed Illogical absences. Illogical absences follow  
485 from the assumption underlying phylogenetic profiling that interacting proteins are often

486 evolutionary conserved. Therefore, it can be considered suspect when a protein  
487 interaction partner is absent. A possible reason could be that the absence is due to gene  
488 prediction, genome annotation or even homology detection errors.

489 We selected the strongest interacting orthologous group pairs by selecting their  
490 phylogenetic profiles with the least cosine distance. This selection removes the complex  
491 interplay between interacting groups of orthologs. For every interacting orthologous group  
492 pair, we calculated the absences of interaction partners in every species. These absences  
493 we termed illogical absences or the IA metric.

494 We can consider the genomes with the most BUSCO absences and illogical  
495 absences as lesser quality genomes. In contrast, we can consider the genomes with the  
496 least BUSCO absences and illogical absences as higher-quality genomes. We selected  
497 50 genomes of lesser-quality and 50 genomes of higher-quality for each of the metrics  
498 and recalculated the protein-protein interaction prediction performance.

## 499 **2. 2. Selecting highly diverse and similar genomes**

500 We calculated the pairwise cosine distance between all species with the presence  
501 and absence profiles of LECA orthologous groups to obtain species sets of maximum  
502 diversity and maximum similarity. We then iterated through the resulting pairwise distance  
503 matrix and selected the maximally distant pairs for the diverse set or minimally distant  
504 pairs for the similar set. Before adding a species of a species pair to a set, we checked  
505 to see if the species also had a distance above a certain arbitrary threshold to the other  
506 species in the growing set (cosine value  $\geq 0.38$  for the dissimilar, cosine value  $\leq 0.58$  for  
507 the similar set). We did this until we obtained the desired amount of 50 genomes per set.

508 The maximum diverse and maximum similar genome sets were each used to recalculate  
509 the protein-protein interaction prediction performance.

### 510 **2.3. Selecting single influential genomes and their combined effect on** 511 **prediction performance**

512 We removed genomes one-by-one from the initial species set of 167 eukaryotes  
513 to see how the different genomes influence the performance of protein interaction  
514 prediction with phylogenetic profiling. We recalculated the performance for each of these  
515 167 sets. The 50 genomes that increased the performance compared to the initial species  
516 set the most when removed from the initial set were labelled as disadvantageous. The 50  
517 genomes that decreased the performance the most when removed from the initial set  
518 were labelled advantageous. For both the disadvantage and advantageous set we  
519 recalculated the protein interaction prediction performance.

## 520 **3. Gene and interactome selection procedures**

521 We compared the results of the orthologous group selection procedures to  
522 randomly selected LECA orthologous groups to exemplify that the differences in  
523 prediction accuracy is not due to random variations in orthologous group composition.  
524 We made a thousand LECA orthologous group sets containing a random selection of  
525 63% of the orthologous groups. We calculated each of these set's protein interaction  
526 prediction performance.



### 527 **3. 1. Selecting orthologous groups**

528 In our initial species set, we used orthologous groups estimated to be in LECA  
529 (Methods section 1.1.). We took the raw output of the orthology inference methods and  
530 filtered out the LECA orthologous groups to get a set that contains post-LECA orthologous  
531 groups. We also recalculated the prediction performance with the raw output of the  
532 orthology prediction methods, which is all inferred orthologous groups.

### 533 **3.2. Selecting different reference interactomes**

534 We compared the five PubMed ID filtered human BioGRID set with the unfiltered  
535 human BioGRID dataset. Every interaction with less than five pubIDs is now included as  
536 well. Removing the five PubMedID filter should indicate how quality filtering of reference  
537 interactions influences prediction performance.

538 We selected next to the human interactions the *Saccharomyces cerevisiae*  
539 BioGRID interaction database (version 3.5.175 July 2019) [39] to analyse the influence  
540 of the reference interactome. We filtered the interactions to keep only the interaction pairs  
541 found in at least five publications (PubMed ID's). We followed the same procedure as with  
542 the human interaction set (Methods section 1.2.).

543 Following this analysis, we hypothesized that the drop in prediction performance  
544 for yeast is caused by the loss of ancestral protein complexes in yeast. To test this, we  
545 chose interacting LECA orthologous groups that contained only human genes (sample  
546 set) and calculated the enrichment to the set with interacting LECA orthologous groups  
547 containing human and yeast genes (population set). We calculated the enrichment using  
548 the following equation:  $\frac{n}{m} \div \frac{k}{q}$ , where n is the total number of genes associated with a GO

549 term (Downloaded GO terms Januari 2021 biomart) in the sample set (overlap),  $m$  is the  
550 total numbers of genes in the sample set,  $k$  is the total number of genes associated with  
551 a specific GO term in the population set, and  $q$  is the total number of genes in the  
552 population set. Since enrichment does not work well for small overlaps, we filtered for a  
553 minimum overlap ( $n$ ) of 3. Enrichment was considered significant for  $p$ -values below 0.01.  
554 Since orthologous groups can contain multiple genes, we randomly selected genes from  
555 an orthologous group to generate a new sample and population sets ten times and  
556 recalculated the enrichment.

## 557 **References**

- 558 1. Burki F, Roger AJ, Brown MW, Simpson AGB. The New Tree of Eukaryotes. Trends  
559 Ecol Evol. 2020;35: 43–55. doi:10.1016/j.tree.2019.08.008
- 560 2. Bork P, Koonin E V. Predicting functions from protein sequences -sequences -  
561 where are the bottlenecks ? Nat Genet. 1998;18: 313–318.
- 562 3. Nagy LG, Merényi Z, Hegedüs B, Bálint B. Novel phylogenetic methods are needed  
563 for understanding gene function in the era of mega-scale genome sequencing.  
564 Nucleic Acids Res. 2020;48: 2209–2219. doi:10.1093/nar/gkz1241
- 565 4. Pellegrini M, Marcotte EM, Thompson MJ, Eisenberg D, Yeates TO. Assigning  
566 protein functions by comparative genome analysis: Protein phylogenetic profiles.  
567 Proc Natl Acad Sci U S A. 1999;96: 4285–4288. doi:10.1073/pnas.96.8.4285
- 568 5. de Wolf B, Oghabian A, Akinyi M V, Hanks S, Tromer EC, van Hooff J, et al.  
569 Chromosomal instability by mutations in a novel specificity factor of the minor  
570 spliceosome. 2020. doi:10.1101/2020.08.06.239418

- 571 6. Moi D, Kilchoer L, Aguilar PS, Dessimoz C. Scalable phylogenetic profiling using  
572 MinHash uncovers likely eukaryotic sexual reproduction genes. Ouzounis CA,  
573 editor. PLOS Comput Biol. 2020;16: e1007553. doi:10.1371/journal.pcbi.1007553
- 574 7. Sherill-Rofe D, Rahat D, Findlay S, Mellul A, Guberman I, Braun M, et al. Mapping  
575 global and local coevolution across 600 species to identify novel homologous  
576 recombination repair genes. Genome Res. 2019;29: 439–448.  
577 doi:10.1101/gr.241414.118
- 578 8. van Dam TJP, Townsend MJ, Turk M, Schlessinger A, Sali A, Field MC, et al.  
579 Evolution of modular intraflagellar transport from a coatomer-like progenitor. Proc  
580 Natl Acad Sci. 2013;110: 6943–6948. doi:10.1073/pnas.1221011110
- 581 9. Dey G, Jaimovich A, Collins SR, Seki A, Meyer T. Systematic Discovery of Human  
582 Gene Function and Principles of Modular Organization through Phylogenetic  
583 Profiling. Cell Rep. 2015;10: 993–1006. doi:10.1016/j.celrep.2015.01.025
- 584 10. van Hooff JJ, Tromer E, van Wijk LM, Snel B, Kops GJ. Evolutionary dynamics of  
585 the kinetochore network in eukaryotes as revealed by comparative genomics.  
586 EMBO Rep. 2017;18: 1559–1571. doi:10.15252/embr.201744102
- 587 11. Li Y, Calvo SE, Gutman R, Liu JS, Mootha VK. Expansion of biological pathways  
588 based on evolutionary inference. Cell. 2014;158: 213–225.  
589 doi:10.1016/j.cell.2014.05.034
- 590 12. Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for comparative  
591 genomics. Genome Biol. 2019;20: 238. doi:10.1186/s13059-019-1832-y
- 592 13. Cosentino S, Iwasaki W. SonicParanoid: Fast, accurate and easy orthology  
593 inference. Bioinformatics. 2018;35: 149–151. doi:10.1093/bioinformatics/bty631

- 594 14. Derelle R, Philippe H, Colbourne JK. Broccoli: combining phylogenetic and network  
595 analyses for orthology assignment. Falush D, editor. *Mol Biol Evol.* 2020;  
596 2019.12.13.875831. doi:10.1093/molbev/msaa159
- 597 15. Snitkin ES, Gustafson AM, Mellor J, Wu J, Delisi C. Comparative assessment of  
598 performance and genome dependence among phylogenetic profiling methods.  
599 *BMC Bioinformatics.* 2006;7: 1–11. doi:10.1186/1471-2105-7-420
- 600 16. Jothi R, Przytycka TM, Aravind L. Discovering functional linkages and  
601 uncharacterized cellular pathways using phylogenetic profile comparisons: A  
602 comprehensive assessment. *BMC Bioinformatics.* 2007;8. doi:10.1186/1471-2105-  
603 8-173
- 604 17. Muley VY, Ranjan A. Effect of reference genome selection on the performance of  
605 computational methods for genome-wide protein-protein interaction prediction.  
606 *PLoS One.* 2012;7. doi:10.1371/journal.pone.0042057
- 607 18. Škunca N, Dessimoz C. Phylogenetic Profiling: How Much Input Data Is Enough?  
608 Escriva H, editor. *PLoS One.* 2015;10: e0114701.  
609 doi:10.1371/journal.pone.0114701
- 610 19. Deutekom ES, Snel B, van Dam TJP. Benchmarking orthology methods using  
611 phylogenetic patterns defined at the base of Eukaryotes. *Brief Bioinform.* 2020;00:  
612 1–9. doi:10.1093/bib/bbaa206
- 613 20. Sun J, Li Y, Zhao Z. Phylogenetic profiles for the prediction of protein-protein  
614 interactions: How to select reference organisms? *Biochem Biophys Res Commun.*  
615 2007;353: 985–991. doi:10.1016/j.bbrc.2006.12.146
- 616 21. Simonsen M, Maetschke SR, Ragan MA. Automatic selection of reference taxa for

- 617 protein-protein interaction prediction with phylogenetic profiling. *Bioinformatics*.  
618 2012;28: 851–857. doi:10.1093/bioinformatics/btr720
- 619 22. Waterhouse RM, Seppey M, Simao FA, Manni M, Ioannidis P, Klioutchnikov G, et  
620 al. BUSCO applications from quality assessments to gene prediction and  
621 phylogenomics. *Mol Biol Evol*. 2018;35: 543–548. doi:10.1093/molbev/msx319
- 622 23. Bloch I, Sherill-Rofe D, Stupp D, Unterman I, Beer H, Sharon E, et al. Optimization  
623 of co-evolution analysis through phylogenetic profiling reveals pathway-specific  
624 signals. *Bioinformatics*. 2020;36: 4116–4125. doi:10.1093/bioinformatics/btaa281
- 625 24. Snel B, Huynen MA. Quantifying modularity in the evolution of biomolecular  
626 systems. *Genome Res*. 2004;14: 391–397. doi:10.1101/gr.1969504
- 627 25. Campillos M, Von Mering C, Jensen LJ, Bork P. Identification and analysis of  
628 evolutionarily cohesive functional modules in protein networks. *Genome Res*.  
629 2006;16: 374–382. doi:10.1101/gr.4336406
- 630 26. Fokkens L, Snel B. Cohesive versus flexible evolution of functional modules in  
631 eukaryotes. *PLoS Comput Biol*. 2009;5. doi:10.1371/journal.pcbi.1000276
- 632 27. von Mering C, Huynen M, Jaeggi D, Schmidt S, Bork P, Snel B. STRING: A  
633 database of predicted functional associations between proteins. *Nucleic Acids Res*.  
634 2003;31: 258–261. doi:10.1093/nar/gkg034
- 635 28. von Mering C, Krause R, Snel B, Cornell M, Oliver SG, Fields S, et al. Comparative  
636 assessment of large-scale data sets of protein-protein interactions. *Nature*.  
637 2002;417: 399–403. doi:10.1038/nature750
- 638 29. Aravind L, Watanabe H, Lipman DJ, Koonin E V. Lineage-specific loss and  
639 divergence of functionally linked genes in eukaryotes. *Proc Natl Acad Sci U S A*.

- 640 2000;97: 11319–11324. doi:10.1073/pnas.200346997
- 641 30. Münsterkötter M, Steinberg G. The fungus *Ustilago maydis* and humans share  
642 disease-related proteins that are not found in *Saccharomyces cerevisiae*. *BMC*  
643 *Genomics*. 2007;8: 1–10. doi:10.1186/1471-2164-8-473
- 644 31. Dujon BA, Louis EJ. Genome diversity and evolution in the budding yeasts  
645 (*Saccharomycotina*). *Genetics*. 2017;206: 717–750.  
646 doi:10.1534/genetics.116.199216
- 647 32. Tabach Y, Billi AC, Hayes GD, Newman MA, Zuk O, Gabel H, et al. Identification  
648 of small RNA pathway genes using patterns of phylogenetic conservation and  
649 divergence. *Nature*. 2013;493: 694–698. doi:10.1038/nature11779
- 650 33. Barker AR, Renzaglia KS, Fry K, Dawe HR. Bioinformatic analysis of ciliary  
651 transition zone proteins reveals insights into the evolution of ciliopathy networks.  
652 *BMC Genomics*. 2014;15: 1–9. doi:10.1186/1471-2164-15-531
- 653 34. Kollmar M, Lbik D, Enge S. Evolution of the eukaryotic ARP2/3 activators of the  
654 WASP family: WASP, WAVE, WASH, and WHAMM, and the proposed new family  
655 members WAWH and WAML. *BMC Res Notes*. 2012;5. doi:10.1186/1756-0500-5-  
656 88
- 657 35. Albalat R, Cañestro C. Evolution by gene loss. *Nat Rev Genet*. 2016;17: 379–391.  
658 doi:10.1038/nrg.2016.39
- 659 36. Rogozin IB, Wolf YI, Badenko VN, Koonin E V. Dollo parsimony and the  
660 reconstruction of genome evolution. In: Albert VA, editor. *Parsimony, Phylogeny*  
661 *and Genomics*. 2006. pp. 1–18. doi:10.1093/acprof
- 662 37. Adl SM, Bass D, Lane CE, Lukeš J, Schoch CL, Smirnov A, et al. Revisions to the

- 663 Classification, Nomenclature, and Diversity of Eukaryotes. *J Eukaryot Microbiol.*  
664 2019;66: 4–119. doi:10.1111/jeu.12691
- 665 38. Stark C, Breitzkreutz B-J, Reguly T, Boucher L, Breitzkreutz A, Tyers M. BioGRID: a  
666 general repository for interaction datasets. *Nucleic Acids Res.* 2006;34: D535–  
667 D539. doi:10.1093/nar/gkj109
- 668 39. Oughtred R, Stark C, Breitzkreutz B-J, Rust J, Boucher L, Chang C, et al. The  
669 BioGRID interaction database: 2019 update. *Nucleic Acids Res.* 2019;47: D529–  
670 D541. doi:10.1093/nar/gky1079
- 671 40. Eddy SR. HMMER. Available: <http://hmmer.org>  
672

## 673 **Supporting information**

674 **S1 Fig. Illogical absences and genome quality selection based on illogical**  
675 **absences using Sonicparanoid inferred orthologous groups (OGs).** **A.** Illogical  
676 absences as a function of retained LECA OGs in different species. We see that where  
677 most Opisthokonta scores similarly low with the BUSCO metric, they score lower with the  
678 IA metric indicating a difference between the two metrics. However, the performance of  
679 genomes selected with both metrics are similar to each other. Filled data points are the  
680 selected genomes for the prediction accuracy calculations. **B.** Receiver-operator Curve  
681 of two species sets ( $n = 50$ ) with the most and least illogical absences. The inset gives  
682 the Area Under the Curve (AUC) values compared with the random backdrop of 1000  
683 random species sets (violin plot) and the initial species set (teal diamond). Human has a

684 perfect score of 0 illogical absences since the interactions are from the human reference  
685 interactome. Therefore, we did not select human for the genome set.

686

687 **S2 Fig. Lesser quality genomes have more impact on protein interaction prediction**  
688 **performance also for Broccoli inferred orthologous groups (OGs). A.** BUSCO and  
689 Illogical absences as a function of retained LECA OGs in different species. Filled data  
690 points are the selected genomes for the prediction accuracy calculations. **B.** Receiver-  
691 operator Curve of two species sets ( $n = 50$ ) with the most and least BUSCO and illogical  
692 absences. The inset gives the Area Under the Curve (AUC) values compared with the  
693 random backdrop of 1000 random species sets (violin plot) and the initial species set (teal  
694 diamond).

695

696 **S3 Fig. Both high and low diversity sets have little impact on protein interaction**  
697 **prediction performance also for Broccoli inferred orthologous groups. A.** The most  
698 similar species form more clusters and are overall more similar to each other. **B.** The most  
699 diverse species show no clustering and are overall less similar to each other. **C.** Receiver-  
700 operator Curve of two species sets ( $n = 50$ ) with the most diverse and most similar  
701 species. The inset gives the Area Under the Curve (AUC) values compared with the  
702 random backdrop of 1000 random species sets (violin plot) and the initial species set  
703 (green diamond).

704

705 **S4 Fig. Influential genomes and their combined effect for Broccoli inferred OGs. A.**  
706 Recalculated Area Under the Curve (AUC) values when a single species is removed from



707 the initial species set. Genomes that increase the AUC value when removed can be  
708 considered disadvantages compared to the initial set when predicting protein interactions  
709 with phylogenetic profiles. Genomes that decrease the AUC value when removed can be  
710 considered advantageous for predicting protein interactions. Top 50 advantageous and  
711 top 50 disadvantageous genomes shown with the black fill in the scatter plot. **B.** Receiver-  
712 operator Curve of two species sets (n=50) with the most advantageous and  
713 disadvantageous genomes. The inset gives the Area Under the Curve (AUC) values  
714 compared with the random backdrop of 1000 random species sets (violin plot) and the  
715 initial species set (green diamond). **C.** Comparison of the counts (histogram) and kernel  
716 density estimates (line plot) of (I) illogical absence ratios (illogical absences divided by  
717 total interaction absences (co-absences + illogical absences)), (II) present interactions,  
718 (III) the cosine distance to human, and (IV) total shared orthologous groups with human.  
719

720 **S5 Fig. Correlations between multiple parameters in the advantageous and**  
721 **disadvantageous genome set.** Given for **A.** Sonicparanoid and **B.** Broccoli inferred  
722 orthologous groups (OGs). From top to bottom (or left to right) the interactions that are  
723 co-absent; illogically absent; and present; the ratio of illogical absences to total absences;  
724 number of OGs shared with the human genome; the cosine distance to the human  
725 genome; LECA OGs loss (Dollo parsimony inferred); species (lineage) specific loss;  
726 (clade) ancestral loss; and the difference in AUC from the initial set AUC when a genome  
727 is removed.  
728

729 **S6 Fig. Entropy of phylogenetic profiles that have interactions.** Given for **A.**  
730 Sonicparanoid and **B.** Broccoli inferred orthologous groups (OGs). From top to bottom,  
731 the entropy is shown in profiles for LECA, post-LECA and all OGs. Median entropy is  
732 presented with a black arrow. Mann-Whitney U test shows significant difference between  
733 distributions of LECA, post-LECA and all OGs, p-value < 0.001.

734

735 **S7 Fig. Orthologous group selection has a large impact on prediction performance**  
736 **also for Broccoli inferred orthologous groups (OGs).** Receiver-operator Curve of  
737 post-LECA orthologous groups and unfiltered orthologous groups. The inset gives the  
738 Area Under the Curve (AUC) values compared with the random backdrop of randomly  
739 selected LECA OGs (violin plot) and the initial species set (green diamond).

740

741 **S8 Fig. Dollo parsimony inferred loss of LECA and post-LECA orthologous groups**  
742 **(OGs).** Given for **A.** Sonicparanoid and **B.** Broccoli inferred OGs. Mann-Whitney U test  
743 shows significant difference between distributions.

744

745 **S9 Fig. Groups of interacting orthologous groups (OGs) where one is in LECA**  
746 **(always the last row in a group subplot) and the others are not.** The profiles are  
747 sorted according to the species tree.

748

749 **S10 Fig. Interactome selection is important for prediction performance. A.** Receiver-  
750 operator Curve of post-LECA and unfiltered orthologous groups (OGs) of Broccoli. The  
751 inset gives the Area Under the Curve (AUC) values compared with the random backdrop

752 of randomly selected LECA OGs (violin plot) and the initial species set (green diamond).

753 **B.** GO-enrichment analysis for genes enriched in interactions present in only human vs.  
754 interactions present in human and yeast. OGs can contain multiple genes. We randomly  
755 selected genes from an OG to generate new sample and population sets 10 times and  
756 recalculated the enrichment (shown by multiple points in the figure rows).

757

758 **S11 Fig. Interactions of human and yeast interactome present in different species**  
759 **(left) and entropy for LECA profiles that have interactions in human and yeast**  
760 **(right).** Given for **A.** Sonicparanoid inferred and **B.** Broccoli inferred orthologous groups  
761 (OGs). Median values are presented with the arrows. Mann-Whitney U test shows  
762 significant difference between distributions.

763

764 **S1 Table. Species table for species used in this study.** Green marked species are the  
765 species that are in the advantageous set, and red marked species in the disadvantageous  
766 set (Sonicparanoid). The measured values are shown in S5 Fig.

767

768 **S2 Table. GO-enrichment table for Sonicparanoid inferred orthologous groups**  
769 **(OGs).** Since there can be multiple genes in an OG, we randomly selected one of the  
770 genes for the GO-enrichment analysis. We did this ten times, creating ten foreground and  
771 background sets (set\_num). These values are shown in Fig 5.

772

773 **S3 Table. GO-enrichment table for Broccoli inferred orthologous groups (OGs).**  
774 Since there can be multiple genes in an OG, we randomly selected one of the genes for

775 the GO-enrichment analysis. We did this ten times, creating ten foreground and  
776 background sets (set\_num). These values are shown in S10 Fig.

777

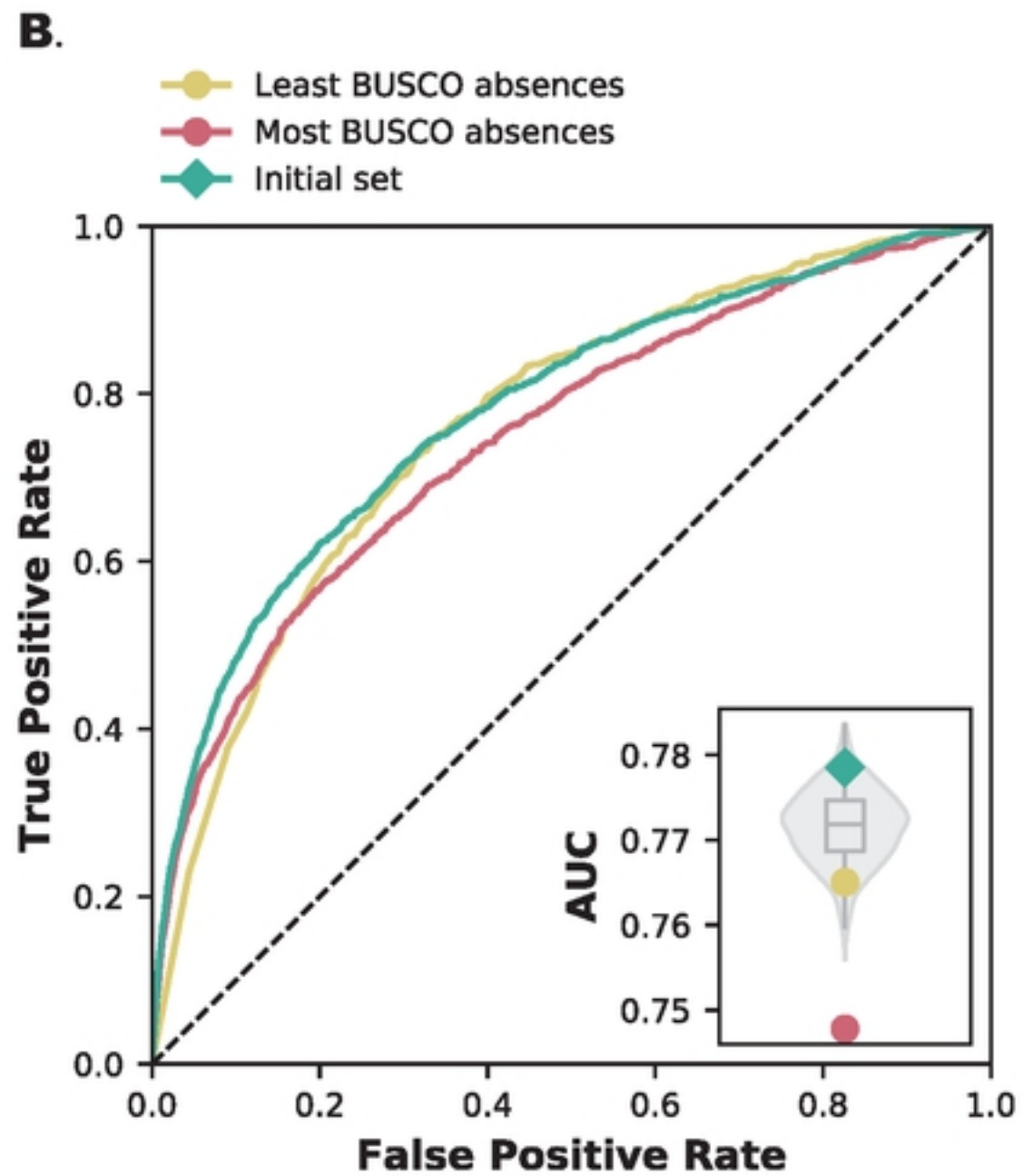
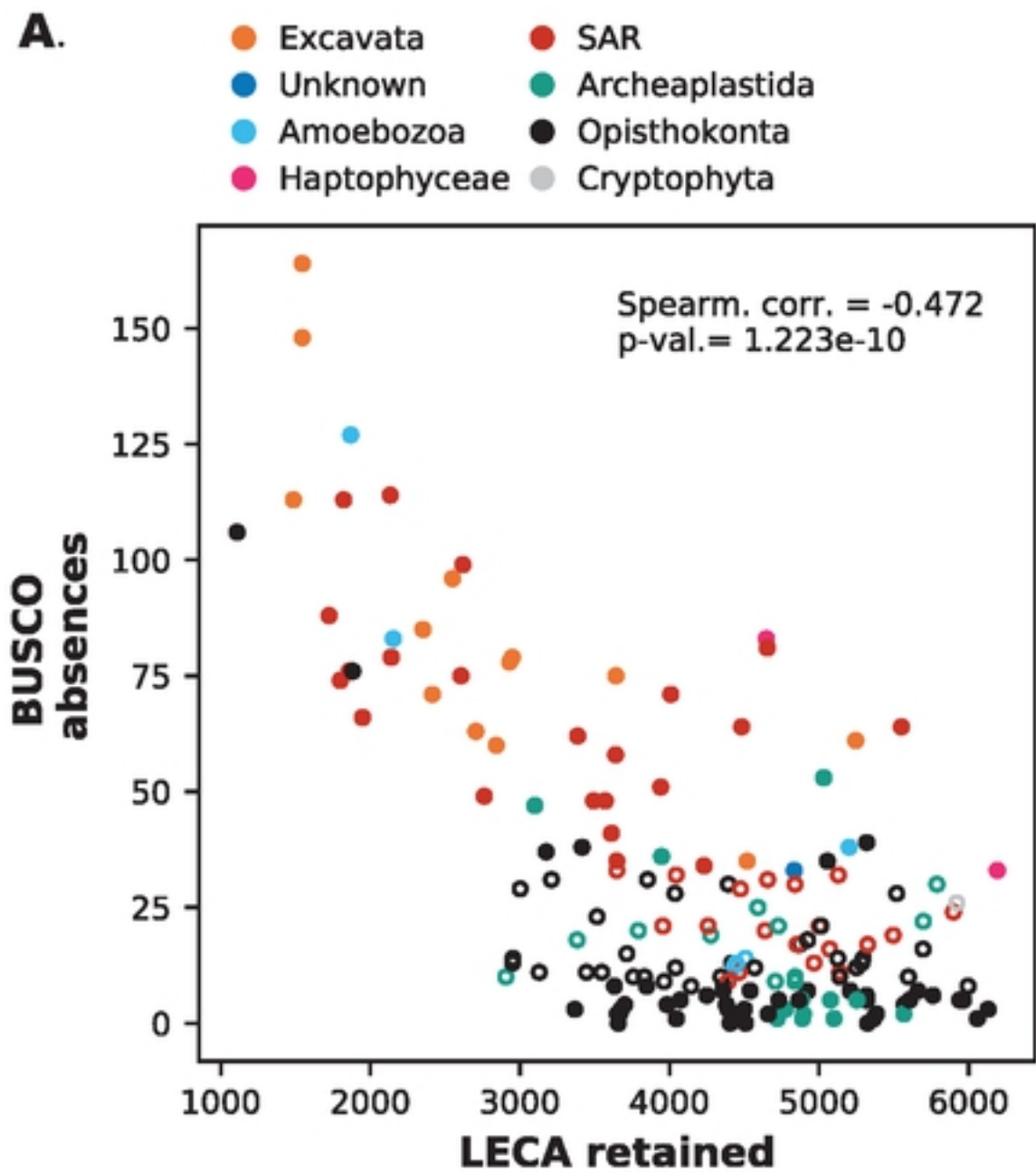
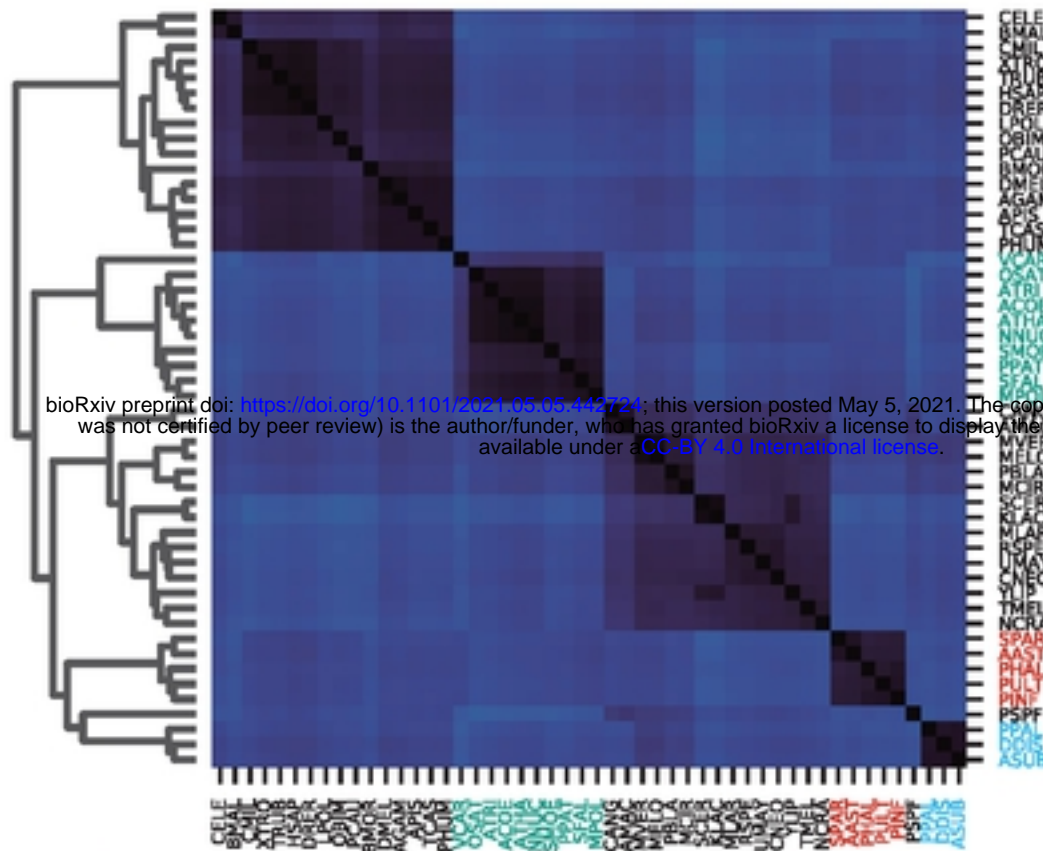


Figure 1

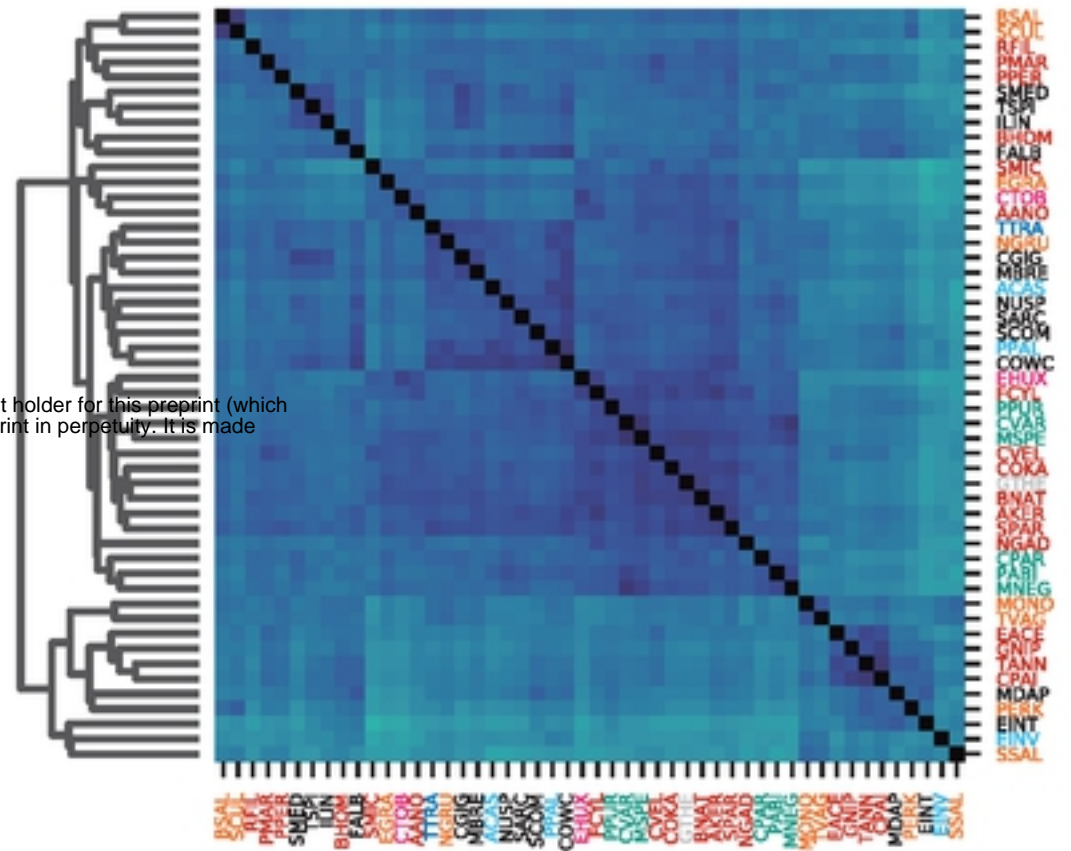
- Excavata
- Unknown
- Amoebozoa
- Haptophyceae
- SAR
- Archeplastida
- Opisthokonta
- Cryptophyta

**A.**

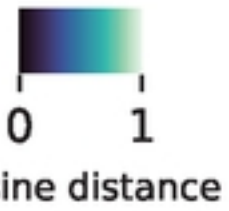


**Most similar species**

**B.**



**Most dissimilar species**



bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.05.445724>; this version posted May 5, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

**C.**

- Most diverse genomes
- Most similar genomes
- ◆ Initial set

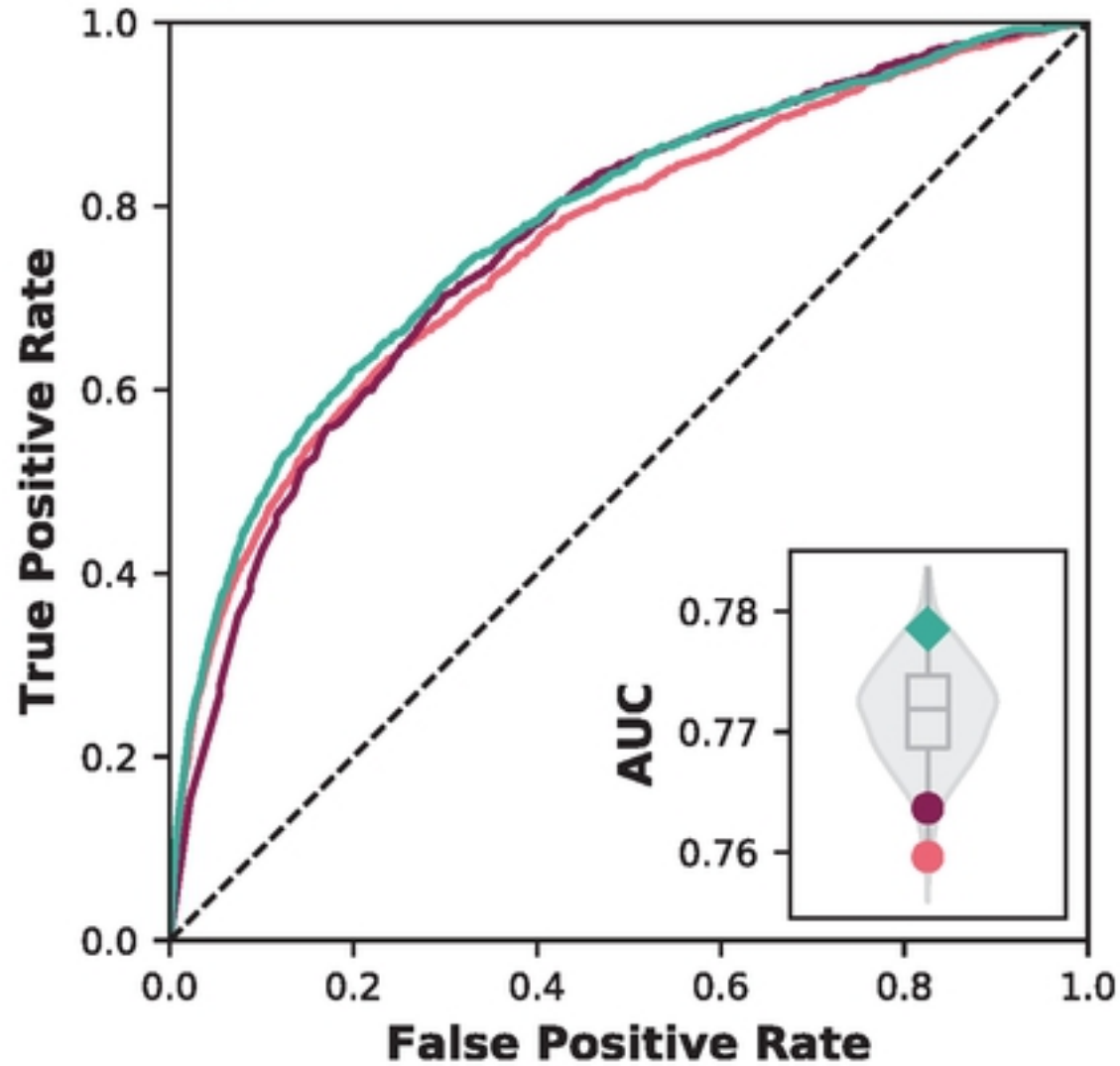


Figure 2



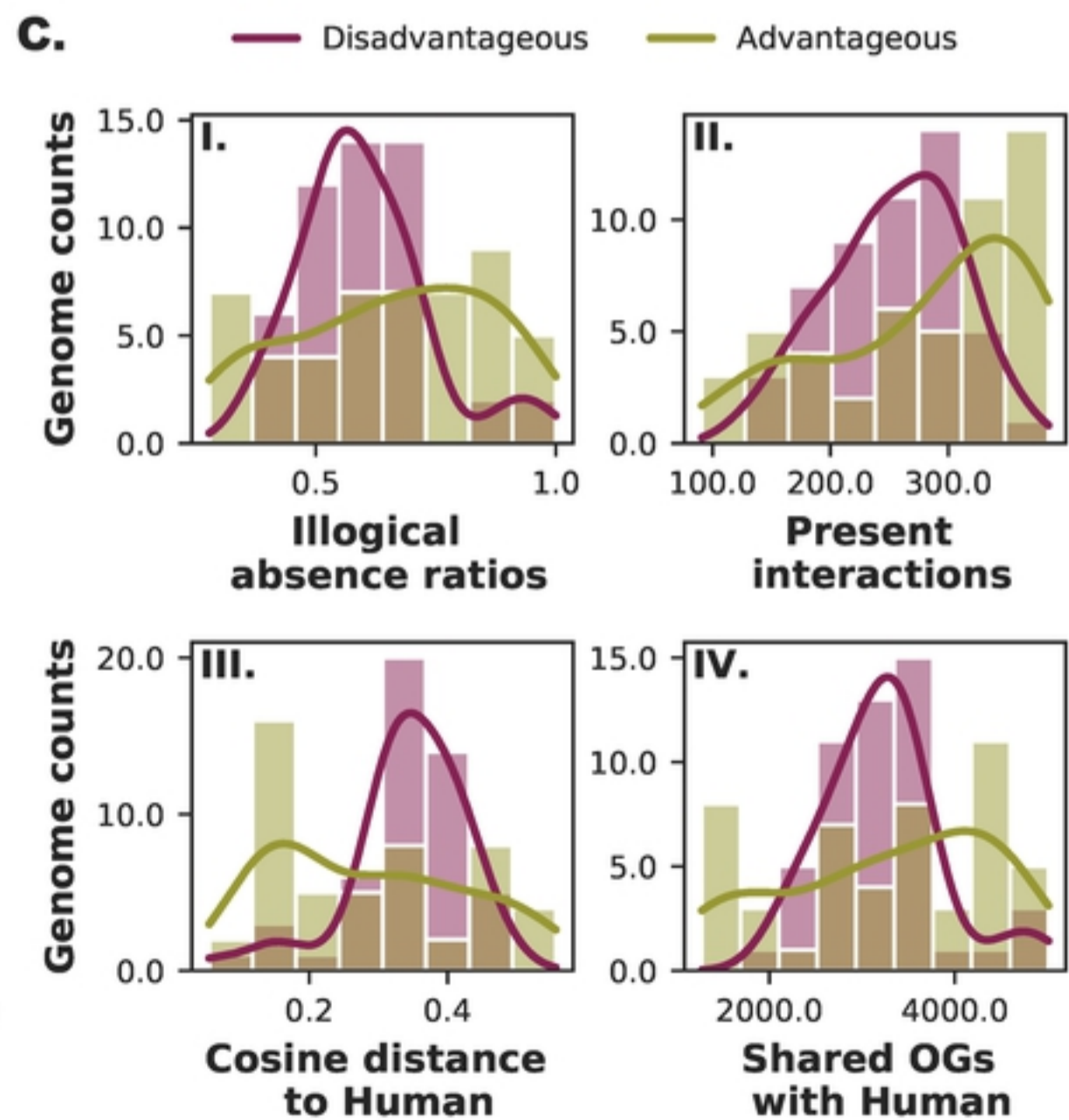
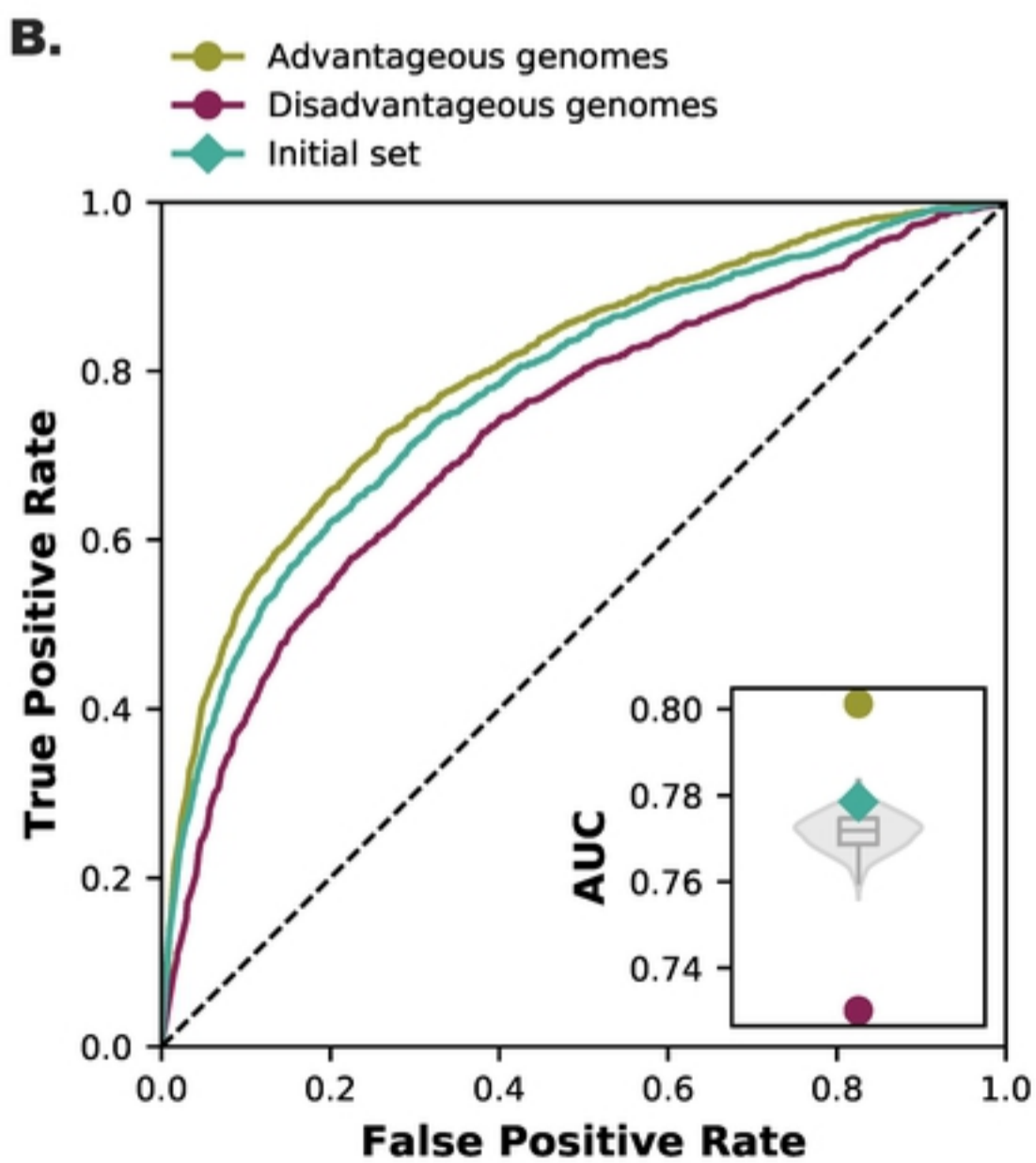
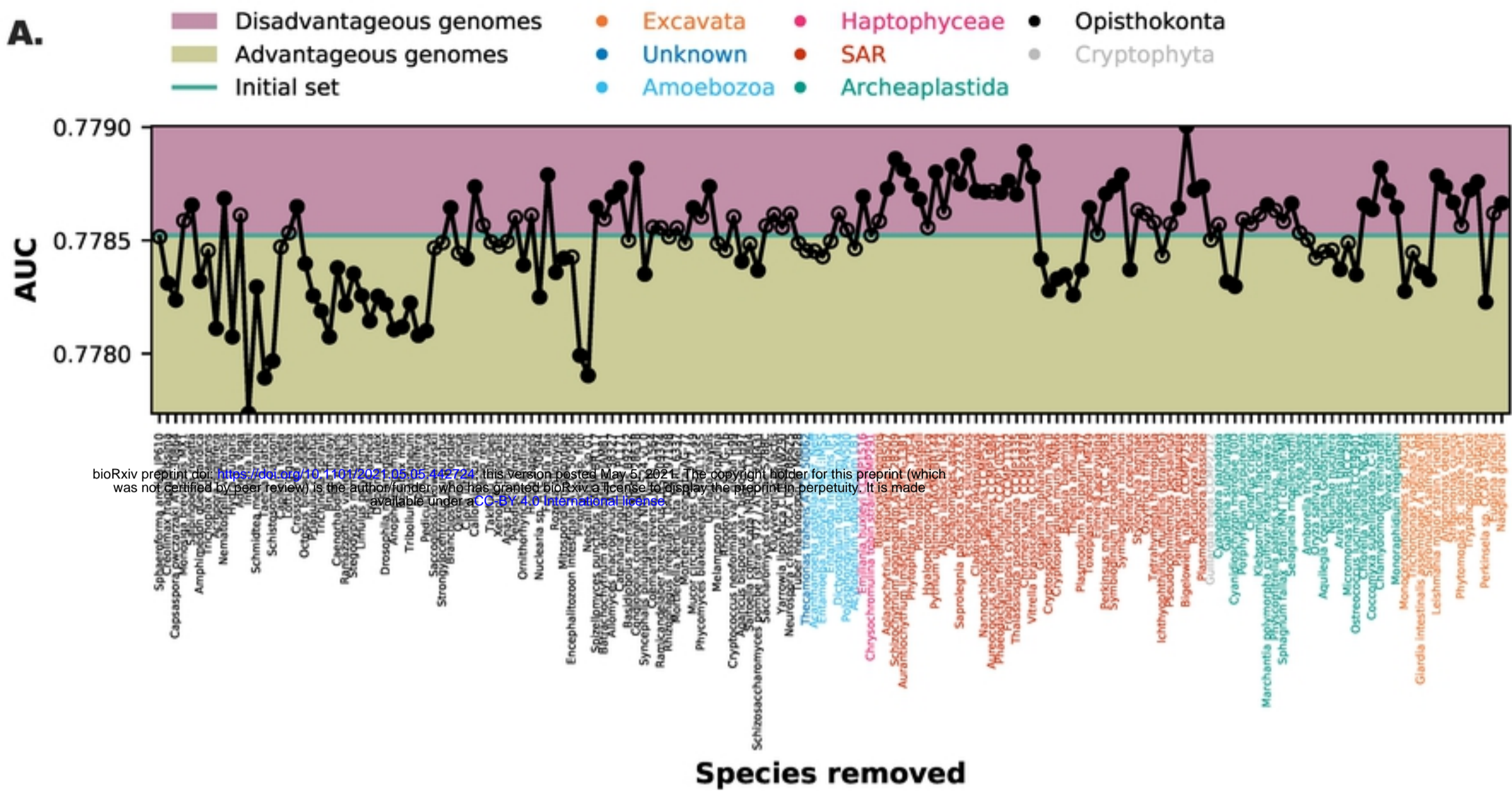


Figure 3

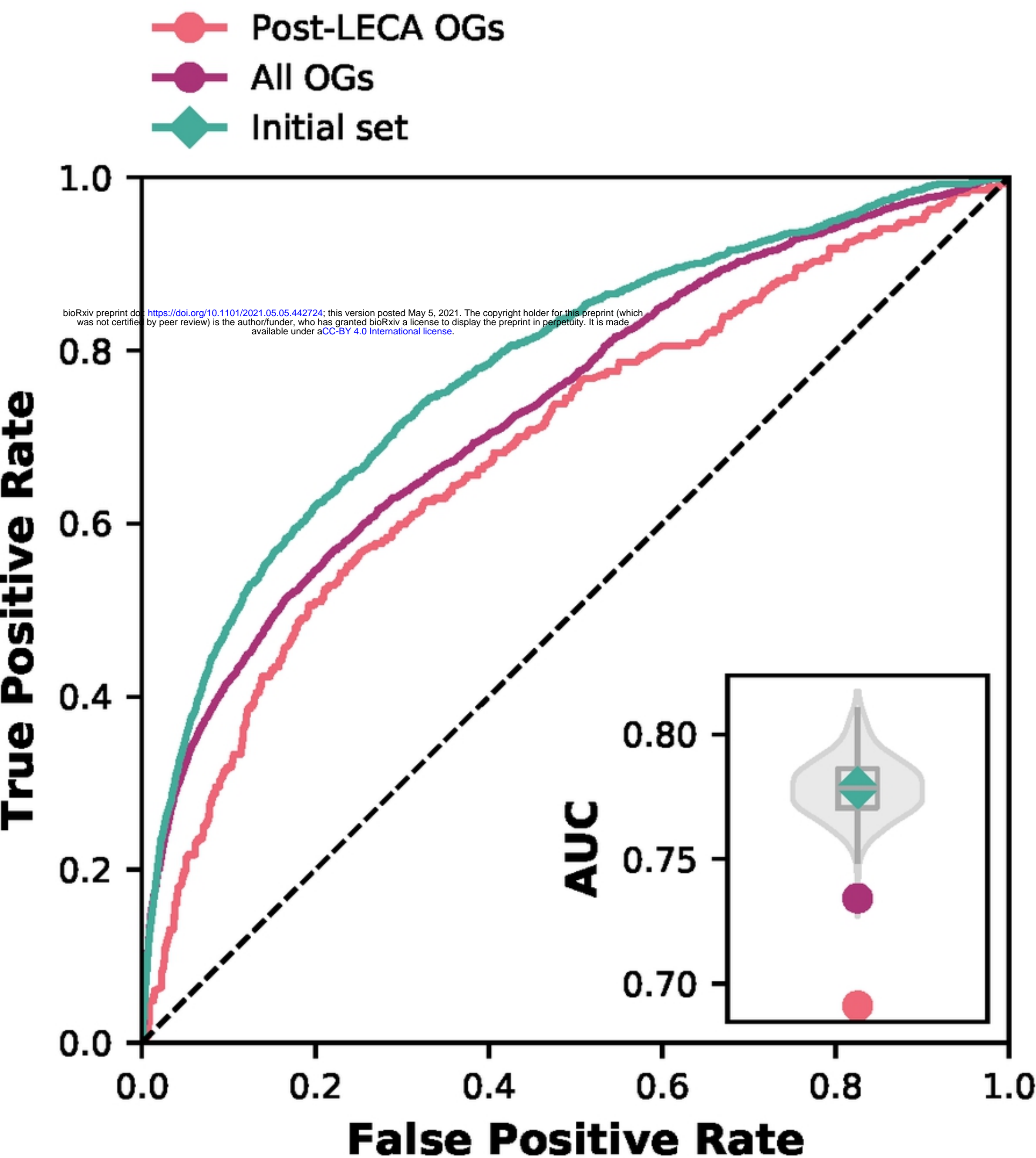


Figure 4



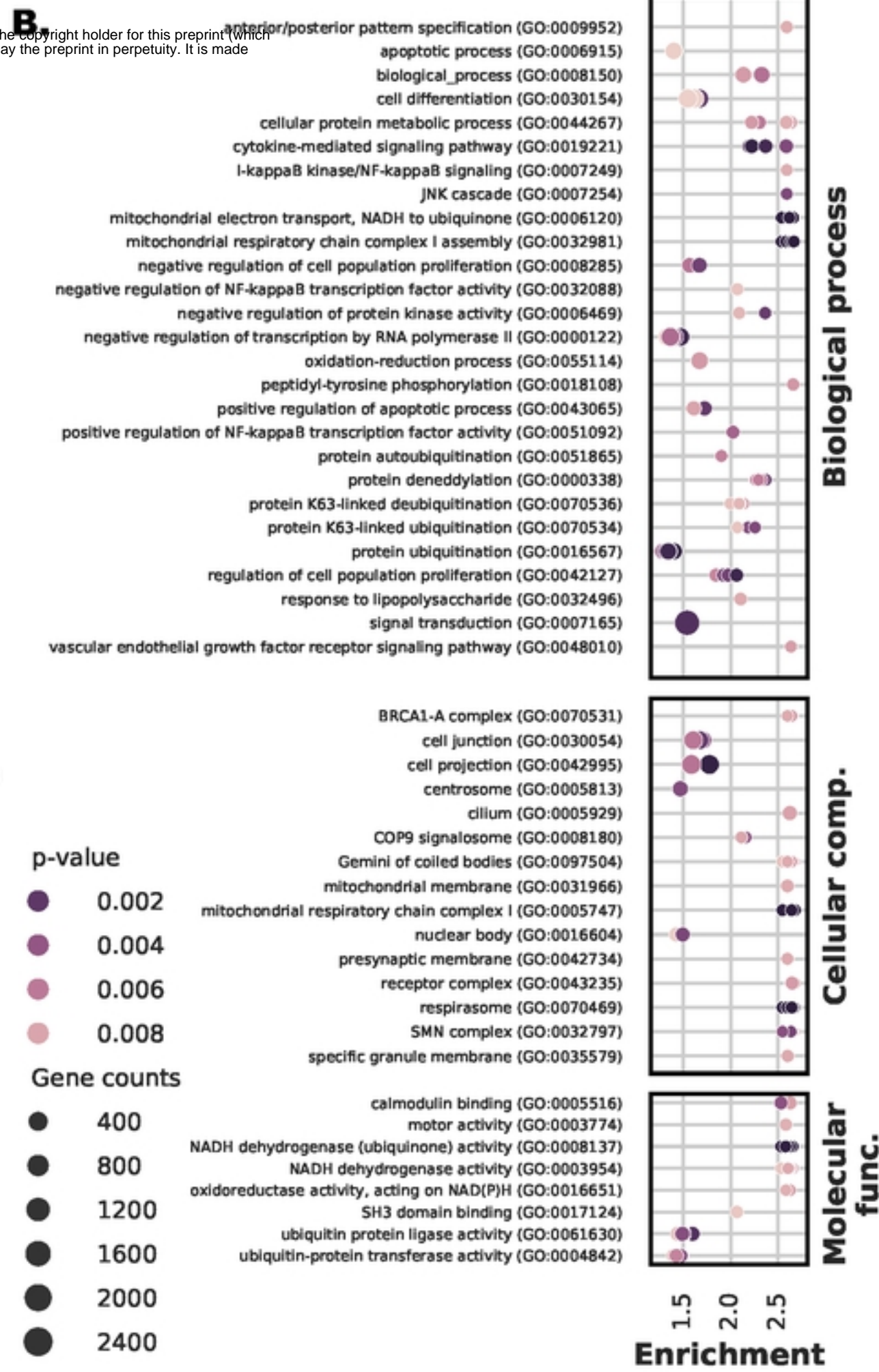
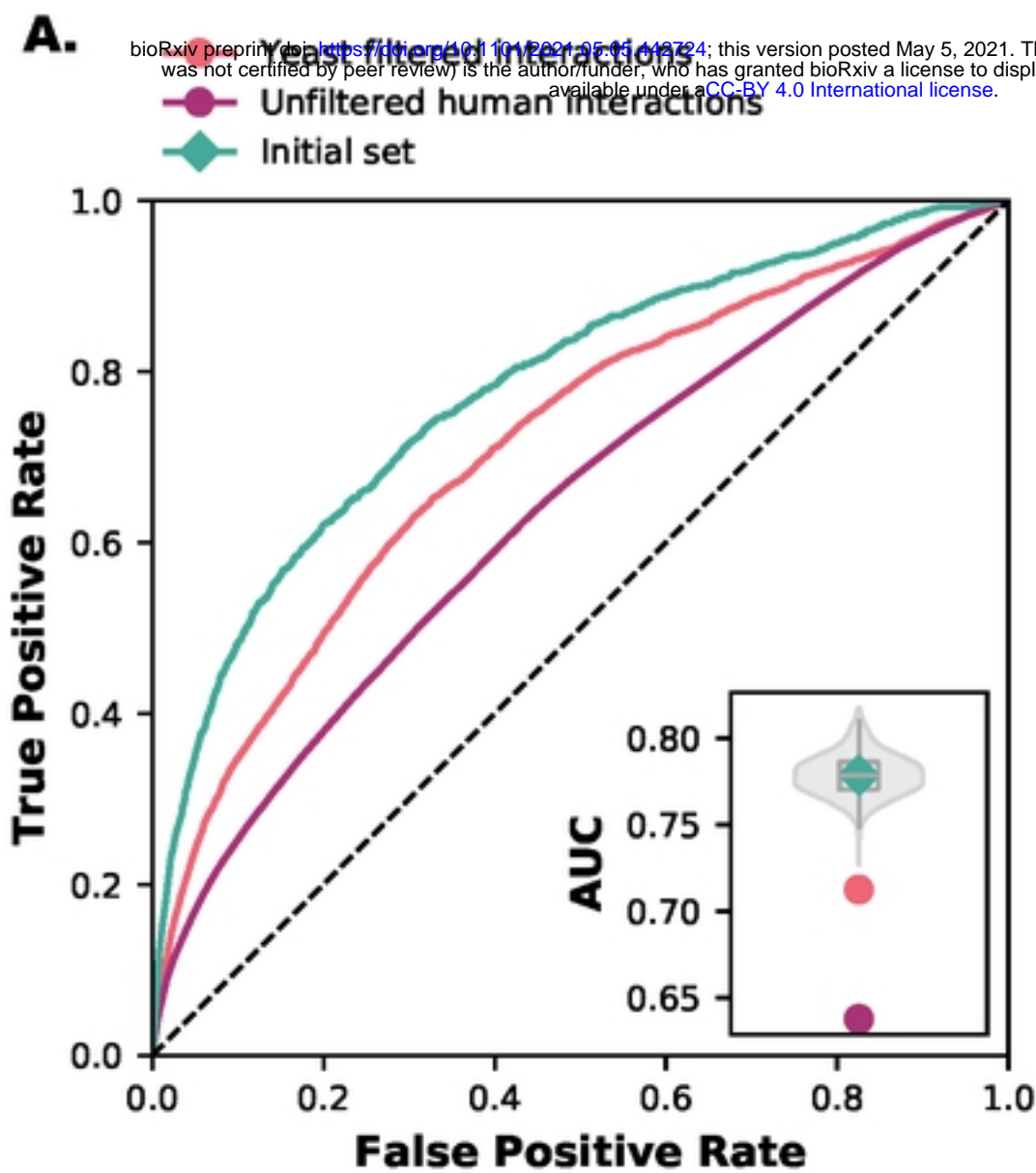


Figure 5