# Mixing genome annotation methods in a comparative analysis inflates the apparent number of lineage-specific genes

Caroline M. Weisman[1] [*], Andrew W. Murray[2], Sean R. Eddy[2,3,4]

[1] Lewis-Sigler Institute for Integrative Genomics, Princeton University, Princeton, NJ, USA [2] Department of Molecular & Cellular Biology, Harvard University, Cambridge MA, USA, [3] Howard Hughes Medical Institute, [4] John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge MA, USA

* Correspondence: cweisman@princeton.edu, @WeismanCara

## Summary

Comparisons of genomes of different species are used to identify lineage-specific genes, those genes that appear unique to one species or clade. Lineage-specific genes are often thought to represent genetic novelty that underlies unique adaptations. Identification of these genes depends not only on genome sequences, but also on inferred gene annotations. Comparative analyses typically use available genomes that have been annotated using different methods, increasing the risk that orthologous DNA sequences may be erroneously annotated as a gene in one species but not another, appearing lineage-specific as a result. To evaluate the impact of such "annotation heterogeneity," we identified four clades of species with sequenced genomes with more than one publicly available gene annotation, allowing us to compare the number of lineage-specific genes inferred when differing annotation methods are used to those resulting when annotation method is uniform across the clade. In these case studies, annotation heterogeneity increases the apparent number of lineage-specific genes by up to 15-fold, suggesting that annotation heterogeneity is a substantial source of potential artifact.

## Introduction

Comparing the genome sequences of different organisms can yield inferences about the genetic basis of the biological differences between them. One such analysis aims to identify genes unique to a particular monophyletic group. Such genes, called "orphan genes" when restricted to one species and "lineage-specific" or "taxonomically-restricted" when restricted to a clade of several species, are interesting from the perspective of genetic and evolutionary novelty. For example, they have been proposed to underlie lineage-specific structural and functional innovations, and to be novel genes that have emerged from noncoding DNA [1-5].

Lineage-specific genes are typically identified by searching for homologs in outgroup species: genes for which homologs cannot be found are considered lineage-specific. Such analyses typically begin not with raw genome sequences, but with particular "annotations" of them: inferences about what genes they encode. Often, only genes included in these annotations are considered in the homology search [2, 6, 7].

44         Previous work has recognized two ways in which errors in genome annotations could
45    produce spurious lineage-specific genes. A real gene could be annotated in the focal lineage, but
46    its homologs incorrectly unannotated in outgroups [8-10]. Conversely, a non-genic sequence
47    could be incorrectly annotated as a gene in the lineage, but correctly omitted in outgroups [11].
48    Such errors could occur even when all genomes in an analysis are consistently annotated by the
49    same annotation methodology, but the potential for error is expected to increase if genomes are
50    annotated by different methods, which use different criteria in determining which sequences are
51    genic.  Because comparative analyses typically depend on publicly available genomes whose
52    annotations come from different authors and sources, such "annotation heterogeneity" is
53    common [12-16]. Many gene annotation methods are in wide use, including custom pipelines at
54    large bioinformatics data providers (NCBI [17], Ensembl [18]), hand-curated model organism
55    annotation (Flybase [19], Wormbase [20]), crowd-sourced annotation (VectorBase [21]), and
56    various software packages (Maker [22], PASA [23]), used independently or in combination, with
57    custom parameters chosen by individual researchers.
58         Here we evaluate the impact of annotation heterogeneity on inferred numbers of lineage-
59    specific protein-coding genes. We identify four clades of species with available genome
60    sequences for which multiple different annotations are publicly available. These enable us to
61    conduct case studies in which we compare the number of lineage-specific genes when all species
62    are annotated with the same method ("uniform annotations") to when they are annotated with
63    different methods ("heterogeneous annotation"). We find that annotation heterogeneity
64    consistently and substantially increases the inferred number of lineage-specific genes. This effect
65    is strongest when all species within the lineage are annotated with one method and all outgroup
66    species with a different one. Our results suggest that annotation heterogeneity can produce many
67    spurious lineage-specific genes, potentially a majority of those found in a study.
68

69 **Results**

70

71 **Identification of clades of sequenced genomes with annotations from two methods**

72

73         To directly compare lineage-specific genes found using uniform annotations and
74    heterogeneous annotations, we manually searched the literature and bioinformatic databases for
75    species groups in which all species were annotated with the same method, and, additionally, the
76    same assembly of each species had been independently annotated with some other method. We
77    used existing annotations from a variety of standard sources instead of generating our own to
78    make results maximally representative of real studies. We identified four groups of five species:
79    cichlids, primates, bats, and rodents (Table 1, Supplemental Table 1). For cichlids and primates,
80    all five species were annotated with the same two methods, whereas for bats and rodents, one
81    method was applied to all five species, and the other available annotation was from three
82    different methods, with each species being annotated by one of the three. Each of these four
83    groups is less than approximately 60 My old.
84

85

86 **Different annotations of the same genome have many proteins unique to each method**
87

88  Spurious lineage-specific genes may result from annotation heterogeneity when different
89  annotation methods differentially annotate homologous sequences. Spurious lineage-specific
90  genes may also result from such erroneous differential annotation even when a single annotation
91  method is used, as sequence differences between the species may alter a given method's
92  determination regarding genic status. To get a sense of how many spurious lineage-specific
93  protein-coding genes annotation heterogeneity *per se* can produce, we compared two protein
94  annotations of the same species to identify proteins appearing to be unique to one of the
95  annotations. Because the underlying genome sequences are identical, any such apparently unique
96  proteins must be spurious, due only to annotation heterogeneity.
97  To mimic a typical analysis, for each species' two annotations, we used BLASTP [24] for
98  all proteins in one annotation to see if a significantly similar (E<0.001) homolog was present in
99  the other annotation. Between 0.6% and 9.7% of proteins in one annotation had no significantly
100  similar sequence in the other (Table 1). Of the 40 (20 pairs) annotations, 19 had over 1000
101  proteins without a significant homolog in the other annotation. In an extreme case of the cichlid
102  *Astatotilapia burtonii*, one annotation (Broad Institute) found 4110 genes that had no significant
103  similarities in the other (NCBI eukaryotic annotation pipeline), and 799 proteins in the NCBI
104  annotation lacked significant similarities in the Broad annotation. These substantial differences
105  between two annotations of one genome illustrate the potential for spurious lineage-specific
106  genes in comparisons of different genomes.
107
**108  Different patterns of annotation heterogeneity may differently affect the inferred number**
**109  of lineage-specific genes**
110
111  When different annotation methods are used for species within an analysis, different
112  patterns in which those methods are arranged on the species topology are possible. These
113  different patterns may differently affect the number of spurious lineage-specific genes produced
114  by annotation heterogeneity. In particular, because a gene is called as "lineage-specific" if no
115  significant homologs are found in any species outside the lineage, we expected that the number
116  of spurious lineage-specific genes would be positively related to the overall degree of difference
117  between the lineage and outgroup annotations.
118  We considered three such patterns. In the first, one annotation method is used for all
119  ingroup species (in the lineage, the gray boxes in the figures), and a different method for all
120  outgroup species (outside the lineage); we refer to this as "phyletic" annotation (Figure 1). In the
121  second, one method is used for all ingroup species, but a mixture of methods is used for the
122  outgroup species; we refer to this as "semi-phyletic" annotation (Figure 2). In the third, a mixture
123  of methods is used for both the ingroup species and the outgroup species; we refer to this as
124  "unpatterned" annotation (Figure 3). We used our four clades to create case studies for each
125  pattern.
126  The differences in annotation methods between ingroups and outgroups are largest for
127  phyletic annotation, intermediate for semi-phyletic annotation, and smallest for unpatterned
128  annotation; we expected the number of spurious lineage-specific genes to scale accordingly.
129
130
**131  Annotating a lineage with one method and outgroups with a different method greatly**
**132  increases the apparent number of lineage-specific genes**

133
134   Phyletic annotation occurs in at least two scenarios. Studies that newly sequence a
135 lineage often use their own method to annotate that lineage, and may then compare it to outgroup
136 annotations from another single source (e.g. Ensembl). Additionally, studies using existing
137 annotations may encounter a correlation between taxon and annotation method because genome
138 sequencing groups (with their annotation teams) often select species taxonomically (e.g. studies
139 of particular taxa, sequencing consortia/database initiatives for particular taxa) [25, 26].
140   We tested the impact of phyletic annotation on the apparent number of lineage-specific
141 genes on two groups of species, where the same genome assembly for every species had been
142 annotated by the same two methods: five cichlids, annotated both by the Broad Institute and
143 NCBI; and five primates, annotated both by Ensembl and NCBI (Supplemental Table 1).
144   For each tree of five species, we exploited the ladder-like topology (Figure 1) of the tree
145 to perform four analyses, comparing each of the four monophyletic groups including the focal
146 species to the remaining outgroups. For each lineage that included the focal species, we
147 conducted a typical analysis of lineage-specific genes by identifying genes in the focal species
148 that have a significantly similar homolog in the deepest rooted member of the ingroup (and thus
149 are "present" in that clade), but lack significant similarity to any protein in any outgroup species
150 in a BLASTP search (Methods). We compared the number of lineage-specific genes found when
151 all species (both ingroups and outgroups) were annotated with the same method to the number
152 found when the annotations for all outgroup species were switched to the other method in a
153 "phyletic" annotation pattern (Figure 1).
154   Heterogeneous annotation consistently caused a large increase of hundreds to thousands
155 of apparent lineage-specific genes, typically about a 4-fold (ranging from 1.4-fold to 15-fold)
156 difference relative to uniform annotation. In all but one of the eight cases in Figure 1, the
157 increase is more than 2-fold, suggesting that the majority of lineage-specific genes inferred in
158 heterogeneous annotations are artifacts of the heterogeneity.
159
160 **Annotating a lineage with one method and outgroups with a mixture of other methods**
161 **increases the apparent number of lineage-specific genes**
162
163   Examples of what we call "semi-phyletic" annotation, where the ingroup is annotated
164 with one method and outgroups with a mixture of methods, are common in the literature on
165 lineage-specific genes [12, 13, 26-31]. This can occur in scenarios similar to phyletic annotation,
166 but where outgroup annotations are available from a mixture of sources (e.g. a combination of
167 Ensembl and NCBI).We created case studies of semi-phyletic annotation using groups of species
168 for which every species had been annotated both by the same method and by one of a mix of
169 other methods: five rodents and five bats (Supplemental Table 1). We repeated the procedure
170 described for phyletic annotation above to compare the number of lineage-specific genes in
171 semi-phyletic annotations to those in uniform annotations (Figure 2).
172   Semi-phyletic annotation heterogeneity caused a smaller but still substantial increase in
173 the number of apparent lineage-specific genes in all lineages in both groups (Figure 2). The
174 magnitude of this effect ranged from 20 to 833 additional lineage-specific genes, corresponding
175 to 1.2-fold to 6-fold increases.
176

**Annotating species with a mixture of methods without taxonomic bias increases the apparent number of lineage-specific genes**

Examples of what we call "unpatterned" annotation, where the annotation method varies within the ingroup as well as the outgroup, are also common in the literature [15, 16, 27, 32-35]. This occurs when studies use existing available annotations for the desired species, which may come from a variety of sources. We created case studies of unpatterned annotation using the same rodent and bat species we used for semi-phyletic annotation (Figure 2), with the difference that we always compared the uniform annotations to the full set of mixed annotations (Figure 3) to produce unpatterned annotation heterogeneity.

Unpatterned annotation heterogeneity usually caused an increase in apparent lineage-specific genes (Figure 3), though the effect was smaller than for phyletic or semi-phyletic annotations. Two cases showed equal numbers or slight decreases, and the other six cases showed increases of 1.1-fold to 5.7-fold; the largest increases were in the cases with a single outgroup species.

**As expected, six-frame translation homology searches dramatically reduce the apparent number of lineage-specific genes**

A homology search in which the query protein is compared directly to a six-frame translation of the target genome does not rely on an annotation of the target species, and so should reduce this source of spurious lineage-specific genes. Such translated searches have previously been shown to reduce the inferred number of lineage-specific genes [8, 9]. In agreement with these expectations, we find that, for all of the lineages described above (depicted in Figures 1-3), a search for the focal species' proteins against six-frame translations of all comparator species genomes dramatically reduces the number of lineage-specific genes: to below the number inferred with uniform annotations, and often to less than one hundred (Supplemental Table 2).

## Discussion

We used six case studies to ask if varying the annotation method across species in a comparative analysis ("annotation heterogeneity") alters the apparent number of lineage-specific genes. We found that switching from uniform to heterogeneous annotations consistently increased the number of genes that were classified as lineage-specific, with increases ranging from tens to thousands of genes, corresponding to increases of up to 15-fold. The largest increases were seen when one annotation method was used for all the ingroup species and another was used for all the outgroup species ("phyletic annotation"). The smallest increases were seen when a mixture of annotation methods were used in both ingroup and outgroup species. Our case studies consist of trees of five species; mixtures of annotations in larger numbers of outgroup and ingroup species may reduce the artifact.

Annotation heterogeneity is common in comparative studies. Our results suggest that the numbers of lineage-specific genes found in these studies may be inflated, especially in "phyletic annotation" cases, and where the number of species compared is small. Annotation heterogeneity

221 may also have consequences that we do not explore here, like producing spurious lineage-
222 specific losses.
223      Recent work from us and others has shown that homology detection failure, in which
224 homology searches fail to detect homologs that are actually present in outgroups, can also
225 produce spurious lineage-specific genes [36, 37]. Previous studies have noted a surprisingly large
226 number of "young" lineage-specific genes found in recently evolved clades [15], which,
227 compared to older lineage-specific genes, are less readily explained by homology detection
228 failure, which is minimized at short evolutionary distances. The results here are all for young
229 (<60 My old) clades, showing that annotation heterogeneity can be a significant source of
230 spurious lineage-specific genes in young clades.
231      In accordance with previous results, we show that annotation heterogeneity artifacts can
232 be reduced by performing homology searches of six-frame translated genomic DNA sequence in
233 search of unannotated homologs in target species. This approach has caveats. At short
234 evolutionary distances, a sequence may be sufficiently similar for successful detection in such a
235 search without having the same coding status as the query; for example, a truly de novo
236 originated gene is expected to have significant nucleotide similarity to a homologous noncoding
237 locus in close outgroup species. This approach also still relies on an accurate annotation of the
238 focal species.
239      When annotation methods disagree, which is correct? Our results do not address this,
240 only demonstrating a consequence of this disagreement. Even homogeneous annotations are
241 imperfect. Of particular concern, methods in general rely on features (homology to known genes,
242 length, expression level, codon optimization) that seem likely to be absent or weaker in newly
243 evolved (*de novo*) genes, and so may fail to identify these genes. We consider annotation
244 accuracy primarily accountable to experimental data. Testing transcription, translation, and
245 function in all species in question is of ultimate importance in accurately identifying lineage-
246 specific genes. In light of our results, we suggest more emphasis on these metrics. In the
247 meantime, the true number of lineage-specific genes remains difficult to ascertain, but better
248 understanding sources of spurious ones helps us constrain it.
249

## Author Contributions

251
252 Conceptualization, C.W.; Formal Analysis, C.W.; Investigation, C.W.; Writing – Original Draft:
253 C.W.; Writing – Review & Editing: C.W., A.W.M, S.R.E; Supervision, A.W.M, S.R.E; Funding
254 Acquisition, S.R.E.
255

## Declaration of interests

257
258 The authors declare no competing interests.
259

## Methods

261

### Identifying lineage-specific proteins

263      For each species group, we defined a protein as specific to a particular lineage if a search
264 using BLASTP [24] version 6.2.0 had no similar protein at a significance threshold of E=0.001

265 in the annotation of any species that was an outgroup to that lineage. We did not require that a
266 protein be present in all members of the lineage to be specific to that lineage: a protein was
267 defined as specific to a lineage based on the most distant species in which it was detected. For
268 example, if a protein in *M. musculus* was detected only in *R. norvegicus*, it was defined as
269 specific to that lineage; if a gene in *M. musculus* was detected in *M. caroli, M. pahari, and R.*
270 *norvegicus,* it was also defined as specific to that same lineage. If a protein was found in the
271 earliest-branching member of the species group, it was considered "conserved" and so not
272 counted as any kind of lineage-specific gene. This way of classifying lineage-specificity coheres
273 with standard practice [6].
274      For the six-frame translated searches, we first generated a six-frame translation of the
275 genome assembly of each species using the 'esl-translate' command in the hmmer easel package,
276 and then used it as the target database in a BLASTP search, as described in the previous
277 paragraph.
278

## Supplemental Information

280

281 Supplemental Table 1: Sources, brief descriptions, and links to protein annotations and genome
282 assemblies used in this study.

283

284 Supplemental Table 2: Results of six-frame translation homology searches. Numbers in the table
285 indicate the inferred number of genes specific to the indicated lineage (corresponding to the four
286 lineages depicted in Figures 1-3) in each of the described taxa.

287

## Data availability

289

290 All raw results summarized in Figures 1-3 are available at
291 https://github.com/caraweisman/Annotation_homology.
292
293

294 **Tables**

295

| Species | Annotation 1 source | No. proteins in annotation 1 | Number/percent of proteins in annotation 1 w/ no homologs in annotation 2 | Annotation 2 source | No. proteins in annotation 2 | Number/percent of proteins in annotation 2 w/ no homologs in annotation 1 |
|---|---|---|---|---|---|---|
| | | | | | | |
| **Cichlid fish** | | | | | | |
| *Metraclimia zebra* | Broad Institute | 51772 | 3592/6.9% | NCBI Eukaryotic annotation pipeline | 40043 | 706/1.8% |
| *Pundamilia nyererei* | Broad Institute | 42152 | 3276/7.7% | NCBI Eukaryotic annotation pipeline | 38583 | 668/1.7% |
| *Astatotilapia burtoni* | Broad Institute | 52845 | 4110/7.8% | NCBI Eukaryotic annotation pipeline | 44653 | 799/1.8% |
| *Neolamprologus brichardi* | Broad Institute | 36873 | 3568/9.7% | NCBI Eukaryotic annotation pipeline | 31372 | 755/2.4% |
| *Oreochromis niloticus* | Broad Institute | 66482 | 4143/6.2% | NCBI Eukaryotic annotation pipeline | 47713 | 765/1.6% |
| | | | | | | |
| **Primates** | | | | | | |
| *Macaca fascicularis* | Ensembl "full genebuild" | 46148 | 1510/3.3% | NCBI Eukaryotic annotation pipeline | 62672 | 1044/1.7% |
| *Macaca nemestrina* | Ensembl "full genebuild" | 46238 | 1295/2.8% | NCBI Eukaryotic annotation pipeline | 66484 | 1623/2.4% |
| *Mandrillus leucophaeus* | Ensembl "full genebuild" | 40903 | 1406/3.4% | NCBI Eukaryotic annotation pipeline | 38336 | 693/1.8% |
| *Rhinopithecus bieti* | Ensembl "full genebuild" | 43730 | 1233/2.8% | NCBI Eukaryotic annotation pipeline | 49595 | 1476/3.0% |
| *Cebus imitator* | Ensembl "full genebuild" | 40677 | 926/1.7% | NCBI Eukaryotic | 55885 | 602/1.5% |

8

preprint doi: https://doi.org/10.1101/2022.01.13.476251; this version posted January 15, 2022. 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-NC-ND 4.0 International license.

| | | | | annotation pipeline | | |
|---|---|---|---|---|---|---|
| | | | | | | |
| | | | | | | |
| **Rodents** | | | | | | |
| *Mus musculus* | Ensembl "full genebuild" | 68381 | 1523/2.2% | NCBI Eukaryotic annotation pipeline | 84985 | 471/0.6% |
| *Mus caroli* | UCSC | 50492 | 1496/3.0% | NCBI Eukaryotic annotation pipeline | 47409 | 590/1.2% |
| *Mus pahari* | UCSC | 50002 | 1596/3.2% | NCBI Eukaryotic annotation pipeline | 42226 | 413/1.0% |
| *Rattus norwegicus* | Ensembl "mixed genebuild" | 29107 | 458/1.6% | NCBI Eukaryotic annotation pipeline | 56110 | 716/1.3% |
| *Peromyscus maniculatus* | Ensembl "full genebuild" | 28866 | 267/0.9% | NCBI Eukaryotic annotation pipeline | 45588 | 517/1.1% |
| | | | | | | |
| | | | | | | |
| **Bats** | | | | | | |
| *Myotis lucifugus* | Ensembl "full genebuild" | 20719 | 197/1.0% | NCBI Eukaryotic annotation pipeline | 43106 | 622/1.4% |
| *Myotis brandtii* | Beijing Genomics Institute | 19484 | 867/4.5% | NCBI Eukaryotic annotation pipeline | 40808 | 1370/3.4% |
| *Myotis myotis* | Bat1K | 46057 | 3448/7.5% | NCBI Eukaryotic annotation pipeline | 61156 | 704/1.2% |
| *Molossus molossus* | Bat1K | 53797 | 3107/5.8% | NCBI Eukaryotic annotation pipeline | 42753 | 486/1.1% |
| *Pteropus alecto* | Beijing Genomics Institute | 19619 | 1338/6.8% | NCBI Eukaryotic annotation pipeline | 39706 | 955/2.4% |

296
297

298    **Table 1:** Genome annotations used in this study. Brief description of annotation source, number
299    of genes in the annotation, and the number and percentage of genes in each annotation with no
300    significant homologs found by a BLASTP search in the other annotation for the given species are
301    listed. Note that, where large differences in the number of proteins included in a pair of
302    annotations occurs, this is often due in part to one annotation including a larger number of
303    different isoforms of the same locus, all or many of which may have significant similarity to the
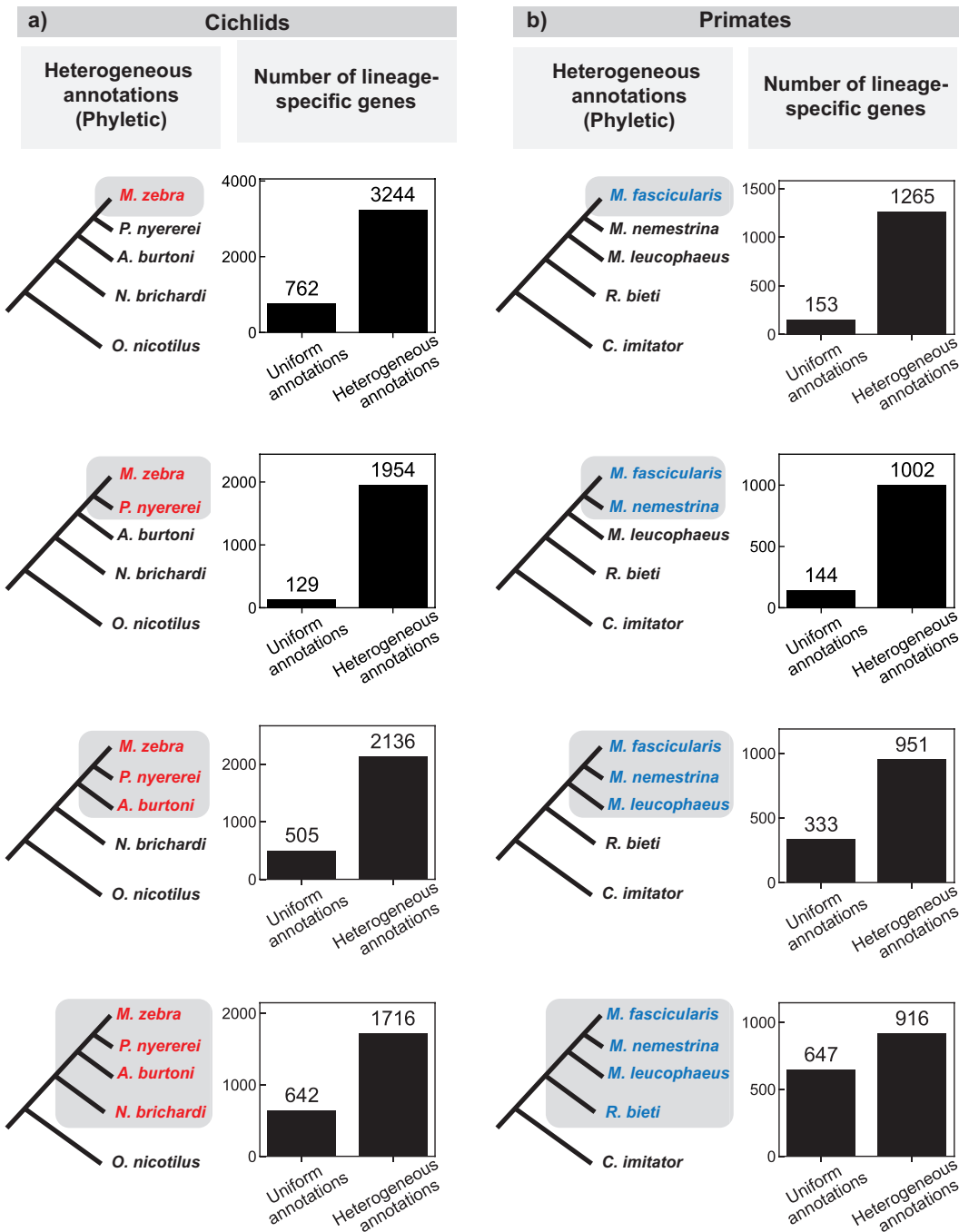304    same protein(s) in the other annotation.
305

306
307 **Figures**
308



309
310
311
312
313

314 **Figure 1: Comparison of the number of lineage-specific genes found using uniform and**
315 **heterogeneous (phyletic) annotations in a) cichlids and b) primates.** The species tree on the
316 left indicates the lineage under consideration (grey shading); different text colors indicate
317 different annotation sources in the heterogeneous annotation analysis (black, NCBI; red, research
318 group at the Broad Institute; blue, Ensembl). A depiction of the uniform annotation pattern, in
319 which all annotations are from NCBI (black), is not shown. Bar graphs indicate the number of
320 genes that appear specific to the lineage shaded on the species tree to the left using either
321 uniform or heterogeneous annotations.
322
323
324
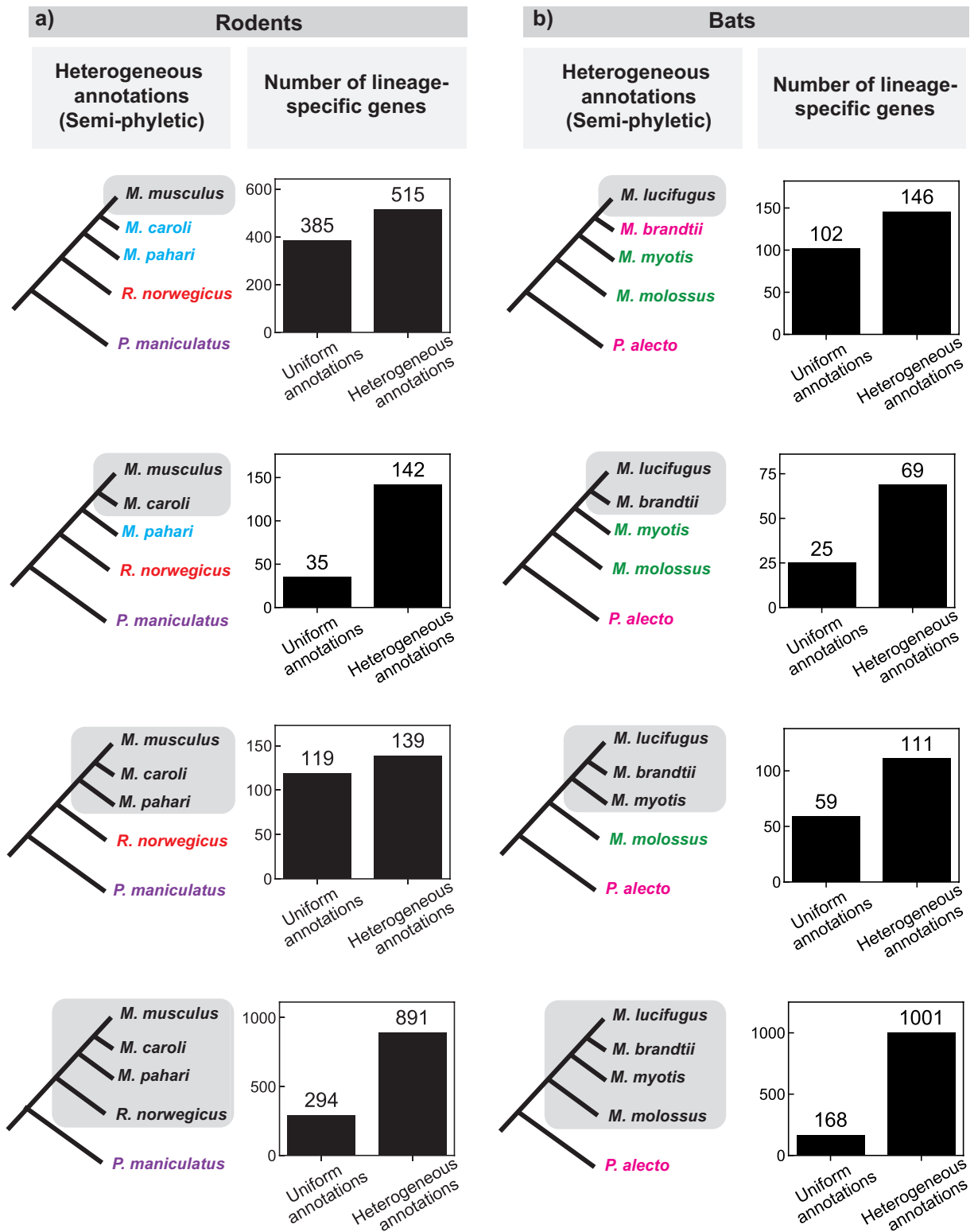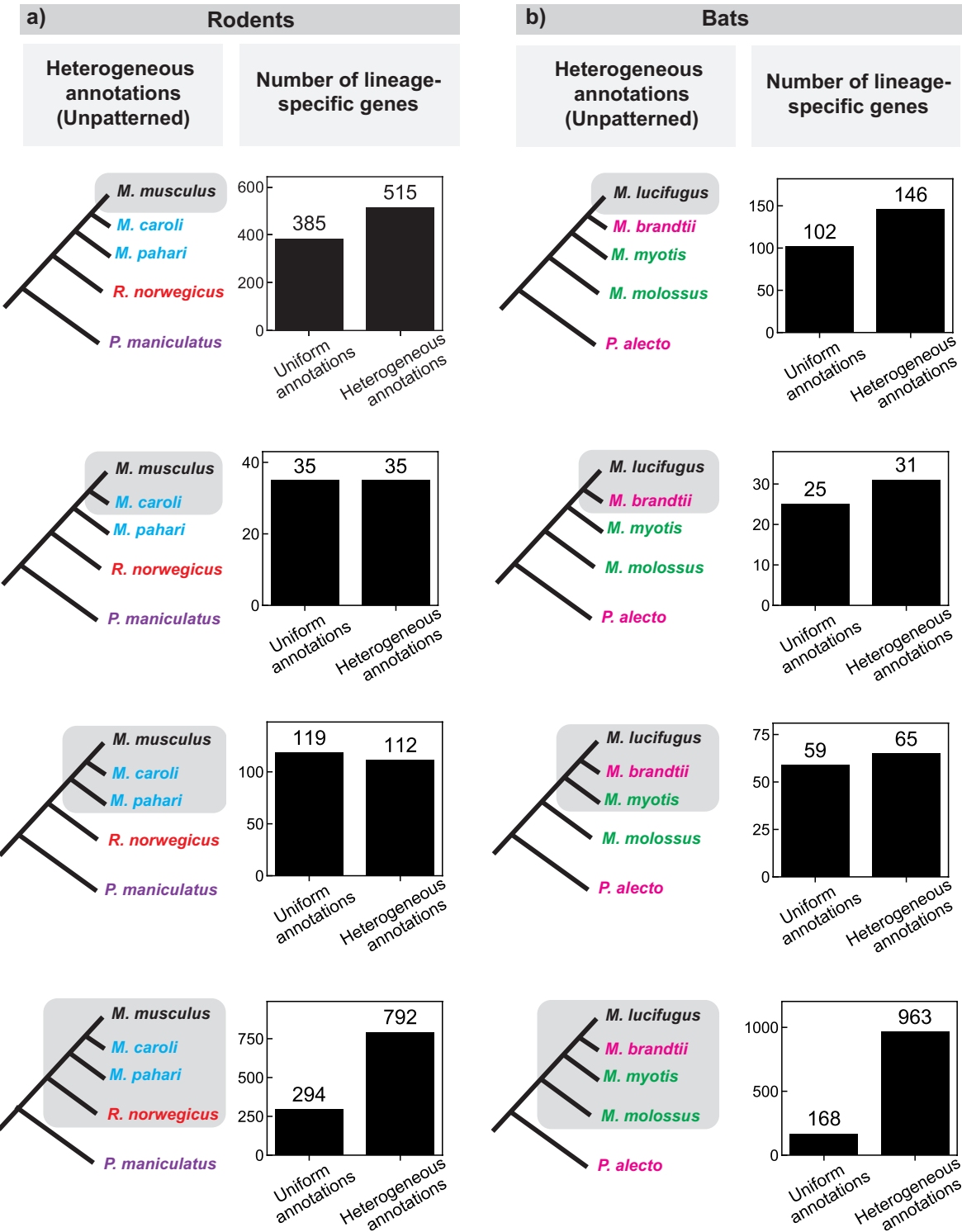325
326

327

328
329
330    **Figure 2: Comparison of the number of lineage-specific genes found using uniform and**
331    **heterogeneous (semi-phyletic) annotations in a) rodents and b) bats.** The species tree on the
332    left indicates the lineage under consideration (grey shading); different text colors indicate
333    different annotation sources in the heterogeneous annotation analysis (black, NCBI; blue, UCSC;
334    red, Ensembl "mixed genebuild"; purple, Ensembl "full genebuild"; green, Bat1k; pink, Beijing
335    Genomics Institute). A depiction of the uniform annotation pattern, in which all annotations are
336    from NCBI (black), is not shown. Bar graphs indicate the number of genes that appear specific to
337    the lineage shaded on the species tree to the left using either uniform or heterogeneous
338    annotations.
339
340
341

342
343
344

345 **Figure 3: Comparison of the number of lineage-specific genes found using uniform and**
346 **heterogeneous (unpatterned) annotations in a) rodents and b) bats.** The species tree on the
347 left indicates the lineage under consideration (grey shading); different text colors indicate
348 different annotation sources in the heterogeneous annotation analysis (black, NCBI; blue, UCSC;
349 red, Ensembl "mixed genebuild"; purple, Ensembl "full genebuild"; green, Bat1k; pink, Beijing
350 Genomics Institute). A depiction of the uniform annotation pattern, in which all annotations are
351 from NCBI (black), is not shown. Bar graphs indicate the number of genes that appear specific to
352 the lineage shaded on the species tree to the left using either uniform or heterogeneous
353 annotations.
354

## References

1. Khalturin K, Hemmrich G, Fraune S, Augustin R, Bosch TC. More than just orphans: are taxonomically-restricted genes important in evolution? Trends in Genetics. 2009;25:404-13.

2. Tautz D, Domazet-Lošo T. The evolutionary origin of orphan genes. Nature Reviews Genetics. 2011;12:692.

3. Wilson G, Bertrand N, Patel Y, Hughes J, Feil E, Field D. Orphans as taxonomically restricted and ecologically important genes. Microbiology. 2005;151:2499-501.

4. McLysaght A, Guerzoni D. New genes from non-coding sequence: the role of de novo protein-coding genes in eukaryotic evolutionary innovation. Philosophical Transactions of the Royal Society B: Biological Sciences. 2015;370:20140332.

5. Tautz D. The discovery of de novo gene evolution. Perspectives in biology and medicine. 2014;57:149-61.

6. Domazet-Lošo T, Brajković J, Tautz D. A phylostratigraphy approach to uncover the genomic history of major adaptations in metazoan lineages. Trends in Genetics. 2007;23:533-9.

7. McLysaght A, Hurst LD. Open questions in the study of de novo genes: what, how and why. Nature Reviews Genetics. 2016;17:567.

8. Basile W, Elofsson A. The number of orphans in yeast and fly is drastically reduced by using combining searches in both proteomes and genomes. bioRxiv. 2017:185983.

9. Casola C. From de novo to "de nono": the majority of novel protein-coding genes identified with phylostratigraphy are old genes or recent duplicates. Genome Biology and Evolution. 2018;10:2906-18.

10. Zile K, Dessimoz C, Wurm Y, Masel J. Only a single taxonomically restricted gene family in the Drosophila melanogaster subgroup can be identified with high confidence. Genome Biology and Evolution. 2020.

11. Wilson BA, Foy SG, Neme R, Masel J. Young genes are highly disordered as predicted by the preadaptation hypothesis of de novo gene birth. Nature ecology & evolution. 2017;1:1-6.

12. Carvunis A-R, Rolland T, Wapinski I, Calderwood MA, Yildirim MA, Simonis N, et al. Proto-genes and de novo gene birth. Nature. 2012;487:370.

13. Vakirlis N, Hebert AS, Opulente DA, Achaz G, Hittinger CT, Fischer G, et al. A molecular portrait of de novo genes in yeasts. Molecular Biology and Evolution. 2017;35:631-45.

14. Bowles AM, Bechtold U, Paps J. The Origin of Land Plants Is Rooted in Two Bursts of Genomic Novelty. Current Biology. 2020.

15. Neme R, Tautz D. Phylogenetic patterns of emergence of new genes support a model of frequent de novo evolution. BMC Genomics. 2013;14:117.

16. Paps J, Holland PW. Reconstruction of the ancestral metazoan genome reveals an increase in genomic novelty. Nature communications. 2018;9:1-8.

17. Thibaud-Nissen F, DiCuccio M, Hlavina W, Kimchi A, Kitts P, Murphy T, et al. P8008 the NCBI eukaryotic genome annotation pipeline. Journal of Animal Science. 2016;94:184-.

18. Howe KL, Achuthan P, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, et al. Ensembl 2021. Nucleic acids research. 2021;49:D884-D91.

19. Drysdale R, Consortium F. FlyBase. Drosophila. 2008:45-59.

397    20.    Howe K, Davis P, Paulini M, Tuli MA, Williams G, Yook K, et al., editors. WormBase:
398    annotating many nematode genomes. Worm; 2012: Taylor & Francis.
399    21.    Giraldo-Calderón GI, Emrich SJ, MacCallum RM, Maslen G, Dialynas E, Topalis P, et al.
400    VectorBase: an updated bioinformatics resource for invertebrate vectors and other organisms
401    related with human diseases. Nucleic acids research. 2015;43:D707-D13.
402    22.    Cantarel BL, Korf I, Robb SM, Parra G, Ross E, Moore B, et al. MAKER: an easy-to-use
403    annotation pipeline designed for emerging model organism genomes. Genome research.
404    2008;18:188-96.
405    23.    Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic
406    gene structure annotation using EVidenceModeler and the Program to Assemble Spliced
407    Alignments. Genome biology. 2008;9:1-22.
408    24.    Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.
409    Journal of molecular biology. 1990;215:403-10.
410    25.    Schmitz JF, Chain FJ, Bornberg-Bauer E. Evolution of novel genes in three-spined
411    stickleback populations. Heredity. 2020;125:50-9.
412    26.    Prabh N, Rödelsperger C. De novo, divergence, and mixed origin contribute to the
413    emergence of orphan genes in pristionchus nematodes. G3: Genes, Genomes, Genetics.
414    2019;9:2277-86.
415    27.    Li Z-W, Chen X, Wu Q, Hagmann J, Han T-S, Zou Y-P, et al. On the origin of de novo genes
416    in Arabidopsis thaliana populations. Genome biology and evolution. 2016;8:2190-202.
417    28.    Heames B, Schmitz J, Bornberg-Bauer E. A continuum of evolving de novo genes drives
418    protein-coding novelty in Drosophila. Journal of molecular evolution. 2020;88:382-98.
419    29.    Aguilera F, McDougall C, Degnan BM. Co-option and de novo gene evolution underlie
420    molluscan shell diversity. Molecular Biology and Evolution. 2017;34:779-92.
421    30.    Huang J, Chen J, Fang G, Pang L, Zhou S, Zhou Y, et al. Two novel venom proteins
422    underlie divergent parasitic strategies between a generalist and a specialist parasite. Nature
423    communications. 2021;12:1-16.
424    31.    Wang Y-W, Hess J, Slot JC, Pringle A. De Novo Gene Birth, Horizontal Gene Transfer, and
425    Gene Duplication as Sources of New Gene Families Associated with the Origin of Symbiosis in
426    Amanita. Genome biology and evolution. 2020;12:2168-82.
427    32.    Papanicolaou A, Schetelig MF, Arensburger P, Atkinson PW, Benoit JB, Bourtzis K, et al.
428    The whole genome sequence of the Mediterranean fruit fly, Ceratitis capitata (Wiedemann),
429    reveals insights into the biology and adaptive evolution of a highly invasive pest species.
430    Genome biology. 2016;17:1-31.
431    33.    Forêt S, Knack B, Houliston E, Momose T, Manuel M, Quéinnec E, et al. New tricks with
432    old genes: the genetic bases of novel cnidarian traits. Trends in Genetics. 2010;26:154-8.
433    34.    Johnson BR, Tsutsui ND. Taxonomically restricted genes are associated with the
434    evolution of sociality in the honey bee. BMC Genomics. 2011;12:164.
435    35.    Donoghue MT, Keshavaiah C, Swamidatta SH, Spillane C. Evolutionary origins of
436    Brassicaceae specific genes in Arabidopsis thaliana. BMC evolutionary biology. 2011;11:1-23.
437    36.    Weisman CM, Murray AW, Eddy SR. Many, but not all, lineage-specific genes can be
438    explained by homology detection failure. PLoS biology. 2020;18:e3000862.

439    37.    Moyers BA, Zhang J. Evaluating phylostratigraphic evidence for widespread de novo
440    gene birth in genome evolution. Molecular biology and evolution. 2016;33:1245-56.
441