

1

## Article: Discoveries

2 Species tree disequilibrium positively misleads models of gene family evolution

3 Authors:

4 M. Elise Lauterbur, Department of Ecology and Evolution, Stony Brook University\*

5 Sarah Heder, Department of Ecology and Evolution, Stony Brook University

6 Laurel R. Yohe, Department of Ecology and Evolution, Stony Brook University\*\*

7 Liliana M. Dávalos, Department of Ecology and Evolution and Consortium for Inter-  
8 Disciplinary Environmental Research, Stony Brook University

9 \* Current address: Department of Ecology and Evolutionary Biology, University of  
10 Arizona

11 \*\* Current address: Department of Geology and Geophysics, Yale University

12 Corresponding Author:

13 M. Elise Lauterbur, [lauterbur@gmail.com](mailto:lauterbur@gmail.com)

## 14 Abstract

15 Gene duplication is a key source of evolutionary innovation, and multigene families  
16 evolve in a birth-death process, continuously duplicating and pseudogenizing through  
17 time. To empirically test hypotheses about adaptive expansion and contraction of  
18 multigene families across species, models infer gene gain and loss in light of speciation  
19 events and these inferred gene family expansions may lead to interpretations of  
20 adaptations in particular lineages. While the relative abundance of a gene subfamily in  
21 the subgenome may reflect its functional importance, tests based on this expectation  
22 can be confounded by the complex relationship between the birth-death process of  
23 gene subfamily evolution and the species phylogeny. Using simulations, we confirmed  
24 tree heterogeneity as a confounding factor in inferring multi-gene adaptation, causing  
25 spurious associations between shifts in birth-death rate and lineages with higher  
26 branching rates. We then used the *olfactory receptor (OR)* repertoire, the largest gene  
27 family in the mammalian genome, of different bat species with divergent diets to test  
28 whether expansions in olfactory receptors are associated with shifts to frugivorous diets.  
29 After accounting for tree heterogeneity, we robustly inferred that certain *OR* subfamilies  
30 exhibited expansions associated with dietary shifts to frugivory. Taken together, these  
31 results suggest ecological correlates of individual *OR* gene subfamilies can be  
32 identified, setting the stage for detailed inquiry into within-subfamily functional  
33 differences.

## 34 Introduction

35 Gene duplication is a key source of evolutionary innovation, facilitating the evolution of  
36 new functions of duplicated genes, and even the acquisition of new and specific  
37 biological roles (Assis and Bachtrog 2013). In contrast with most protein-coding regions  
38 of the genome for which duplication is relatively rare, certain regions duplicate  
39 frequently (Nei 1969) and, released from purifying selection, mutation pseudogenizes  
40 some of the copies over time (Nei and Hughes 1991) generating multigene families.  
41 Unlike single copy genes, multigene families evolve through a process of duplication  
42 (birth) and pseudogenization and deletion (death). Assuming multiple copies as a  
43 starting point, an equilibrium between births and deaths will result in turnover in the  
44 composition and identity of the genes (Han et al. 2013), independent of gene identity  
45 and function. Since multigene families evolving through the birth-death process encode  
46 for important functions such as chemosensation (Nei et al. 1997), immune defense (Nei  
47 et al. 1997), and even development (Nei and Rooney 2005), their evolutionary rates and  
48 shifts in duplication or loss rates are of broad interest in molecular evolution, particularly  
49 when shifts correspond to major ecological transitions.

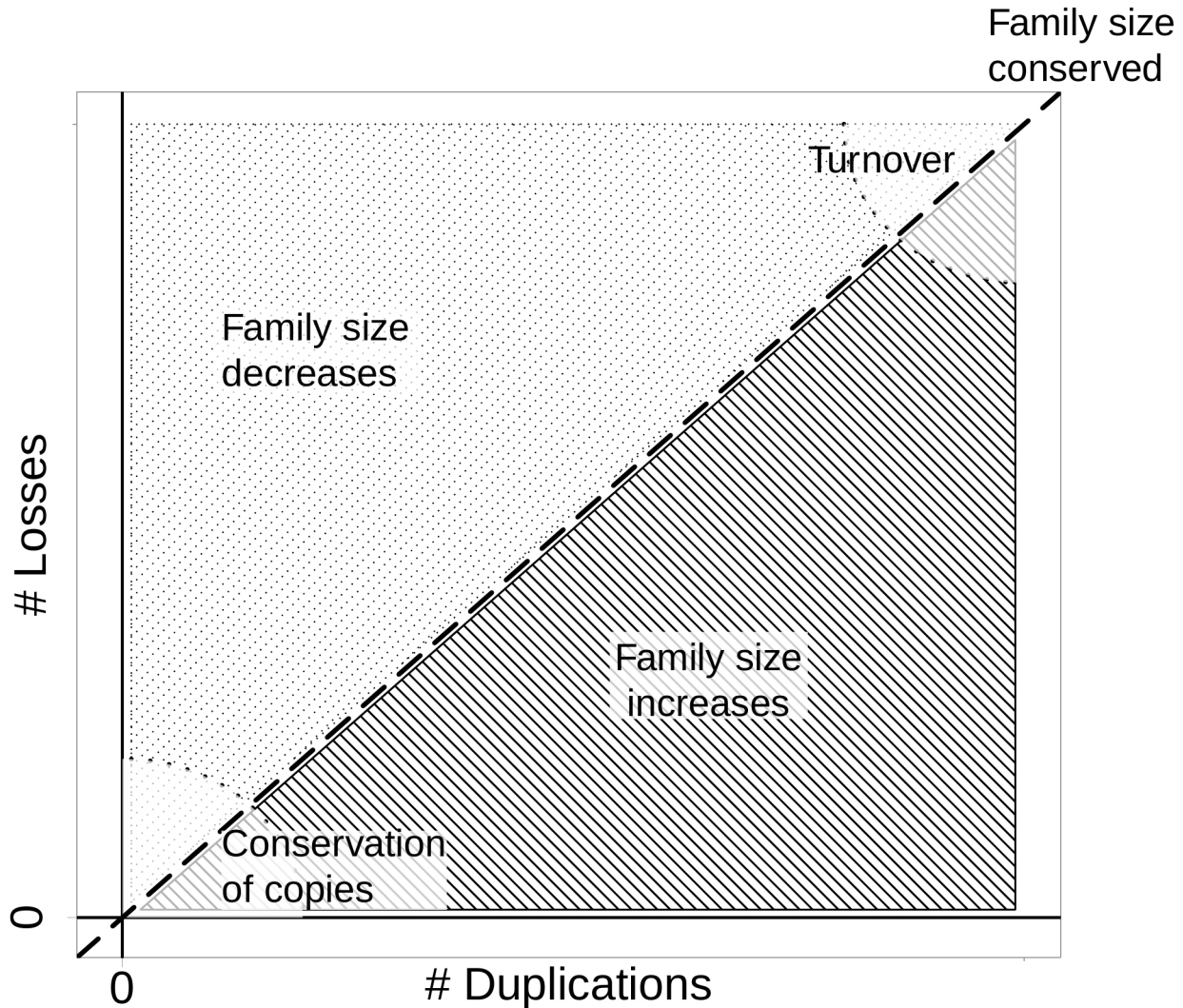
50 There are two main approaches to infer changes in duplication and loss: (1) estimating  
51 the rates of birth and death along branches of the species tree inferred from the number  
52 of gene copies present in each species and (2) and gene tree-species tree  
53 reconciliation to infer incongruences that suggest a duplication or loss event (Yohe et al.  
54 2019). In the former approach, multi-gene family evolution models typically  
55 parameterize a probability distribution over the counts of the number of genes within the  
56 family from the tips of the species tree (Bie et al. 2006). This parameterization thus  
57 integrates the birth-death model and the structure of the phylogeny, making it possible  
58 to test for equilibrium in the birth-death model for a gene family, as well as pinpointing  
59 the specific branch on the tree responsible for disequilibrium, if found. In similar  
60 integrated models, however, disequilibrium can arise from heterogeneity in the birth-  
61 death process (the intended hypothesis test), in the branching process of the phylogeny  
62 (not intended), or both, thus confounding tree and gene subfamily disequilibria. First

63 identified in binary state-dependent speciation (-SSE) models (Rabosky and Goldberg  
64 2015), this problem has been found to be more general (Rojas et al. 2018) when there  
65 is heterogeneity in the branching process of the phylogeny (Beaulieu and O'Meara  
66 2016). In other words, a rejection of the null emerges not from the evolution of the trait  
67 of interest, but from rejecting the embedded assumption that a single lineage branching  
68 rate applies to the entire phylogeny when it does not (Beaulieu and O'Meara 2016).  
69 Although hitherto unexplored for inference of multi-gene family evolution, previous  
70 findings for -SSE models suggest caution in applying hypothesis tests using birth-death  
71 models to phylogenies with known heterogeneity in their branching pattern.

72 To determine how species tree heterogeneity affects birth-death models for multigene  
73 families, particularly when linked to ecological diversity, three components are needed:  
74 (1) a diverse multigene family across species with divergent ecologies; (2) shifts in  
75 ecological function that may relate to gene family evolution; and (3) a system with  
76 strong shifts in diversification rates that may lead to tree heterogeneity. Olfactory  
77 receptor (*OR*) proteins are encoded by precisely such a diverse gene family, comprising  
78 the largest protein-coding fraction of a given mammalian genome. Although rich OR  
79 repertoires have been recorded across vertebrates (Vandeweghe et al. 2016), the  
80 greatest OR diversity is found in mammals, with hundreds of *olfactory receptor (OR)*  
81 genes encoded in tandem arrays of similar copies (Niimura and Nei 2003). If  
82 mammalian ecology changes such that the sense of smell is no longer usable, as  
83 among many aquatic mammals, disequilibrium ensues with higher rates of loss. In an  
84 extreme case, odontocete whales — dolphins, porpoises, and sperm whales among  
85 others — lack an olfactory bulb, as well as other morphological and molecular correlates  
86 of olfaction, and a high proportion of their olfactory subgenome is pseudogenized (74-  
87 100%) compared to terrestrial relatives (McGowen et al. 2014). In contrast, when  
88 function is important, selection will act against the random pseudogenization or deletion  
89 of copies, resulting in neofunctionalization of the duplicated gene or the retention of  
90 more functional gene copies and evidenced by expansion of specific gene subfamilies.  
91 As a result, comparing the relative sizes of gene subfamilies across species and  
92 inferring the corresponding duplications and losses can yield evidence for selection on

93 particular gene subfamilies, or shifts in selection therein (figure 1).

94 We use the *OR* gene family in bats, the second most diverse mammalian order with  
95 diverse ecological specializations previously linked to *OR* evolution, to test the  
96 association between divergent ecologies and olfactory receptors and evaluate model  
97 performance. Bats (Mammalia: Chiroptera) are an ideal system in which to test this -  
98 while most bats are insectivorous and strongly rely on echolocation to find food  
99 resources, several lineages of bats have independently evolved to feed on plants  
100 (Jones et al. 2005). Behavioral studies suggest they use scents to find fruiting trees  
101 (Korine and Kalko 2005). Previous analyses of the bat olfactory subgenome related  
102 specific gene subfamilies to the evolution of frugivory (Hayden et al. 2014), however  
103 tree heterogeneity was unexplored as a source of error. Since then, heterogeneity in  
104 diversification across the entire bat phylogeny was traced to a single branch that also  
105 corresponds to a shift toward a primarily frugivorous diet (Shi and Rabosky 2015), and  
106 misleading effects of this heterogeneity on -SSE analyses have been described (Rojas  
107 et al. 2018).



108 Figure 1. Relationship between the duplications, losses and multi-gene family size. The  
109 dashed line indicates equilibrium between losses and duplications, turnover indicates  
110 changes in the identity of genes over time resulting from high duplication and loss rates.

111 Given the potential confounding effects of tree heterogeneity on birth-death equilibrium,  
112 we used simulations to determine if such an effect was present in models integrating the  
113 birth-death process with the phylogeny. We also used two approaches, both based on a  
114 combination of gene and species trees, to infer adaptation in olfactory receptors to the  
115 new frugivorous diet: Poisson mixed models that estimate rates of birth and death

116 (Sackton et al. 2017), and phylogenetic instability analyses that uses gene-tree/species-  
117 tree reconciliation (Curran et al. 2018). We hypothesize that one or more *OR* gene  
118 families involved in shifts to frugivory have expanded in number relative to other *OR*  
119 gene families across bats. Changes linked to frugivory could then be identified as a shift  
120 in duplication or loss rate (increased duplication or decreased loss in frugivorous vs.  
121 non-frugivorous bats), or unusual discordance between these gene trees and the  
122 species tree. The results confirmed tree heterogeneity as a confounding factor in testing  
123 multi-gene adaptation, identified gene subfamilies experiencing higher turnover, and  
124 uncovered high discordance linked to the evolution of frugivory.

## 125 Results

### 126 Influence of topology on CAFE results

127 We randomly permuted observed olfactory receptor (*OR*) gene copy numbers on  
128 inferred, non-Yule species trees of Yangochiroptera and Yinpterochiroptera, comparing  
129 single- $\lambda$  (single-rate) and two- $\lambda$  (two-rate) models. This revealed that CAFE has a false  
130 positive rate on the yangochiropteran tree between 26% and 69%, and between 0% and  
131 76% on the yinpterochiropteran tree (table 1). When branch lengths were Yule-  
132 transformed, the false positive rates for the yangochiropteran tree increased for all but  
133 one *OR* gene subfamily (gene subfamily 4) to between 44% and 91%. The false positive  
134 rates for the Yule-transformed yinpterochiropteran tree decreased to 0% - 0.1%.

135 Table 1. Proportion of false positives by tree and gene subfamily.

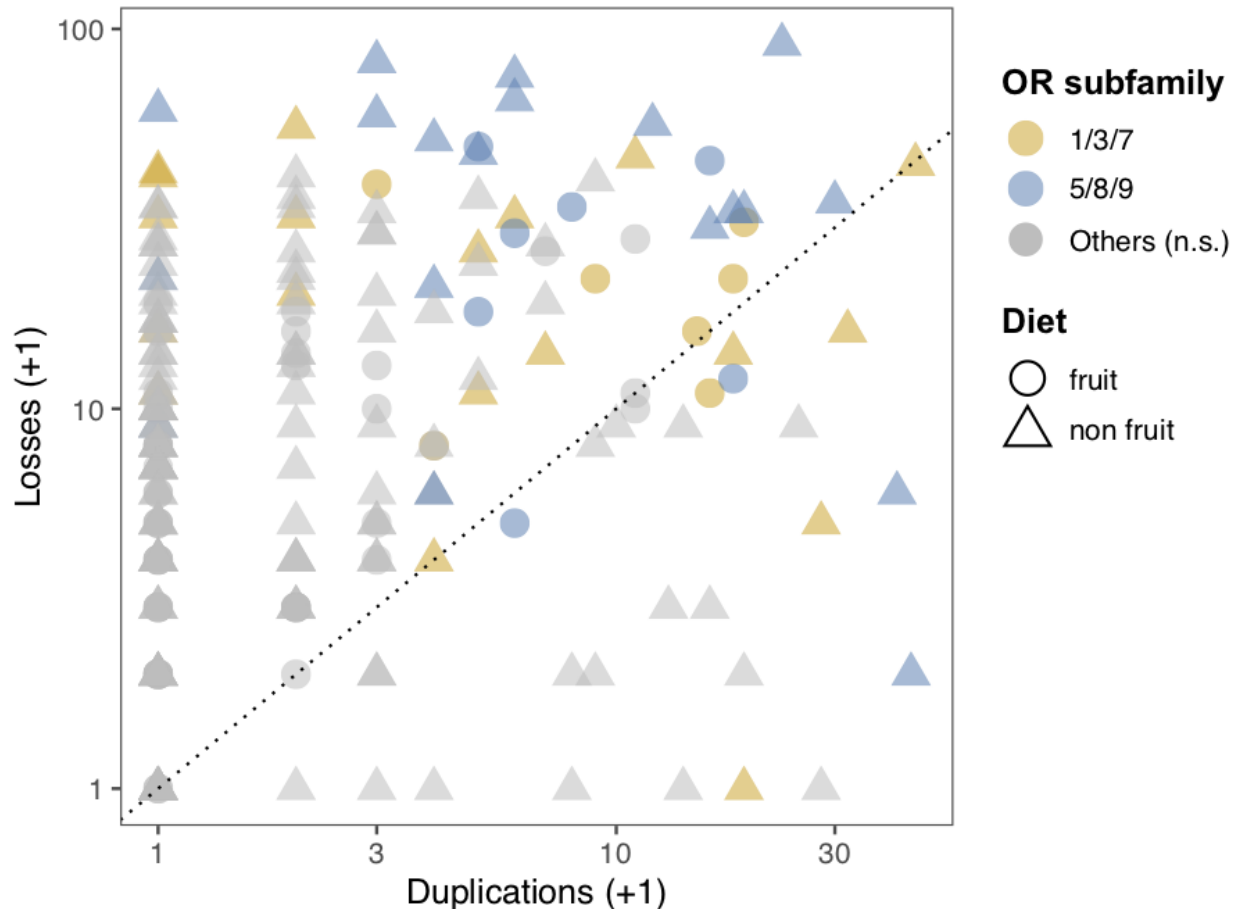
Gene	Yangochiroptera	Yinpterochiroptera	Yule-transformed Yangochiroptera	Yule-transformed Yinpterochiroptera
1/3/7	0.36	0	0.88	0
2/13	0.28	0.37	0.55	0
4	0.46	0.74	0.44	0.001
5/8/9	0.26	0.22	0.86	0

6	0.57	0.70	0.90	0
11	0.69	0.74	0.91	0
52	0.65	0.76	0.57	0

136 Duplication and Loss Events (Notung)

137 We used Notung to estimate the number of gene duplication and loss events for each  
138 OR gene subfamily. The number of duplication events inferred per OR gene subfamily  
139 per species or internal branch ranged from 0 to 44 (maximum in OR gene subfamily  
140 1/3/7 in *Anoura geoffroyi*), and the number of loss events ranged from 0 to 90  
141 (maximum in OR gene subfamily 5/8/9 in *Desmodus rotundus*). Most branches had  
142 relatively more loss events than duplication events, implying a decrease in gene family  
143 size (figure 3).





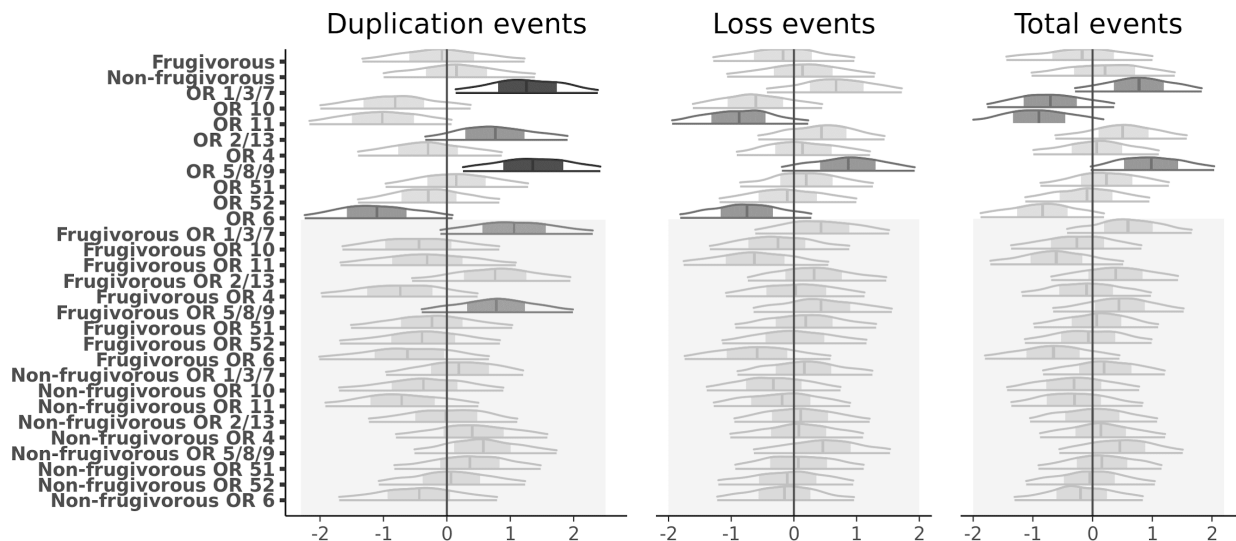
144 Figure 3. The number of loss events vs the number of duplication events for each  
145 olfactory receptor gene subfamily in each phyllostomid species. Triangles represent  
146 non-frugivorous species (animalivores and omnivores), circles represent frugivorous  
147 species. Colored dots correspond to *OR* gene families in which the number of  
148 duplication events was found to be significantly different between frugivorous and non-  
149 frugivorous species.

150 Poisson overdispersion model

151 Two *OR* gene subfamilies were identified as having more duplication events based on  
152 their posterior predictive intervals excluding 0, *OR* 1/3/7 and *OR* 5/8/9. *OR* 1/3/7 was  
153 associated with a mean 1.26X increase in duplication events, and *OR* 5/8/9 was  
154 associated with a mean 1.34X increase in duplication events. However, diet did not  
155 predict duplication events (figure 4). There is some evidence for increases in duplication  
156 events in two combinations of diet and gene subfamily: While the 90% credible intervals

157 of *OR* 2/13, *OR* 6, and the interactions of *OR* 1/3/7 and *OR* 5/8/9 with frugivory  
158 included 0, their 70% credible intervals did not (figure 4).

159 None of the three sets of predictors tested (*OR* gene family, diet, and the interactions of  
160 *OR* gene family and diet) predicted either loss events, or the total number of duplication/  
161 loss events per species based on the 90% credible intervals. However, at 70% credible  
162 intervals, *OR* 11 and *OR* 6 had fewer loss events, and *OR* 5/8/9 more loss events  
163 (figure 4). *OR* 10 and *OR* 11 had fewer total duplication/loss events, and *OR* 1/3/7 and  
164 *OR* 5/8/9 had more total duplication/loss events (figure 4).

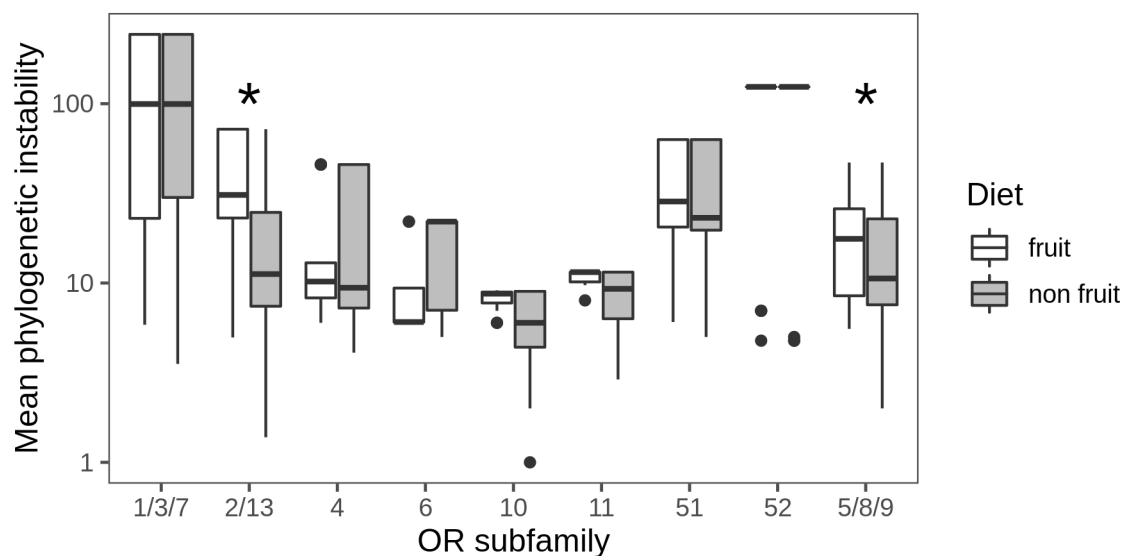


165 Figure 4. The effect of diet (frugivorous vs non-frugivorous), *OR* gene family, and *OR*  
166 gene family in only frugivorous or non-frugivorous species on the number of duplication  
167 events (first panel), the number of loss events (second panel), and the total number of  
168 duplication and loss events per branch (third panel). Distributions outline the 90%  
169 credible intervals of the coefficient of the effect of each predictor on the total number of  
170 duplication events per branch, with the 50% credible intervals shaded. Distributions in  
171 dark grey have 90% credible intervals that do not include 0, distributions in medium grey  
172 have 70% credible intervals that do not include 0, and distributions in light grey have  
173 70% credible intervals that do include 0. Grey boxes highlight the diet x *OR* family  
174 interaction terms.

175 Phylogenetic instability

176 Three *OR* gene families, *OR* 5/8/9, *OR* 2/13, and *OR* 6, had significant differences in  
177 instability scores between frugivores and non-frugivores based on a Mann-Whitney U  
178 test (*OR* 5/8/9:  $Z = 3.77$ ,  $p = 0.000165$ ; *OR* 2/13:  $Z = 7.36$ ,  $p = 1.84 \times 10^{-13}$ ; *OR* 6:  $Z = -$   
179  $2.47$ ,  $p = 0.0135$ ) before correcting for multiple comparisons, but only two (*OR* 5/8/9 and  
180 *OR* 2/13) were significant after a Hommel correction (figure 5). Frugivorous bats had  
181 significantly higher rank instability scores than non-frugivorous bats for both *OR* 5/8/9  
182 and *OR* 2/13. While *OR* 6 was not significant after correcting for multiple comparisons,  
183 there was a trend toward frugivorous bats having lower instability scores than non-  
184 frugivorous bats (Sup. figure 1A).

185 *OR* 5/8/9 and *OR* 2/13 had significant differences in instability between frugivorous and  
186 non-frugivorous species, and were the two *OR* subfamilies with the most orthologs  
187 (Sup. figure 1B, C). *OR* 5/8/9 contained 580 genes and was sorted into 80 groups, with  
188 instability scores ranging from 1.00 to 46.97. *OR* 2/13 contained 239 genes and was  
189 sorted into 27 instability groups, with scores ranging from 1.38 to 72.18. *OR* 6, however,  
190 was a relatively smaller subfamily, containing only 35 genes which were sorted into 7  
191 instability groups. *OR* 10 is an example of a subfamily with little variation between diets.  
192 There are fewer orthologs in the subfamily, and each instability group only had between  
193 one and five genes in it (Sup. figure 1D).



194 Figure 5. Average instability scores of frugivorous (white) vs. non-frugivorous (grey) bat  
195 species in each *OR* gene family. Asterisks designate *OR* subfamilies with significant  
196 differences in rank instability scores between frugivorous and non-frugivorous species.

## 197 Discussion

198 We show that species tree heterogeneity misleads inference of gene duplication and  
199 loss rates, with implications for tests of ecological associations (e.g., Chang and Duda  
200 2012; Dahan et al. 2015; Ramasamy et al. 2016). Although the influence of tree  
201 heterogeneity on analyses of speciation rates has been demonstrated (Rabosky and  
202 Goldberg 2015; Rojas et al. 2018), the conflation of tree heterogeneity with rate shifts in  
203 multi-gene family evolution had not been previously demonstrated. We show that the  
204 most commonly used method for modeling selection on gene copy number, CAFE (Bie  
205 et al. 2006; Han et al. 2013), is prone to error with unbalanced species trees, with  
206 implications for analyses of many empirical data sets. This type of error, however, can  
207 be overcome by adopting tree transformations, or using alternative methods, as shown  
208 here. Therefore, whenever tree heterogeneity is suspected, simulations such as those  
209 presented here can help determine what methods are appropriate for testing links  
210 between rates of multi-gene family evolution and specific clades or branches.

211 Using simulations we find CAFE has an unacceptably high false positive rate in non-  
212 Yule trees, making it inappropriate for analyses of many empirical systems. CAFE  
213 provides an elegant model comparison framework for testing for duplication (and/or  
214 loss) rate shifts at specific hypothesized nodes of a species phylogeny, and is widely  
215 applied for this purpose across many systems. However, its framework assumes a Yule  
216 tree, a premise rarely tested for empirical phylogenies. The high false positive rate we  
217 show is the result of heterogeneity in the tree shape (non-Yule branching patterns and  
218 branch lengths) being ascribed to a duplication or loss rate shift. In effect, CAFE is  
219 detecting a deviation from a neutral, single-rate model, but in a non-Yule tree that  
220 deviation may come from tree shape instead of from rate shifts in gene family evolution.

221 While the spurious association between a non-null model and a factor being tested (in

222 this case, duplication/loss rate shifts) was first demonstrated with binary state-speciation  
223 and extinction models (BiSSE) (Rabosky and Goldberg 2015), this challenge is a more  
224 general consequence of the integration of the species phylogeny with rates  
225 superimposed therein, and has also been found with continuous data (Harvey and  
226 Rabosky 2018; Rojas et al. 2018). Solutions and alternatives have been formulated for  
227 SSE-based models (e.g., (Beaulieu and O'Meara 2016; Rabosky and Huang 2016;  
228 Harvey and Rabosky 2018; May and Moore 2019), but are unavailable for this newly-  
229 identified instance of tree imbalance producing spurious results. Our simulations show  
230 transforming the species tree to meet the Yule assumption can overcome the conflation  
231 of heterogeneity in the species tree with gene family rate shifts. This, however, must be  
232 tested for particular cases. Methods that do not make assumptions about the tree shape  
233 are available as alternatives.

234 In an empirical application, and using methods that are not confounded by the effects of  
235 tree heterogeneity on birth-death equilibrium, we show that despite an overall decrease  
236 in most olfactory receptor (*OR*) gene subfamily sizes across phyllostomid bat species,  
237 certain *OR* subfamilies have expansions associated with a dietary shift to frugivory.  
238 These results contradict some of the findings of a previous study associating *OR*  
239 subfamily size changes with a shift to frugivory in phyllostomid bats (Hayden et al.  
240 2014), suggesting that collapsing the gene subfamily data into principal components  
241 allows tree-wide instability in a gene subfamily to be misinterpreted as having ecological  
242 associations. Our analyses overcome this shortcoming by using reconciliation to  
243 analyze rate shifts directly, revealing in greater depth in the dynamics of duplication and  
244 loss events associated with frugivory.

245 By using alternative methods not dependent on tree topology, we identified an increase  
246 in duplication events in *OR* subfamily 5/8/9 as linked with the shift to frugivory in  
247 phyllostomid bats. Previous work also found a significant increase in duplication events  
248 in this subfamily (Hayden et al. 2014), but for non-frugivorous phyllostomids instead. We  
249 explain this contradiction through the interaction between duplication events and  
250 methods of analysis. The number of duplication events in *OR* subfamily 5/8/9 suggests

251 expansion across phyllostomids, frugivorous and non-frugivorous (figure 3 this paper,  
252 Hayden et al. figure 4). Using methods that directly examine the interaction of the  
253 number of duplication events and dietary category (Poisson overdispersion) and the  
254 deviation of duplication patterns according to diet (phylogenetic instability), we were  
255 able to measure the association of duplication events and diet despite an overall  
256 increase in the number of duplication events in *OR* subfamily 5/8/9 across the tree.

257 Analyses by Hayden et al. 2014 added even more abstraction by collapsing *OR*  
258 subfamily variation into principal components. This approach compressed the data in  
259 ways that could prove consequential. First, while the principal components can relate  
260 back to specific gene subfamilies, variation across multiple subfamilies influences their  
261 value, making it difficult to discover subfamily correlates of ecology. Second, there was  
262 no explicit model of ecological trait evolution and comparisons of the principal  
263 components using ANOVA were limited to the tips and secondarily inferred to be  
264 localized to internal branches. This highlights the importance of methods that use the  
265 number of duplication events directly, as is evident in the resulting instability in the *OR*  
266 subfamily 5/8/9 gene tree (Sup. figure 1C).

267 Based on its expansion within frugivorous phyllostomids, elements of *OR* subfamily  
268 5/8/9 may therefore specifically respond to volatile organic compounds (VOCs) released  
269 by ripening fruits. Indeed, expansions of this subfamily have recently been linked to  
270 herbivory in mammals (Hughes et al. 2018), lending confidence to the conclusion that  
271 *OR* subfamily 5/8/9 is involved in detecting plant VOCs. It has been previously  
272 hypothesized that the size of an *OR* repertoire could influence olfactory sensitivity {rtf  
273 (Rouquier et al. 2000; but see Wackermannová et al., 2016), and some evidence  
274 suggests this may be the case (Laska and Shepherd 2007; Rizvanovic et al. 2013).  
275 These additional copies of subfamily 5/8/9 genes could increase a bat's perception of  
276 fruit VOCs. However, this could be either because of increased sensitivity to a single  
277 VOC as a result of having more receptors for that specific VOC, or increased  
278 discrimination between many similar VOCs, resulting from slight functional differences  
279 between *OR* gene copies within the subfamily (Laska and Galizia 2001; Rizvanovic et

280 al. 2013). To discern whether sensitivity or discrimination is responsible for this  
281 expansion, individual-level differences in *OR* gene subfamily repertoires and associated  
282 VOC preferences should be integrated with models of gene subfamily birth/death  
283 dynamics.

284 *OR* subfamily 2/13 was also previously found to be important in the shift to frugivory  
285 within phyllostomids (Hayden et al. 2014), and this is corroborated by phylogenetic  
286 instability scores in this study (figure 5). In contrast, the number of duplication events  
287 within this subfamily was not higher in frugivorous species, as determined by posterior  
288 predictive intervals (figure 4). While related, the number of duplication events and  
289 phylogenetic instability cannot be expected to always yield consistent results.  
290 Incongruence (instability) could increase at a branch without quite pushing the number  
291 of events at the ancestral branch toward significance. Alternatively, an expansion in this  
292 subfamily could be associated with individual branches only partially coincident with a  
293 shift to frugivory, instead of the single branch associated with frugivory itself.

294 While *OR* subfamily 1/3/7 was previously found to be important in the shift to frugivory  
295 within phyllostomids (Hayden et al. 2014), our results suggest this subfamily is highly  
296 unstable across the whole phyllostomid tree. Within fruit-specialized phyllostomids,  
297 there is a residual trend toward more duplications, but this subfamily experiences high  
298 turnover across the whole tree. Thus, its expansion may be related to either an inherent  
299 genetic mechanism, or an ecological factor affecting all of Phyllostomidae, rather than  
300 the shift to frugivory within the family.

301 We have shown that expansions of at least one *OR* subfamily, 5/8/9, are associated  
302 with a shift to frugivory in phyllostomid bats. Nonetheless, we have also shown that  
303 expansions in *OR* subfamilies are not always associated with an important ecological  
304 change. For example, the results of the Poisson overdispersion model show that about  
305 half of the subfamilies tested may simply have more duplication events than expected  
306 based on the phylogeny, regardless of diet. This means other factors can confound  
307 these analyses, some of which will often be unaccounted for. By combining a method



308 that measures variation in the number of duplication(/loss) events (Poisson  
309 overdispersion model) with one that measures deviation from the expected number of  
310 copies (phylogenetic instability), we are able to identify *OR* subfamilies for which  
311 conflicting results (e.g., *OR* subfamily 2/13) point to confounding factor(s). Changes in  
312 the per-lineage numbers of duplication and loss events not accounted for in the shift to  
313 frugivory may also be driven by other behavioral and ecological factors (e.g., mate  
314 identification, habitat type), or changes in the rates of duplication events for specific  
315 genomic regions.

316 Within *OR* subfamily 5/8/9, and in other *OR* subfamilies found to be associated with  
317 ecological shifts in other systems, it is now important to identify whether the subfamily  
318 expansion is beneficial because of increased sensitivity or increased discrimination.  
319 That is, do more copies of an *OR* gene increase the likelihood that the animal will detect  
320 a scent when there are few molecules, or increase the likelihood that the animal will be  
321 able to sort out a particular scent of interest from many similar ones? Determining  
322 functional differences between *OR* genes within a subfamily, by comparing amino acid  
323 substitutions with different physiochemical properties, or by *in vivo* or *in vitro* response  
324 experiments, will be the next step to making this distinction.

325 Testing the adaptive significance of gene subfamily size or changes in gene duplication  
326 or loss rates is complex and potentially confounded by the relationship between the  
327 birth-death process of gene subfamily evolution and the species phylogeny. This study  
328 confirms that, as with state-dependent speciation models, species tree heterogeneity is  
329 a confounding factor in inferring multi-gene adaptation, and can cause spurious  
330 associations between shifts in birth-death rate and species ecology. Since testing the  
331 ecological associations of gene subfamily size and/or duplication and loss rates is a  
332 necessary step in determining their adaptive significance, it is thus imperative to use  
333 methods that are not confounded by the effects of tree heterogeneity on birth-death  
334 equilibrium.



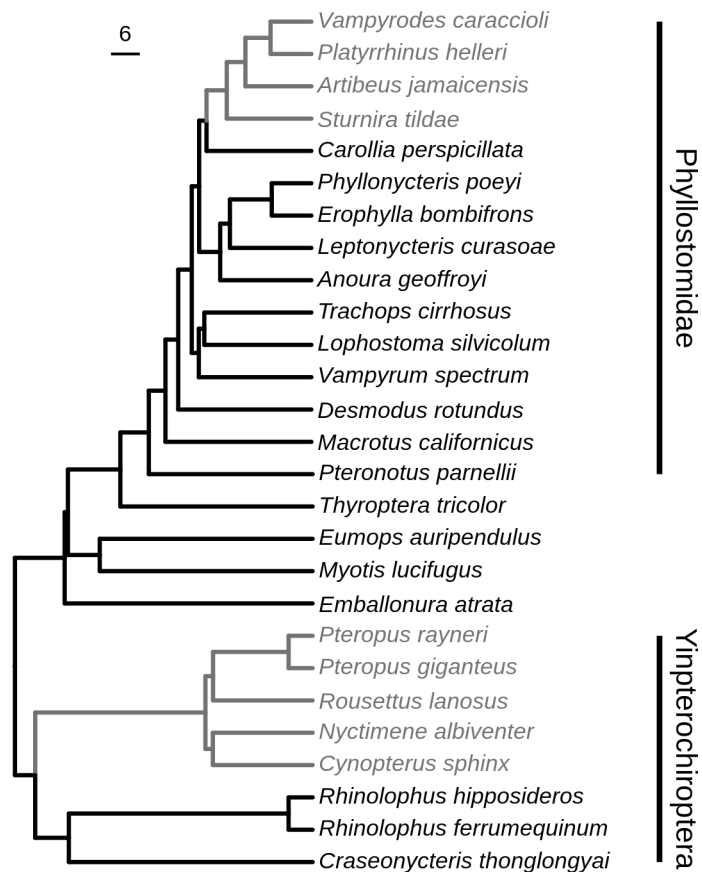
## 335 Methods

### 336 Sequence data and copy number

337 Aligned consensus sequences of bat *OR* gene families and copy number counts for  
338 nine *OR* gene subfamilies (*OR* 1/3/7, *OR* 2/13, *OR* 4, *OR* 5/8/9, *OR* 6, *OR* 10, *OR* 11,  
339 *OR* 51, and *OR* 52) in 14 species of Phyllostomids and eight species of  
340 Yinpterochiropterans were obtained from Hayden et al. (2014) (figure 2).

### 341 Influence of topology

342 In order to test the variation in duplication rates in gene families, CAFE assumes a birth-  
343 death (Yule) species tree (Han et al. 2013). Many species trees, however, do not fit this  
344 assumption. To test the performance on real phylogenetic trees, we ascertained the  
345 false positive rate of changes in the duplication rate  $\lambda$  associated with a shift to frugivory  
346 in bats. We used both actual species trees and those trees transformed to meet Yule  
347 branching time expectations, using copy number counts from Hayden et al. (2014). We  
348 ran CAFE independently for *OR* gene families in two bat clades, family Phyllostomidae  
349 (poor fit to Yule expectations) and suborder Yinpterochiroptera (better fit to Yule  
350 expectations) (figure 2).



351 Figure 2. Bat species and clades used in the analysis of topology. Species with  
352 branches highlighted in grey are frugivorous.

353 For each gene subfamily and species tree, copy number counts from Hayden et al.  
354 (2014) were randomized across the tips in 100 permutations, and CAFE was run on  
355 each permutation. This randomization preserved observed magnitudes and relative  
356 copy number counts, but associated each with a different tip, thus removing the  
357 influence and signal of selection on copy numbers corresponding to ecological traits.

358 Each set of CAFE runs included both a single- $\lambda$  and a two- $\lambda$  model, with the shift in  $\lambda$   
359 corresponding to the shift to frugivory. Models were compared via likelihood ratio test  
360 using CAFE's built-in *lhtest* command, which creates a null distribution of likelihood

361 ratios for the null (single- $\lambda$ ) hypothesis. Any significant likelihood ratio test, indicating a  
362 better fit to the data for the two- $\lambda$  model, was thus a false positive for selection. A false  
363 positive implies spurious influence of tree shape on the CAFE test, especially if false  
364 positives are more prevalent in trees with a poorer fit to a Yule model.

365 To further discern the influence of non-Yule branch length vs. topology on the false  
366 positive rate, we transformed branch lengths of both trees to meet Yule expectations by  
367 fitting a Yule model to each tree with `yule()`, applying Yule expected branching times  
368 simulated to fit that model with `sim.bdtree()`, and applying those branching times to each  
369 original tree with `compute.brtime()` from the `ape` (Paradis et al. 2004) and `geiger`  
370 (Harmon et al. 2008) packages in R (R Core Team 2018).

#### 371 Inferring Gene Trees and Estimating Duplication and Loss Events

372 We used ModelOMatic (Whelan et al. 2015) to determine the best substitution model for  
373 each *OR* gene subfamily, and used `garli` (Zwickl 2006) to infer gene trees from the  
374 alignments provided by Hayden et al. (2014). We then ran Notung (Chen et al. 2000)  
375 with default duplication (1.5) and loss (1.0) costs to reconcile the resulting gene and  
376 species trees and obtain estimates of the number and location of duplication and loss  
377 events.

#### 378 Poisson overdispersion model

379 To determine the influence of diet on the duplication and loss rates of each *OR* gene  
380 family, we built a Poisson generalized mixed model of diet, *OR* gene subfamily, and the  
381 interaction effect between diet and *OR* gene subfamily, based on the model from  
382 Sackton et al. (2017). Although the Sackton et al. (2017) implementation used a  
383 maximum likelihood implementation of hierarchical mixed Poisson regressions, here we  
384 apply Bayesian method to account for overdispersion in the counts. The model was run  
385 using `MCMCglmm` in the `MCMCglmm` package (Hadfield 2010) including a branch  
386 length offset as a fixed effect to control for species relationships, and all other  
387 characters as random effects after Gelman and Hill (2006).

388 We ran three separate models, one with the number of duplication events per branch

389 and tip, one with the number of loss events, and one with the total number of events as  
390 the response variable. The bayesplot package (Gabry et al. 2019) was used to visualize  
391 posterior probability distributions.

### 392 Phylogenetic instability

393 We determined phylogenetic instability for all nine *OR* gene subfamilies in  
394 Phyllostomids (figure 2) using MIPhy (Minimizing Instability in Phylogenetics) (Curran et  
395 al. 2018) and gene and species trees from Hayden et al. (2014). MIPhy uses patterns of  
396 duplications and losses in gene trees to identify groups of genes that are under  
397 selection based on deviations from expected baseline patterns. MIPhy quantifies  
398 instability in each gene, and assigns each group a score based on how its pattern of  
399 gene events differs from what would be expected if the relationships in the gene tree  
400 mirrored those of the species tree. A higher instability score indicates that the gene tree  
401 is discordant with the species tree, and suggests that it is under positive or negative  
402 selection.

403 To quantify the phylogenetic instability for each of the nine *OR* gene families, we used  
404 the MIPhy online web tool with default weights (Curran et al. 2018). The species tree  
405 and the list of *OR* genes corresponding to each species were compiled into an  
406 information file compatible with MIPhy.

407 To determine if the resulting instability scores were more likely to be associated with  
408 frugivorous species than expected, we performed a 2-tailed Wilcoxon-Mann-Whitney  
409 test with species classified as frugivorous or non-frugivorous, using the `wilcoxon_test`  
410 function in the `coin` packages to account for ties (Hothorn et al. 2019 Aug). No  
411 correction for species relatedness was performed because the instability scores already  
412 take this into account. A greater score of the genes in an *OR* subfamily would suggest  
413 that the family was under selective pressure during the transition to frugivory, causing it  
414 to deviate from the expected pattern of duplications and losses based on the gene tree.  
415 We applied the Hommel correction for multiple comparisons.

## 416 Acknowledgements

417 We thank Katie Martin and Sharlene Santana for comments and feedback on the work.  
418 The authors would like to thank Stony Brook Research Computing and  
419 Cyberinfrastructure, and the Institute for Advanced Computational Science at Stony  
420 Brook University for access to the high-performance SeaWulf computing system, which  
421 was made possible by a \$1.4M National Science Foundation grant (#1531492). This  
422 work was supported in part by the National Science Foundation (DEB 1442142, DEB  
423 1456455, and DEB 1838273 to L.M.D., DEB 1701414 to L.M.D. and L.R.Y., PRFB  
424 181203 to L.R.Y.), and the Tinker Foundation and American Association of University  
425 Women fellowships to M.E.L.

426 References

- 427 Assis R, Bachtrog D. 2013. Neofunctionalization of young duplicate genes in  
428 *Drosophila*. *Proc Natl Acad Sci*. 110(43):17409–17414.  
429 doi:10.1073/pnas.1313759110.
- 430 Beaulieu JM, O’Meara BC. 2016. Detecting Hidden Diversification Shifts in Models of  
431 Trait-Dependent Speciation and Extinction. *Syst Biol*. 65(4):583–601.  
432 doi:10.1093/sysbio/syw022.
- 433 Bie TD, Cristianini N, Demuth JP, Hahn MW. 2006. CAFE : a computational tool for the  
434 study of gene family evolution. 22(10):1269–1271.  
435 doi:10.1093/bioinformatics/btl097.
- 436 Chang D, Duda TF. 2012. Extensive and Continuous Duplication Facilitates Rapid  
437 Evolution and Diversification of Gene Families. *Mol Biol Evol*. 29(8):2019–2029.  
438 doi:10.1093/molbev/mss068.
- 439 Chen K, Durand D, Farach-Colton M. 2000. NOTUNG: A Program for Dating Gene  
440 Duplications and Optimizing Gene Family Trees. *J Comput Biol*. 7(3–4):429–447.  
441 doi:10.1089/106652700750050871.
- 442 Curran DM, Gilleard JS, Wasmuth JD. 2018. MIPhy: identify and quantify rapidly  
443 evolving members of large gene families. *PeerJ*. 6:e4873.  
444 doi:10.7717/peerj.4873.
- 445 Dahan RA, Duncan RP, Wilson AC, Dávalos LM. 2015. Amino acid transporter  
446 expansions associated with the evolution of obligate endosymbiosis in sap-  
447 feeding insects (Hemiptera: sternorrhyncha). *BMC Evol Biol*. 15(1):52.  
448 doi:10.1186/s12862-015-0315-3.
- 449 Gabry J, Simpson D, Vehtari A, Betancourt M, Gelman A. 2019. Visualization in  
450 Bayesian workflow. *J R Stat Soc Ser A Stat Soc*. 182(2):389–402.  
451 doi:10.1111/rssa.12378.
- 452 Gelman A, Hill J. 2006. *Data Analysis Using Regression and Multilevel/Hierarchical*  
453 *Models*. Cambridge University Press (Analytical Methods for Social Research).
- 454 Hadfield JD. 2010. MCMC Methods for Multi-Response Generalized Linear Mixed  
455 Models: The **MCMCglmm** R Package. *J Stat Softw*. 33(2).

- 456           doi:10.18637/jss.v033.i02. [accessed 2019 Sep 7].  
457           <http://www.jstatsoft.org/v33/i02/>.
- 458   Han MV, Thomas GWC, Lugo-Martinez J, Hahn MW. 2013. Estimating gene gain and  
459           loss rates in the presence of error in genome assembly and annotation using  
460           CAFE 3. *Mol Biol Evol.* 30(8):1987–97. doi:10.1093/molbev/mst100.
- 461   Harmon LJ, Weir JT, Brock CD, Glor RE, Challenger W. 2008. GEIGER: investigating  
462           evolutionary radiations. *Bioinformatics.* 24(1):129–131.  
463           doi:10.1093/bioinformatics/btm538.
- 464   Harvey MG, Rabosky DL. 2018. Continuous traits and speciation rates: Alternatives to  
465           state-dependent diversification models. Cooper N, editor. *Methods Ecol Evol.*  
466           9(4):984–993. doi:10.1111/2041-210X.12949.
- 467   Hayden S, Bekaert M, Goodbla A, Murphy WJ, Dávalos LM, Teeling EC. 2014. A  
468           Cluster of Olfactory Receptor Genes Linked to Frugivory in Bats. *Mol Biol Evol.*  
469           31(4):917–927. doi:10.1093/molbev/msu043.
- 470   Hothorn T, Winell H, Hornik K, van de Wiel MA, Zeileis A. 2019 Aug. Conditional  
471           Inference Procedures in a Permutation Test Framework. CRAN.
- 472   Hughes GM, Boston ESM, Finarelli JA, Murphy WJ, Higgins DG, Teeling EC. 2018.  
473           The Birth and Death of Olfactory Receptor Gene Families in Mammalian Niche  
474           Adaptation. Satta Y, editor. *Mol Biol Evol.* 35(6):1390–1406. doi:10.1093/molbev/  
475           msy028.
- 476   Jones KE, Bininda-Emonds ORP, Gittleman JL. 2005. Bats, Clocks, and Rocks:  
477           Diversification Patterns in Chiroptera. *Evolution.* 59(10):2243–2255.  
478           doi:10.1111/j.0014-3820.2005.tb00932.x.
- 479   Laska M, Galizia CG. 2001. Enantioselectivity of Odor Perception in Honeybees (*Apis*  
480           *mellifera carnica*). *Behav Neurosci.* 115(3):632–639.
- 481   Laska M, Shepherd GM. 2007. Olfactory discrimination ability of CD-1 mice for a large  
482           array of enantiomers. *Neuroscience.* 144(1):295–301.  
483           doi:10.1016/j.neuroscience.2006.08.063.
- 484   May MR, Moore BR. 2019. A Bayesian Approach for Inferring the Impact of a Discrete  
485           Character on Rates of Continuous-Character Evolution in the Presence of  
486           Background-Rate Variation. *Syst Biol.* syz069.

- 487 McGowen MR, Gatesy J, Wildman DE. 2014. Molecular evolution tracks  
488 macroevolutionary transitions in Cetacea. *Trends Ecol Evol.* 29(6):336–346.  
489 doi:10.1016/j.tree.2014.04.001.
- 490 Nei M. 1969. Gene Duplication and Nucleotide Substitution in Evolution. *Nature.*  
491 221(5175):40–42. doi:10.1038/221040a0.
- 492 Nei M, Gu X, Sitnikova T. 1997. Evolution by the birth-and-death process in multigene  
493 families of the vertebrate immune system. *Proc Natl Acad Sci.* 94(15):7799–  
494 7806. doi:10.1073/pnas.94.15.7799.
- 495 Nei M, Hughes AL. 1991. Polymorphism and evolution of the major histocompatibility  
496 complex loci in mammals. In: Selander R, Clark A, Whittam T, editors. *Evolution*  
497 *at the Molecular Level.* Sunderland, MA: Sinauer Associates. p. 222–247.  
498 [accessed 2019 Dec 5].  
499 <http://www.personal.psu.edu/nxm2/1991%20Publications/1991-nei-hughes2.pdf>.
- 500 Nei M, Rooney AP. 2005. Concerted and Birth-and-Death Evolution of Multigene  
501 Families. *Annu Rev Genet.* 39(1):121–152.  
502 doi:10.1146/annurev.genet.39.073003.112240.
- 503 Paradis E, Claude J, Strimmer K. 2004. APE: Analyses of Phylogenetics and Evolution  
504 in R language. *Bioinformatics.* 20(2):289–290. doi:10.1093/bioinformatics/btg412.
- 505 R Core Team. 2018. R: A language and environment for statistical computing. R Found  
506 Stat Comput Vienna Austria. <https://www.R-project.org>.
- 507 Rabosky DL, Goldberg EE. 2015. Model Inadequacy and Mistaken Inferences of Trait-  
508 Dependent Speciation. *Syst Biol.* 64(2):340–355. doi:10.1093/sysbio/syu131.
- 509 Rabosky DL, Huang H. 2016. A Robust Semi-Parametric Test for Detecting Trait-  
510 Dependent Diversification. *Syst Biol.* 65(2):181–193. doi:10.1093/sysbio/syv066.
- 511 Ramasamy S, Ometto L, Crava CM, Revadi S, Kaur R, Horner DS, Pisani D, Dekker T,  
512 Anfora G, Rota-Stabelli O. 2016. The Evolution of Olfactory Gene Families in  
513 *Drosophila* and the Genomic Basis of chemical-Ecological Adaptation in  
514 *Drosophila suzukii*. *Genome Biol Evol.* 8(8):2297–2311.  
515 doi:10.1093/gbe/evw160.
- 516 Rizvanovic A, Amundin M, Laska M. 2013. Olfactory Discrimination Ability of Asian  
517 Elephants (*Elephas maximus*) for Structurally Related Odorants. *Chem Senses.*



- 518 38(2):107–118. doi:10.1093/chemse/bjs097.
- 519 Rojas D, Ramos Pereira MJ, Fonseca C, Dávalos LM. 2018. Eating down the food  
520 chain: generalism is not an evolutionary dead end for herbivores. Harmon L,  
521 editor. *Ecol Lett.* 21(3):402–410. doi:10.1111/ele.12911.
- 522 Rouquier S, Blancher A, Giorgi D. 2000. The olfactory receptor gene repertoire in  
523 primates and mouse: Evidence for reduction of the functional fraction in primates.  
524 *Proc Natl Acad Sci.* 97(6):2870–2874. doi:10.1073/pnas.040580197.
- 525 Sackton TB, Lazzaro BP, Clark AG. 2017. Rapid expansion of immune-related gene  
526 families in the house fly, *Musca domestica*. *Mol Biol Evol.*  
527 doi:10.1093/molbev/msw285. [accessed 2017 Mar 17].  
528 <https://academic.oup.com/mbe/article-lookup/doi/10.1093/molbev/msw285>.
- 529 Vandewege MW, Mangum SF, Gabaldón T, Castoe TA, Ray DA, Hoffmann FG. 2016.  
530 Contrasting patterns of evolutionary diversification in the olfactory repertoires of  
531 reptile and bird genomes. *Genome Biol Evol.* doi:10.1093/gbe/evw013.  
532 [accessed 2019 Dec 5].  
533 <https://academic.oup.com/gbe/article-lookup/doi/10.1093/gbe/evw013>.
- 534 Wackermannová M, Pinc L, Jebavý L. 2016. Olfactory Sensitivity in Mammalian  
535 Species. *Physiol Res.* 65:369–390.
- 536 Whelan S, Allen JE, Blackburne BP, Talavera D. 2015. ModelOMatic: Fast and  
537 Automated Model Selection between RY, Nucleotide, Amino Acid, and Codon  
538 Substitution Models. *Syst Biol.* 64(1):42–55. doi:10.1093/sysbio/syu062.
- 539 Yohe LR, Liu L, Dávalos LM, Liberles DA. 2019. Protocols for the Molecular  
540 Evolutionary Analysis of Membrane Protein Gene Duplicates. In: Sikosek T,  
541 editor. *Computational Methods in Protein Evolution*. Vol. 1851. New York, NY:  
542 Springer New York. p. 49–62. [accessed 2019 Mar 14].  
543 [http://link.springer.com/10.1007/978-1-4939-8736-8\\_3](http://link.springer.com/10.1007/978-1-4939-8736-8_3).
- 544 Zwickl DJ. 2006. Genetic algorithm approaches for the phylogenetic analysis of large  
545 biological sequence datasets under the maximum likelihood criterion  
546 [Dissertation]. The University of Texas at Austin.