

Version dated: August 10, 2017

GHOST: INFERRING HETEROTACHOUS EVOLUTION

GHOST: Recovering Historical Signal from Heterotachously-evolved Sequence Alignments

STEPHEN M CROTTY^{†*1,2,3}, BUI QUANG MINH^{†1}, NIGEL G BEAN^{2,3}, BARBARA R HOLLAND⁴, JONATHAN TUKE^{2,3}, LARS S JERMIIN^{5,6}, ARNDT VON HAESELER^{1,7}

¹*Center for Integrative Bioinformatics Vienna, Max F. Perutz Laboratories, University of Vienna and Medical University of Vienna, Vienna, Austria*

²*School of Mathematical Sciences, University of Adelaide, Adelaide, SA 5005, Australia.*

³*ARC Centre of Excellence for Mathematical and Statistical Frontiers, The University of Adelaide, Adelaide, SA, Australia.*

⁴*School of Physical Sciences, University of Tasmania, Hobart, TAS 7001, Australia.*

⁵*CSIRO Land & Water, Black Mountain Laboratories, Canberra, ACT 2601, Australia.*

⁶*Research School of Biology, Australian National University, Canberra, ACT 2601, Australia*

⁷*Bioinformatics & Computational Biology, Faculty of Computer Science, University of Vienna, Vienna, Austria.*

19 † Joint first authors (these authors contributed equally to the work)

20 *Corresponding author: Stephen Crotty, Center for Integrative Bioinformatics

21 Vienna, Max F. Perutz Laboratories, University of Vienna and Medical University

22 of Vienna, Vienna, Austria; E-mail: stephen.crotty@univie.ac.at

23 *Abstract.*— Molecular sequence data that have evolved under the influence of
24 heterotachous evolutionary processes are known to mislead phylogenetic inference.
25 We introduce the General Heterogeneous evolution On a Single Topology (GHOST)
26 model of sequence evolution, implemented under a maximum-likelihood framework
27 in the phylogenetic program IQ-TREE. Extensive simulations show that the
28 GHOST model can accurately recover the tree topology, branch lengths,
29 substitution rate and base frequency parameters from heterotachously-evolved
30 sequences. We apply our model to an electric fish dataset and identify a subtle
31 component of the historical signal, linked to the previously established convergent
32 evolution of the electric organ in two geographically distinct lineages of electric fish.
33 We compare the GHOST model to the partition model and show that, owing to the
34 minimization of model constraints, the GHOST model is able to offer unique
35 biological insights when applied to empirical data.

36 Keywords: Phylogenetics, heterotachy, mixture model, maximum likelihood,
37 convergent evolution

38 The success and reliability of model-based phylogenetic inference methods
39 are limited by the adequacy of the models that are assumed to approximate the
40 evolutionary process. Homogeneous evolutionary models have long been recognised
41 as inadequate since the rate of evolution is known to vary across sites (Fitch and
42 Margoliash, 1967; Holmquist et al., 1983) and across lineages (Baele et al., 2006;
43 Lopez et al., 2002; Wu and Susko, 2011; Jayaswal et al., 2014). There are many
44 models that have been proposed to compensate for rate heterogeneity across sites.
45 The classical example is the discrete Γ model (Yang, 1994), which allows different
46 classes of variable sites to have their rates drawn from a Γ distribution. More
47 recently, Kalyaanamoorthy et al. (2017) relaxed the requirement for the rates of the
48 classes to fit a Γ distribution, implementing a probability distribution-free
49 rates-across-sites model. However, these models still assume that the substitution
50 rate for each site is constant across all lineages. This is too restrictive; biologically
51 speaking it is not hard to accept that evolutionary processes can be both lineage
52 and time dependent. In the context of a phylogenetic tree this manifests as
53 lineage-specific shifts in evolutionary rate, coined heterotachy (Philippe and Lopez,
54 2001; Lopez et al., 2002), resulting in sequences that cannot be characterised as
55 having evolved according to a single set of branch lengths and substitution
56 model.

57 The effect of heterotachy on phylogenetic inference was thrust into the
58 spotlight by Kolaczkowski and Thornton (K&T) (2004). They used a simulation
59 study to show that heterotachously-evolved sequences could mislead the popular
60 inference methods of maximum-likelihood (ML) and Bayesian Markov Chain
61 Monte-Carlo (BMCMC) to a greater extent than maximum parsimony (MP). Their
62 findings were controversial and were widely challenged on the grounds that the
63 simulations captured only a special case of heterotachy (Gadagkar and Kumar,
64 2005; Philippe et al., 2005; Spencer et al., 2005; Steel, 2005), and more general
65 studies of heterotachy concluded that ML performed at least as well as, and in
66 most cases better than, MP (Gadagkar and Kumar, 2005; Spencer et al., 2005).
67 Valid as these criticisms may have been, the key issue that K&T's study brought to
68 light stood firm - heterotachy was a primary source of model misspecification and
69 the models and methods of the time were ill-equipped to deal with it. The main
70 impediment to the development of models that can accommodate
71 heterotachously-evolved sequences has been the computational expense. Models
72 that account for heterogeneity of rates of change across sites can be integrated
73 relatively cheaply, but modeling heterotachy is not so simple. One approach has
74 been to use partition models (Lanfear et al., 2012), which require the data to be
75 partitioned *a priori*. The analysis then proceeds by inferring separate branch length

76 and model parameters for each partition. Sequence data is commonly partitioned
77 based on genes and/or codon position. However, the inherent assumption of such a
78 partitioning scheme is that heterotachy only occurs between partitions, not within
79 each partition. This may not be a valid assumption, so the requirement to partition
80 the data in advance of the analysis is a possible source of model misspecification.
81 Another approach has been to use mixture models, in which the likelihood of the
82 data at each site in the alignment is calculated as a weighted sum across multiple
83 classes (see Pagel and Meade (2005) for a detailed description of phylogenetic
84 mixture models). The most common approaches can be referred to as mixed
85 substitution rate (MSR) models (Lartillot and Philippe, 2004; Pagel and Meade,
86 2004), whereby each class has its own substitution rate matrix; and mixed branch
87 length (MBL) models (Kolaczkowski and Thornton, 2004; Meade and Pagel, 2008),
88 whereby each class has its own set of branch lengths on the tree. As a consequence
89 of their parameter rich nature, these models have all been implemented only within
90 a Bayesian framework. Wu and Susko (2009) proposed a general framework for
91 heterotachy, encompassing both mixed substitution rate and mixed branch length
92 models as special cases. Another example is the CAT models of Lartillot and
93 Philippe (2004), which have been widely used (Whelan and Halanaych (2017) and
94 references therein). Whelan and Halanaych (2017) carried out extensive simulation

95 and empirical studies comparing the performance of the CAT models to partition
96 models. They concluded that despite their additional complexity and associated
97 increase in runtime, the CAT models generally perform no better than partition
98 models. They also found that when new mixture models are introduced in the
99 literature their performance is not always assessed against the current popular
100 methods for phylogenetic analysis, such as partition models.

101 As a consequence of their varied nature, mixture models require many
102 parameters and the associated computational expense has thus far impeded their
103 implementation in a ML framework. The issue of computational expense is an ever
104 diminishing one; as computing power increases and algorithmic architecture
105 improves, the opportunity to employ more and more complex models of sequence
106 evolution does also. We introduce the General Heterogeneous evolution On a Single
107 Topology (GHOST) model for ML inference. The GHOST model combines features
108 of both MSR and MBL models. It consists of a number of classes, all evolving on
109 the same tree topology. For each class the branch lengths, nucleotide or amino-acid
110 frequencies, substitution rates and class weight are parameters to be inferred. It
111 minimises the number of assumptions that must be made *a priori* by inferring all
112 parameters directly from the data. Therefore, GHOST is free of the artificial
113 constraints common in other models, often included for computational expedience

114 rather than biological relevance. This means that the GHOST model has the
115 necessary freedom to extract any historical signals present in the data. We provide
116 an easy to use implementation of the GHOST model in the phylogenetic program
117 IQ-TREE (Nguyen et al., 2015), the first mixture model of comparable flexibility to
118 be made available in a ML framework.

119 METHODS AND MATERIALS

120 *Model Description*

121 The GHOST model consists of m classes and one tree topology, T , common to all
122 classes. All other parameters are inferred separately for each class. For the j^{th} class
123 we define λ_j as the set of branch lengths on T ; \mathbf{R}_j , the relative substitution rate
124 parameters; \mathbf{F}_j , the set of nucleotide or amino acid frequencies; and w_j , the class
125 weight ($w_j > 0, \sum w_j = 1$). Given a multiple sequence alignment (MSA), A , we
126 define L_{ij} as the likelihood of the data observed at the i^{th} site in A under the j^{th}
127 class of the GHOST model. L_{ij} is computed using Felsenstein's pruning algorithm
128 (Felsenstein, 1981). The likelihood of the i^{th} site, L_i , is then given by the weighted
129 sum of the L_{ij} over all j :

$$L_i = \sum_{j=1}^m w_j L_{ij}(T, \lambda_j, \mathbf{R}_j, \mathbf{F}_j).$$

130 Therefore, if S contains N sites (length of the alignment), the full
131 log-likelihood, ℓ , is given by:

$$\ell = \sum_{i=1}^N \log \left(\sum_{j=1}^m w_j L_{ij}(T, \lambda_j, \mathbf{R}_j, \mathbf{F}_j) \right).$$

132 We make use of the existing parameter optimisation algorithms within
133 IQ-TREE, extending it where necessary, to incorporate parameter estimation
134 across the m classes.

135 *Model Parameter Estimation for a Fixed Tree, T*

136 Let $\Theta = \{w_1, \dots, w_m, \lambda_1, \dots, \lambda_m, \mathbf{R}_1, \dots, \mathbf{R}_m, \mathbf{F}_1, \dots, \mathbf{F}_m\}$ denote the GHOST
137 model parameters (*i.e.*, class weights, branch lengths, relative substitution rates,
138 and nucleotide or amino-acid frequencies) for each of the m classes. To estimate all
139 parameters for a tree T we employ an expectation-maximization (EM) algorithm

140 (Dempster et al., 1977; Wang et al., 2008). We initialize Θ with all $\hat{\mathbf{R}}_j = \mathbf{1}$ in each
141 class, uniform nucleotide or amino-acid frequencies $\hat{\mathbf{F}}_j$ (i.e., the Jukes-Cantor
142 model), and \hat{w}_j and $\hat{\lambda}_j$ obtained by parsimonious branch lengths rescaled by a
143 discrete, distribution-free rates-across-sites model (Kalyaanamoorthy et al., 2017)
144 with m categories. This becomes the current estimate $\hat{\Theta}$. The EM algorithm
145 iteratively performs an expectation (E) step and a maximization (M) step to
146 update the current estimate until a (local) maximum likelihood is reached.

147 *E-step.*— For each site i and class j compute the posterior probability \hat{p}_{ij} of site i
148 belonging to class j based on the current estimate $\hat{\Theta}$:

$$\hat{p}_{ij} = \frac{\hat{w}_j L_{ij}(T, \hat{\lambda}_j, \hat{\mathbf{R}}_j, \hat{\mathbf{F}}_j)}{\sum_{k=1}^m \hat{w}_k L_{ik}(T, \hat{\lambda}_k, \hat{\mathbf{R}}_k, \hat{\mathbf{F}}_k)}.$$

149 *M-step.*— For each class j , maximize the log-likelihood function:

$$\ell_j = \sum_{i=1}^N \hat{p}_{ij} \log \left(L_{ij}(T, \lambda_j, \mathbf{R}_j, \mathbf{F}_j) \right)$$

150 to obtain the next $\hat{\lambda}_j^{NEW}$, \hat{R}_j^{NEW} , \hat{F}_j^{NEW} . This can be done with standard
151 phylogenetic optimization routines for each class.

152 Finally, the weights are updated by:

$$\hat{w}_j^{NEW} = \frac{1}{N} \sum_{i=1}^N \hat{p}_{ij}.$$

153 That is, the new weight for class j is the mean posterior probability of each
154 site belonging to class j . This completes the proposal of the new estimate $\hat{\Theta}^{NEW}$.
155 If $\ell(\hat{\Theta}^{NEW}) > \ell(\hat{\Theta}) + \epsilon$ (where ϵ is a user-defined tolerance, $\epsilon = 0.01$ by default),
156 then $\hat{\Theta}$ is replaced by $\hat{\Theta}^{NEW}$ and the E and M steps are repeated. Otherwise, the
157 EM algorithm finishes.

158 An auxiliary benefit of the ML implementation of the GHOST model in
159 IQ-TREE is that once the EM-algorithm has converged, we can soft-classify sites
160 according to their probability of belonging to a particular class. Post convergence,
161 the final values of p_{ij} can be directly interpreted as the probability that the i^{th} site
162 in the alignment belongs to the j^{th} class. This classification can be used to identify
163 sites in the alignment that belong with high probability to a particular class of
164 interest.

Software

165

166 The GHOST model has been implemented in IQ-TREE (Nguyen et al., 2015)
167 (<http://www.iqtree.org>), the first model of this type and complexity to be made
168 available in a ML framework. The GHOST model can be run with both nucleotide
169 and amino acid sequences. The GHOST model is executed in IQ-TREE v1.6 by
170 augmenting the model argument as shown below. For example if one wants to fit a
171 four-class GHOST model in conjunction with the GTR model of evolution to
172 sequences contained in `data.fst`, one would use the following command:

```
173     iqtree -s data.fst -m GTR+H4
```

174 By default the above command will infer only one set of equilibrium base
175 frequencies and apply these to all classes. To infer separate equilibrium base
176 frequencies for each class then we must add the `+FO` option:

```
177     iqtree -s data.fst -m GTR+FO+H4
```

178 The above command implements the linked version of the GHOST model.
179 This means that only one set of GTR rate parameters will be inferred and applied
180 to all classes. If one wishes to infer separate GTR rate parameters for each class
181 then the unlinked version is required:

```
182     iqtree -s data.fst -m GTR+FO*H4
```


201 carried out by K&T.

202 *12-taxon simulations.*— The replication of the K&T simulations focused on
203 recovering tree topology only. However, the GHOST model is parameter rich and
204 naturally the validation process must address its ability to accurately recover
205 branch lengths and model parameters. We constructed independent sets of
206 parameters for two classes on a randomly generated 12-taxon tree using the GTR
207 model of evolution. For each class the branch lengths were drawn randomly from an
208 exponential distribution with a mean of 0.1. When specifying a GTR rate matrix
209 in *Seq-Gen*, the G \leftrightarrow T substitution rate is fixed at 1 and all other substitution rates
210 are expressed relatively. Within each class, the five relative substitution rates were
211 drawn randomly from a uniform distribution between 0.5 and 5. The four base
212 frequencies for each class were assigned a minimum of 0.1, with the remainder
213 allocated proportionally by scaling a normalised set of four observations from a
214 uniform distribution. From these two classes MSAs were constructed (again using
215 *Seq-Gen*) by varying the weight of each class. The weight of Class 1, w_1 , was varied
216 from 0.2 to 0.8 in increments of 0.05 and at each increment 20 separate MSAs were
217 simulated. Each MSA was constructed by concatenating two independently
218 simulated sets of sequences, the first of length $10000 \times w_1$ simulated using the Class
219 1 parameters, and the second of length $10000 \times (1 - w_1)$ simulated using the Class

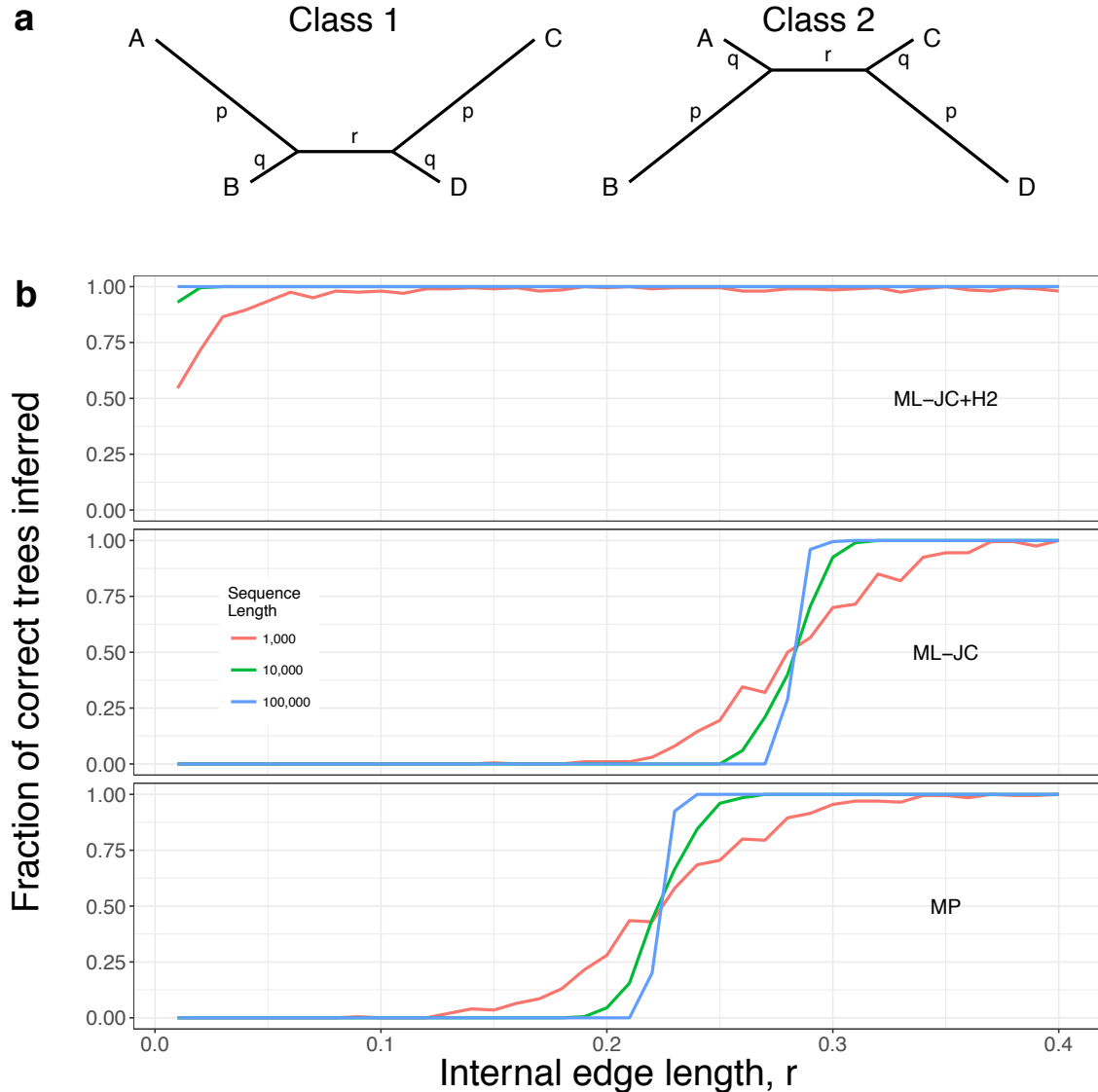


Figure 1: Replication of the simulations of Kolaczkowski & Thornton. (a) We simulated DNA sequences on two symmetric, 4-taxon trees of identical topology using the Jukes-Cantor (JC) model of evolution (Jukes and Cantor, 1969). The branch lengths were constructed such that each tree comprised of two non-sister long branches and two non-sister short branches. Thus each tree was susceptible to long branch attraction (Felsenstein, 1978) (LBA). Importantly, the LBA artefact in both trees was complementary - the bias was in the direction of the AC|BD tree. (b) Performance of maximum likelihood (ML) using a JC, two-class mixture model (ML-JC+H2), ML using a single-class JC model (ML-JC) and maximum parsimony (MP) for data generated under strong heterotachy, $p=0.75$ and $q=0.05$. The length of the internal branch, r , was varied between 0.01 and 0.4 with 200 replicates at each value of r . ML-JC+H2 was able to reliably recover the tree topology for this data even when the internal branch is very short.

220 2 parameters. We used IQ-TREE to infer parameters from each MSA under a
221 GHOST model with two GTR classes (GTR+FO*H2). We also inferred parameters
222 from each MSA under a GTR edge-unlinked partition model.

223 *Parameter recovery: metrics.*— The recovery accuracy of base frequency and
224 relative rate parameters for the 12-taxon simulations was measured by calculating
225 the mean absolute difference between the inferred and true parameters. The
226 accuracy of branch length estimates was assessed using the branch score metric, BS
227 (Kuhner and Felsenstein, 1994). One challenge in assessing accuracy of branch
228 length recovery is that BS is an absolute distance metric. Therefore, we established
229 a frame of reference so that we could assess whether the results obtained are
230 suitably close to the truth or not. To do this we made use of the estimates under
231 the edge-unlinked partition model as a baseline. The fundamental difference
232 between the partition model and the GHOST model is that the partition model has
233 *a priori* knowledge of which sites in the alignment belong to which class. This
234 means that in effect (and excluding the possibility of inferring the incorrect
235 topology) the results of the partition model are identical to those that would be
236 obtained by fitting GTR models to the Class 1 and Class 2 sequences
237 independently. Thus we can consider the accuracy of the partition model as a
238 benchmark.

239 *Convergent Evolution of the $Na_v1.4a$ Gene Among Teleosts*

240 To investigate the performance of the GHOST model using real data we applied it
241 to a sequence alignment (2178 bp) taken from the coding region of a sodium channel
242 gene, $Na_v1.4a$, for 11 teleost species. We used Akaike's Information Criterion (AIC)
243 (Akaike, 1974) to determine the model of sequence evolution and number of classes
244 that provided the best fit to the data. We also used PartitionFinder (Lanfear et al.,
245 2012) and IQ-TREE to fit the best edge-unlinked partition model to the alignment.
246 The data was partitioned based on codon position.

247 **RESULTS & DISCUSSION**

248 *Validation - $K&T$ Simulations*

249 *Experiment 1.*— We fixed $p = 0.75$ and $q = 0.05$ (see Fig. 1a) and varied the
250 internal branch length, r , on the interval $[0.01, 0.4]$ in increments of 0.01. For each
251 value of r , 200 simulated MSAs were constructed by concatenating two
252 sub-alignments of equal length, one simulated on each of the trees in Figure 1a. We
253 carried out phylogenetic inference on each MSA using MP, ML-JC and
254 ML-JC+H2. The experiment was repeated for sequence lengths of 1,000, 10,000
255 and 100,000 base pairs. The results are shown in Figure 1b. We found that both

256 ML-JC and MP were misled when r was short, but as r increased MP recovered
257 before ML. For a sequence length of 100kb, MP was misled to some extent for
258 $r < 0.24$ and ML-JC was misled for $r < 0.3$. These findings mirrored those of K&T
259 precisely. However, the ML-JC+H2 model however was never misled. Figure 1b
260 shows that given sufficient sequence length, the ML-JC+H2 model inferred the
261 correct topology from the heterogeneous sequences 100% of the time with r as low
262 as 0.01. Our results clearly demonstrate that the ML-JC+H2 model can correctly
263 infer the tree topology when ML-JC and MP both are misled by the heterotachous
264 nature of the data.

265 *Experiment 2.*— We tested nine different combinations of $p \in \{0.3, 0.5, 0.7\}$ and
266 $q \in \{0.001, 0.1, 0.2, 0.3, 0.4\}$ (see Fig. 1a). For each of the three methods/models
267 (MP, ML-JC and ML-JC+H2) and at each combination of p and q we determined
268 the smallest value of r (subject to the minimum $r = 0.001$), denoted BL_{50} by K&T,
269 such that the correct topology was returned at least 50% of the time. The results
270 (Fig. 2) indicate that ML-JC+H2 comprehensively outperformed the two
271 alternatives, with the difference most apparent when the influence of heterotachy
272 was strongest (most notably when p is large and q is small). Again the results we
273 observed for MP and ML-JC closely emulated the findings of K&T.

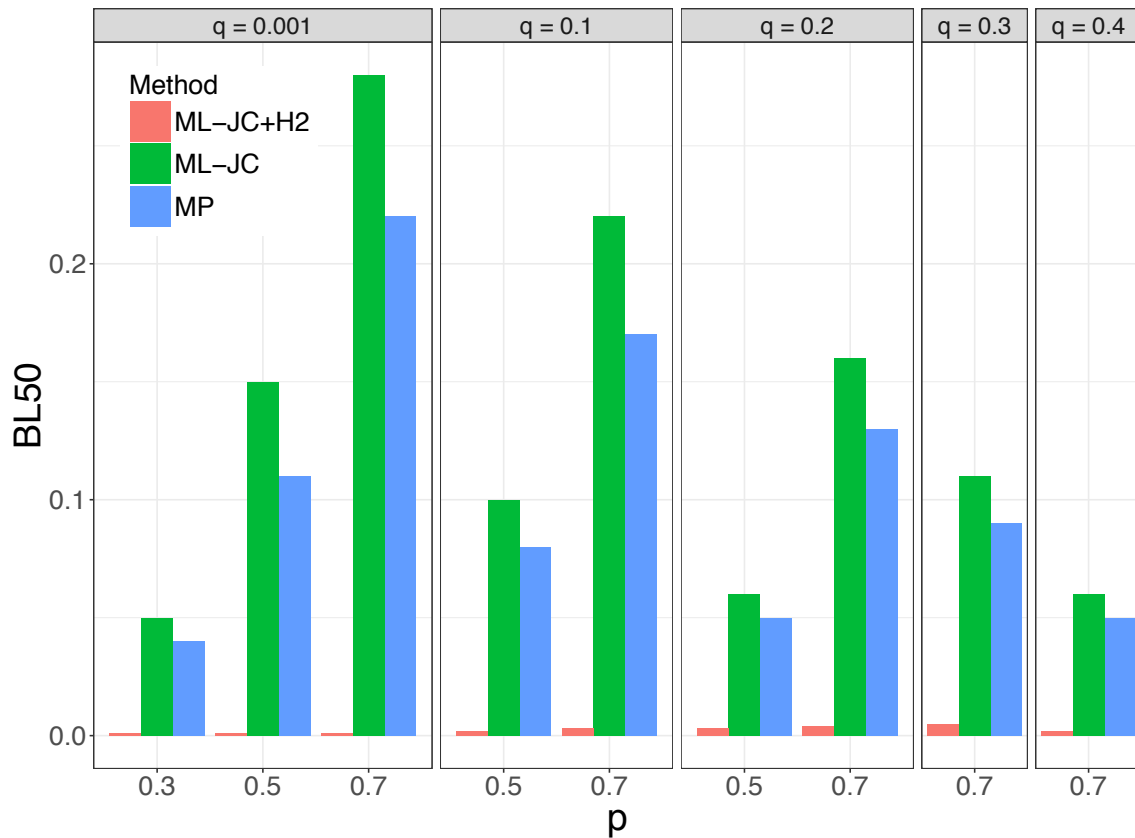


Figure 2: The ML-JC+H2 model clearly outperforms MP and ML-JC over the range of heterotachous conditions tested by K&T. They introduced the BL_{50} measure as the minimum internal branch length required for the method to recover the correct tree topology at least 50% of the time. Small values of BL_{50} indicate that the model is less likely to be misled by the heterotachous nature of the data.

274 *Experiment 3.*— We tested the impact of varying the weight, w , of each class in the
275 simulated MSAs for a variety of branch length combinations. Initially p and q (see
276 Fig. 1a) were fixed at 0.75 and 0.05 respectively, with $r \in \{0.05, 0.15, 0.25\}$ and
277 $w \in \{0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.99\}$. The process was then
278 repeated, this time with p and r fixed at 0.75 and 0.15 respectively, with
279 $q \in \{0.05, 0.15, 0.25\}$ and w as before. Sequence length was held fixed throughout
280 at 100,000bp and 200 replicates were simulated at each combination of branch
281 lengths and weight. We found that for almost all branch length combinations
282 ML-JC+H2 was able to recover the correct topology for all replicates. In the entire
283 experiment, only one dataset (out of 13,200) returned the incorrect topology. The
284 results of K&T indicate that ML-JC could not reliably recover the correct topology
285 for all weights for any of the branch length combinations.

286 The good performance of the GHOST model over the three K&T
287 experiments should be expected in some sense, as ML-JC+H2 enjoys significant
288 advantage over the two alternatives. It is in no way misspecified, having the
289 freedom to fit two classes evolved under the JC substitution model, precisely the
290 conditions used to simulate the data. Conversely, ML-JC has only a single class
291 and therefore is subject to model misspecification. No single set of branch lengths
292 can reproduce the signal present in the simulated alignments. MP is obviously not

293 subject to model misspecification as the method is non-parametric, but it is subject
294 to the long-established artefact of long branch attraction (LBA) (Felsenstein, 1978).
295 Felsenstein showed that having long non-sister branches separated by a relatively
296 short internal branch can result in MP incorrectly inferring the long branches as
297 sisters. Figure 1a shows the two trees used for the classes in the mixture, both
298 sharing the same AB|CD topology. The Class 1 tree has long terminal branches on
299 the A and C lineages, therefore the LBA artefact leads MP to incorrectly favour
300 the AC|BD topology. The Class 2 tree is in a sense the symmetric opposite of the
301 Class 1 tree, it has long terminal edges on the B and D lineages so the result is the
302 same: LBA leads MP to incorrectly infer the AC|BD topology.

303 Therefore the successful replication of the K&T simulations is a necessary
304 but not sufficient condition for the GHOST model's endorsement. It indicates that
305 the implementation of the GHOST model within IQ-TREE's algorithm structure
306 has been successful, but these simulations are on only four taxa and use the most
307 simple model of sequence evolution. Moreover, they only focus on recovering
308 correct tree topology and not inferring branch length parameters.

309 *12-taxon simulations.*— We simulated heterotachously-evolved MSAs on a random
310 12-taxon tree topology under a GTR+FO*H2 model. Using the true GTR+FO*H2

311 model, IQ-TREE accurately recovered the correct tree topology in all 260
312 simulated datasets. Figure 3 shows the performance of the GHOST model in
313 recovering the various tree and model parameters for Class 1 of the simulated data.
314 The analogous plots for Class 2 can be found in Supplementary Figures S1 - S4.
315 The results of the 12-taxon simulations clearly show that under the GTR+FO*H2
316 model IQ-TREE recovered the base frequencies, relative rate parameters and
317 weights to a high degree of accuracy for both classes. With respect to the branch
318 score (BS) (Figs. 3c and S3), we see that the GHOST model again performs very
319 well. The mean BS for the GHOST model approaches that obtained by the
320 partition model as class weight (and therefore share of sequence length in the
321 mixture) increases. This is a very impressive result, given that the partition model
322 enjoys the significant advantage of having full knowledge of which sites were
323 simulated under which class. A BS of zero would imply that the true simulation
324 parameters were inferred for every simulated alignment. Thus, the magnitude of
325 the BS for the partition model can be thought of as a measure of the stochastic
326 simulation error. The difference between the BS for the GHOST and partition
327 models can then be considered the error attributable to losing the knowledge of the
328 partitioning scheme. Clearly this error is negligible in comparison to the simulation
329 error. In Figure 3c, when $w_1 > 0.5$ (or equivalently Fig. S3 when $w_1 < 0.5$), the

330 clear overlap of the error bars (which represent ± 2 standard errors of the mean)
331 suggests that the trees inferred by the GHOST model are not significantly different
332 from those inferred by the partition model. This is a promising result, as in
333 empirical data any partitioning of the MSA is based on assumptions, and therefore
334 introduces a significant potential source of model misspecification. The GHOST
335 model can be applied without any such assumptions.

336 To demonstrate the ability of the GHOST model to provide meaningful
337 information about which sites might belong to which class, we performed a soft
338 classification on one of the MSAs generated for the 12-taxon simulations. For
339 simplicity we have chosen an MSA where Class 1 and Class 2 are of equal weight.
340 Figure 4 clearly indicates, as one would expect, that the probability of a site
341 belonging to Class 1 is generally higher for those sites that were simulated under
342 the Class 1 parameters. However, given the stochastic element of the simulations,
343 there are some sites simulated under the Class 2 parameters that are classified as
344 having a higher probability of evolving under Class 1, and *vice versa*. For this
345 reason we never attempt to hard classify specific sites to a particular class. Rather
346 we consider a specific site's probability distribution of evolving under all of the
347 classes.

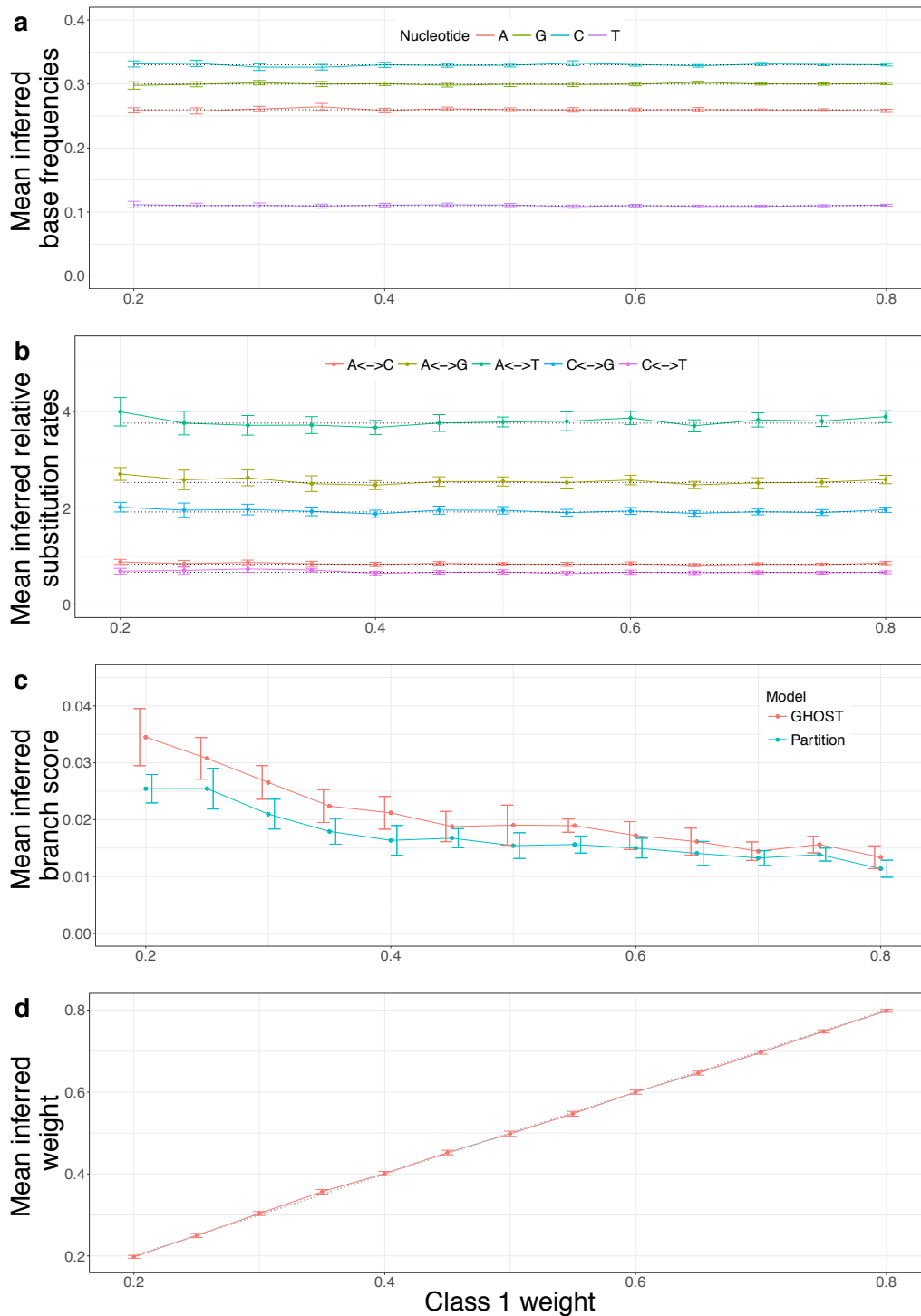


Figure 3: 12-taxon simulations - Class 1 inferred parameters vs Class 1 weight. The data points indicate the mean value of the inferred parameter or statistic, the error bars represent ± 2 standard errors of the mean. Dotted lines represent the true parameter value used for data simulation. (a) Base frequencies (b) Relative substitution rates (c) Branch score (BS) for both the GHOST and partition models (d) Inferred Class 1 weight.

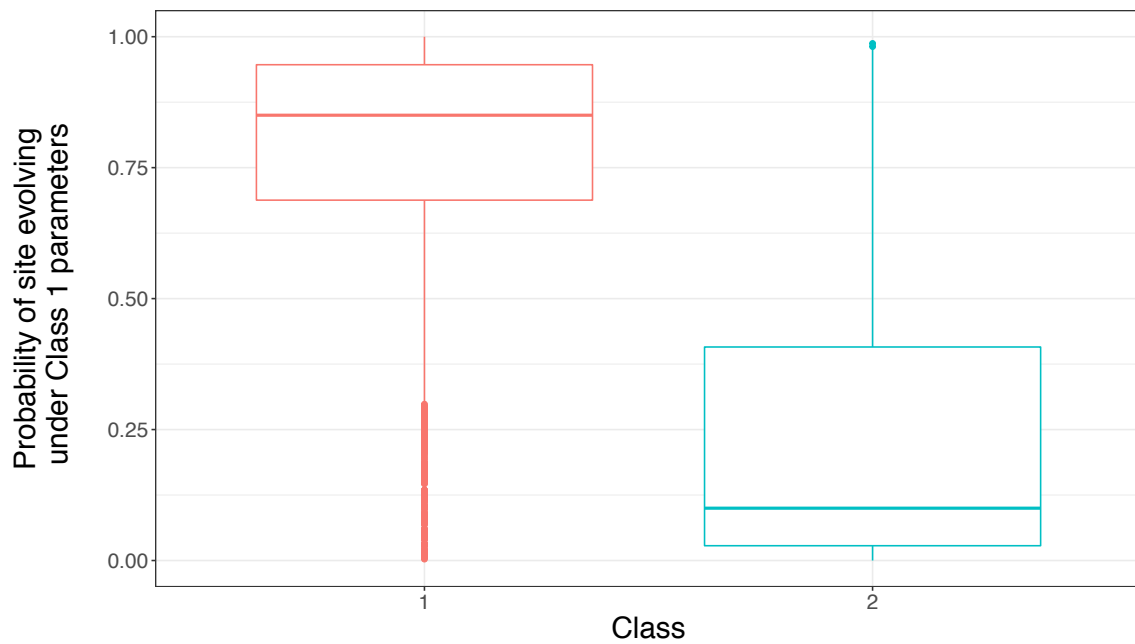


Figure 4: Soft classification of sites to classes - the probability of a site belonging to Class 1 is shown on the y-axis, the two Classes are shown on the x-axis. The boxplots clearly show that in general, sites generated under Class 1 parameters tend to have a higher probability of belonging to Class 1 than sites generated under Class 2.

348 *Convergent Evolution of the $Na_v1.4a$ Gene Among Teleosts*

349 To investigate the performance of the GHOST model using empirical data we
350 applied it to the coding region of a sodium channel gene, $Na_v1.4a$, for 11 teleost
351 species. Zakon et al. (2006) demonstrated the role of this gene in the convergent
352 evolution of the electric organ amongst electric fish species from South America
353 and Africa. AIC determined that GTR+FO*H4 provided the best fit between tree,
354 model and data (Supplementary Fig. S5). The trees inferred by the GHOST model
355 can be found in Figure 5. We then partitioned the electric fish sequence alignment
356 into three partitions, based on codon position (CP). PartitionFinder suggested
357 GTR+FO+G4 (GTR with inferred equilibrium base frequencies plus discrete Γ
358 with four classes) for both the CP1 and CP2 partitions, and GTR+FO+I+G4
359 (same as above but with the inclusion of an invariable sites class) for the CP3
360 partition. We used IQ-TREE to run the partition model with the models indicated
361 by PartitionFinder. The trees inferred by the partition model can be found in
362 Figure 6.

363 We labelled the four classes inferred by the GHOST model in order of
364 increasing total tree length (TTL): the ‘Conserved Class’ ($TTL_{Cons}=0.23$), the
365 ‘Convergent Class’ ($TTL_{Conv}=0.99$), ‘Fast-evolving Class A’ ($TTL_{FEA}=4.06$) and
366 ‘Fast-evolving Class B’ ($TTL_{FEB}=4.18$). Of particular interest is the Convergent

367 Class, so named as it corresponds well to Zakon *et al.*'s (2006) hypothesis of
368 convergent evolution of $Na_v1.4a$ among the South American and African electric
369 fish clades. The convergent class tree displays much more evolution in the electric
370 rather than the non-electric fish lineages (Fig. 7). This is indicative of either a
371 relaxation of purifying selection pressure, an introduction of positive selection
372 pressure or a combination of both. The notable exception is the Brown Ghost
373 Knifefish, which appears relatively conserved. The Brown Ghost Knifefish is unique
374 amongst the other electric fish in the dataset, in that its electric organ has evolved
375 from neural rather than muscle tissue. Consequently in the Brown Ghost Knifefish
376 the $Na_v1.4a$ gene is still expressed in muscle, just as it is in the non-electric fish.
377 The clear distinction in terminal edge length between the Brown Ghost Knifefish
378 and the other electric fishes is obvious and compelling. It provides strong evidence
379 that the GHOST model has indeed identified a subtle component of the historical
380 signal related to the convergent evolution of $Na_v1.4a$, as opposed to returning an
381 arbitrary combination of numerical parameters that happen to maximize the
382 likelihood function. The ability of the GHOST model to isolate such a small
383 component of the signal (the inferred weight of the convergent class being 0.13, the
384 smallest of the 4 classes) is most encouraging. Furthermore, we can expect that the
385 sites belonging with high probability to the convergent class are likely to have been

386 influential in the functional development of the electric organ.

387 *Soft classification of sites to classes.*— The soft classification of sites to classes
388 facilitates the prospective identification of functionally important sites in an
389 alignment. Zakon et al. (2006) report several amino acid sites from the dataset that
390 are influential in the inactivation of the sodium channel, a process critical to
391 electric organ pulse duration. Figure 8a shows that these sites generally have a
392 higher than average probability of belonging to the convergent class in at least one
393 codon position. For example, at amino acid site 647 an otherwise conserved proline
394 (codon CCN) is replaced by a valine (GTN) in the Pintailed Knifefish and a
395 cysteine (TGY) in the Electric Eel. Unique substitutions at codon positions 1 and 2
396 are necessary for both of these amino acid replacements and we find these two sites
397 have a very high probability of belonging to the convergent class. With this result
398 in mind, for each amino acid we summed the probability of codon positions 1 and 2
399 belonging to the Convergent Class. Figure 8b shows the results for the eight amino
400 acid sites with the highest score. Comparing the magnitude of these bars with those
401 of the amino acids in Figure 8a (which are known to be functionally important),
402 one can suspect that these amino acids might also be critical to the operation of
403 the sodium channel gene. Given that there are many other sites in the alignment
404 with a high probability of belonging to the convergent class, one can envisage the

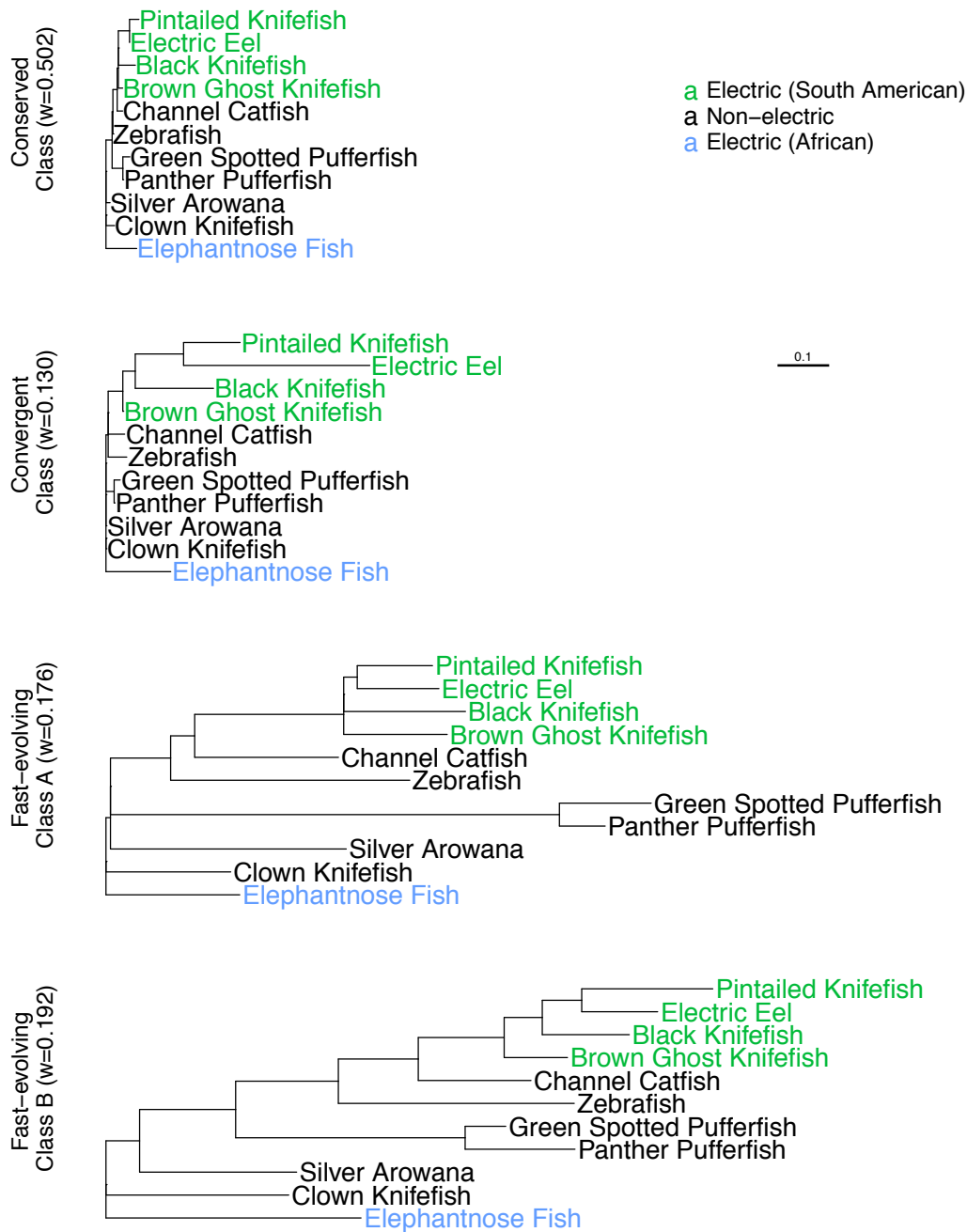


Figure 5: The four trees inferred under the General Time Reversible, four class mixture model (GTR+FO*H4) for the electric fish data. We can clearly see the variability of the branch lengths among the four classes. The classes are displayed in order of increasing tree size, as determined by the sum of the branch lengths. We refer to this as the total tree length (TTL): $TTL_{Cons} = 0.23$, $TTL_{Conv} = 0.99$, $TTL_{FEA} = 4.06$ and $TTL_{FEB} = 4.18$.

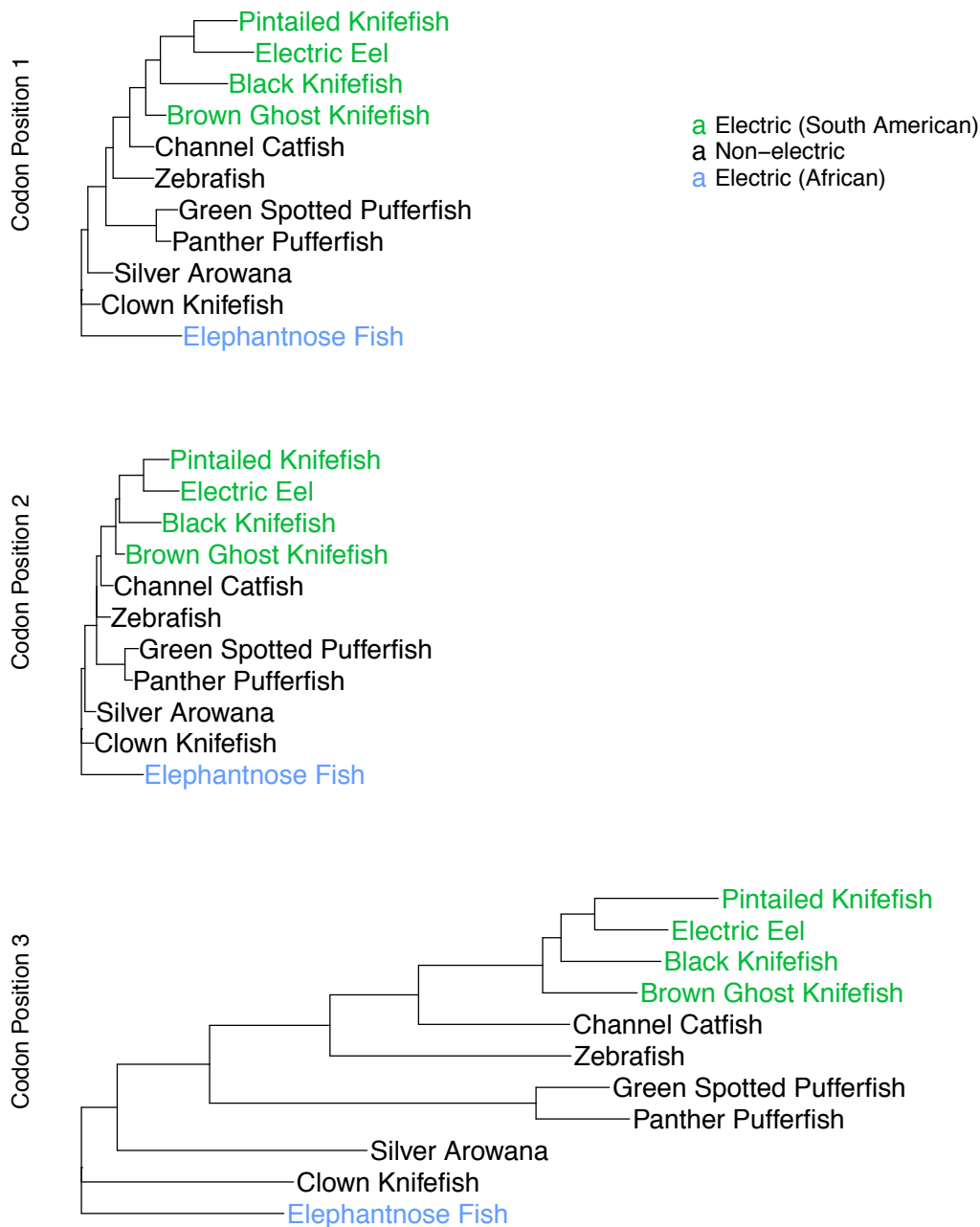


Figure 6: The three trees inferred under the edge-unlinked partition model for the electric fish dataset, with the alignment partitioned based on codon position (CP). The CP1 and CP2 partitions used a GTR+FO+G model, while the CP3 partition used a GTR+FO+I+G model.

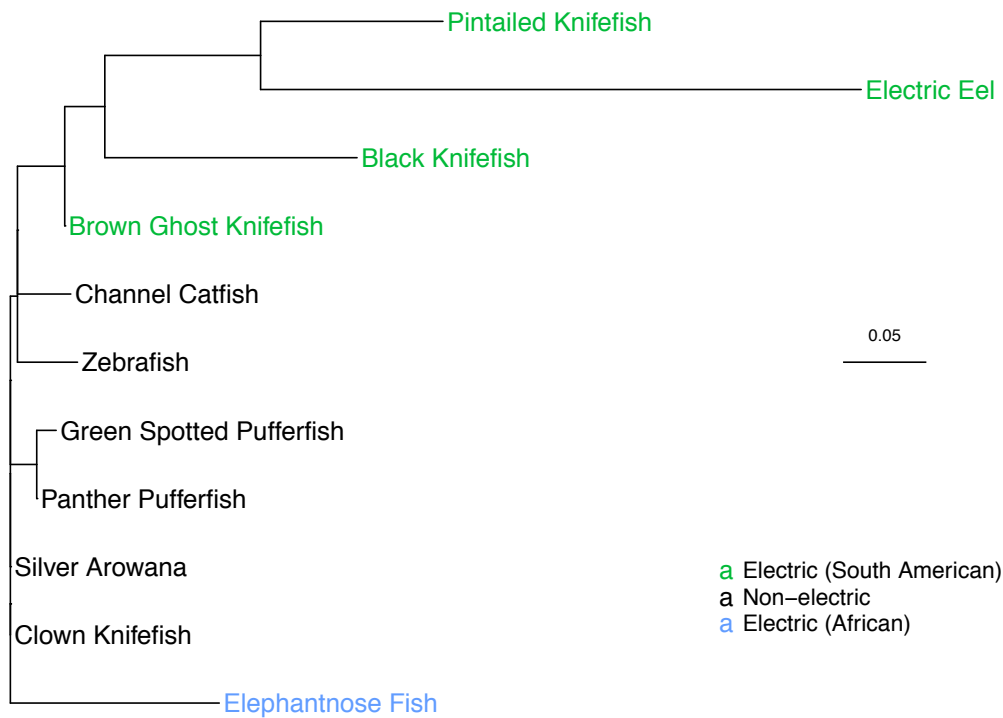


Figure 7: The convergent class inferred by ML-GTR+FO*H4. The 11 fish species comprised four South American electric fish (green), one African electric fish (blue), and six non-electric fish (black) from various locations. The tree for this class shows that in comparison to the electric fish, the non-electric species are relatively conserved.

405 GHOST model helping to identify sites of potential functional importance in an
406 alignment, thereby focusing the experimental work of biologists.

407 In addition to providing insight on an individual site basis, the soft
408 classification can also help to inform us about the nature of the classes themselves.
409 Summing the weighted TTLs for each of the inferred classes results in an estimated
410 1.766 substitutions per site under the inferred model. Table 1 reports the
411 contributions to this figure, stratified by codon position and class. If class
412 membership and codon position were independent attributes of each site then we
413 should expect the contribution of each codon position to be approximately one
414 third for each class. This is not what we observe. Overall we can see that sites in
415 CP1(23%) and CP2 (16%) contribute only 39% of the total of 1.766 substitutions
416 per site. However, within the Conserved and Convergent Classes, sites in CP1 and
417 CP2 are responsible for 90% and 76% of their contribution respectively. This would
418 suggest that a comparatively larger proportion of the substitutions attributed to
419 these classes are non-synonymous: resulting in amino acid replacements that
420 influence the fitness of the organism. We can therefore conclude that even though
421 the Conserved and Convergent Classes are smallest (as determined by substitutions
422 per site), they appear to be the primary catalyst of evolution via natural selection
423 within *Na_v1.4a* amongst these species.

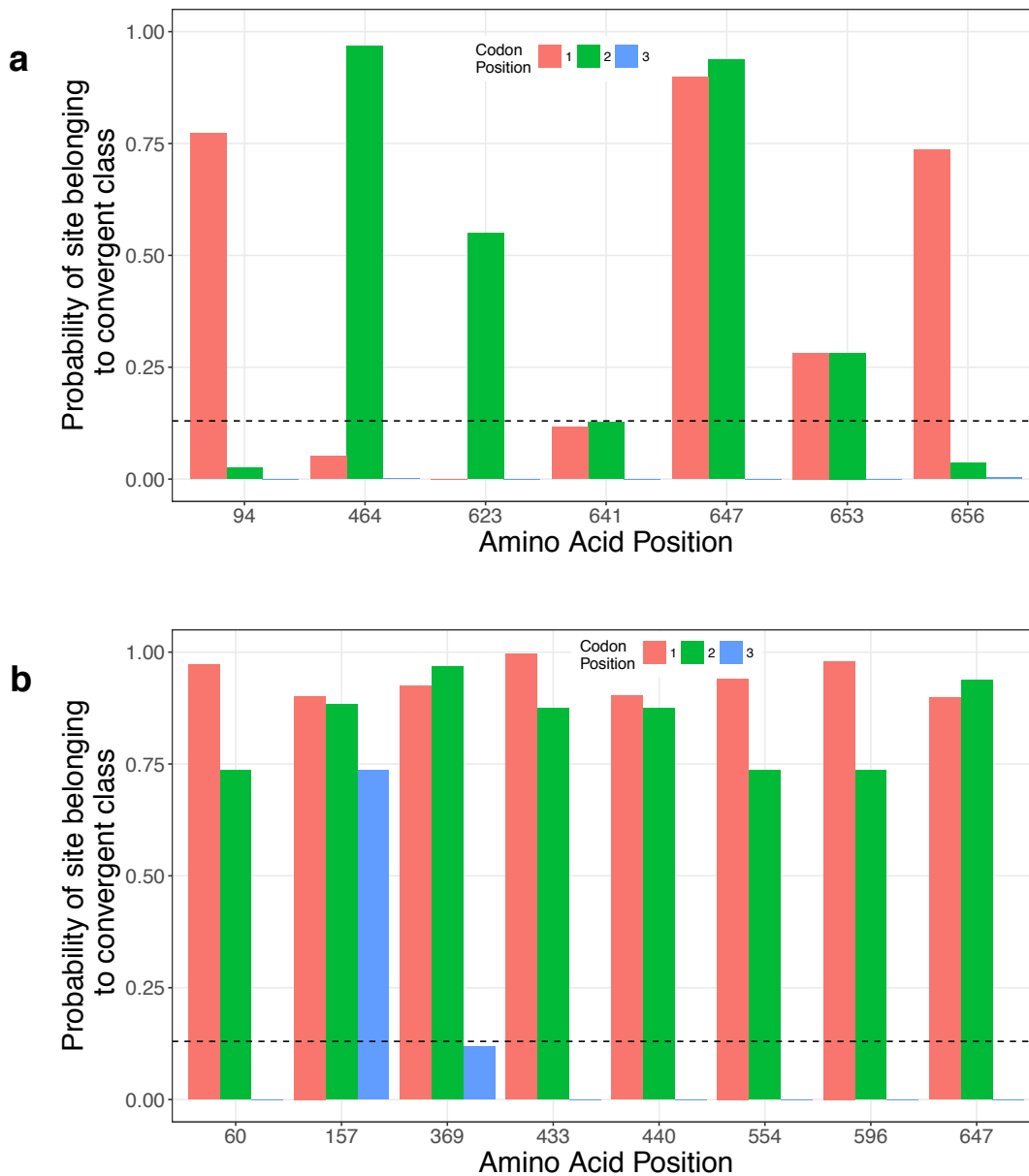


Figure 8: Probability of sites belonging to the convergent class by codon position. (a) The amino acid positions selected correspond with those identified by Zakon et al. (2006) as being functionally important to the inactivation of the Na^+ gene. The horizontal dotted line at 0.13 represents the average probability of belonging to the convergent class over all sites in the alignment. (b) The amino acid positions selected correspond to those with the highest probability of belonging to the convergent class, summed across the first two codon positions.

Class	CP1	CP2	CP3	Subs/site
Conserved	0.049 (41%)	0.058 (49%)	0.012 (10%)	0.119
Convergent	0.051 (40%)	0.047 (36%)	0.031 (24%)	0.129
Fast-evolving A	0.135 (19%)	0.076 (11%)	0.504 (70%)	0.715
Fast-evolving B	0.175 (22%)	0.100 (12%)	0.528 (66%)	0.803
All Classes	0.410 (23%)	0.280 (16%)	1.076 (61%)	1.766

Table 1: Expected number of substitutions per site (bold), weighted by class and separated by codon position (CP). For each inferred class, the expected substitutions per site are calculated by multiplying the total tree length (TTL) by the class weight. The CP1, CP2 and CP3 columns show the contribution to these figures from only the sites within each CP. The grand total indicates that under the parameters inferred by ML-GTR+H4 we would expect 1.766 nucleotide substitutions per site. We can then see, for example, that the Convergent Class is responsible for 0.129 of these substitutions per site. Finally, of the 0.129 substitutions per site attributable to the Convergent Class, 0.051 (or 40%) is the contribution from sites in CP1, 0.047 (36%) is the contribution from sites in CP2 and 0.031 (24%) is the contribution from sites in CP3.

424 *Comparison to the Partition Model.*— It is apparent upon examination of the trees
425 in Figure 6 that the evidence of convergent evolution highlighted by the GHOST
426 model (Fig. 7) has not been recovered by the partition model. None of the three
427 trees in Figure 6 have the distinctive pattern, whereby the majority of the total
428 tree length is associated with the electric fish species (with the exception of the
429 Brown Ghost Knifefish). The reason that the partition model failed to recover this
430 signal is clear when considering the contribution of each CP to the Convergent
431 Class. Table 1 indicates that the substitutions associated with the Convergent
432 Class are attributable to CP1 sites (40%), CP2 sites (36%) and CP3 sites (24%).
433 The partition model constrains the analysis, such that sites in different CPs are
434 modeled independent of each other. It is impossible for a model constrained in such
435 a way to recover the convergent evolution signal, or any other signal whose
436 components are distributed across multiple partitions. The decision to partition the
437 data based on codon position may make sense superficially, but in doing so the
438 analysis is constrained and the results are compromised. We no longer have the
439 ability to uncover the evolutionary stories concealed within the data. We can only
440 hope to obtain those stories that happen not to conflict with the assumptions and
441 constraints that have been placed on the analysis *a priori*. Minimizing these
442 assumptions and constraints where possible, while computationally expensive, is

443 necessary in order to illuminate the evolutionary history without distorting it in
444 the process.

445 *On the Identifiability of the GHOST Model*

446 An ongoing concern with regard to parameter-rich mixture models has been
447 whether or not they are identifiable. There are several examples of theoretically
448 non-identifiable mixture models in the literature (Matsen and Steel, 2007;
449 Štefankovič and Vigoda, 2007b). These examples have inspired much theoretical
450 work on the identifiability or otherwise of different types of phylogenetic mixture
451 models (Allman and Rhodes, 2006; Štefankovič and Vigoda, 2007a; Allman et al.,
452 2008; Allman and Rhodes, 2008; Allman et al., 2011). Of particular interest to the
453 current study, Allman et al. (2011) showed that for a single topology, four taxa,
454 two-class mixture under the JC model, only the tree topology is identifiable but not
455 the branch lengths. This provides a theoretical justification for the procedure
456 carried out by K&T (and replicated here), measuring performance of the models
457 based only on recovery of the topology and paying no attention to recovery of
458 branch length parameters. With regard to the identifiability of the GHOST model
459 more generally, we rely on a result from Rhodes and Sullivant (2012). They
460 established an upper bound on the number of classes for which tree topology,

461 branch lengths and model parameters are identifiable, as a function of the number
462 of character states and the number of taxa. For the simulations we carry out in the
463 current study, with 12 taxa and four character states, the model is identifiable up
464 to a maximum of 15 classes. In the case of the electric fish dataset, with four
465 character states and only 11 taxa, the model is identifiable up to 11 classes.
466 However, there is a technical caveat. The result is shown based on assuming a
467 general Markov model across the tree. There are specific choices of parameters that
468 can result in non-identifiability, but these are of little concern in practical data
469 analysis. Problems arise only when the parameters selected collapse the parameter
470 space to some lower dimension. For example, we could fit the GTR model but if we
471 chose parameters such that all base frequencies were equal and all substitution
472 rates were equal then we are in fact using a JC model, and identifiability may be
473 compromised. However, these technical examples of non-identifiability are not
474 relevant in practice, as in the absence of any constraints there is no likelihood of
475 inferring parameters that collapse the parameter space in such a way.

476 CONCLUSION

477 Heterotachy has been somewhat of an Achilles heel for ML since K&T published
478 their study. The implementation of the GHOST model in IQ-TREE represents a

479 positive advance for ML based phylogenetic inference. Through minimization of
480 model assumptions the GHOST model offers significant flexibility to infer
481 heterotachous evolutionary processes, illuminating historical signals that might
482 otherwise remain hidden. The GHOST model seems well suited to the analysis of
483 phylogenomic datasets, commonly used to address deep phylogenetic questions.
484 While we only present the method and one single-gene empirical example in the
485 current paper, forthcoming empirical studies will compare the performance of the
486 GHOST model to currently popular phylogenomic analysis tools, such as partition
487 and CAT models. One can also envisage many other potential uses for the GHOST
488 model. It could be applied to datasets for which the topology is poorly supported
489 or disputed. It could also provide more accurate parameter estimates, leading to
490 sounder divergence date estimation. The model provides intuitive, biologically
491 meaningful visualizations of the different evolutionary pressures that act on a group
492 of taxa. Structural biologists may find it useful for highlighting functionally
493 important areas within proteins. We have demonstrated its use as a method for
494 identifying changes in selection pressure, as well as bringing to light evidence of
495 convergent evolution. Similarly, one can envisage the GHOST model illuminating
496 the subtle evolutionary relationships between hosts and parasites, disease and
497 immune cells, or the countless evolutionary arms races that are observed

498 throughout the natural world.

499

ACKNOWLEDGEMENTS

500 The authors would like to thank Elizabeth Allman and John Rhodes for helpful
501 discussion about the manuscript.

502 B.Q.M. and A.v.H were supported by the Austrian Science Fund (FWF
503 I-2805-B29).

504 *COMPETING FINANCIAL INTERESTS.*— The authors declare no competing
505 financial interests.

REFERENCES

506

507 Akaike, H. (1974). A new look at the statistical model identification. *IEEE*
508 *Transactions on Automatic Control*, 19(6):716–723.

509 Allman, E. S., Ané, C., and Rhodes, J. A. (2008). Identifiability of a Markovian
510 model of molecular evolution with gamma-distributed rates. *Advances in*
511 *Applied Probability*, pages 229–249.

512 Allman, E. S., Petrovic, S., Rhodes, J. A., and Sullivant, S. (2011). Identifiability
513 of two-tree mixtures for group-based models. *IEEE/ACM Transactions on*
514 *Computational Biology and Bioinformatics (TCBB)*, 8(3):710–722.

515 Allman, E. S. and Rhodes, J. A. (2006). The identifiability of tree topology for
516 phylogenetic models, including covarion and mixture models. *Journal of*
517 *Computational Biology*, 13(5):1101–1113.

518 Allman, E. S. and Rhodes, J. A. (2008). Identifying evolutionary trees and
519 substitution parameters for the general Markov model with invariable sites.
520 *Mathematical Biosciences*, 211(1):18–33.

521 Baele, G., Raes, J., Van de Peer, Y., and Vansteelandt, S. (2006). An improved
522 statistical method for detecting heterotachy in nucleotide sequences. *Molecular*
523 *Biology and Evolution*, 23(7):1397–1405.

- 524 Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from
525 incomplete data via the EM algorithm. *Journal of the Royal Statistical*
526 *Society, Series B*, pages 1–38.
- 527 Felsenstein, J. (1978). Cases in which parsimony or compatibility methods will be
528 positively misleading. *Systematic Biology*, 27(4):401–410.
- 529 Felsenstein, J. (1981). Evolutionary trees from DNA sequences: a maximum
530 likelihood approach. *Journal of Molecular Evolution*, 17(6):368–376.
- 531 Fitch, W. M. and Margoliash, E. (1967). A method for estimating the number of
532 invariant amino acid coding positions in a gene using cytochrome *c* as a model
533 case. *Biochemical Genetics*, 1(1):65–71.
- 534 Gadagkar, S. R. and Kumar, S. (2005). Maximum likelihood outperforms
535 maximum parsimony even when evolutionary rates are heterotachous.
536 *Molecular Biology and Evolution*, 22(11):2139–2141.
- 537 Holmquist, R., Goodman, M., Conroy, T., and Czelusniak, J. (1983). The spatial
538 distribution of fixed mutations within genes coding for proteins. *Journal of*
539 *Molecular Evolution*, 19(6):437–448.
- 540 Jayaswal, V., Wong, T. K., Robinson, J., Poladian, L., and Jermin, L. S. (2014).
541 Mixture models of nucleotide sequence evolution that account for

- 542 heterogeneity in the substitution process across sites and across lineages.
543 *Systematic Biology*, 63(5):726–742.
- 544 Jukes, T. and Cantor, C. (1969). Evolution of protein molecules. In Munro H.N.
545 *Mammalian Protein Metabolism*, pages 21–123, New York: Academic Press.
- 546 Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K., von Haeseler, A., and Jermini,
547 L. S. (2017). Modelfinder: fast model selection for accurate phylogenetic
548 estimates. *Nature Methods*, 14(6):587–589.
- 549 Kolaczkowski, B. and Thornton, J. W. (2004). Performance of maximum
550 parsimony and likelihood phylogenetics when evolution is heterogeneous.
551 *Nature*, 431(7011):980–984.
- 552 Kuhner, M. K. and Felsenstein, J. (1994). A simulation comparison of phylogeny
553 algorithms under equal and unequal evolutionary rates. *Molecular Biology and*
554 *Evolution*, 11(3):459–468.
- 555 Lanfear, R., Calcott, B., Ho, S. Y., and Guindon, S. (2012). PartitionFinder:
556 combined selection of partitioning schemes and substitution models for
557 phylogenetic analyses. *Molecular Biology and Evolution*, 29(6):1695–1701.
- 558 Lartillot, N. and Philippe, H. (2004). A Bayesian mixture model for across-site

- 559 heterogeneities in the amino-acid replacement process. *Molecular Biology and*
560 *Evolution*, 21(6):1095–1109.
- 561 Lopez, P., Casane, D., and Philippe, H. (2002). Heterotachy, an important process
562 of protein evolution. *Molecular Biology and Evolution*, 19(1):1–7.
- 563 Matsen, F. A. and Steel, M. (2007). Phylogenetic mixtures on a single tree can
564 mimic a tree of another topology. *Systematic Biology*, 56(5):767–775.
- 565 Meade, A. and Pagel, M. (2008). A phylogenetic mixture model for heterotachy. In
566 *Evolutionary Biology from Concept to Application*, pages 29–41. Springer.
- 567 Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., and Minh, B. Q. (2015).
568 IQ-TREE: a fast and effective stochastic algorithm for estimating
569 maximum-likelihood phylogenies. *Molecular Biology and Evolution*,
570 32(1):268–274.
- 571 Pagel, M. and Meade, A. (2004). A phylogenetic mixture model for detecting
572 pattern-heterogeneity in gene sequence or character-state data. *Systematic*
573 *Biology*, 53(4):571–581.
- 574 Pagel, M. and Meade, A. (2005). Mixture models in phylogenetic inference.
575 *Mathematics of Evolution and Phylogeny*, pages 121–142.

- 576 Philippe, H. and Lopez, P. (2001). On the conservation of protein sequences in
577 evolution. *Trends in Biochemical Sciences*, 26(7):414–416.
- 578 Philippe, H., Zhou, Y., Brinkmann, H., Rodrigue, N., and Delsuc, F. (2005).
579 Heterotachy and long-branch attraction in phylogenetics. *BMC Evolutionary*
580 *Biology*, 5(1):50.
- 581 Rambaut, A. and Grassly, N. C. (1997). Seq-Gen: an application for the Monte
582 Carlo simulation of DNA sequence evolution along phylogenetic trees.
583 *Computer Applications in the Biosciences: CABIOS*, 13(3):235–238.
- 584 Rhodes, J. A. and Sullivant, S. (2012). Identifiability of large phylogenetic mixture
585 models. *Bulletin of Mathematical Biology*, 74(1):212–231.
- 586 Spencer, M., Susko, E., and Roger, A. J. (2005). Likelihood, parsimony, and
587 heterogeneous evolution. *Molecular Biology and Evolution*, 22(5):1161–1164.
- 588 Steel, M. (2005). Should phylogenetic models be trying to fit an elephant? *Trends*
589 *in Genetics*, 21(6):307–309.
- 590 Štefankovič, D. and Vigoda, E. (2007a). Phylogeny of mixture models: Robustness
591 of maximum likelihood and non-identifiable distributions. *Journal of*
592 *Computational Biology*, 14(2):156–189.

- 593 Štefankovič, D. and Vigoda, E. (2007b). Pitfalls of heterogeneous processes for
594 phylogenetic reconstruction. *Systematic Biology*, 56(1):113–124.
- 595 Wang, H.-C., Li, K., Susko, E., and Roger, A. J. (2008). A class frequency mixture
596 model that adjusts for site-specific amino acid frequencies and improves
597 inference of protein phylogeny. *BMC Evolutionary Biology*, 8(1):331.
- 598 Whelan, N. V. and Halanych, K. M. (2017). Who let the CAT out of the bag?
599 Accurately dealing with substitutional heterogeneity in phylogenomic analyses.
600 *Systematic Biology*, 66(2):232–255.
- 601 Wu, J. and Susko, E. (2009). General heterotachy and distance method
602 adjustments. *Molecular Biology and Evolution*, 26(12):2689–2697.
- 603 Wu, J. and Susko, E. (2011). A test for heterotachy using multiple pairs of
604 sequences. *Molecular Biology and Evolution*, 28(5):1661–1673.
- 605 Yang, Z. (1994). Maximum likelihood phylogenetic estimation from DNA sequences
606 with variable rates over sites: approximate methods. *Journal of Molecular*
607 *Evolution*, 39(3):306–314.
- 608 Zakon, H. H., Lu, Y., Zwickl, D. J., and Hillis, D. M. (2006). Sodium channel genes
609 and the evolution of diversity in communication signals of electric fishes:

610 convergent molecular evolution. *Proceedings of the National Academy of*
611 *Sciences of the United States of America*, 103(10):3675–3680.