

Version dated: December 7, 2018

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 JERMIN^{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.*

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

*Corresponding author: Stephen Crotty, Center for Integrative Bioinformatics

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

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

[‡] Current address: Research School of Biology, Australian National University,

Canberra, ACT 2601, Australia

25 *Abstract.*— Molecular sequence data that have evolved under the influence of
 26 heterotachous evolutionary processes are known to mislead phylogenetic inference.
 27 We introduce the General Heterogeneous evolution On a Single Topology (GHOST)
 28 model of sequence evolution, implemented under a maximum-likelihood framework
 29 in the phylogenetic program IQ-TREE (<http://www.iqtree.org>). Simulations show
 30 that using the GHOST model, IQ-TREE can accurately recover the tree topology,
 31 branch lengths and substitution model parameters from heterotachously-evolved
 32 sequences. We develop a model selection algorithm based on simulation results, and
 33 investigate the performance of the GHOST model on empirical data by sampling
 34 phylogenomic alignments of varying lengths from a plastome alignment. We then
 35 carry out inference under the GHOST model on a phylogenomic dataset composed
 36 of 248 genes from 16 taxa, where we find the GHOST model concurs with the
 37 currently accepted view, placing turtles as a sister lineage of archosaurs, in contrast
 38 to results obtained using traditional variable rates-across-sites models. Finally, we
 39 apply the model to a dataset composed of a sodium channel gene of 11 fish taxa,
 40 finding that the GHOST model is able to infer a subtle component of the historical
 41 signal, linked to the previously established convergent evolution of the electric
 42 organ in two geographically distinct lineages of electric fish. We compare inference
 43 under the GHOST model to partitioning by codon position and show that, owing

44 to the minimization of model constraints, the GHOST model is able to offer unique
45 biological insights when applied to empirical data.

46 Keywords: Phylogenetics, heterotachy, mixture model, maximum likelihood,
47 convergent evolution

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

66 The effect of heterotachy on phylogenetic inference was thrust into the

67 spotlight by Kolaczkowski and Thornton (K&T) (2004). They used a simulation
68 study to show that heterotachously-evolved sequences could mislead the popular
69 inference methods of maximum-likelihood (ML) and Bayesian Markov Chain
70 Monte-Carlo (BMCMC) to a greater extent than maximum parsimony (MP). Their
71 findings were controversial and were widely challenged on the grounds that the
72 simulations captured only a special case of heterotachy (Gadagkar and Kumar,
73 2005; Philippe et al., 2005; Spencer et al., 2005; Steel, 2005), and more general
74 studies of heterotachy concluded that ML performed at least as well as, and in
75 most cases better than, MP (Gadagkar and Kumar, 2005; Spencer et al., 2005).
76 Valid as these criticisms may have been, the key issue that the K&T study brought
77 to light stood firm - heterotachy was a primary source of model misspecification
78 and the models and methods of the time were ill-equipped to deal with it.

79 The main impediment to the development of models that can accommodate
80 heterotachously-evolved sequences has been the computational expense. Models
81 that account for heterogeneity of rates of change across sites can be integrated
82 relatively cheaply, but modeling heterotachy is not so simple. One approach has
83 been covarion (COV) models (Fitch and Markowitz, 1970). Tuffley and Steel
84 (1998) described a model in which sites could switch between variable and
85 invariable states in different lineages. All variable sites in the model shared a

86 common substitution model and rate. This model was gradually extended (Galtier,
87 2001; Huelsenbeck, 2002), eventually reaching its most complex form in which sites
88 can switch along lineages between a number of different rates as well as an
89 invariable state (Wang et al., 2007).

90 Another approach has been to use partition models (Lanfear et al., 2012),
91 which require the data to be partitioned *a priori*. The analysis then proceeds by
92 inferring separate branch length and model parameters for each partition. Sequence
93 data are commonly partitioned based on genes and/or codon position. However,
94 the inherent assumption of this approach is that heterotachy only occurs between
95 partitions, not within each partition. This may not be a valid assumption, so the
96 requirement to partition the data in advance of the analysis is a possible source of
97 model misspecification.

98 An alternative approach has been to use mixture models, in which the
99 likelihood of the data at each site in the alignment is calculated as a weighted sum
100 across multiple classes (see Pagel and Meade (2005) for a detailed description of
101 phylogenetic mixture models). The most common approaches can be referred to as
102 mixed substitution rate (MSR) models (Lartillot and Philippe, 2004; Pagel and
103 Meade, 2004; Foster, 2004), whereby each class has its own substitution rate
104 matrix; and mixed branch length (MBL) models (Kolaczkowski and Thornton,

2004; Meade and Pagel, 2008), whereby each class has its own set of branch lengths on the tree. Hybrid versions of these models have also been proposed, such as the HAL-HAS model of Jayaswal et al. (2014). Zhou et al. (2007) compared a covarion model to an MBL model, finding the covarion model to be more efficient at handling heterotachy. They did however conclude that both methods warranted further exploration, going on to propose the covarion mixture model (CM) (Zhou et al., 2010), which incorporates covarion parameters that vary across sites. As a consequence of their parameter rich nature, these models have all been implemented only within a Bayesian framework. Wu and Susko (2009) proposed a general framework for heterotachy, encompassing both MSR and MBL models as special cases. Another example is the CAT models of Lartillot and Philippe (2004), which have been widely used (Whelan and Halanaych (2017) and references therein). Whelan and Halanaych (2017) carried out extensive simulation and empirical studies comparing the performance of the CAT models to partition models. They concluded that despite their additional complexity and associated increase in runtime, the CAT models generally perform no better than partition models. They also lamented that when new mixture models are introduced in the literature their performance is not always assessed against the current popular methods for phylogenetic analysis, such as partition models.

As a consequence of their varied nature, mixture models require many parameters and the associated computational expense has thus far impeded their implementation in a ML framework. The issue of computational expense is an ever diminishing one; as computing power increases and algorithmic architecture improves, the opportunity to employ more and more complex models of sequence evolution does also. We introduce the General Heterogeneous evolution On a Single Topology (GHOST) model for ML inference. The GHOST model combines features of both MSR and MBL models. It consists of a number of classes, all evolving on the same tree topology. For each class, the branch lengths, nucleotide or amino-acid frequencies, substitution rates and class weight are all parameters to be inferred. The motivation behind this modelling philosophy is the desire to minimise assumptions that might lead to model misspecification. Although the cost of this philosophy, in terms of model complexity and the associated risk of over-parameterisation, is not to be ignored, by refraining from placing strict constraints on the inference we allow the opportunity to recover new, and perhaps surprising, historical signals from the data. We provide an easy-to-use, ML implementation of the GHOST model in the phylogenetic program IQ-TREE (Nguyen et al., 2015) (<http://www.iqtree.org>), the first mixture model of comparable flexibility to be made available in a ML framework.

METHODS AND MATERIALS

Model Description

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

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

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

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

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

Model Parameter Estimation for a Fixed Tree, T

Given a fixed tree topology, T , let

$\Theta = \{w_1, \dots, w_m, \lambda_1, \dots, \lambda_m, R_1, \dots, R_m, F_1, \dots, F_m\}$ denote the GHOST model parameters (*i.e.*, class weights, branch lengths, relative substitution rates, and nucleotide or amino-acid frequencies) for each of the m classes. To estimate all parameters for a tree, T , we employ the expectation-maximization (EM) algorithm (Dempster et al., 1977; Wang et al., 2008). We initialize Θ with all $\hat{R}_j = \mathbf{1}$ in each class, uniform nucleotide or amino-acid frequencies \hat{F}_j (*i.e.*, the Jukes-Cantor model), and \hat{w}_j and $\hat{\lambda}_j$ obtained by parsimonious branch lengths rescaled according to the rate parameters of a discrete, PDF rate model (Kalyaanamoorthy et al., 2017) with m categories. This becomes the current estimate $\hat{\Theta}$. The EM algorithm iteratively performs an expectation (E) step and a maximization (M) step to update the current estimate until an optimum in likelihood is reached.

Derivation of Expectation-Maximization algorithm

The premise underlying the GHOST model is that each site evolved according to

just one of the m classes, however we do not have any information about which sites belong to which class. We define $\mathbf{c} = \{c_1, c_2, \dots, c_N\}$ as a vector that maps the N sites to one of the m classes. The EM algorithm works by formulating an expression for the expected value of our objective function and then maximizing that expectation. In the context of GHOST, we can restate the likelihood equation as follows:

$$\ell = \sum_{j=1}^m \sum_{i=1}^N I\{c_i = j\} \log \left(L_{ij}(T, \boldsymbol{\lambda}_j, \mathbf{R}_j, \mathbf{F}_j) \right),$$

where $I\{c_i = j\}$ is an indicator function that is equal to 1 when the class of the i^{th} site is equal to j , and 0 otherwise. Taking the expectation of this expression yields:

$$E[\ell] = \sum_{j=1}^m \sum_{i=1}^N E[I\{c_i = j\}] \log \left(L_{ij}(T, \boldsymbol{\lambda}_j, \mathbf{R}_j, \mathbf{F}_j) \right),$$

E-step.— In the context of the GHOST mixture model, the goal of the E-step is to evaluate the quantity $E[I\{c_i = j\}]$ for a fixed set of tree and model parameters. An intuitive interpretation of the expected value of this indicator function, is that it is simply the probability that a given site i belongs to a given class j . For simplicity,

we define this quantity as \hat{p}_{ij} and evaluate it using a simple application of Bayes Theorem. Given the current parameter estimates, we can calculate \hat{p}_{ij} as follows:

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

M-step.— The goal of the M-step is then to update the parameter estimates to maximize the expected likelihood, fixing the \hat{p}_{ij} that were calculated during the E-step. For each class j , we maximize the expectation of the log-likelihood function:

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

to obtain the next $\hat{\lambda}_j^{NEW}$, $\hat{\mathbf{R}}_j^{NEW}$, $\hat{\mathbf{F}}_j^{NEW}$. Within IQ-TREE, $\hat{\lambda}_j^{NEW}$ is obtained via Newton-Raphson optimization, while $\hat{\mathbf{R}}_j^{NEW}$ and $\hat{\mathbf{F}}_j^{NEW}$ are estimated by the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm (Fletcher, 2013). Finally, the weights are updated by:

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

198 That is, the new weight for class j is the mean posterior probability of each site
 199 belonging to class j . This completes the proposal of the new estimate $\hat{\Theta}^{NEW}$. If
 200 $\ell(\hat{\Theta}^{NEW}) > \ell(\hat{\Theta}) + \epsilon$ (where ϵ is a user-defined tolerance, $\epsilon = 0.01$ by default), then
 201 $\hat{\Theta}$ is replaced by $\hat{\Theta}^{NEW}$ and the E and M steps are repeated. Otherwise, the EM
 202 algorithm finishes.

203 An auxiliary benefit of the ML implementation of the GHOST model in
 204 IQ-TREE is that once the EM-algorithm has converged, we can soft-classify sites to
 205 classes, according to their probability of belonging to a particular class. This
 206 classification can be used to identify sites in the alignment that belong with high
 207 probability to a particular class of interest.

208 *Tree search*

209 The tree search algorithm in IQ-TREE (Nguyen et al., 2015) is based on the
 210 construction of a candidate tree set. Trees from the candidate tree set are
 211 rearranged by Nearest Neighbour Interchange (NNI) to explore the tree space. This
 212 algorithm was tested extensively during the ML implementation of the GHOST
 213 model and two significant changes to the heuristic were required:

- 214 1. In the original implementation of IQ-TREE, after each NNI is performed, a
 215 single branch length optimization step for the five branches adjacent to the

216 NNI is carried out. We found that this amount of branch length optimization
217 was insufficient for the GHOST model. Instead, IQ-TREE now performs m
218 branch length optimization steps after each NNI (where m is the number of
219 classes in the GHOST model).

220 2. Prior to the ML implementation of the GHOST model, IQ-TREE only fully
221 optimised the model parameters of the best tree in the candidate tree set.
222 During the ML implementation of the GHOST model we found that this
223 technique proved to provide too much of an advantage to the current best
224 tree. The algorithm was modified such that when the GHOST model is used,
225 all trees in the set of candidate trees are fully optimized.

226 *Software*

227 The GHOST model has been implemented in IQ-TREE (Nguyen et al., 2015)
228 (<http://www.iqtree.org>). A list of commands for use of the GHOST model in
229 IQ-TREE can be found in the supplementary material (available on the Dryad data
230 repository, doi: TBA).

231 *Testing of the GHOST Model in IQ-TREE*

232 We tested the efficacy of the ML implementation of the GHOST model in

233 IQ-TREE by carrying out three separate simulation studies. The first study was a
 234 replication of the simulations carried out by Kolaczkowski and Thornton (2004),
 235 focusing on IQ-TREE’s ability to recover the correct tree topology from
 236 heterotachously-evolved data on quartet trees. We found that using the GHOST
 237 model, IQ-TREE was able to reliably recover the simulation parameters in all
 238 cases. The methods, results (Supplementary Figs. S1 and S2) and discussion of
 239 these simulations can be found in the supplementary material. The second study
 240 used 12-taxon trees and focused on IQ-TREE’s ability to recover branch length and
 241 substitution model parameters from heterotachously-evolved data. The third study
 242 used 32-taxon trees and focused on the establishment of a sound model selection
 243 procedure, specifically to determine the correct number of classes from simulated
 244 alignments. Finally, we investigated the effect of using the incorrect number of
 245 classes on topological accuracy.

246 *12-taxon simulations.*— The replication of the Kolaczkowski and Thornton
 247 simulations focused on recovering tree topology only. However, the GHOST model
 248 is parameter rich and naturally the implementation process must assess the ability
 249 of IQ-TREE to accurately recover branch lengths and model parameters under the
 250 GHOST model. We constructed independent sets of parameters for two classes on
 251 a randomly generated 12-taxon tree using the GTR model of evolution. For each

class the branch lengths were drawn randomly from an exponential distribution with a mean of 0.1. We then used *Seq-Gen* (Rambaut and Grassly, 1997) to simulate MSAs. When specifying a GTR rate matrix in *Seq-Gen*, the G \leftrightarrow T substitution rate is fixed at 1 and all other substitution rates are expressed relatively. Within each class, the five relative substitution rates were drawn randomly from a uniform distribution between 0.5 and 5. The four base frequencies for each class were assigned a minimum of 0.1, with the remainder allocated proportionally by scaling a normalised set of four observations from a uniform (0, 1) distribution. From these two classes MSAs were constructed by varying the weight of each class. The weight of Class 1, w_1 , was varied from 0.2 to 0.8 in increments of 0.05 and at each increment 20 separate MSAs were simulated. Each MSA was constructed by concatenating two independently simulated sets of sequences, the first of length $10,000 \times w_1$ simulated using the Class 1 parameters, and the second of length $10,000 \times (1 - w_1)$ simulated using the Class 2 parameters. We used IQ-TREE to infer parameters from each MSA under a GHOST model with two GTR classes (GTR+FO*H2). We also inferred parameters from each MSA under a partitioned GTR model, where the branch length parameters were unlinked (*i.e.*, estimated separately for each partition). We also repeated the procedure with a range of shorter sequence lengths: 100, 500, 1000, and 5000 nucleotides.

271 The accuracy of inferred base frequency and relative rate parameters for the
 272 12-taxon simulations was measured by calculating the mean absolute difference
 273 between the inferred and true parameters. The accuracy of branch length estimates
 274 was assessed using the branch score metric, BS (Kuhner and Felsenstein, 1994). In
 275 order to assess the accuracy of branch length recovery we needed to establish a
 276 frame of reference to gauge whether the results obtained are suitably close to the
 277 truth or not. To do this we made use of the estimates under the branch-unlinked
 278 partition model as a baseline. The fundamental difference between the partition
 279 model and the GHOST model is that the partition model has *a priori* knowledge of
 280 which sites in the alignment belong to which class. This means that in effect (and
 281 excluding the possibility of inferring the incorrect topology) the results of the
 282 partition model are identical to those that would be obtained by fitting GTR
 283 models to the Class 1 and Class 2 sequences independently. Naturally we cannot
 284 expect that the GHOST model can perform better than this, so we can consider
 285 the accuracy of the partition model as a benchmark.

286 *Model selection*

287 *32-taxon simulations.*— In order for the GHOST model to be used on empirical
 288 sequence alignments we must have some method of model selection, in particular

selecting the appropriate number of classes. Information criterion methods such as Akaike’s Information Criterion (AIC) (Akaike, 1974) or Bayesian Information Criterion (BIC) (Schwarz et al., 1978) are commonly used in phylogenetics for this purpose, so we carried out simulations to establish whether these criteria could accurately predict the correct number of classes that generated the alignment. We also investigated the influence of the number of classes inferred on topological accuracy. We generated 300 heterotachous sequence alignments for each of $m = 2, 3$ and 4 classes. Each alignment was 10,000 bp long, contained 32 taxa and used the GTR model of sequence evolution for each class. The weight of each class, w_i , was held fixed at $\frac{1}{m}$. For each alignment, the model parameters for each of the m classes were generated as in the 12-taxon simulations. For each alignment, a ‘base set’ of branch lengths, $\boldsymbol{\lambda}$, was generated randomly from an exponential distribution with a mean of 0.1. The branch length parameters for the m classes were then generated as follows:

1. For the i^{th} class, a vector of random variables (of same length as $\boldsymbol{\lambda}$), \mathbf{s}_i , was drawn from a uniform distribution on $(0, 1)$.
2. For the i^{th} class, a class scaling factor, α_i , was drawn from a uniform distribution on $(0, 1)$.
3. Finally, an overall scaling factor, β , was calculated to ensure that the

308 weighted total tree length (TTL) of the m classes was equivalent to the TTL
309 of the ‘base set’:

$$\beta = \frac{\sum \lambda}{\sum_{i=1}^m w_i \alpha_i \mathbf{s}_i \lambda}$$

310 4. The branch length vectors for the i^{th} class were then given by:

$$\lambda_{Ci} = \beta \alpha_i \mathbf{s}_i \lambda$$

311 For the i^{th} class, we used *Seq-Gen* to simulate a sequence alignment of length
312 $10,000 \times w_i$ bp. These were concatenated together to form the heterotachous
313 alignment. This procedure was repeated to generate 300 heterotachous sequence
314 alignments for each of $m \in \{2, 3, 4\}$.

315 For each of the 900 simulated alignments we used IQ-TREE to fit GHOST
316 models with $1, 2, 3, \dots, 8$ classes. For each alignment, we used AIC and BIC (where
317 we used sequence length as the proxy for n in the BIC formula) to determine the
318 number of classes that provided the best fit between tree, model and data. In order
319 to select between models when AIC and BIC do not agree, we examined the

320 inferred parameters when too many classes were used to see if we could recognise
321 characteristic signs of overfitting. We also investigated the influence of the inferred
322 number of classes on the topological accuracy, as measured by the Robinson-Foulds
323 (RF) distance (Robinson and Foulds, 1981). Finally, we investigated the
324 computation time required for IQ-TREE to arrive at ML estimates under the
325 GHOST model, as a function of the number of classes in the model (Supplementary
326 Figure S3).

327 *Plastome alignments.*— In order to investigate the variability in the number of
328 classes recommended by AIC and BIC on empirical alignments, we created separate
329 empirical alignments by subsampling from a plastome alignment, taken from Yan
330 et al. (2017), which consisted of 66 genes for 26 species. We discarded all genes
331 shorter than 1000 bp, leaving a total of 20 genes. From these 20 genes, we
332 randomly sampled 20 groups of 1, 3, 5, 10 and 15 genes to create a total of 100
333 separate alignments. We then fitted GHOST models with increasing number of
334 classes to each alignment to determine the number of classes that provided the best
335 fit according to both AIC and BIC.

336 *Placement of Turtles Among Archosaurs*

337 One can think of the linked version of the GHOST model in terms of the discrete Γ

model, with the removal of some constraints. The linked GHOST model does not require the classes to be of equal weight, nor does it impose that the branch lengths between classes are correlated. The PDF rate model can be thought of as an intermediate step between the discrete Γ and the linked GHOST models. To demonstrate the effect of relaxing these constraints we applied 4-class discrete Γ , PDF rate and linked GHOST models to a phylogenomic alignment consisting of 248 genes (187,026 bp) for 16 taxa. The alignment was taken from Chiari et al. (2012), in which they concluded that turtles were a sister group to birds and crocodiles, as opposed to crocodiles only.

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

We applied the GHOST model to a sequence alignment (2178 bp) taken from the coding region of a sodium channel gene, $Na_v1.4a$, for 11 teleost species.

Model selection is the first challenge when using the GHOST model on an empirical alignment. We tested a wide variety of substitution models, as shown in Supplementary Figure S4. Starting with the two-class GHOST model, we used IQ-TREE to optimise the likelihood of the data under each substitution model. Subsequently, we repeated the process with up to a maximum of six classes. For each run we used the unlinked version of the GHOST model, so that each class had

its own set of branch lengths, base frequencies and substitution model parameters inferred. We then applied our information theory-based model selection criteria to determine the substitution model and number of classes that provided the best fit. For the best GHOST model, we also tested the linked versions to evaluate whether inferring model parameters individually for each class was necessary. Finally, we found the best PDF rate model (Kalyaanamoorthy et al., 2017) and compared that to the best GHOST model based on AIC.

In order to compare the GHOST model to alternative current phylogenetic methods, we also used IQ-TREE to fit a branch-unlinked partition model. The electric fish alignment was split into three partitions, based on codon structure. We then used PartitionFinder (Lanfear et al., 2012) to evaluate the best substitution models to use on each partition. Finally, IQ-TREE was used to fit the best branch-unlinked partition model to the alignment, using the models of sequence evolution suggested by PartitionFinder.

RESULTS & DISCUSSION

12-taxon simulations

We simulated heterotachously-evolved MSAs of varying lengths (100, 500, 1,000, 5,000 and 10,000 bp) on a random 12-taxon tree topology, with two classes evolving

374 according to the GTR model of evolution. Figure 1 shows the performance of the
375 GHOST model in recovering the various tree and model parameters for Class 1 of
376 the 10,000 bp simulated alignments. The analogous plots for all sequence lengths
377 and both classes can be found in Supplementary Figures S5- S12. The results of
378 the 12-taxon simulations show that under the GTR+FO*H2 model IQ-TREE
379 recovered the base frequencies, relative rate parameters and weights to a high
380 degree of accuracy for both classes. With respect to the branch score (BS) (Fig. 1c
381 and Supplementary Figs. S9 and S10), we see that the GHOST model again
382 performs well. The mean BS for the GHOST model approaches that obtained by
383 the partition model as class weight (and therefore share of sequence length in the
384 mixture) increases, despite the partition model enjoying the advantage of having
385 full knowledge of which sites were simulated under which class. A BS of zero would
386 imply that the true simulation parameters were inferred for every simulated
387 alignment. Thus, the magnitude of the BS for the partition model can be thought
388 of as a measure of the stochastic simulation error. The difference between the BS
389 for the GHOST and partition models can then be considered the error attributable
390 to losing the knowledge of the partitioning scheme. This error appears negligible in
391 comparison to the simulation error. In Figure 1c, when $w_1 > 0.5$, the overlap of the
392 error bars (which represent ± 2 standard errors of the mean) suggests that the trees

393 inferred by the GHOST model are not significantly different from those inferred by
394 the partition model. This is a promising result, as in empirical data any
395 partitioning of the MSA is based on assumptions, and therefore introduces a
396 potential source of model misspecification. The GHOST model can be applied
397 without any such assumptions.

398 To demonstrate the ability of the GHOST model to provide meaningful
399 information about which sites might belong to which class, we performed a soft
400 classification on one of the MSAs generated for the 12-taxon simulations. For
401 simplicity we have chosen an MSA where Class 1 and Class 2 are of equal weight.
402 Supplementary Figure S13 indicates that the probability of a site belonging to
403 Class 1 is generally higher for those sites that were simulated under the Class 1
404 parameters. However, given the stochastic element of the simulations, there are
405 some sites simulated under the Class 2 parameters that are classified as having a
406 higher probability of evolving under Class 1, and *vice versa*. For this reason, we
407 never attempt to ‘hard classify’ the sites, that is, allocating specific sites to a
408 particular class with absolute certainty. Rather, we ‘soft classify’ the sites, that is,
409 we consider that a site belongs to all classes, according to its probability
410 distribution of evolving under each class.

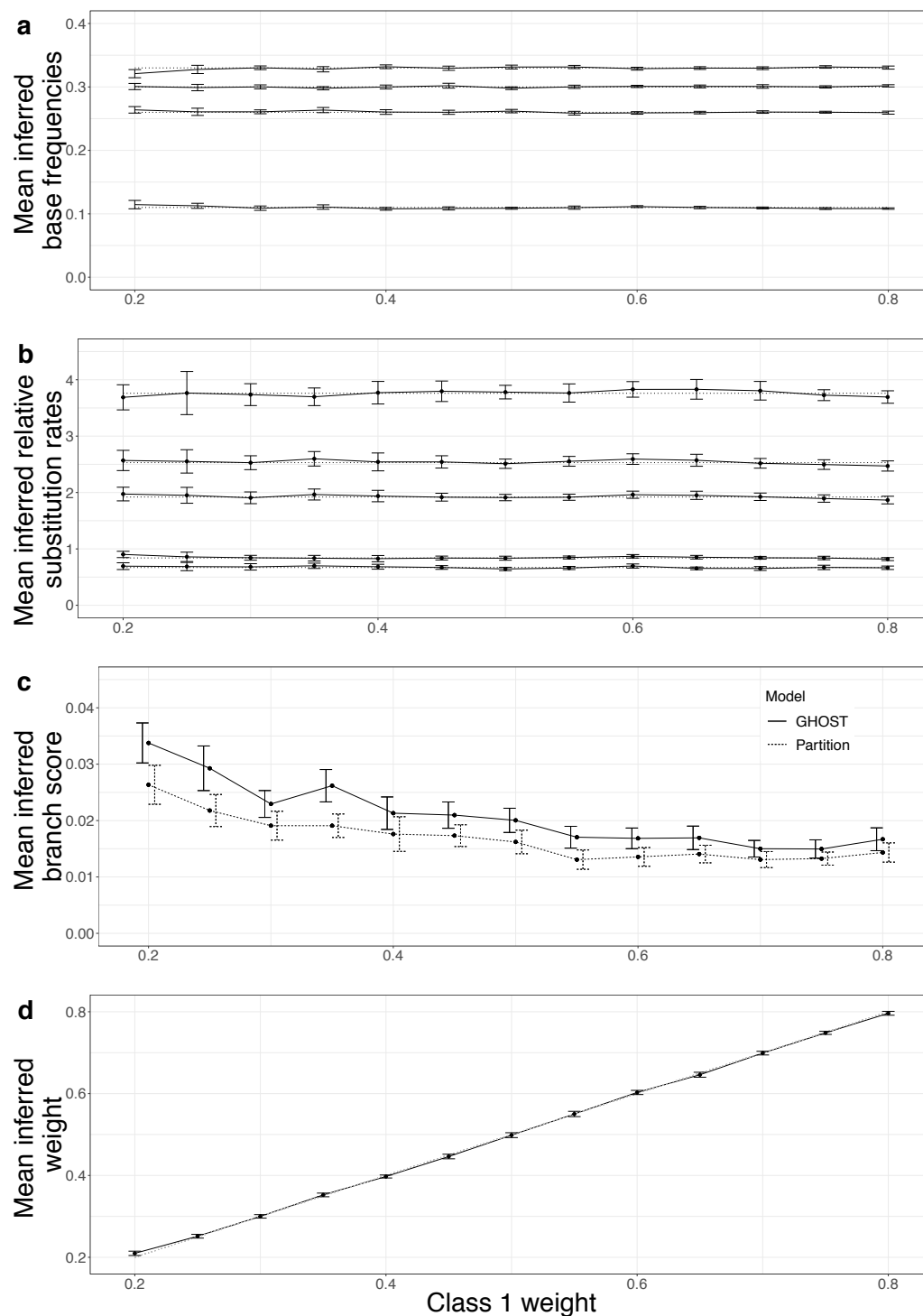


Figure 1: 12-taxon simulations, 10,000 bp alignments - Class 1 inferred parameters vs Class 1 weight. The data points indicate the mean value of the inferred parameter or statistic, the error bars represent ± 2 standard errors of the mean. Dotted lines represent the true parameter value used for data simulation. (a) Base frequencies (b) Relative substitution rates (c) Branch score (BS) for both the GHOST and partition models (d) Inferred Class 1 weight.

411 *The effect of sequence length.*— An important consideration when employing
 412 parameter rich models is the amount of information in the alignment. Estimating
 413 many parameters from an insufficient amount of information will result in
 414 unreliable parameter estimates. Supplementary Figure S14 shows that the GHOST
 415 model and the partition model recover the correct tree topology at similar rates.
 416 For simulated alignments of length 100 bp, tree inference was poor for both
 417 GHOST (30.8% inferred trees correct) and the partition model (33.5%). This
 418 failing is quickly remedied by increasing sequence length, with topological accuracy
 419 for both models greater than 90% for 500 bp alignments. When looking at
 420 parameter inference we see a similar story. Supplementary Figures S5 to S12 show
 421 the progressive improvement in the accuracy of inferred parameters as sequence
 422 length is increased. For sequence lengths of 100 bp the parameter estimates are
 423 completely unreliable, as is the inferred topology. This is not surprising given the
 424 dearth of information on which to base the inference. As sequence length increases
 425 so does the strength of the phylogenetic signal from each class. At 500 and 1,000
 426 bp, the estimates are reasonably close to the true values but still exhibit a
 427 moderate level of variance. For 5,000 and 10,000 bp the parameter estimates are
 428 very close to the true values and with little variance. These 12-taxon, 2-class
 429 simulations have a total of 59 free parameters to be estimated. Based on these

430 results it would seem prudent when applying the GHOST model to empirical
431 datasets to ensure a minimum of $10k$ sites in the alignment, where k is the number
432 of free parameters under the proposed model.

433 *Model Selection*

434 *32-taxon simulations.*— The primary purpose of the 32-taxon simulations was to
435 establish a sound model selection technique to allow the GHOST model to be
436 applied to empirical alignments with confidence. Information theory methods such
437 as AIC and BIC are typically used by phylogeneticists to choose amongst models.
438 How these two methods perform on complex mixture models such as GHOST is
439 unclear. Zhou et al. (2007) found that when applied to models with high numbers
440 of parameters, AIC tended to overfit the data (inclusion of parameters is penalised
441 too lightly) whereas BIC tended to underfit the data (inclusion of parameters is
442 penalised too heavily). Dziak et al. (2018) counsel that while information criteria
443 are useful guides, they do have their limitations, and so nuance and judgment
444 remain important elements in the model selection process.

445 For each of the 900 simulated alignments (300 for each $m \in \{2, 3, 4\}$, 10,000
446 bp long), we used AIC and BIC to determine the optimal number of classes for
447 IQ-TREE to infer under the GHOST model. The results are summarised in Table

	Inferred number of classes									
m	1	2	3	4	5	6	7	8	Total	
AIC	2	0	285	14	1	0	0	0	0	300
	3	0	0	292	6	2	0	0	0	300
	4	0	0	0	298	2	0	0	0	300
BIC	2	0	300	0	0	0	0	0	0	300
	3	0	5	295	0	0	0	0	0	300
	4	0	0	30	270	0	0	0	0	300

Table 1: 32-taxon simulations, model selection using AIC or BIC. For each of the 900 simulated alignments (300 for each $m \in \{2, 3, 4\}$, where m is the true number of simulated classes), we used AIC and BIC to determine the optimal number of classes to infer under the GHOST model.

1. AIC selects the correct number of classes in 95% of cases for $m = 2$, always
 erring on the side of overfitting. As m increases, the accuracy of AIC rises to more
 than 99% for $m = 4$. BIC selects the correct number of classes 100% of the time for
 $m = 2$, but the accuracy of BIC decreases as m increases, dropping to 90% for
 $m = 4$. Conversely to AIC and in line with expectations based on the literature,
 BIC always erred on the side of underfitting.

Plastome alignments.— The results of the 32-taxon simulations discussed above
 indicate that BIC and AIC agree on the number of classes in the vast majority of

cases, so there is little ambiguity in the model selection process. However, this may not be the case in empirical alignments. We subsampled genes from a phylogenomic alignment to create 100 different alignments, 20 each of single-gene, 3-gene, 5-gene, 10-gene and 15-gene alignments. Supplementary Figure S15 shows the level of variability between the number of classes recommended by BIC and AIC. It is apparent that the broad agreement between BIC and AIC when applied to simulated alignments is not mirrored in empirical data. One reason for this might be that when applied to the simulated alignments, the true model is available as one of the candidate models and so both criteria tend to select this model or something quite close to it. This is obviously not the case for empirical data, and so this may explain why we see considerably more variation in the results between the criteria. Regardless, it does highlight that when applying the GHOST model to empirical alignments, choosing the number of classes requires a more nuanced approach to be developed.

Choosing the number of classes.— Model selection can be thought of as a trade-off between bias (the chosen model has too few parameters to adequately represent the underlying evolutionary processes) and variance (the model has too many parameters to provide stable parameter estimates) (Burnham and Anderson, 2003; Posada and Buckley, 2004). Given that the primary motivation behind the

development of the GHOST model was the minimization of model misspecification, we should prefer modest overfitting to modest underfitting. A model that has too many classes has the advantage that the true model is nested within it, and therefore the true parameters remain recoverable, albeit with some undesirable redundancy. Conversely, a model with too few classes must merge at a minimum two classes into one, and therefore the true parameters are not recoverable. Thus, we can respectively consider the BIC and AIC-based optimal number of classes as a lower and upper bound on the number of classes in the best-fit GHOST model. The challenge is to find a way to sensibly choose the optimal number of classes between these bounds.

Intuitively, there does not seem to be any way to predict the effect of underfitting (fitting less classes than was used to generate the data) on the inferred parameters. However, the same is not true of overfitting. If we fit too many classes then we may expect one of two things to happen:

1. We will recover the true branch lengths, model parameters and weights for the correct number of classes, with any remaining classes having weight very close to zero.
2. We may have two or more inferred classes in which the inferred branch lengths and model parameters are very similar to each other, with the sum of

494 their weights being approximately equal to the weight of a single true class.

495 To investigate whether the above propositions hold, we examined the results of the

496 32-taxon, 4-class simulations. We compared results when the correct, 4-class

497 GHOST model was used vs. the overfit, 5-class GHOST model. Each of the 300

498 alignments were simulated from 4 classes of equal weight, so we decided to look at

499 the weight of the inferred classes as a proxy for successful recovery of the true

500 classes. It would seem an unlikely coincidence if we inferred four classes of

501 approximately equal weight without them closely resembling the true classes in

502 terms of branch lengths and model parameters. Supplementary Figure S16 shows

503 the variability in inferred class weights when the correct number of classes

504 (GTR+FO*H4) is used. With the exception of a few outliers, all the inferred

505 classes have a weight close to the simulated value (0.25). If we add an extra class

506 such that the model is overfitted (GTR+FO*H5) as in Supplementary Figure S17 ,

507 we see that the three largest inferred classes have weight close to the true value of

508 0.25, with the remaining weight split between the two smallest classes. The

509 rightmost box in Figure S17 shows the sum of the weights for the two smallest

510 classes, which is again consistently close to 0.25. These observations are consistent

511 with the propositions outlined above, in which three of the four simulated classes

512 are inferred with reasonable accuracy, while the fourth simulated class is split into

513 two inferred classes. To further explore this hypothesis, we checked if the branch
 514 lengths inferred in the smallest class were more strongly correlated with the branch
 515 lengths of the second smallest inferred class than those of the three largest classes.
 516 We expect this effect to be stronger as the weight of the smallest inferred class
 517 increases. Supplementary Figure S18 shows two matrices, displaying the correlation
 518 between the branch lengths inferred by the five classes. The classes are ordered by
 519 weight, C1 referring to the largest class and C5 referring to the smallest.
 520 Supplementary Figure S18 (a) shows the correlation for alignments in which the
 521 inferred weight of C5 was less than 0.05, and we see that for these cases (138 of 300
 522 alignments) the branch lengths of C5 and C4 have a correlation of 0.32; whereas
 523 (b) shows the correlation among those alignments for which it was greater than
 524 0.05, and we see that for these cases (162 of 300 alignments) the correlation is much
 525 higher at 0.79. Based on this evidence, we can highlight two characterisitic signs of
 526 overfitting, that users can check for when choosing the number of classes to fit to
 527 empirical data:

- 528 1. One or more of the inferred classes has a negligible weight in comparison with
 529 the other classes.
- 530 2. The trees of two or more of the inferred classes show strong similarities in
 531 terms of branch lengths.

532 It is possible then to recommend the following approach when selecting the
533 number of classes to fit to empirical alignments:

- 534 1. Calculate the maximum number of classes that is reasonable to use for a
535 given alignment. The criterion that should be used is that the number of free
536 parameters in the model must be no more than $\frac{N}{10}$, where N is the length of
537 the alignment. Call this number U .
- 538 2. Without exceeding U , find the optimal number of classes as judged by both
539 BIC (call this number L) and AIC (if this is less than U , then update U
540 accordingly). We now have a lower (L) and upper (U) bound on the number
541 of classes that should be considered as potentially providing the best fit.
- 542 3. Examine the class weights and trees inferred by the GHOST model with U
543 classes. If none of the class weights are negligible and the trees are all
544 reasonably distinct then accept U as the optimal number of classes for the
545 GHOST model.
- 546 4. If signs of overfitting are present, examine the class weights and trees for the
547 GHOST model with $U - 1$ classes. If no overfitting is present then accept this
548 number of classes as optimal.
- 549 5. Repeat Step 4, continually removing classes until no signs of overfitting are

550 present, or until L is reached.

551 The forthcoming discussion of the convergent evolution of the $\text{Na}_v1.4\text{a}$ is an
552 example of an empirical alignment in which AIC appears to give a reasonable
553 number of classes, with no signs of overfitting present. A counter example is
554 provided Crotty et al. (2018), where AIC is found to overfit the data whereas BIC
555 offers a more reasonable fit.

556 *Impact of model misspecification.*— While we consider the model selection
557 procedure outlined above to be reasonable, it must be remembered that it is not
558 deterministic and it was developed based on the performance of the GHOST model
559 on simulated alignments. We must therefore recognise the potential for
560 over/underfitting to occur in practice with empirical alignments, and assess the
561 potential impact of such errors. To do so we used the 32-taxon simulations to
562 investigate the effect of choosing the wrong number of classes on IQ-TREE’s ability
563 to infer the correct topology under the GHOST model. We calculated the RF
564 distance between the trees used for simulation and those inferred by IQ-TREE.
565 Figure 2 displays the mean RF distance as a function of the number of classes in
566 the fitted model, expressed relative to m , the true number of classes used to
567 simulate the alignments. As we should expect, for all values of m the mean RF

distance is minimized when m classes are inferred. However, the mean RF distance increases much faster in the presence of underfitting than it does in the presence of overfitting. This finding supports the use of the top-down approach when choosing the number of classes as described above, as any errors will tend to be on the side of overfitting rather than underfitting. Detailed summaries of the distribution of RF distances are given in Supplementary Tables S1 , S2 and S3 for $m = 2, 3$ and 4, respectively.

Placement of Turtles Among Archosaurs

The placement of turtles in the phylogenetic tree of amniotes has been controversial, due in part to their morphological peculiarities. It is currently accepted that turtles are a sister lineage to archosaurs (birds and crocodiles), as opposed to crocodiles alone. Chiari et al. (2012) assembled and analyzed a 248-gene, 187,026 nucleotide alignment of 16 taxa, concluding that the tendency to place turtles as sister to crocodiles was a phylogenetic artefact caused by saturation at codon position 3 sites. They found the preferred grouping of turtles as sister to archosaurs was returned when the alignment was partitioned by codon position or when only codon position 1 and 2 sites were included. Among the models that returned the non-preferred topology was the GTR+G, with four rate categories. To

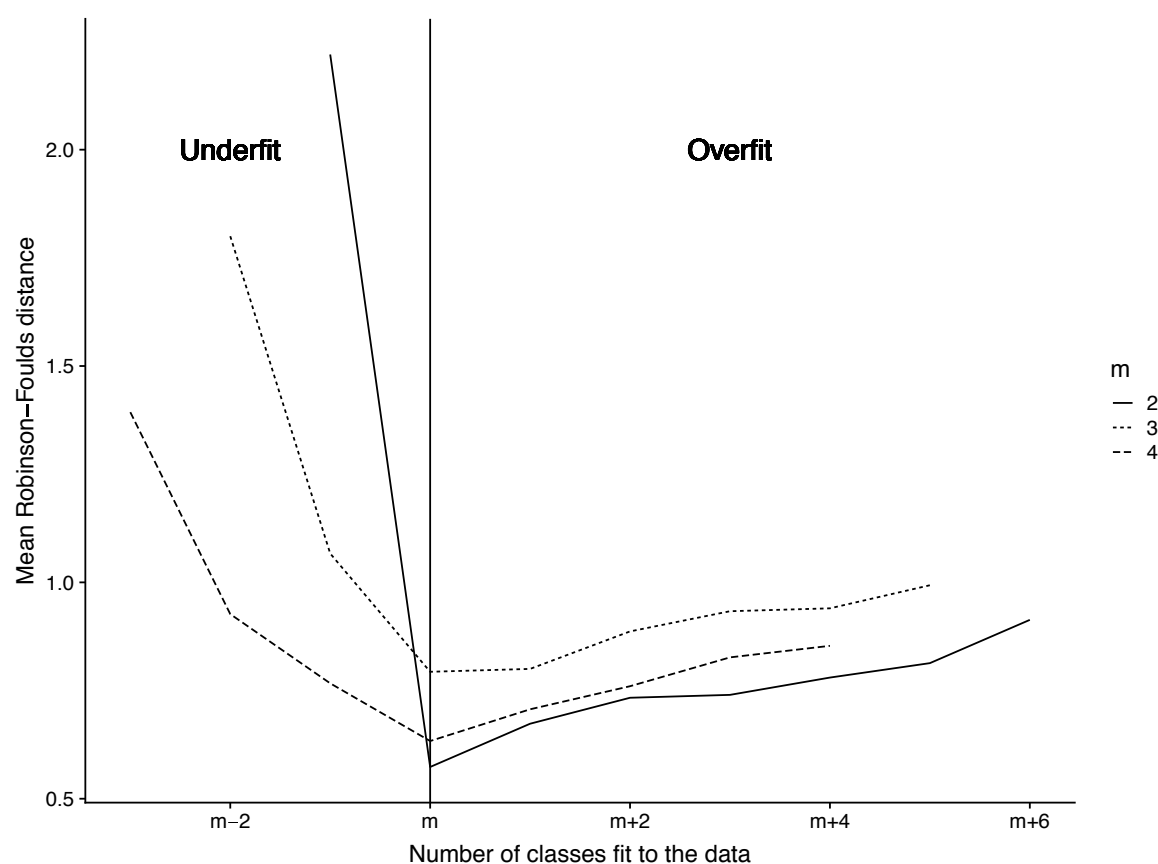


Figure 2: 32-taxon simulations, effect of under/overfitting on topological accuracy, for the 900 simulated alignments (300 each for $m \in \{2, 3, 4\}$). The y-axis displays the mean RF distance between the inferred trees and the trees used to simulate the alignment. The x-axis shows the number of classes used for the inference, expressed relative to m , the true number of classes used to simulate the alignments.

586 examine the influence on this result of the restrictions imposed by the discrete Γ
 587 model, we tested the discrete Γ , the PDF rate model and the GHOST model on the
 588 same alignment. In order to ensure a fair comparison all models used four classes
 589 (as in Chiari et al. (2012) and the linked version of the GHOST model was used.
 590 Supplementary Table S4 indicates that the GHOST model proved superior in terms
 591 of both AIC and BIC. The resulting tree topologies can be found in Figure 3,
 592 showing that the discrete Γ and PDF rate models returned the turtles and
 593 crocodiles grouping, whereas the GHOST model returned the turtles and archosaurs
 594 grouping. Therefore, the GHOST model is not misled by the saturation found at
 595 codon position 3 sites, whereas the discrete Γ and PDF rate models are.

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

597 *Model selection and interpretation.*— To investigate its performance on empirical
 598 data, we applied the GHOST model to the coding region of a sodium channel gene,
 599 $Na_v1.4a$, for 11 teleost species. Zakon et al. (2006) demonstrated the role of this
 600 gene in the convergent evolution of the electric organ amongst electric fish species
 601 from South America and Africa. AIC determined that GTR+FO*H4 (AIC=27602)
 602 provided the best fit between tree, model and data (Supplementary Fig. S4).
 603 Conversely, BIC determined that GTR+FO*H2 provided the best fit. Examining

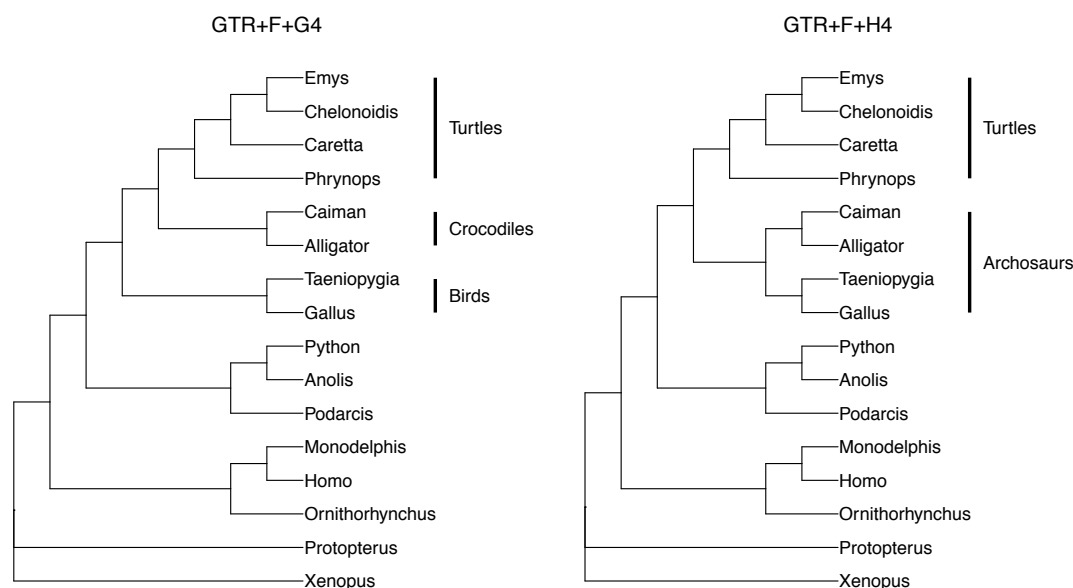


Figure 3: Turtle alignment - The two different topologies obtained from the turtle alignment. The topology on the left is returned by the 4-class discrete Γ and PDF rate models and places turtles as sister to crocodiles. The topology on the right is returned by the 4-class unlinked GHOST model and places turtles as sister to archosaurs (crocodiles and birds).

the class weights and trees (Figure 4) inferred by GTR+FO*H4 indicates that all classes have non-negligible weight (minimum class weight is 0.13) and all four trees appear reasonably distinct. Thus, we conclude that there are no obvious signs of overfitting present, and we accept four classes as optimal for this alignment. We also tested the empirical base frequencies version (GTR+F*H4, AIC=27749) and linked substitution rates version (GTR+FO+H4, AIC=27860). Each of these models returned a significantly higher AIC value, indicating that the unlinked version provided the best fit. We then tested the PDF rate model, finding that the best such model had six classes (GTR+FO+R6), but still a much higher AIC (27813) than that of the GTR+FO*H4 model. In order to confirm the stability of the parameter estimates we repeated the analysis using the GTR+FO*H4 model 100 times.

We then partitioned the electric fish sequence alignment into three partitions, based on codon position (CP). PartitionFinder suggested GTR+FO+G4 (GTR with inferred equilibrium base frequencies plus discrete Γ with four classes) for both the CP1 and CP2 partitions, and GTR+FO+I+G4 (same as above but with the inclusion of an invariable sites class) for the CP3 partition. We used IQ-TREE to run the codon partition model with the models indicated by PartitionFinder. The trees inferred by the partition model can be found in

623 Supplementary Figure S19.

624 *Interpretation of results.*— We labelled the four classes inferred by IQ-TREE under
625 the GTR+FO*H4 model in order of increasing TTL: the ‘Conserved Class’
626 ($TTL_{Cons}=0.23$), the ‘Convergent Class’ ($TTL_{Conv}=0.99$), ‘Fast-evolving Class A’
627 ($TTL_{FEA}=4.06$) and ‘Fast-evolving Class B’ ($TTL_{FEB}=4.18$). Of particular
628 interest is the Convergent Class, so named as it corresponds well to Zakon *et al.* ’s
629 (2006) hypothesis of convergent evolution of $Na_v1.4a$ among the South American
630 and African electric fish clades. They explained that the $Na_v1.4a$ gene arose from a
631 single gene duplication event which occurred in a species ancestral to all 11 fish
632 species in the alignment, and was historically expressed in muscle tissue. They then
633 show that the gene is now expressed in the electric organ of all but one of the
634 electric fish species in both the South American and African electric fishes, but
635 obviously not in the non-electric fishes. Since these lineages constitute two separate
636 clades, one conclusion that can be drawn is that this morphological trait evolved
637 twice independently, once in the South American clade and once in the African
638 clade. Hence, this appears to be an interesting example of convergent evolution
639 (convergent at the morphological level, but not necessarily at the molecular level).
640 The inferred tree associated with the Convergent Class displays much more
641 evolution in the electric rather than the non-electric fish lineages (Fig. 5). This is

indicative of either a relaxation of purifying selection pressure, an introduction of positive selection pressure or a combination of both. The notable exception is the Brown Ghost Knifefish, which appears relatively conserved. The Brown Ghost Knifefish is unique amongst the electric fish in the dataset, in that its electric organ has evolved from neural rather than muscle tissue. Consequently, in the Brown Ghost Knifefish the *Na_v1.4a* gene is still expressed in muscle, just as it is in the non-electric fish. The distinction in terminal branch length between the Brown Ghost Knifefish and the other electric fishes offers compelling evidence that the GHOST model has identified a subtle component of the historical signal related to the convergent evolution of *Na_v1.4a*, as opposed to returning a somewhat arbitrary combination of numerical parameters that happen to maximize the likelihood function. To further verify that this conclusion was justified, we examined the trees inferred under the GTR+FO*H5 and GTR+FO*H6 models. If a convergent evolution signal is indeed present in the alignment then it should also be revealed under these models. Supplementary Figures S20 and S21 show the trees inferred by the five and six class model respectively. The third class in each Figure appears to capture a similar signal to that captured by the Convergent Class of the four class model. The ability of the GHOST model to isolate such a small component of the signal (the inferred weight of the convergent class being 0.13, the smallest of the 4

661 classes) is encouraging. Furthermore, we can hypothesize that the sites belonging
662 with high probability to the Convergent Class are likely to have been influential in
663 the functional development of the electric organ.

664 *Soft classification of sites to classes.*— The soft classification of sites to classes
665 facilitates the prospective identification of functionally important sites in an
666 alignment. Zakon et al. (2006) report several amino-acid sites from the dataset that
667 are influential in the inactivation of the sodium channel, a process critical to
668 electric organ pulse duration. Figure 6a shows that these sites generally have a
669 higher than average probability of belonging to the convergent class in at least one
670 codon position. For example, at position 647, an otherwise conserved proline
671 (codon CCN) is replaced by a valine (GTN) in the Pintailed Knifefish and a
672 cysteine (TGY) in the Electric Eel. Unique substitutions at codon positions 1 and
673 2 are necessary for both of these amino-acid replacements, and we find these two
674 sites have a very high probability of belonging to the convergent class. With this
675 result in mind, for each amino acid we summed the probability of codon positions 1
676 and 2 belonging to the Convergent Class. Figure 6b shows the results for the eight
677 amino-acid sites with the highest score. Comparing the magnitude of these bars
678 with those of the amino-acid sites in Figure 6a (which are known to be functionally
679 important), one is led to suspect that these amino acids might also be critical to

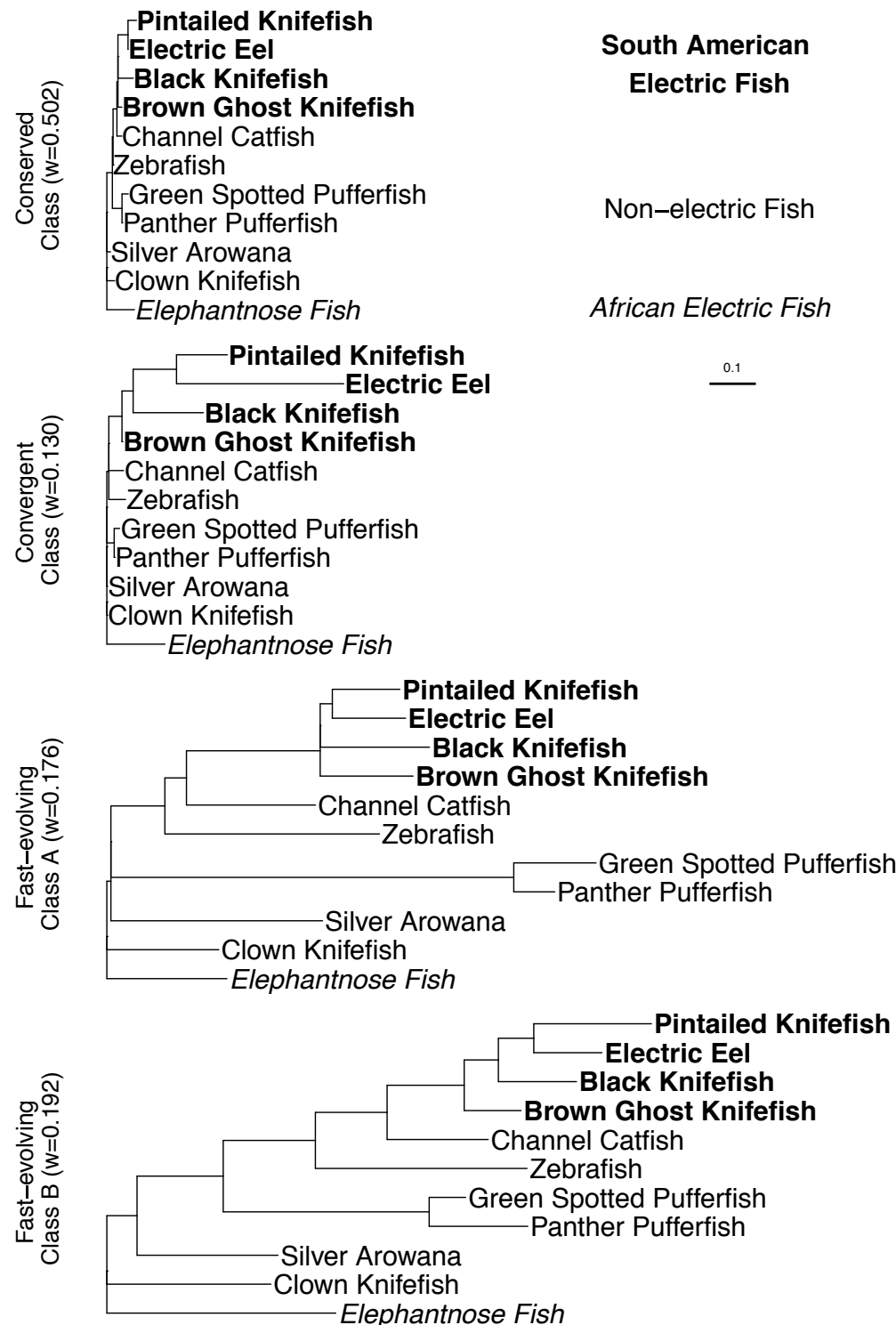


Figure 4: The four trees inferred under the General Time Reversible, four-class mixture model (GTR+FO*H4) for the electric fish data. The classes are displayed in order of increasing tree size, as determined by the sum of the branch lengths. We refer to this as the total tree length (TTL): $TTL_{Cons} = 0.23$, $TTL_{Conv} = 0.99$, $TTL_{FEA} = 4.06$ and $TTL_{FEB} = 4.18$.

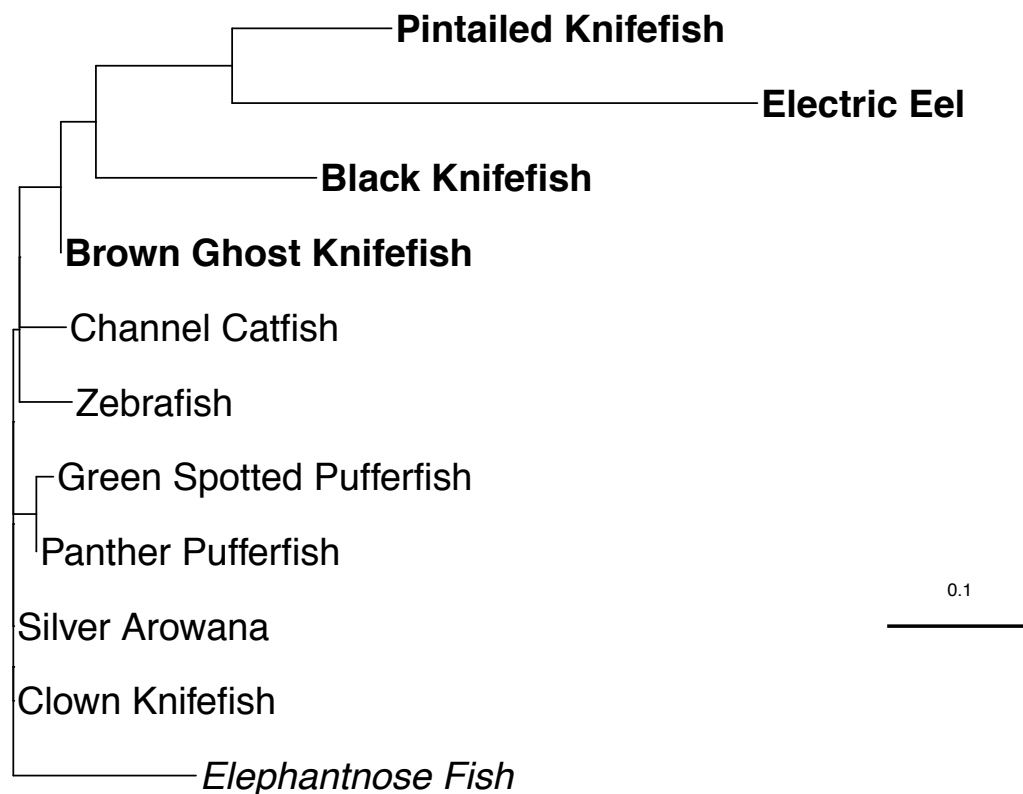


Figure 5: The convergent class inferred by ML-GTR+FO*H4. The 11 fish species comprised four South American electric fish (bold), one African electric fish (italics), and six non-electric fish (normal font) from various locations. The tree for this class shows that in comparison to the electric fish, the non-electric species are relatively conserved.

the operation of the sodium channel gene. Given that there are many other sites in the alignment with a high probability of belonging to the convergent class, one can envisage the GHOST model helping to identify sites of potential functional importance in an alignment, thereby focusing the experimental work of biologists.

In addition to providing insight on an individual site basis, the soft classification can also help to inform us about the nature of the classes themselves. Generally speaking, the branch lengths inferred by IQ-TREE can be interpreted as the expected number of substitutions per site. Therefore, summing the weighted TTLs for each of the inferred classes tells us that we expect 1.766 substitutions per site under the inferred model. Table 2 reports the contributions to this figure, stratified by codon position and class. If class membership and codon position were independent attributes of each site, then we should expect the contribution of each codon position to be approximately one third for each class. This is not what we observe. Overall we can see that sites in CP1(23%) and CP2 (16%) contribute only 39% of the total of 1.766 substitutions per site. However, within the Conserved and Convergent Classes, sites in CP1 and CP2 are responsible for 90% and 76% of their contribution respectively. This would suggest that a comparatively larger proportion of the substitutions attributed to these classes are non-synonymous:

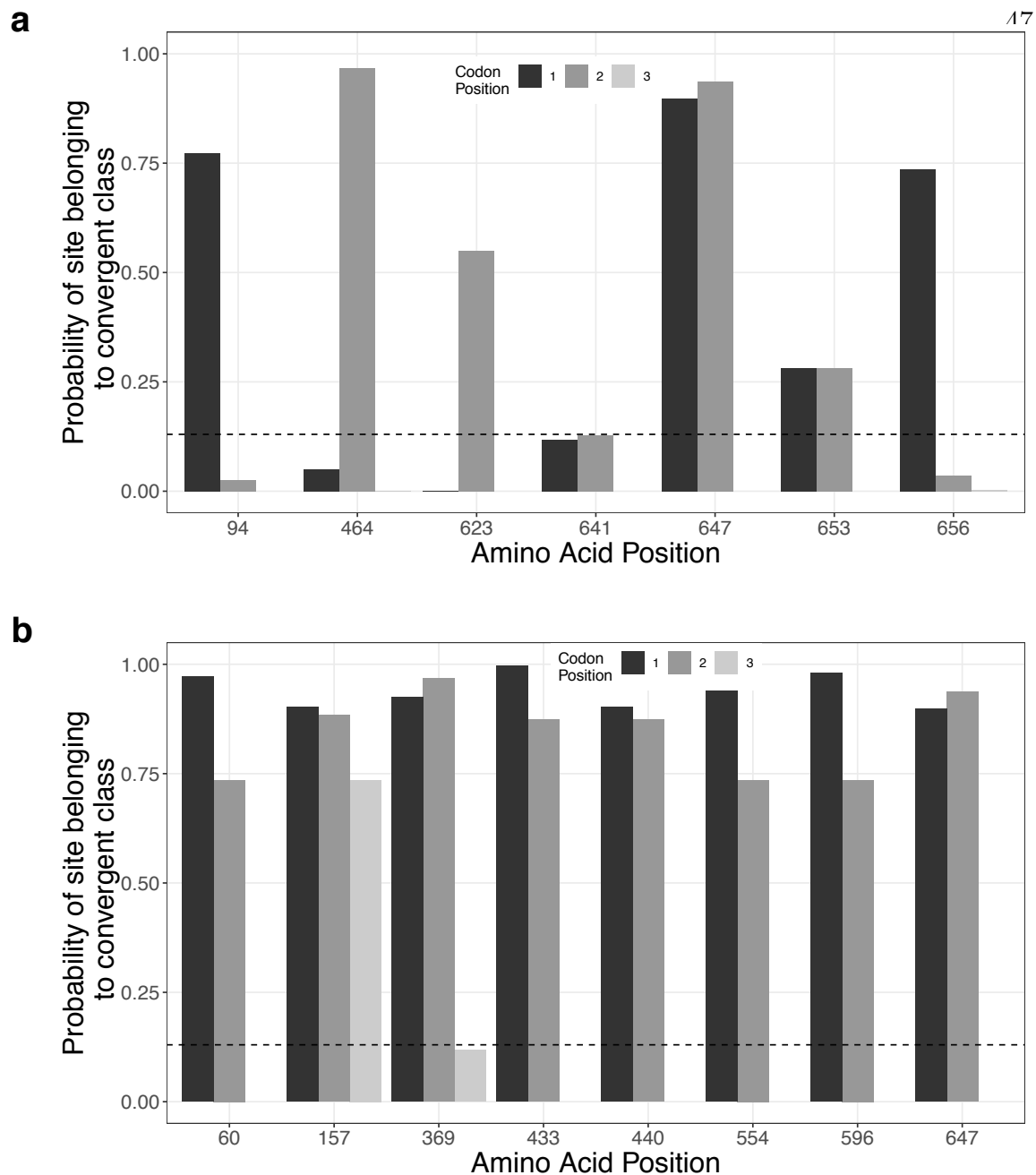


Figure 6: Probability of sites belonging to the convergent class by codon position. (a) The amino-acid positions selected correspond with those identified by Zakon et al. (2006) as being functionally important to the inactivation of the Na^+ channel gene. The horizontal dotted line at 0.13 represents 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 2: Expected number of substitutions per site (bold) for the electric fish alignment, 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.

699 resulting in amino-acid replacements that influence the fitness of the organism. We
700 can therefore conclude that even though the Conserved and Convergent Classes are
701 smallest (as determined by substitutions per site), they appear to be the primary
702 catalyst of evolution via natural selection within $Na_v 1.4a$ amongst these
703 species.

704 *Comparison to the Partition Model.*— It is apparent upon examination of the trees

in Supplementary Figure S19 that the evidence of convergent evolution highlighted by the GHOST model (Fig. 5) has not been recovered by the codon-based partition model. None of the three trees in Supplementary Figure S19 have the distinctive pattern, whereby the majority of the total tree length is associated with the electric fish species (with the exception of the Brown Ghost Knifefish). The reason that the partition model was unable to recover this signal has to do with the relative contribution of sites from each CP to the Convergent Class. Table 2 indicates the extent to which the substitutions associated with the Convergent Class are attributable to CP1 sites (40%), CP2 sites (36%) and CP3 sites (24%). The partition model constrains the analysis, such that sites in different CPs are modeled independent of each other. It is impossible for a model constrained in such a way to effectively recover the convergent evolution signal because the signal is distributed across all three partitions. The decision to partition the data based on codon position may make sense superficially, but in doing so the analysis is constrained and the results are compromised. We no longer have the ability to uncover the evolutionary stories concealed within the data. We can only hope to obtain those stories that happen not to conflict with the assumptions and constraints that have been placed on the analysis *a priori*. Minimizing these assumptions and constraints where possible, while computationally expensive, is necessary in order to illuminate

the evolutionary history without distorting it in the process.

On the Identifiability of the GHOST Model

An ongoing concern regarding parameter-rich mixture models has been whether or not they are identifiable. There are several examples of theoretically non-identifiable mixture models in the literature (Matsen and Steel, 2007; Štefankovič and Vigoda, 2007b). These examples have inspired much theoretical work on the identifiability or otherwise of different types of phylogenetic mixture models (Allman and Rhodes, 2006; Štefankovič and Vigoda, 2007a; Allman et al., 2008; Allman and Rhodes, 2008; Steel, 2010; Allman et al., 2011). Of particular interest to the current study, Allman et al. (2011) showed that for a single topology, four taxa, two-class mixture under the JC model, only the tree topology is identifiable but not the branch lengths. This provides a theoretical justification for the procedure carried out by K&T (and replicated here), measuring performance of the models based only on recovery of the topology and paying no attention to recovery of branch length parameters. With regard to the identifiability of the GHOST model more generally, we rely on a result from Rhodes and Sullivan (2012). They established an upper bound on the number of classes, m , for which tree topology, branch lengths and model parameters are identifiable, as a function

of the number of character states, κ , and the number of taxa, n :

$$m < \kappa^{\lceil \frac{n}{4} \rceil - 1}$$

For the simulations we carry out in the current study, with 12 taxa and four character states, the model is identifiable up to a maximum of 16 classes. For 32 taxa and four character states, the model is identifiable up to a maximum of 16,384 classes. In the case of the electric fish dataset, with four character states and only 11 taxa, the model is identifiable up to 16 classes. However, there is a technical caveat. The result is shown based on assuming a general Markov model across the tree. There are specific choices of parameters that can result in non-identifiability, but these are of little concern in practical data analysis. Problems arise only when the parameters selected collapse the parameter space to some lower dimension. For example, we could fit the GTR model but if we chose parameters such that all base frequencies were equal and all substitution rates were equal then we are in fact using a JC model, and identifiability may be compromised. However, these technical examples of non-identifiability are not relevant in practice, as in the absence of any constraints there is no reasonable chance of inferring parameters that collapse the parameter space in such a way.

CONCLUSION

Heterotachy has been somewhat of an Achilles heel for ML since K&T published their study. The ML implementation of the GHOST model in IQ-TREE represents a positive advance for ML-based phylogenetic inference. Through minimization of model assumptions, the GHOST model offers significant advantages and flexibility to infer heterotachous evolutionary processes, illuminating historical signals that might otherwise remain hidden. Owing to the diversity of selective pressures acting on different genes, the GHOST model seems well suited to the analysis of phylogenomic datasets (albeit with the limitation of being constrained to a single tree topology), commonly used to address deep phylogenetic questions. Forthcoming empirical studies will further compare the performance of the GHOST model to currently popular phylogenomic analysis tools, such as partition and CAT models. In addition, further simulation studies will help to better establish practical limitations to the use of the GHOST model, in terms of number of taxa, number of classes, sequence length and computation time. Many opportunities for refining or extending the GHOST model also present themselves. It could be used in conjunction with partition models, to account for heterotachy within partitions; if the data suggests correlation of some branch lengths across classes, then these could be linked to decrease the parameter space; it could form the basis of a test

777 for heterotachy itself, by comparing results obtained under the GHOST model to
 778 those obtained under a discrete Γ or PDF rate model. One can also envisage many
 779 other potential applications for the GHOST model. It may be insightful when
 780 applied to datasets for which the topology is poorly supported or disputed. It could
 781 also provide more accurate parameter estimates, leading to sounder divergence date
 782 estimation. The model provides intuitive, biologically meaningful visualizations of
 783 the different evolutionary pressures that act on a group of taxa. Structural
 784 biologists may find it useful for highlighting functionally important areas within
 785 proteins. We have demonstrated its use as a method for identifying changes in
 786 selection pressure, as well as bringing to light evidence of convergent evolution.
 787 Similarly, one can envisage the GHOST model illuminating the subtle evolutionary
 788 relationships between hosts and parasites, disease and immune cells, or the
 789 countless evolutionary arms races that are observed throughout the natural
 790 world.

SUPPLEMENTARY MATERIAL

Supplementary material, including further simulation results, figures and IQ-TREE command line instructions can be found in the Dryad data repository (doi TBA).

ACKNOWLEDGEMENTS

The authors would like to thank Elizabeth Allman, John Rhodes and Edward Susko for helpful discussions about the manuscript.

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

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

REFERENCES

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723.
- Allman, E. S., Ané, C., and Rhodes, J. A. (2008). Identifiability of a Markovian model of molecular evolution with gamma-distributed rates. *Advances in Applied Probability*, pages 229–249.
- Allman, E. S., Petrovic, S., Rhodes, J. A., and Sullivant, S. (2011). Identifiability of two-tree mixtures for group-based models. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, 8(3):710–722.
- Allman, E. S. and Rhodes, J. A. (2006). The identifiability of tree topology for phylogenetic models, including covarion and mixture models. *Journal of Computational Biology*, 13(5):1101–1113.
- Allman, E. S. and Rhodes, J. A. (2008). Identifying evolutionary trees and substitution parameters for the general Markov model with invariable sites. *Mathematical Biosciences*, 211(1):18–33.
- Baele, G., Raes, J., Van de Peer, Y., and Vansteelandt, S. (2006). An improved statistical method for detecting heterotachy in nucleotide sequences. *Molecular Biology and Evolution*, 23(7):1397–1405.

- 820 Burnham, K. P. and Anderson, D. R. (2003). *Model selection and multimodel*
821 *inference: a practical information-theoretic approach*. Springer Science &
822 Business Media.
- 823 Chiari, Y., Cahais, V., Galtier, N., and Delsuc, F. (2012). Phylogenomic analyses
824 support the position of turtles as the sister group of birds and crocodiles
825 (archosauria). *BMC Biology*, 10(1):65.
- 826 Crotty, S. M., Rohrlach, A. B., Ndunguru, J., and Boykin, L. M. (2018).
827 Characterising genetic diversity in cassava brown streak virus. *bioRxiv*, page
828 455303.
- 829 Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from
830 incomplete data via the EM algorithm. *Journal of the Royal Statistical*
831 *Society, Series B*, pages 1–38.
- 832 Dziak, J. J., Coffman, D. L., Lanza, S. T., Li, R., and Jermin, L. S. (2018).
833 Sensitivity and specificity of information criteria. *bioRxiv*, page 449751.
- 834 Felsenstein, J. (1981). Evolutionary trees from DNA sequences: a maximum
835 likelihood approach. *Journal of Molecular Evolution*, 17(6):368–376.
- 836 Fitch, W. M. and Margoliash, E. (1967). A method for estimating the number of

- invariant amino acid coding positions in a gene using cytochrome *c* as a model case. *Biochemical Genetics*, 1(1):65–71.
- Fitch, W. M. and Markowitz, E. (1970). An improved method for determining codon variability in a gene and its application to the rate of fixation of mutations in evolution. *Biochemical genetics*, 4(5):579–593.
- Fletcher, R. (2013). *Practical methods of optimization*. John Wiley & Sons.
- Foster, P. G. (2004). Modeling compositional heterogeneity. *Systematic Biology*, 53(3):485–495.
- Gadagkar, S. R. and Kumar, S. (2005). Maximum likelihood outperforms maximum parsimony even when evolutionary rates are heterotachous. *Molecular Biology and Evolution*, 22(11):2139–2141.
- Galtier, N. (2001). Maximum-likelihood phylogenetic analysis under a covarion-like model. *Molecular Biology and Evolution*, 18(5):866–873.
- Holmquist, R., Goodman, M., Conroy, T., and Czelusniak, J. (1983). The spatial distribution of fixed mutations within genes coding for proteins. *Journal of Molecular Evolution*, 19(6):437–448.
- Huelsenbeck, J. P. (2002). Testing a covariotide model of dna substitution. *Molecular Biology and Evolution*, 19(5):698–707.

- 855 Jayaswal, V., Wong, T. K., Robinson, J., Poladian, L., and Jermini, L. S. (2014).
856 Mixture models of nucleotide sequence evolution that account for
857 heterogeneity in the substitution process across sites and across lineages.
858 *Systematic Biology*, 63(5):726–742.
- 859 Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K., von Haeseler, A., and Jermini,
860 L. S. (2017). Modelfinder: fast model selection for accurate phylogenetic
861 estimates. *Nature Methods*, 14(6):587–589.
- 862 Kolaczkowski, B. and Thornton, J. W. (2004). Performance of maximum
863 parsimony and likelihood phylogenetics when evolution is heterogeneous.
864 *Nature*, 431(7011):980–984.
- 865 Kuhner, M. K. and Felsenstein, J. (1994). A simulation comparison of phylogeny
866 algorithms under equal and unequal evolutionary rates. *Molecular Biology and*
867 *Evolution*, 11(3):459–468.
- 868 Lanfear, R., Calcott, B., Ho, S. Y., and Guindon, S. (2012). PartitionFinder:
869 combined selection of partitioning schemes and substitution models for
870 phylogenetic analyses. *Molecular Biology and Evolution*, 29(6):1695–1701.
- 871 Lartillot, N. and Philippe, H. (2004). A Bayesian mixture model for across-site
872 heterogeneities in the amino-acid replacement process. *Molecular Biology and*

- 873 *Evolution*, 21(6):1095–1109.
- 874 Lopez, P., Casane, D., and Philippe, H. (2002). Heterotachy, an important process
875 of protein evolution. *Molecular Biology and Evolution*, 19(1):1–7.
- 876 Matsen, F. A. and Steel, M. (2007). Phylogenetic mixtures on a single tree can
877 mimic a tree of another topology. *Systematic Biology*, 56(5):767–775.
- 878 Meade, A. and Pagel, M. (2008). A phylogenetic mixture model for heterotachy. In
879 Pontarotti, P., editor, *Evolutionary Biology from Concept to Application*, pages
880 29–41. Springer.
- 881 Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., and Minh, B. Q. (2015).
882 IQ-TREE: a fast and effective stochastic algorithm for estimating
883 maximum-likelihood phylogenies. *Molecular Biology and Evolution*,
884 32(1):268–274.
- 885 Pagel, M. and Meade, A. (2004). A phylogenetic mixture model for detecting
886 pattern-heterogeneity in gene sequence or character-state data. *Systematic
887 Biology*, 53(4):571–581.
- 888 Pagel, M. and Meade, A. (2005). Mixture models in phylogenetic inference. In
889 Gascuel, O., editor, *Mathematics of Evolution and Phylogeny*, pages 121–142.
890 Oxford University Press Oxford, United Kingdom.

- Philippe, H. and Lopez, P. (2001). On the conservation of protein sequences in evolution. *Trends in Biochemical Sciences*, 26(7):414–416.
- Philippe, H., Zhou, Y., Brinkmann, H., Rodrigue, N., and Delsuc, F. (2005). Heterotachy and long-branch attraction in phylogenetics. *BMC Evolutionary Biology*, 5(1):50.
- Posada, D. and Buckley, T. R. (2004). Model selection and model averaging in phylogenetics: advantages of Akaike Information Criterion and Bayesian approaches over likelihood ratio tests. *Systematic Biology*, 53(5):793–808.
- Rambaut, A. and Grassly, N. C. (1997). Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Computer Applications in the Biosciences: CABIOS*, 13(3):235–238.
- Rhodes, J. A. and Sullivant, S. (2012). Identifiability of large phylogenetic mixture models. *Bulletin of Mathematical Biology*, 74(1):212–231.
- Robinson, D. F. and Foulds, L. R. (1981). Comparison of phylogenetic trees. *Mathematical Biosciences*, 53(1):131–147.
- Schwarz, G. et al. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464.

- 908 Spencer, M., Susko, E., and Roger, A. J. (2005). Likelihood, parsimony, and
909 heterogeneous evolution. *Molecular Biology and Evolution*, 22(5):1161–1164.
- 910 Steel, M. (2005). Should phylogenetic models be trying to fit an elephant? *Trends*
911 *in Genetics*, 21(6):307–309.
- 912 Steel, M. (2010). Can we avoid “SIN” in the house of “No Common Mechanism”?
913 *Systematic Biology*, 60(1):96–109.
- 914 Štefankovič, D. and Vigoda, E. (2007a). Phylogeny of mixture models: Robustness
915 of maximum likelihood and non-identifiable distributions. *Journal of*
916 *Computational Biology*, 14(2):156–189.
- 917 Štefankovič, D. and Vigoda, E. (2007b). Pitfalls of heterogeneous processes for
918 phylogenetic reconstruction. *Systematic Biology*, 56(1):113–124.
- 919 Tuffley, C. and Steel, M. (1998). Modeling the covarion hypothesis of nucleotide
920 substitution. *Mathematical Biosciences*, 147(1):63–91.
- 921 Wang, H.-C., Li, K., Susko, E., and Roger, A. J. (2008). A class frequency mixture
922 model that adjusts for site-specific amino acid frequencies and improves
923 inference of protein phylogeny. *BMC Evolutionary Biology*, 8(1):331.
- 924 Wang, H.-C., Spencer, M., Susko, E., and Roger, A. J. (2007). Testing for

- 925 covarion-like evolution in protein sequences. *Molecular Biology and Evolution*,
926 24(1):294–305.
- 927 Whelan, N. V. and Halanych, K. M. (2017). Who let the CAT out of the bag?
928 Accurately dealing with substitutional heterogeneity in phylogenomic analyses.
929 *Systematic Biology*, 66(2):232–255.
- 930 Wu, J. and Susko, E. (2009). General heterotachy and distance method
931 adjustments. *Molecular Biology and Evolution*, 26(12):2689–2697.
- 932 Wu, J. and Susko, E. (2011). A test for heterotachy using multiple pairs of
933 sequences. *Molecular Biology and Evolution*, 28(5):1661–1673.
- 934 Yan, M., Moore, M. J., Meng, A., Yao, X., and Wang, H. (2017). The first
935 complete plastome sequence of the basal asterid family styracaceae (Ericales)
936 reveals a large inversion. *Plant Systematics and Evolution*, 303(1):61–70.
- 937 Yang, Z. (1994). Maximum likelihood phylogenetic estimation from DNA sequences
938 with variable rates over sites: approximate methods. *Journal of Molecular*
939 *Evolution*, 39(3):306–314.
- 940 Zakon, H. H., Lu, Y., Zwickl, D. J., and Hillis, D. M. (2006). Sodium channel genes
941 and the evolution of diversity in communication signals of electric fishes:

- 942 convergent molecular evolution. *Proceedings of the National Academy of*
943 *Sciences of the United States of America*, 103(10):3675–3680.
- 944 Zhou, Y., Brinkmann, H., Rodrigue, N., Lartillot, N., and Philippe, H. (2010). A
945 Dirichlet process covarion mixture model and its assessments using posterior
946 predictive discrepancy tests. *Molecular Biology and Evolution*, 27(2):371–384.
- 947 Zhou, Y., Rodrigue, N., Lartillot, N., and Philippe, H. (2007). Evaluation of the
948 models handling heterotachy in phylogenetic inference. *BMC Evolutionary*
949 *Biology*, 7(1):206.