Nucleotide Substitutions during Speciation may Explain Substitution Rate Variation

Janzen, Thijs; Bokma, Folmer; Etienne, Rampal S

[Link to publication in University of Groningen/UMCG research database](#)

The data/supplementary material associated with this paper are available for review via

Dryad. The following is a temporary direct download link. Please copy and paste it directly

into a web browser to download the data files to your computer (unfortunately this may

not work as a link to click on)


https://datadryad.org/stash/share/5cEzwaEQalH1eIOl0xDZiOZjfEvuAEm1n1jpt4Ijw5U

1    Nucleotide substitutions during speciation may explain substitution rate variation

2

3    Thijs Janzen[1*], Folmer Bokma[2], Rampal S. Etienne[1]

4

5    [1] Groningen Institute for Evolutionary Life Sciences, University of Groningen, Box

6    11103, 9700 CC Groningen, The Netherlands

7    [2] Center for Ecological and Evolutionary Synthesis (CEES) Department of
8    BioSciences, University of Oslo, PO Box 1066, Blindern, 0316 Oslo, Norway
9

10

11    * corresponding author: t.janzen@rug.nl

12

13    RUNNING TITLE: Nucleotide substitutions during speciation

15    ABSTRACT

16    Although molecular mechanisms associated with the generation of mutations are

17    highly conserved across taxa, there is widespread variation in mutation rates

18    between evolutionary lineages. When phylogenies are reconstructed based on

19    nucleotide sequences, such variation is typically accounted for by the assumption

20    of a relaxed molecular clock, which is a statistical distribution of mutation rates

21    without much underlying biological mechanism. Here, we propose that variation

22    in accumulated mutations may be partly explained by an elevated mutation rate

23    during speciation. Using simulations, we show how shifting mutations from

24    branches to speciation events impacts inference of branching times in

25    phylogenetic reconstruction. Furthermore, the resulting nucleotide alignments

26    are better described by a relaxed than by a strict molecular clock. Thus, elevated

27    mutation rates during speciation potentially explain part of the variation in

28    substitution rates that is observed across the tree of life.

29

30    Keywords: molecular clock, speciation, phylogenetic reconstruction, substitution

31    rate variation

32

33

34

35

37     INTRODUCTION

38

39          Phenotypic diversification occurs at a higher rate in some clades than in

40     others (Simpson 1945; van Valen 1985; Ricklefs 2006; Rabosky et al. 2007;

41     Jansson and Davies 2008) and similarly, there is substantial variation across

42     evolutionary lineages in the rate of molecular evolution (King and Wilson 1975),

43     such as that of nucleotide sequences (Nabholz et al. 2008; Bromham 2011; Dowle

44     et al. 2013; Sung et al. 2016).  As a consequence, studies attempting to reconstruct

45     the phylogeny of a clade often find that the sequence data do not support the

46     assumption of a strict molecular clock, i.e. constant substitution rates across

47     lineages. For such cases, phylogenetic inference software allows one to use a

48     relaxed molecular clock (Drummond et al. 2006; Lepage et al. 2007), which

49     assumes that the substitution rate varies between lineages according to a

50     statistical distribution such as a gamma or lognormal distribution. However, the

51     relaxed molecular clock thus introduces at least one additional degree of freedom,

52     namely the variance of the distribution of substitution rates (although some argue

53     that an uncorrelated relaxed clock in effect adds one additional degree of freedom

54     *per branch* (Dornburg et al. 2012; Bromham 2019; Zhang and Drummond 2020;

55     Douglas et al. 2021). Moreover, the relaxed clock is a rather ad-hoc solution with

56     little underlying biological reasoning (but see Lartillot and Poujol 2014; Lartillot

57     et al. 2016; Saclier et al. 2018).

58

59          A first formal test to detect the impact of speciation on sequence evolution

60     was formulated by Avise and Ayala (Avise and Ayala 1975, 1976), who

61     distinguished gradual evolution from "punctuated equilibria" by comparing

62    sequence evolution in species-rich and species-poor clades. Whereas Avise and

63    Ayala found no evidence for increased sequence evolution in species-rich clades,

64    others did, in tetrapods (Mindell et al. 1989, 1990), sauropsids (Eo and DeWoody

65    2010) and angiosperms (Duchene and Bromham 2013; Bromham et al. 2015).

66    Furthermore, substitution rates have been found to be positively associated with

67    diversification rates (Fontanillas et al. 2007; Eo and DeWoody 2010; Lanfear et al.

68    2010a, 2010b; Ezard et al. 2013, but see Goldie et al. 2011).

69           Several biological processes acting at speciation could lead to accelerated

70    sequence evolution, including, but not limited to founder effects, bottlenecks,

71    inbreeding, hybridization, selection for an increased mutation rate, divergent

72    selection and local adaptation (Venditti and Pagel 2010). Here, we explore how

73    such processes driving sequence evolution during speciation events might affect

74    phylogenetic reconstruction; we posit that differences in apparent substitution

75    rates between lineages are due to processes acting exclusively or predominantly

76    during speciation. Due to (effectively) random extinction of lineages, different

77    branches of a reconstructed phylogeny will differ in how often they experienced

78    such short episodes of accelerated substitution rates, resulting in differences in

79    apparent substitutions rates along these branches. Our approach is two-fold: first,

80    we explore whether inclusion of substitutions during speciation affects

81    phylogenetic inference, and, if so, which aspects of the inferred phylogenetic tree

82    are affected. Second, we explore whether substitutions during speciation can

83    explain variation in estimated substitution rates.

85    METHODS

86         We propose a model where substitutions accumulate not only along the

87    branches of a phylogeny, but also at speciation events, including not only the

88    internal nodes of the phylogeny but also those pruned from the phylogeny by

89    extinction. We first make the standard assumption that gradual sequence

90    evolution along a phylogenetic branch can be modeled as a time-homogeneous

91    Markov process with substitution matrix:

92    $$\mathbf{Q} = \begin{bmatrix} -\mu_{AT} - \mu_{AC} - \mu_{AG} & \mu_{AT} & \mu_{AC} & \mu_{AG} \\ \mu_{TA} & -\mu_{TA} - \mu_{TC} - \mu_{TG} & \mu_{TC} & \mu_{TG} \\ \mu_{CA} & \mu_{CT} & -\mu_{CA} - \mu_{CT} - \mu_{CG} & \mu_{CG} \\ \mu_{GA} & \mu_{GT} & \mu_{GC} & -\mu_{GA} - \mu_{GT} - \mu_{GC} \end{bmatrix}$$

93    where $\mu_{ij}$ denotes the mutation rate from nucleotide $i$ to nucleotide $j$. The

94    transition probabilities of nucleotide substitutions after time $t$ of gradual

95    sequence evolution are then given by the matrix

96                               $$\mathbf{P_a}(t) = \exp(\mathbf{Q}t)$$

97    where the subscript **a** indicates anagenetic change, i.e. gradual accumulation of

98    substitutions over time. This matrix can be multiplied with an initial probability

99    vector at time $t = 0$ to yield the probabilities for each of the four nucleotides at

100   time $t$.

101   In addition to gradual sequence evolution over time, we assume that sequences

102   may change rapidly during speciation. We can thus assume another matrix $\mathbf{P_c}$

103   (subscript c for "cladogenetic") that describes the nucleotide transition

104   probabilities during a single speciation event.

105

106    The processes that may accelerate sequence evolution during speciation, such as

107    founder effects, bottlenecks, inbreeding, hybridization, and adaptation to novel

108    environments, may well result in different kinds of substitutions than those that

109    take place over time in established species. However, for mathematical

110    convenience we will here assume that we can write:

111 $$\mathbf{P}_c = \exp\left(\mathbf{Q}\tau\right)$$

112    where $\tau$ is a parameter that measures the effect of substitutions during speciation.

113    In other words, we assume that nucleotide sequence evolution is only accelerated

114    during speciation events, but not qualitatively altered: the $\mu_{ij}$ used in $\mathbf{P}_c$ must be

115    identical to those used in $\mathbf{P_a}$. The acceleration is then measured by parameter $\tau$: a

116    single speciation event causes as much sequence evolution as $\tau$ years of gradual

117    evolution over time within each lineage. Thus, larger values of $\tau$ correspond to a

118    larger experienced effect at the nodes, similar to sequence evolution along a

119    branch of length $\tau$. For $\tau = 0$, $\mathbf{P}_c$ becomes the identity matrix, and our model

120    reduces to the standard model of sequence evolution that only assumes

121    substitutions along phylogenetic branches. Important to note here is that both

122    daughter lineages resulting from a speciation event experience substitutions

123    independently (see the Supplement for a model where substitutions in both

124    daughter lineages are dependent on each other). Furthermore, we emphasize that

125    we assume the speciation process to happen in a similar fashion across a tree,

126    assuming an identical $\tau$ for all nodes in the tree. Later versions of the model could

127    potentially relax this assumption, provided independent information about

128    speciation dynamics.

129

130    For simplicity we assume in our simulations that sequence evolution can be

131    modeled as a Jukes-Cantor process (Jukes and Cantor 1969), for which **Q** is given

132    by:

133

$$\mathbf{Q} = \begin{matrix} -\dfrac{3\mu}{4} & \dfrac{\mu}{4} & \dfrac{\mu}{4} & \dfrac{\mu}{4} \\[2ex] \dfrac{\mu}{4} & -\dfrac{3\mu}{4} & \dfrac{\mu}{4} & \dfrac{\mu}{4} \\[2ex] \dfrac{\mu}{4} & \dfrac{\mu}{4} & -\dfrac{3\mu}{4} & \dfrac{\mu}{4} \\[2ex] \dfrac{\mu}{4} & \dfrac{\mu}{4} & \dfrac{\mu}{4} & -\dfrac{3\mu}{4} \end{matrix}$$

134    Several existing software packages (e.g. the R package phangorn, Schliep 2011; the

135    python module pyvolve, Spielman & Wilke 2015; the R package phylosim, Sipos *et*

136    *al.* 2011),  provide algorithms to simulate sequence evolution along the branches

137    of the phylogeny, given a rooted phylogeny and a root sequence (e.g. some

138    arbitrary sequence assumed to represent the ancestral sequence), by applying the

139    transition matrix sequentially along the phylogenetic tree. Here, we extend this

140    methodology to also include substitutions accumulated at the nodes of the

141    phylogeny. We implemented this in the R package 'nodeSub', available via

142    https://CRAN.R-project.org/package=nodeSub.

144     *Testing the impact of node substitution models using simulations*

145        To identify the amount of error in phylogenetic inference caused by

146     assuming a (relaxed) molecular clock when substitutions actually arise (in part)

147     during speciation, we simulated sequence evolution on known trees and then

148     reconstructed the phylogeny from the simulated sequences, assuming strict and

149     relaxed molecular clocks. We simulated sequence evolution with the node

150     substitution model introduced above, with various degrees of sequence

151     accumulation at the nodes of the tree ($\tau$), and with various extinction rates. We

152     then compared the resulting trees with the original true tree using a number of

153     statistics: the gamma statistic (Pybus and Harvey 2000), the beta statistic (Aldous

154     2001), the mean branch length (Faith 1992; Clarke and Warwick 2001), crown

155     age, the normalized Lineages Through Time (nLTT) statistic (Janzen et al. 2015)

156     and the Jenson-Shannon distance metric comparing the Laplacian spectrum

157     (Lewitus and Morlon 2016).

158        Phylogenetic reconstruction was performed with BEAST2 (Bouckaert et al.

159     2019) using the R package *babette* (Bilderbeek and Etienne 2018). BEAST2

160     inference was performed using default priors (see the Supplementary information

161     for an example XML file), with a birth-death prior as tree prior (or a Yule prior if

162     the extinction rate was zero), the Jukes-Cantor nucleotide substitution model, and

163     a strict or relaxed clock model. The BEAST chain was run for 10 million steps,

164     whilst sampling a tree every 5000 steps. After completion, the first 10% of the

165     chain was discarded as burn-in.

166

168    *Assessing error in phylogenetic reconstruction: the twin tree*

169    Errors observed when comparing with the true tree include both errors

170    incurred by the node substitution model chosen, and errors accumulated in the

171    phylogenetic inference process even when the models used in inference are

172    identical to those generating the data (e.g. stochasticity in substitution

173    accumulation, stochasticity in phylogenetic tree creation). Furthermore,

174    additional effects arising during alignment simulation might affect our findings,

175    such as the impact of parameter values (sequence length, substitution rate, birth

176    rate, death rate), and of multiple substitutions at the same site (the node-density-

177    effect) as well as potential biases or interactions between summary statistics. To

178    correct for these effects, so as to isolate the error induced by using a node

179    substitution model from other sources of error, we inferred a phylogenetic tree

180    for a twin alignment (*sensu* Bilderbeek, Laudanno & Etienne 2020). This twin

181    alignment has exactly the same number of accumulated substitutions as the

182    original alignment. The total number of substitutions is tracked during simulation

183    of the substitution model, and not just the resulting number of variable sites in the

184    alignment. The twin alignment is based on the same true tree, but instead of using

185    a node substitution model to generate the alignment, it results from using either a

186    strict-clock or relaxed-clock substitution model. Using this twin alignment, we

187    performed phylogenetic reconstruction with BEAST2 as for the original

188    alignment, and estimated the same summary statistics for the posterior

189    distribution of trees. The error introduced by the node substitution model is then

190    the difference between the error of the node substitution posterior and the error

191    in the twin posterior. In summary, we use this twin approach as a control

192    treatment, in order to correct for all potential sources of additional error other

193    than that of our proposed substitution model.

194

195    *Obtaining a twin alignment*

196        We generated a twin alignment conditional on a phylogeny, a node

197    substitution model, and a mutation rate. Because an alignment generated using a

198    node substitution model (with $\tau > 0$) has accumulated substitutions at the nodes

199    in addition to those along the branches, the overall number of substitutions

200    accumulated is higher than for an alignment simulated using the same mutation

201    rate and a model with substitutions only on the branches. Thus, in order to

202    generate a *twin* alignment that contains the same amount of information

203    (substitutions) we increased the mutation rate. We did this by calculating the

204    estimated time spent at the nodes, relative to the time spent on the branches, and

205    using this as an estimate of the expected fraction of the number of substitutions

206    on the nodes, relative to the number of substitutions on the branches, assuming

207    that substitutions accumulate at the same rate on both branches and nodes. That

208    is, the mutation rate used in generating the twin alignment is calculated as:

209    $$\mu_{twin} = \mu \left( 1 + \frac{\tau(2N + H)}{\sum t_{branch}} \right) \ (1)$$

210    where $\mu$ is the mutation rate used in the node substitution model, $\tau$ is the time

211    spent on the node, $N$ is the number of internal nodes in the tree, $H$ is the number

212    of hidden nodes in the tree and $\sum t_{branch}$ is the total branch length of the tree. The

213    factor *2N* arises from the independent accumulation of substitutions during a

214    node substitution event for both daughter lineages.

215   During simulation of node substitution alignments, we kept track of the

216   substitutions accumulated at each node and branch, which allowed us to directly

217   measure the contribution of substitutions accumulated at the nodes (i.e.

218   $\tau(2N + H)$) relative to those accumulated at the branches (i.e.. $\sum t_{branch}$) This

219   provided us with an estimate of $\mu_{twin}$, and with an estimate of the total number of

220   substitutions arising during simulation of the alignment. We then used the

221   obtained estimate for $\mu_{twin}$ to generate _twin_ alignments, again tracking all

222   substitutions, until we obtained an alignment exactly matching the number of

223   accumulated substitutions of the alignment simulated with the node substitution

224   model.

225

226   *Node-density-effect*

227        The method we used to simulate substitutions along branches (and nodes)

228   ignores repeated mutations at the same site, which may lead to a node-density-

229   effect. Because the node-density-effect can mask the effect of node substitutions,

230   we made sure in two distinct ways that our results are not affected by this effect.

231   Firstly, by using a *twin* alignment, any resulting node-density-effects are mirrored

232   in the *twin* alignment as well, ensuring that any additional errors picked up do not

233   reflect errors induced by the node-density-effect. Secondly, we repeated our

234   analysis using a different simulation method that explicitly tracks repeated

235   mutations for the Jukes-Cantor model (see Supplementary Information for details

236   and results). We find that this more explicit simulation method yielded virtually

237   identical results to the more general approach described in the main text.

238

239   *Simulation settings*

240   We generated birth-death trees with varying degrees of extinction rate *d* in

241 [0, 0.1, 0.3, 0.5] and a single speciation rate of *b* = 1. Trees were simulated

242 conditional on 100 tips, using the function `sim.bd.taxa` from the R package

243 TreeSim (Stadler 2011). Across all settings, we simulated sequences of 10 kb, with

244 $\mu$ = 0.001.

245

246 *Varying the time spent on the nodes relative to the crown age*

247   We varied $\tau$ in [0, 0.01, 0.05, 0.1, 0.2, 0.4] times the crown age (e.g. when

248 the crown age of the simulated tree is 3MY, $\tau$ = [0, 0.03, 0.15, 0.3, 0.6, 1.2] MY).

249 Again, for each combination of $\tau$ and extinction (*d* = [0, 0.1, 0.3, 0.5]) we simulated

250 100 trees and for each tree we generated one node substitution alignment and one

251 *twin* alignment.

252

253 *The impact of tree balance*

254   In unbalanced trees, some terminal branches are connected by many more

255 past branching events to the root of the tree than are other terminal branches.

256 Hence, we expect that balance of a tree might have a substantial effect on the error

257 in phylogenetic inference: less balanced trees are expected to have higher error.

258 To test this, we compared fully balanced ($\beta$ = 10.0) with extremely unbalanced

259 "caterpillar" trees ($\beta$ = -2). We did so by simulating the branching times of a birth-

260 death tree, and assigning these to a fully balanced or fully unbalanced topology.

261 Thus, the only difference between the trees is the topology. Then, for both the

262 balanced and unbalanced tree a node substitution alignment was generated, with

263 the same number of total substitutions, and setting $\tau$ as a function of crown age.

264    As an extra check, we also generated a node substitution alignment for the original

265    birth-death tree from which the branching times were used. For all three

266    alignments we inferred a phylogenetic tree as in the other scenarios, and

267    compared the error in phylogenetic inference. For caterpillar trees with extreme

268    unbalance we were unable to calculate the Laplacian Spectrum, hence we omitted

269    the Laplacian Spectrum summary statistic in this analysis.

270

271    *Support for strict and relaxed clock models*

272            To test whether the alignment originating from a process with node

273    substitutions was better described by a relaxed than by a fixed clock model, we

274    repeated the analysis, but now inferred the marginal likelihood of the relaxed

275    clock and strict clock models using the "NS" package for BEAST2, which applies

276    Nested Sampling to obtain the marginal posterior likelihood for both models

277    (Russel et al. 2019). We used the function 'bbt_run' from the *babette* package

278    (Bilderbeek    and    Etienne    2018)    in    combination    with    the    function

279    'create_ns_mcmc' from the *beautier* package (Bilderbeek and Etienne 2018).

280    This performs a Nested Sampling MCMC run using BEAST2 (an example XML file

281    outlining the default settings used can be found in the Supplementary

282    information), which runs until convergence is detected. Then, we converted the

283    obtained marginal likelihoods to a relative weight for each model (by dividing

284    both marginal likelihoods by their sum), which allows for comparison of posterior

285    support for each model across parameter settings and trees.

286

287    *Empirical example*

288    As an illustration  of the impact of node substitutions on a real phylogeny (rather

289    than a simulated one), we applied our model to an empirical dataset which is

290    feasible under the assumption that there is no extinction (see below). The dataset

291    consists of sequence data (Ast 2001; Fitch et al. 2006) of 35 species of Australian

292    monitor lizards, of the family of *Varanidae*, which covers all known species of

293    *Varanidae* occurring in the Indo-Australian realm. For each species, mitochondrial

294    DNA was retrieved from GenBank, consisting of ND4, 16S and CO1 genes.

295    Sequences were aligned using the "—auto" setting for mafft (Katoh and Standley

296    2013), and concatenated for ease of use. Assuming a substitution rate of $3.35*10^{-9}$

297    per site (Eo and DeWoody 2010), we inferred a Maximum Likelihood tree from

298    the alignment, using the R package phangorn (Schliep 2011), assuming a

299    Generalised Time Reversible (GTR) model of substitution. This yielded a reference

300    tree, assuming no node substitutions.

301    Then, we made use of a new feature of phangorn (version 2.7.1.2, added upon our

302    request) which allows for the incorporation of node substitutions under the

303    assumption that there is no extinction, because then all nodes where node

304    substitutions occur are observable in the tree, and the branches connected to

305    these nodes can all be extended by a length of $\tau$. Thus, in this new version of

306    phangorn one can specify a value for $\tau$, and compute the tree likelihood (i.e. the

307    probability of the alignment given the tree and the substitution model

308    parameters) for this value. We explored the tree likelihood for values of $\tau$ ranging

309    from $10^{-4}$ to 1 MY, for the Jukes-Cantor (JC)  and GTR  substitution models.

310

312    RESULTS

313    *Summary statistics*

314         We compared summary statistics of trees inferred from alignments using

315    the node substitution model, with summary statistics of twin trees inferred from

316    alignments with identical information content, but generated without the node

317    substitution model (e.g. with only substitutions along the branches, and a fixed

318    clock rate). We find that summary statistics that are influenced by branching times

319    are affected (Figure 1, gamma, nLTT statistic, mean branch length and crown age).

320    For these summary statistics, we find an increased difference with increasing $\tau$.

321    The impact of extinction seems to be limited, as the error in these summary

322    statistics remains around the same level, regardless of the extinction rate used.

323
324

326    *The impact of tree balance*

327    Tree balance clearly influences the sensitivity of inference to node substitutions

328    (Figure 2). The inference error is larger for unbalanced trees, again only for the

329    gamma and the nLTT statistic. Fully balanced trees show slightly less error than

330    birth-death trees. Overall, all three types of trees show an increased error when

331    alignments are generated using the node substitution model. Errors are

332    particularly large for the beta statistic, but that is expected because it measures

333    topological features of the tree that we modified artificially.

335    *Support for strict and relaxed clock models*

336         We compared the relative support for each model, reflected by the relative

337    weight of the marginal likelihood. With an increasing amount of time spent at the

338    nodes $\tau$, the median weight of the relaxed clock model increases for the node

339    substitution alignment, with generally (across extinction rates) a higher weight

340    than the strict clock model for values of $\tau$ that are equal or larger than 0.1 times

341    the crown age (Figure 3). For the twin alignment, the strict clock model is

342    preferred, as expected, because this is the generating model.

343    For low values of $\tau$ (smaller than 0.1 times the crown age), we do not find any

344    effect of the balance of the tree on the marginal likelihood of the relaxed clock

345    model (Figure 4), in line with our finding above. However, for intermediate values

346    of $\tau$ (0.1 and 0.2), we find that unbalanced trees tend to have a higher marginal

347    weight for the relaxed clock model. For high values of $\tau$ (0.4), we find that the

348    marginal weight for the relaxed clock model is always higher, regardless of the

349    balance of the tree.

350

351    *Empirical example*

352    We first verified that extinction was low by fitting a birth-death model to the

353    Maximum Likelihood tree inferred without node substitutions. Here, we found an

354    estimate for *d/b* of 0 (95% CI: [-1.65, 0.24), and for *b – d* of 0.013 (95% CI: [0.0095,

355    0.01889], which together indicate that extinction is low indeed. This provides

356    justification for using the likelihood computations in the new version of phangorn

357    which assumes that the extinction rate is zero.

358    Next, we inferred $\tau$ and found a non-zero estimate for $\tau$ of 0.74MY when using the

359    JC model, and 2.53MY when using the GTR model (Figure 5, A & D). Comparing the

360   resulting trees for these ML estimates, we find that the crown age of the tree is

361   inferred to be much lower. Without node substitutions, the crown age is estimated

362   to be 48.22 MY for the JC model and 46.98 MY for the GTR model. When including

363   node substitutions, the crown age shifts to 33.9 MY for the JC model and 34.1MY

364   for the GTR model. Rescaling of the trees relative to the crown age (Figure 5 C &

365   E) shows that including node substitutions does not merely rescale all branching

366   points proportional to the newly inferred crown age, but that the relative

367   positions of the different branching points shift as well.

368

370    DISCUSSION

371        We have shown that an increased substitution rate during speciation events

372    potentially provides a mechanistic explanation of variation in substitution rates

373    across the branches of phylogenetic trees. Trees inferred from alignments

374    generated with this substitution model differ substantially from trees inferred

375    from alignments generated with a standard substitution model, especially

376    concerning branching times. Furthermore, we find that this new substitution

377    model can potentially explain widespread support for relaxed molecular clocks.

378        If sequence evolution mainly occurs during speciation, this would lead to a

379    correlation between species richness and substitution rate. However, this

380    correlation could also be an artifact of phylogenetic reconstruction known as the

381    node-density-effect (Fitch and Bruschi 1987; Fitch and Beintema 1990). The node-

382    density-effect reflects the inability to detect multiple mutations occurring at the

383    same site, thus causing an underestimate of the true branch length, especially for

384    longer branches where the probability of multiple mutations occurring at the

385    same site is higher. Because species-rich parts of phylogenies tend to have shorter

386    branches, sequence evolution in these species-rich parts is less underestimated

387    than in species-poor parts, causing a correlation between the number of observed

388    substitutions and species diversity. Pagel et al. tested for the impact of speciation

389    events, and of the node-density-effect in 122 phylogenies, spanning 4 taxa (Pagel

390    et al. 2006). Using previously demonstrated methodology to detect the node-

391    density effect (Webster et al. 2003; Venditti et al. 2006), they showed that in 57 of

392    the 122 examined phylogenies, they could detect a signature of increased

393    sequence evolution during speciation events.  However, this was the result of the

394    node-density effect in 22 out of these 57 phylogenies. Here, disentangling

395    sequence evolution during speciation from confounding factors such as the node-

396    density effect, but also stochasticity in tree simulation, stochasticity during

397    alignment simulation and error or bias in tree inference, has proven to be a non-

398    trivial endeavor. In order to assess the impact of node substitutions, we therefore

399    separated error due to assuming an alternative substitution model from the errors

400    introduced by the factors mentioned above. To do so we extended the twinning

401    approach (introduced by Bilderbeek, Laudanno and Etienne (2020)) to assess the

402    impact of choosing a different tree prior to explore the impact of a different

403    substitution model. The twinning approach succeeds by replicating the chosen

404    analysis pipeline, but using *control* data that have been generated using known

405    models and priors. The impact of the node substitution model then follows from

406    the difference between results obtained with the node substitution model and

407    results obtained with the *twin* (control) pipeline: errors are then due to model

408    misspecification, and not stochastic uncertainty produced by the analysis pipeline.

409    Our results show thus, that when we correct for the background effects of

410    (amongst other factors) the node-density effect, we observe strong effects of node

411    substitutions. However, we expect that for small values of $\tau$, the impact of node

412    substitutions might become comparable to the node-density effect, and

413    disentangling these sources of substitution rate variation might become difficult.

414        One might expect that a high extinction rate, by elevating numbers of hidden

415    nodes, would lead to a greater impact of node substitutions. It may therefore be

416    counterintuitive that in our simulation study we did not find such an effect of

417    higher rates of extinction. However, we conditioned our alignments on the same

418    total number of substitutions, to ensure that alignments with and without node

419    substitutions    contained    the    same    information    content.    Thus,    with    higher

420    extinction and hence more hidden nodes, relatively fewer substitutions occur on

421    the observed nodes. Because the number of hidden nodes is proportional to

422    branch length (Eq.2), the number of hidden nodes is interpreted as substitutions

423    on the branches. Potentially, this provides a way to distinguish between

424    phylogenetic models: although every constant-rate birth-death model has a

425    corresponding zero-extinction model with a time-varying speciation rate that

426    yields the same probability of the reconstructed tree (Nee et al. 1994; Louca and

427    Pennell 2020), the resulting alignments under the node substitution model will

428    not be similar. Because the birth-death tree includes extinction events,

429    substitution patterns will be different from those of the tree generated with the

430    time-varying speciation rate model.

431    Distinguishing phylogenetic models will become more feasible if some of the

432    simplifying assumptions made here are relaxed. The model we propose here takes

433    the simplest form, assuming a Jukes-Cantor (Jukes and Cantor 1969) substitution

434    matrix, identical substitution rates, identical substitution matrices between nodes

435    and branches, and constant birth-death rates over time. These assumptions were

436    made as a most basic starting point, but can be relaxed in future analyses, for

437    instance by introducing a different substitution matrix at the nodes, or by studying

438    the effect of node substitutions on trees that are generated by diversity-dependent

439    speciation rates (Etienne et al. 2012). By starting with the most tractable version

440    of the node substitution model, we have provided a first proof of concept of the

441    potential impact of node substitutions without overcomplicating matters.

442    Previous methods have applied rather ad-hoc corrections to account for

443    differences in substitution rates across different branches in the same phylogeny,

444    typically referred to as the 'relaxed clock' approach. These methods provide

445    satisfying statistical solutions to account for variation in substitution rates, but

446    refrain from providing biological explanations for this observed phenomenon. The

447    node substitution model we introduce here provides this explanation: branches

448    that have accumulated a number of 'hidden' branching events, e.g. speciation

449    events of species that have subsequently gone extinct, have a higher number of

450    accumulated substitution events during these 'hidden' speciation events. When

451    we compared the marginal likelihood of the relaxed clock model versus the strict

452    clock model for alignments generated with the node substitution model, we found

453    that marginal likelihoods for the relaxed clock model are much higher. This

454    indicates that our proposed process of accumulating substitutions during

455    speciation events can generate patterns in the alignment that are picked up by

456    phylogenetic methods as evidence for a relaxed clock model, without actually

457    using a relaxed clock model.

458          The notion of accelerated evolution during speciation events ties in

459    closely with the theory of punctuated equilibrium; where Eldredge and Gould

460    (1972) proposed that evolution perhaps is not a gradual process, but rather a

461    process with distinct bursts of phenotypic and morphological change. Their

462    theory was influenced by ideas like Lerner's "genetic homeostasis" (Lerner 1954),

463    which had earlier inspired Mayr (1954) to suggest that the formation of new

464    species involves "genetic revolutions". Our framework provides a step towards

465    being able to test this notion, where information on the estimated fraction of

466    substitutions accumulated at the nodes can directly inform us about whether the

467    majority of substitutions is accumulated over long periods of time in established

468    lineages (e.g. along branches), or during speciation (e.g. at the nodes).

469        To infer whether node substitutions really occur, we should fit the node

470    substitution model to empirical sequence alignments, and find a nonzero estimate

471    for $\tau$. However, the computation of the likelihood of our model (and estimation of

472    associated $\tau$ values), is non-trivial because it requires integration across the

473    enormous state space of complete trees (trees including extinct species). Manceau

474    *et al.* (2020) have taken a first step towards formulating such a likelihood. They

475    introduced an alternative solution for punctuated equilibrium-like patterns in

476    molecular evolution through the implementation of spikes of substitution, e.g.

477    moments in time at which there is an increased substitution rate. They let these

478    moments occur at speciation events, and also model the probability of such an

479    event happening at a speciation event (rather than assuming that they always

480    occur, as we did here). However, they have to assume both the topology and

481    branching times to be fixed. We have provided an alternative inference approach

482    that does not require topology or branching times to be fixed, but assumes

483    extinction to be zero. The absence of extinction greatly reduces computational

484    complexity, and allows us to use Maximum Likelihood to infer the most likely tree,

485    using the R package *phangorn* (Schliep 2011). We inferred a phylogenetic tree via

486    maximum likelihood for 35 species of *Varanidae* and recovered a non-zero

487    estimate for $\tau$. Furthermore, we found that the resulting tree was substantially

488    different from a tree with $\tau = 0$; not only were the crown age and branching times

489    drastically different, the relative position of the branching times was also affected.

490    As expected from our simulation results, topology of the tree was not affected.

491        In order to be able to infer the phylogenetic tree, we had to make several

492    restricting assumptions. Firstly, as stated above, we had to assume extinction to

493    be zero. This ignores any effects that hidden nodes might have. Yet, it seems

494  unlikely that in the 40 million years since the origination of the clade of *Varanidae*,

495  no extinctions took place. Secondly, we were limited to only using a strict clock

496  (other clocks are not yet incorporated in phangorn). Future work could explore

497  how incorporation of a relaxed clock in the maximum likelihood framework we

498  used impacts our findings, particularly whether using a relaxed clock could

499  mitigate some of the differences we recovered.

500  The present study aims to demonstrate that substitutions accumulated

501  during speciation might explain the prevalence of the relaxed molecular clock in

502  phylogenetic analysis. We found that substitutions during speciation may

503  profoundly affect phylogenetic inference: if node substitutions are not taken into

504  account, branching times tend to be overestimated, even when a relaxed clock is

505  used to counteract the effect of "hidden nodes". This suggests that incorporation

506  of a node substitution model may improve phylogenetic inference.

507  With our introduction of the node substitution model, we hope to stimulate

508  discussion on the biological explanation of variation in substitution rates within

509  and across phylogenies. Furthermore, we hope to have set a first step in improving

510  our understanding of this variation, and improving phylogenetic inference as a

511  whole.

512

513  CONFLICT OF INTEREST

514  The authors have no conflicts of interest to report.

515

525    Data Availability

526    R code to simulate the node substitution model has been made available as an R

527    package called 'nodeSub', and can be found here: https://CRAN.R-

528    project.org/package=nodeSub. All code used in simulations, and scripts

529    used to visualize obtained results, are available on dryad via:

530    https://doi.org/10.5061/dryad.t1g1jwt1x .

531    REFERENCES

532    Aldous D. 2001. Stochastic models and descriptive statistics for phylogenetic

533        trees, from Yule to today. Stat. Sci. 16:23–34.

534    Ast J.C. 2001. Mitochondrial DNA evidence and evolution in Varanoidea

535        (Squamata). Cladistics. 17:211–226.

536    Avise J.C., Ayala F.J. 1975. Genetic change and rates of cladogenesis. Genetics.

537        81:757–773.

538    Avise J.C., Ayala F.J. 1976. Genetic Differentiation in Speciose Versus Depauperate

539        Phylads: Evidence from the California Minnows. Evolution. 30:46.

540    Bilderbeek R.J.C., Etienne R.S. 2018. babette : BEAUti 2 , BEAST2 and Tracer for R.

541        2018:2034–2040.

542    Bilderbeek R.J.C., Laudanno G., Etienne R.S. 2020. Quantifying the impact of an

543        inference model in Bayesian phylogenetics. Methods Ecol. Evol.:1–8.

544    Bouckaert R., Vaughan T.G., Barido-Sottani J., Duchene S., Fourment M.,

545        Gavryuskina A., Heled J., Jones G., Kuhnert D., De Maio N., Matschiner M., K.

546        Mendes F., Muller N.F., Ogilvie H.A., du Plessis L., Popinga A., Rambaut A.,

547        Rasmussen D., Siveroni I., Suchard M.A., Wu C.H., Xie D., Zhang C., Stadler T.,

548        Drummond A.J. 2019. BEAST 2 .5 : An advanced software platform for

549        Bayesian evolutionary analysis. PLoS Comput. Biol. 15:e1006650.

550    Bromham L. 2011. The genome as a life-history character: Why rate of molecular

551        evolution varies between mammal species. Philos. Trans. R. Soc. B Biol. Sci.

552        366:2503–2513.

553    Bromham L. 2019. Six Impossible Things before Breakfast: Assumptions, Models,

554        and Belief in Molecular Dating. Trends Ecol. Evol. 34:474–486.

555    Bromham L., Hua X., Lanfear R., Cowman P.F. 2015. Exploring the relationships

556        between mutation rates, life history, genome size, environment, and species

557        richness in flowering plants. Am. Nat. 185:508–524.

558   Clarke K., Warwick R. 2001. A further biodiversity index applicable to species

559        lists: variation in taxonomic distinctness. Mar. Ecol. Prog. Ser. 216:265–278.

560   Dornburg A., Brandley M.C., McGowen M.R., Near T.J. 2012. Relaxed clocks and

561        inferences of heterogeneous patterns of nucleotide substitution and

562        divergence time estimates across whales and dolphins (Mammalia:

563        Cetacea). Mol. Biol. Evol. 29:721–736.

564   Douglas J., Zhang R., Bouckaert R. 2021. Adaptive dating and fast proposals:

565        Revisiting the phylogenetic relaxed clock model. PLOS Comput. Biol.

566        17:e1008322.

567   Dowle E.J., Morgan-Richards M., Trewick S.A. 2013. Molecular evolution and the

568        latitudinal biodiversity gradient. Heredity (Edinb). 110:501–510.

569   Drummond A.J., Ho S.Y.W., Phillips M.J., Rambaut A. 2006. Relaxed phylogenetics

570        and dating with confidence. PLoS Biol. 4:e88.

571   Duchene D., Bromham L. 2013. Rates of molecular evolution and diversification

572        in plants: Chloroplast substitution rates correlate with species-richness in

573        the Proteaceae. BMC Evol. Biol. 13:1.

574   Eldredge N., Gould S.J. 1972. Punctuated equilibria: an alternative to phyletic

575        gradualism. Models In Paleobiology. Freeman Cooper and Co. p. 82–115.

576   Eo S.H., DeWoody J.A. 2010. Evolutionary rates of mitochondrial genomes

577        correspond to diversification rates and to contemporary species richness in

578        birds and reptiles. Proc. R. Soc. B Biol. Sci. 277:3587–3592.

579   Etienne R.S., Haegeman B., Stadler T., Aze T., Pearson P.N., Purvis A., Phillimore

580        A.B. 2012. Diversity-dependence brings molecular phylogenies closer to

581    agreement with the fossil record. Proc. R. Soc. B Biol. Sci. 279:1300–1309.

582    Ezard T.H.G., Thomas G.H., Purvis A. 2013. Inclusion of a near-complete fossil

583        record reveals speciation-related molecular evolution. Methods Ecol. Evol.

584        4:745–753.

585    Faith D.P. 1992. Conservation evaluation and phylogenetic diversity. Biol.

586        Conserv. 61:1–10.

587    Fitch A.J., Goodman A.E., Donnellan S.C. 2006. A molecular phylogeny of the

588        Australian monitor lizards (Squamata : Varanidae) inferred from

589        mitochondrial DNA sequences. Aust. J. Zool. 54:253–269.

590    Fitch W.M., Beintema J.J. 1990. Correcting parsimonious trees for unseen

591        nucleotide substitutions: The effect of dense branching as exemplified by

592        ribonuclease. Mol. Biol. Evol. 7:438–443.

593    Fitch W.M., Bruschi M. 1987. The evolution of prokaryotic ferredoxins--with a

594        general method correcting for unobserved substitutions in less branched

595        lineages. Mol. Biol. Evol. 4:381–394.

596    Fontanillas E., Welch J.J., Thomas J.A., Bromham L. 2007. The influence of body

597        size and net diversification rate on molecular evolution during the radiation

598        of animal phyla. BMC Evol. Biol. 7:1–12.

599    Goldie X., Lanfear R., Bromham L. 2011. Diversification and the rate of molecular

600        evolution: No evidence of a link in mammals. BMC Evol. Biol. 11:286.

601    Jansson R., Davies T.J. 2008. Global variation in diversification rates of flowering

602        plants: Energy vs. climate change. Ecol. Lett. 11:173–183.

603    Janzen T., Höhna S., Etienne R.S.R.S. 2015. Approximate Bayesian Computation of

604        diversification rates from molecular phylogenies: introducing a new

605        efficient summary statistic, the nLTT. Methods Ecol. Evol. 6:566–575.

606    Jukes T., Cantor C. 1969. Evolution of protein molecules. Mamm. Protein Metab.

607        21.

608    Katoh K., Standley D.M. 2013. MAFFT multiple sequence alignment software

609        version 7: Improvements in performance and usability. Mol. Biol. Evol.

610        30:772–780.

611    King M.C., Wilson A.C. 1975. Evolution at two levels in humans and chimpanzees.

612        Science. 188:107–116.

613    Lanfear R., Ho S.Y.W., Love D., Bromham L. 2010a. Mutation rate is linked to

614        diversification in birds. Proc. Natl. Acad. Sci. U. S. A. 107:20423–20428.

615    Lanfear R., Welch J.J., Bromham L. 2010b. Watching the clock: Studying variation

616        in rates of molecular evolution between species. Trends Ecol. Evol. 25:495–

617        503.

618    Lartillot N., Phillips M.J., Ronquist F. 2016. A mixed relaxed clock model. Philos.

619        Trans. R. Soc. B Biol. Sci. 371.

620    Lartillot N., Poujol R. 2014. Correlated evolution of substitution rates and

621        quantitative traits. .

622    Lepage T., Bryant D., Philippe H., Lartillot N. 2007. A general comparison of

623        relaxed molecular clock models. Mol. Biol. Evol. 24:2669–80.

624    Lerner I.M. 1954. Genetic Homeostasis. .

625    Lewitus E., Morlon H. 2016. Characterizing and comparing phylogenies from

626        their laplacian spectrum. Syst. Biol. 65:495–507.

627    Louca S., Pennell M.W. 2020. Extant timetrees are consistent with a myriad of

628        diversification histories. Nature. 580:502–505.

629    Manceau M., Marin J., Morlon H., Lambert A. 2020. Model-based inference of

630        punctuated molecular evolution. Mol. Biol. Evol.:msaa144.

Systematic Biology

631    Mayr E. 1954. Change of genetic environment and evolution. Evol. as a

632        Process.:157–180.

633    Mindell D.P., Sites J.W., Graur D. 1989. Speciational Evolution: a Phylogenetic

634        Test With Allozymes in Sceloporus (Reptilia). Cladistics. 5:49–61.

635    Mindell D.P., Sites J.W., Graur D. 1990. Mode of allozyme evolution: Increased

636        genetic distance associated with speciation events. J. Evol. Biol. 3:125–131.

637    Nabholz B., Glémin S., Galtier N. 2008. Strong variations of mitochondrial

638        mutation rate across mammals - The longevity hypothesis. Mol. Biol. Evol.

639        25:120–130.

640    Nee S., May R.M., Harvey P.H., Trans P., Lond R.S. 1994. The reconstructed

641        evolutionary process. Philos. Trans. R. Soc. Lond. B. Biol. Sci. 344:305–11.

642    Pagel M., Venditti C., Meade A. 2006. Large punctuational contribution of

643        speciation to evolutionary divergence at the molecular level. Science.

644        314:119–121.

645    Pybus O., Harvey P. 2000. Testing macro–evolutionary models using incomplete

646        molecular phylogenies. Proc. R. Soc. B Biol. Sci. 267:2267–72.

647    Rabosky D.L., Donnellan S.C., Talaba A.L., Lovette I.J. 2007. Exceptional among-

648        lineage variation in diversification rates during the radiation of Australia's

649        most diverse vertebrate clade. Proc. R. Soc. B Biol. Sci. 274:2915–2923.

650    Ricklefs R.E. 2006. Global variation in the diversification rate of passerine birds.

651        Ecology. 87:2468–2478.

652    Russel P.M., Brewer B.J., Klaere S., Bouckaert R.R. 2019. Model Selection and

653        Parameter Inference in Phylogenetics Using Nested Sampling. Syst. Biol.

654        68:219–233.

655    Saclier N., François C.M., Konecny-Dupre L., Lartillot N., Guéguen L., Duret L.,

656    Malard F., Douady C.J., Lefébure T. 2018. Life history traits impact the

657        nuclear rate of substitution but not the mitochondrial rate in isopods. Mol.

658        Biol. Evol. 35:2900–2912.

659    Schliep K.P. 2011. phangorn : phylogenetic analysis in R. Bioinformatics. 27:592–

660        593.

661    Simpson G.G. 1945. SECTION OF BIOLOGY: Tempo and Mode in Evolution. Trans.

662        N. Y. Acad. Sci. 8:45–60.

663    Sipos B., Massingham T., Jordan G.E., Goldman N. 2011. PhyloSim - Monte Carlo

664        simulation of sequence evolution in the R statistical computing

665        environment. .

666    Spielman S.J., Wilke C.O. 2015. Pyvolve : A Flexible Python Module for Simulating

667        Sequences along Phylogenies. :1–7.

668    Stadler T. 2011. Simulating trees with a fixed number of extant species. Syst. Biol.

669        60:676–684.

670    Sung W., Ackerman M.S., Dillon M.M., Platt T.G., Fuqua C., Cooper V.S., Lynch M.

671        2016. Evolution of the insertion-deletion mutation rate across the tree of

672        life. G3 Genes, Genomes, Genet. 6:2583–2591.

673    van Valen L.M. 1985. Why and how do mammals evolve unusually rapidly. Evol.

674        Theory. 7:127–132.

675    Venditti C., Meade A., Pagel M. 2006. Detecting the node-density artifact in

676        phylogeny reconstruction. Syst. Biol. 55:637–643.

677    Venditti C., Pagel M. 2010. Speciation as an active force in promoting genetic

678        evolution. Trends Ecol. Evol. 25:14–20.

679    Webster A.J., Payne R.J.H., Pagel M. 2003. Molecular phylogenies link rates of

680        evolution and speciation. Science. 301:478.

681     Zhang R., Drummond A. 2020. Improving the performance of Bayesian

682         phylogenetic inference under relaxed clock models. BMC Evol. Biol. 20:1–28.

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706     Figure Legends

707     Figure 1. Difference in summary statistic values for trees inferred from an

708         alignment generated with node substitutions, and twin trees that were

709         inferred from an alignment generated without node substitutions, both

710         compared with the summary statistics of the true tree. We explored $\tau$ (the

711         amount of time spent on each node) as a fraction of crown age (horizontal

712         axis), and the impact of extinction (d, columns). The summary statistics are

713         the beta and gamma statistic, Laplacian spectrum, mean branch length,

714         nLTT statistic, and crown age. The figure shows that with increasing $\tau$, trees

715         inferred from an alignment generated with node substitutions show larger

716         differences with the true tree than trees inferred from an alignment

717         generated without node substititions. Differences with the true tree are

718         larger for trees inferred using the strict clock model than for those using the

719         relaxed clock model, but only for the alignment generated with node

720         substitutions.

721     Figure 2.  Effect of the node substitution model for phylogenies differing in

722         balance A) example plots of a randomly generated birth-death tree (top), a

723         fully balanced tree generated using the same branching times as the birth-

724         death tree (middle) and a very unbalanced tree generated using the same

725         branching times as the birth-death tree (bottom). Shown are trees with 20

726         tips for illustrative purposes, but results in B are from trees with 100 tips. B)

727         Difference in summary statistic with the true birth-death tree for

728         phylogenetic trees inferred from alignments generated using the node

729         substitution model on either balanced, unbalanced or random trees. We

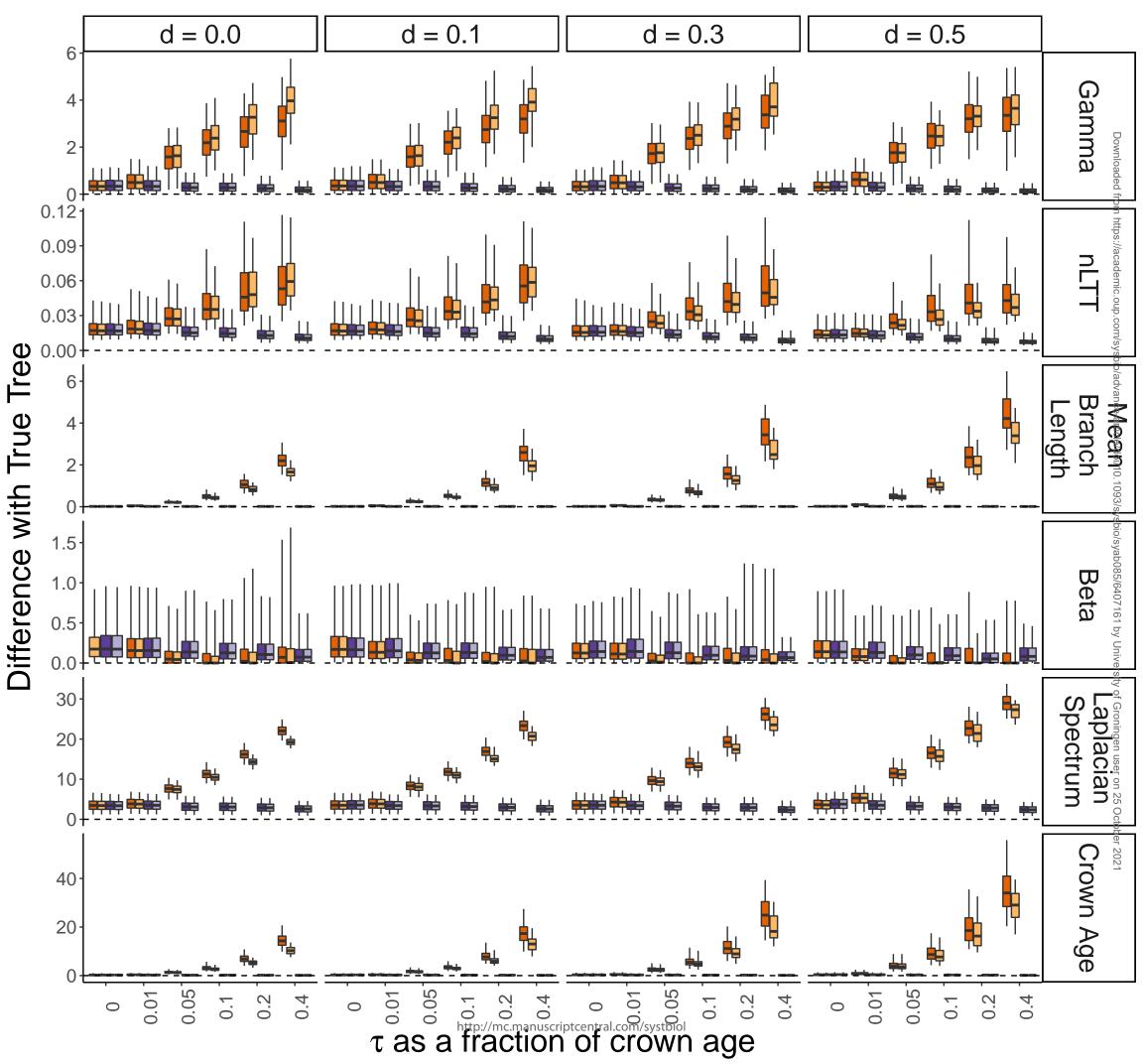730         explore $\tau$ as a fraction of crown age (horizontal axis), and the impact of
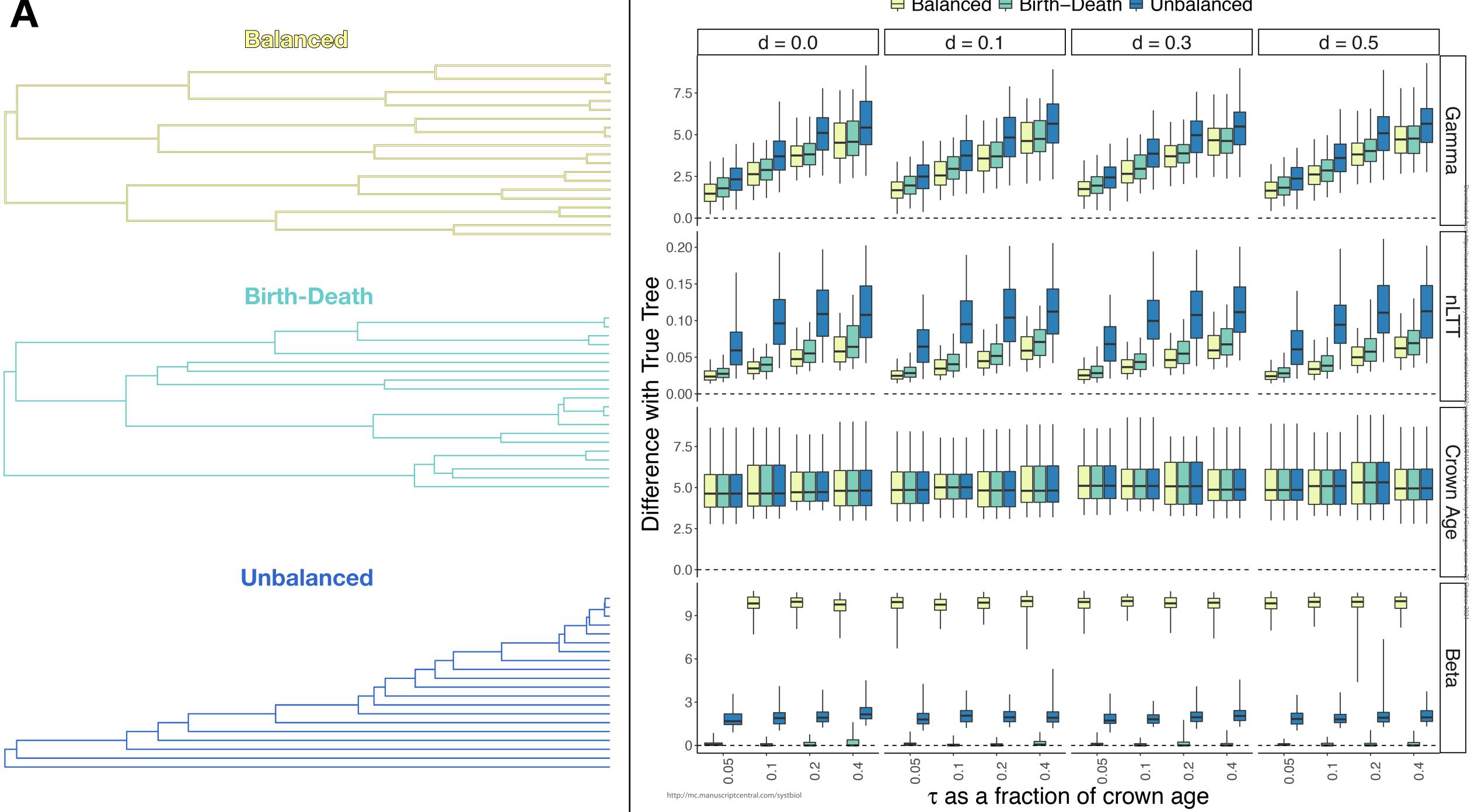
731     extinction (d, columns). The dotted line indicates zero difference with the

732     true tree. The summary statistics are the beta and gamma statistic, nLTT

733     statistic and tree height. Balanced and birth-death trees tend to have similar

734     inferred error, whereas unbalanced trees differ strongly, with a much larger

735     error for the gamma and nLTT statistic.

736  Figure 3. Marginal likelihood weight of the relaxed and strict clock model for

737     varying time spent on the nodes ($\tau$), where $\tau$ ⬚is chosen as fraction of the

738     crown age. Alignments generated with a node substitution model (top row)

739     are compared with alignments generated without node substitutions

740     (bottom row). Per parameter combination, 100 replicate trees were

741     analyzed. Because many dots are plotted on top of each other, we use solid

742     lines to indicate the best fitting locally estimated scatterplot smoothing

743     (LOESS), and the 95% Confidence interval (grey shaded area) of the LOESS

744     curve. As the time spent on the nodes increases, posterior support for the

745     relaxed clock model increases, but only if the alignment was generated with

746     a node substitution model.

747  Figure 4. Marginal likelihood weight of the relaxed clock model for trees of

748     varying balance, split out across different extinction rates (d = [0, 0.1, 0.3,

749     0.5]) and time spent on the nodes ($\tau$), where $\tau$ ⬚is chosen as fraction of the

750     crown age

751  (e.g. $\tau$ = 0.1 reflects a node time of 10% of the crown age). Per parameter

752     combination, Solid lines indicate the best fitting linear regression and the

753     95% confidence interval (grey shaded area) of regression. With increasing

754     values of $\tau$, the relative weight of the relaxed clock model becomes larger.

755     For smaller values of $\tau$, the relative weight of the relaxed clock model is

756    negatively correlated with the balance of the tree, with unbalanced trees

757    having a higher relative weight.

758  Figure 5. Results applying our node substitution model to an alignment

759    consisting of all 35 species of Varanidae occuring in the Indo-Australian

760    realm, assuming no extinction. Likelihood profiles with respect to $\tau$ of the JC

761    (A) and the GTR model (D) are shown. Figures B and E show the inferred

762    trees for both $\tau = 0$ and for the Maximum Likelihood value of $\tau$, for the JC

763    and GTR substitution models respectively. Figures C and F show these same

764    inferred trees, but here the branching times have been rescaled with respect
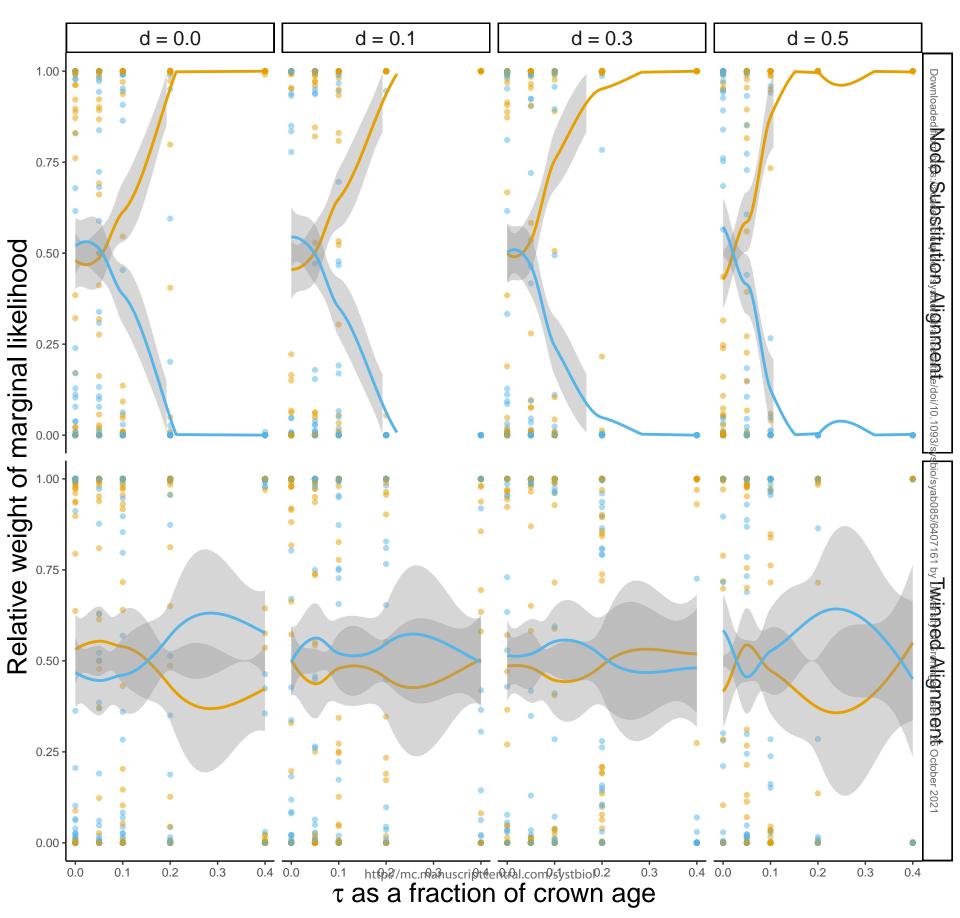
765    to the crown age.

766

767

768

769