1

2

3    **Model misspecification confounds the estimation of rates, and exaggerates their time**
4                                                   **dependency**

5

6

7

8

9            Brent C. Emerson[1,2*], Diego F. Alvarado-Serrano[3] and Michael J. Hickerson[3,4,5]

10

11

12    1. Island Ecology and Evolution Research Group, Instituto de Productos Naturales y Agrobiología

13    (IPNA-CSIC), C/Astrofísico Francisco Sánchez 3, La Laguna, Tenerife, Canary Islands, 38206,

14    Spain.

15    2. School of Biological Sciences, University of East Anglia, Norwich Research Park, Norwich NR4

16    7TJ, UK.

17    3. Biology Department, City College of New York, New York, NY, 10031, USA.

18    4. The Graduate Center, City University of New York, New York, NY, 10016, USA.

19    5. Division of Invertebrate Zoology, American Museum of Natural History, New York, NY 10024,

20    USA.

21

22    * Contact author: bemerson@ipna.csic.es

23

24

25 **Abstract**

26

27 While welcoming the comment of Ho *et al.* (2015), we find little that undermines the strength of

28 our criticism, and it would appear they have misunderstood our central argument. Here we respond

29 with the purpose of reiterating that we are (i) generally critical of much of the evidence presented in

30 support of the time dependent molecular rate (TDMR) hypothesis, and (ii) specifically critical of

31 estimates of $\mu$ derived from tip-dated sequences that exaggerate the importance of purifying

32 selection as an explanation for TDMR over extended timescales. In response to assertions put

33 forward by Ho *et al.* (2015), we use panmictic coalescent simulations of temporal data to explore a

34 fundamental assumption for tip-dated tree shape and associated mutation rate estimates, and the

35 appropriateness and utility of the date-randomisation test. The results reveal problems for the joint

36 estimation of tree topology, effective population size and $\mu$ with tip-dated sequences using BEAST.

37 Given the simulations, BEAST consistently obtains incorrect topological tree structures that are

38 consistent with the substantial overestimation of $\mu$ and under-estimation of effective population

39 size. Data generated from lower effective population sizes were less likely to fail the date-

40 randomisation test yet still resulted in substantially upwardly biased estimates of rates, bringing

41 previous estimates of $\mu$ from temporally sampled DNA sequences into question. We find that our

42 general criticisms of both the hypothesis of time-dependent molecular evolution, and Bayesian

43 methods to estimate $\mu$ from temporally sampled DNA sequences, are further reinforced.

44      **Introduction**

45

46      In their opening paragraph, and then repeated within their comment, Ho *et al.* (2015) state that we

47      (Emerson & Hickerson 2015) "claim that there is a lack of support for a time-dependent pattern in

48      molecular rate estimates". This is not correct. What we argue for, both in our original paper and

49      here, is that (i) there is a lack support for the inferred magnitude of TDMR patterns, and that (ii)

50      explanations of purifying selection over extended timescales to reconcile differences between

51      spontaneous $\mu$ and phylogenetic estimates of $\mu$ have been greatly exaggerated, largely because of

52      issues with biased rate estimates derived from ancient DNA (aDNA) analyses. Neither in this

53      response, nor in our original article, do we deny there to be evidence for time dependent patterns for

54      molecular rate estimates. Nor do we deny that purifying selection will lead to lower values for

55      spontaneous $\mu$. What we argue for in our original article (Emerson & Hickerson 2015), but

56      apparently misunderstood by Ho *et al.* (2015), is that the support for purifying selection

57      underpinning these observed patterns is greatly overstated when most of the observed changes in

58      estimates of $\mu$ can be explained as methodological artifacts. Purifying selection will lead to lower

59      values for spontaneous $\mu$. This is a truism that we have recognised previously (Emerson 2007).

60      However, the assumption of Ho *et al.* (2015) that pattern is evidence for process exaggerates both

61      the inferred extent of and timescale for rate reduction due to purifying selection. This is our central

62      argument and cause for concern.

63

64

65      **Evidence for pattern is not evidence of process**

66

67      A substantial part of the comment of Ho *et al.* (2015) is devoted to presenting many examples of

68      evidence for time-dependent rate estimates, although for nuclear data, Ho *et al.* (2015) acknowledge

69      that there is no strong evidence for such a pattern. As stated above, we are not in denial of the many

70      published estimates supporting the pattern for mtDNA, and as such our position is somewhat

71      misrepresented by Ho *et al.* (2015). It is important to point out that, if a pattern can be explained by

72      something other than the hypothesis (the hypothesis here being purifying selection), then the pattern

73      itself cannot be used as evidence in support of the hypothesis. In this context, the examples

74      presented by Ho *et al.* (2015) do not in themselves contradict the points raised in Emerson &

75      Hickerson (2015), as these may be subject to the methodological issues raised in our original article.

76      Indeed, some of the examples where we highlight methodological issues (e.g. *Caenorhabditis*

77      *elegans*) are presented again by Ho *et al.* (2015) as supporting the hypothesis of time-dependent

78  molecular evolution without further discussion of the concerns we raised. We focus the remainder

79  of this response on specific points within the comment of Ho *et al.* (2015), where we feel they may

80  have either failed to provide an adequate response, or misrepresented our work, when discussing the

81  evidence for the hypothesis that purifying selection is the driver of TDMR estimates.

82

83

84  **Adélie penguin data**

85

86  In our original article (Emerson & Hickerson 2015) we pointed out that, in contradiction to the

87  TDMR hypothesis (i.e. the hypothesis that molecular rate estimates decrease toward the past as a

88  consequence of purifying selection) mean pedigree-based estimates of the mutation rate of

89  mitochondrial DNA in Adélie penguins are lower than those inferred from aDNA. In response to

90  this, Ho *et al.* (2015) make two points. They first suggest that the non-reporting of 95% credibility

91  intervals may somehow limit the significance of our observation, and further claim there to be

92  substantial overlap in the 95% credibility intervals between aDNA estimates and the pedigree

93  estimate. They then state that we acknowledged that both the pedigree rate and aDNA rate estimates

94  "greatly" exceed those inferred from fossil-calibrated analyses of birds. The first point is incorrect,

95  and thus misrepresents our original work (Emerson & Hickerson 2015), as the 95% CI of one of the

96  three published aDNA estimates of $\mu$ does not overlap with the pedigree-derived estimate of $\mu$. The

97  second point requires further context (see below) to understand the extent to which both pedigree

98  and aDNA rates for Adélie penguins can be compared to a phylogenetic rate.

99  With regard to the first point, we stated in our original work (Emerson & Hickerson 2015)

100  that the Adélie aDNA rate estimate of Ho *et al.* (2007a) is *significantly* higher than the pedigree

101  rate. Thus, in contrary to the claim of Ho *et al.* (2015), there is no overlap among their 95%

102  credibility intervals. We do not deny that the 95% credibility intervals of the aDNA rate estimates

103  of Lambert *et al.* (2002) and Millar *et al.* (2008), which have a lower mean value than that of Ho *et*

104  *al.* (2007a), overlap with the pedigree rate. However, this should not be seen as somehow

105  undermining the discrepancy between these two aDNA rate estimates and the pedigree rate in a

106  field (TDMR) where trends in mean values are frequently reported as support for the hypothesis.

107  With regard to the second point, we recognize that the mean values for all aDNA rate

108  estimates and the pedigree-derived rate estimate of $\mu$ are higher than the bird phylogenetic

109  divergence rate of 0.208 mutations/site/Myr presented by Shields and Wilson (1987) that has been

110  used in previous comparisons (e.g. Lambert *et al.* 2002; Millar *et al.* 2008). However, there are

111  several features of this phylogenetic rate estimate that limit its use for comparative purposes.

112    Firstly, it is not a general bird rate estimate, it is an estimate derived from the analysis of 5 species

113    of geese. A difference between a phylogenetically derived mutation rate for geese, and aDNA or

114    pedigree-derived rates for penguins may equally be explainable by fundamental differences

115    between these very different, phylogenetically distant taxonomic groups. Secondly, the

116    phylogenetic rate is probably underestimated, as recognised by Shields and Wilson (1987), due to

117    the difficulty of estimating genetic divergences from restriction fragment analysis.

118

119

120    **Comparing pedigree-derived rate estimates with phylogenetic rate estimates**

121

122    We have previously pointed out, using *Caenorhabditis elegans* as an example, that a mutation

123    accumulation line or pedigree-derived estimate of $\mu$ for a given taxa can only be considered high if

124    it exceeds a taxonomically relevant phylogenetic rate (Emerson & Hickerson 2015). We provide an

125    additional example of this problem above, with the inappropriate comparison of Adélie penguin

126    pedigree and aDNA-derived estimates of $\mu$ with a phylogenetic estimate of $\mu$ derived from geese.

127    Rather than providing suitable comparisons within their reply, Ho *et al.* (2015) continue to cite the

128    spontaneous mutation rate for *C. elegans* (Denver *et al.* 2004), as well as *Drosophila melanagoster*

129    (Keightley et al. 2014) and *Heliconius melpomene* (Keightley *et al.* 2015), as being higher than

130    "corresponding phylogenetic estimates". There are no phylogenetic estimates within the response of

131    Ho *et al.* (2015), nor within the original articles, with the exception of Keightley *et al.* (2015), who

132    note that applying their spontaneous mutation rate to estimate the age of the *Heliconius* suggests

133    that the fossil-calibrated age for the genus is approximately correct. The spontaneous rate is

134    however higher than the fossil rate, and as pointed out by Keightley *et al.* (2015), further work is

135    needed to reconcile the two estimates. But the difference itself is not evidence for the TDMR

136    hypothesis when alternative equally plausible explanations exist. For example, a difference could

137    arise because (1) the data sets being compared are very different (whole genome vs a non-random

138    set of protein coding genes), or (2) only secondary calibration points were used for the phylogeny

139    (i.e. there are no fossil Heliconiini). But let's assume the difference is real. What does it tell us? It

140    tells us that purifying selection results in the underestimation of spontaneous $\mu$ when using a

141    phylogenetic calibration. What it does not tell us is the timescale over which this occurs, and thus

142    such data is uninformative about the timescale for the TDMR hypothesis.

143

144

145    **Estimates of $\mu$ from temporally sampled DNA, and their lack of validation**

146

147 Ho *et al.* (2015) take issue with our claim that, while many studies have produced estimates of $\mu$

148 from aDNA, none have provided validation of their estimates independently of the Bayesian

149 implementation within BEAST (Drummond *et al.* 2012) from which they were derived (Emerson &

150 Hickerson 2015). To support that we are "demonstrably wrong", they cite two tests to evaluate the

151 information content of time-structured data. However, these either have not provided, or do not

152 provide, independent estimates of $\mu$. The first of these, the regression of tree height against

153 sampling time of Fitch *et al.* (1991) can, with some caveats, be used to estimate $\mu$ but has not, to

154 our knowledge, ever been used to validate a Bayesian estimate of $\mu$. The second test cited by Ho *et*

155 *al.* (2015), that of Ramsden *et al.* (2009), which has been further developed by Duchêne *et al.*

156 (2015), is not independent. It is a test of information content, where the Bayesian estimate of $\mu$ is

157 compared to the distribution of $\mu$ estimated when dates are randomised across the tree. Thus, our

158 original assertion still stands - Bayesian estimates of $\mu$ have yet to be independently validated.

159

160

161 **Measurably evolving populations, date-randomisation and $\mu$**

162

163 Ho *et al.* (2015) provide a summary of the date-randomisation test, first presented by Ramsden *et*

164 *al.* (2009) to test for sufficient signal within temporally sampled DNA data sets to estimate $\mu$ and

165 divergence dates. It is important to consider what the 95% credibility interval of the date-

166 randomised rate estimate represents. Ho *et al.* (2015) correctly point out that the two data sets

167 presented in the schematic trees in Fig. 2 of Emerson & Hickerson (2015) would yield positive and

168 misleading estimates of $\mu$. We agree with this, but we do not agree with their conclusion that both

169 data sets do not represent "measurably evolving populations". On the contrary, both data sets do

170 represent measurably evolving populations. The definition of genetic change in populations used by

171 Ho *et al.* (2015) and elsewhere (e.g. Drummond *et al.* 2003; Ewing *et al.* 2004) is of mutation

172 between sampling time points. However, it has been long understood that genetic change in

173 populations involves changes in allele frequencies under the dynamic between mutation, selection

174 and drift (Hartl & Clark 2007), and it is important to clarify that the mutation rate $\mu$ is the rate of

175 mutation along any branch of a sampled gene genealogy, rather than being the rate of new

176 mutations within a population or rate of mutational turnover between sampling time points. For

177 example, due to the coalescent process, the vast majority of mutations between two temporally

178 different samples can often occur at times older than either of the samples. As recognised by Ho *et*

179 *al.* (2015), the sampling scenarios in panels C and D of Fig. 2 (Emerson & Hickerson 2015) will

180    yield non-zero estimates of $\mu$. Ho *et al.* (2015) also suggest that both data sets would fail the date-

181    randomisation test of Ramsden *et al.* (2009). We agree that they probably would fail (although that

182    can only be assessed by direct analysis). However, from this point we disagree with Ho *et al.*

183    (2015), and the accepted interpretation of the date-randomisation test - that if the empirical estimate

184    exceeds the 95% confidence intervals from the randomised distribution, then the empirical value is

185    a reliable estimate of $\mu$.

186         Regardless of whether a dataset passes the randomization test or not, estimates of $\mu$ from

187    temporally sampled data using BEAST may be overestimated because of other population genetic

188    (drift and the coalescent) and sampling processes, as well as phylogenetic constraints that BEAST

189    imposes on temporally sampled data (Box 1). Citing Duchêne *et al.* (2015), Ho *et al.* (2015) point

190    out that data sets that fail the date-randomisation test tend to yield overestimates of $\mu$, which could

191    be taken to suggest that data sets that pass the test provide meaningful approximations of $\mu$. This is

192    not the case. A careful examination of Duchêne *et al.* (2015) reveals that data sets can pass the test

193    *and* yield significant overestimates of $\mu$, where the the 95% confidence interval of the estimate does

194    not include $\mu$. In fact, the parameter space within which both the estimation of $\mu$ is correct, and the

195    test is passed, is limited (Fig. 1 of Duchêne *et al.* 2015). The take home point is that passing the

196    date-randomisation test is not validation for an estimation of $\mu$ using the BEAST temporally

197    sampled model. To more fully explore this dynamic, we have conducted coalescent simulations of

198    temporally sampled data, matching parameters commonly associated with ancient mtDNA data, and

199    show that BEAST can systematically overestimate $\mu$ given temporally sampled data due to incorrect

200    topological estimates that arise from constraining tip dates (Box 1).

201

202

203    **TDMR for some genomes, and not for others?**

204

205    Ho *et al.* (2015) suggest that there is scant evidence for an observed TDMR pattern in nuclear

206    genomes. It will be interesting to see what is learned from new genomic data as it emerges,

207    although it is worth pointing out that much of this observed discrepancy between nDNA and

208    mtDNA evaporates if the studies using tip-dating methods with ancient mitochondrial DNA are

209    confirmed to be the non-trivial overestimates as suggested from our simulation-based exploration.

210    Furthermore, their assertion that "unfortunately, there remains considerable uncertainty about

211    nuclear mutation rates in humans", is vague and misleading, as the various papers show strong

212    evidence that there is genetic variation for the mutation rate and that paternal age can drive

213    differences in mutation rates (e.g. Scally & Durbin 2012; Thomas & Hahn 2014). It also seems

214   somewhat incongruous for Ho *et al.* (2015) to criticise us for reporting short-term estimates of $\mu$ for

215   nuclear data, while they themselves report such data when they believe it to support their argument

216   (e.g. Denver *et al.* 2004; Keightley *et al.* 2014; Keightley *et al.* 2015, but see comments above).

217

218

219   **Bison data and the Bayesian estimation of $\mu$ from temporally sampled DNA**

220

221   Ho *et al.* (2015) cast doubt on two aspects of our reanalysis of the *Bison bison* data first published

222   by Shapiro *et al.* (2004) and reanalysed by Ho *et al.* (2015). Their concerns regarding the impact of

223   fixing effective populations size are vague and misleading, as they seem to suggest that there are

224   "other parameters" in the cataclysmic demographic model that might somehow explain our results.

225   As we have made all our input files publicly available, it is not clear why Ho *et al.* (2015) do not

226   quantitatively assess their concern. A reanalysis exploring their parameters of concern would

227   suffice. We therefore see nothing in the argument of Ho *et al.* (2015) regarding the fixing of modern

228   effective population size for *B. bison*, that explains our results.

229          With regard to their other doubt, Ho *et al.* (2015) state that fixing the root age of the analysis

230   explains our result because "removing the sequences from older samples to reduce the sampling

231   window preferentially removes older branches in the gene tree". In doing so, Ho *et al.* (2015)

232   assume a correlation between DNA sequence sampling time, and the coalescence time of the

233   sampled sequence, which is in stark contrast to expectations under the standard Kingman coalescent

234   for a single panmictic population without size change or subdivision (Tajima 1983). When we

235   examined this assumption of Ho *et al.* (2015) it was apparent that, when compared to an

236   unconstrained tree of the *B. bison* data, constraining the tree with tip dates positively contributes to

237   such a correlation. The maximum clade credibility tree for the *B. bison* data with tip date constraints

238   is topologically very different from the unconstrained tree, with DNA sequences of older age

239   branching more basally within the tip date-constrained tree (Appendix S1, Supporting Information).

240   As an explanation for this, we can only conclude that enforcing tip dates as a constraint contributes

241   to the overestimation of $\mu$, due to additional mutation change in the tree required to accommodate

242   topological difference. We further explore these issues using coalescent simulations of temporally

243   sampled data under a single panmictic population and find that indeed BEAST tends to incorrectly

244   misestimate the gene genealogies as well as consistently overestimate $\mu$ given the sample size and

245   temporal distribution of tips of the *B. bison* data (Box 1). Our analyses (Box 1) call into question all

246   previous estimates of $\mu$ from tip-dated sequences using BEAST.

247    Ho *et al.* (2015) seem to be dismissive of their *B. bison* data, suggesting it to be small by

248    current measures. It is in fact among the biggest data sets that have been analysed to date, providing

249    an apparently compelling example of significance with respect to the date-randomisation test (Ho *et*

250    *al.* 2011). Their argument that bigger data sets for a greater variety of genes will yield more

251    decisive results will only be realised if the concerns we raise both here and in Emerson & Hickerson

252    (2015) are taken on board. There are clear and identifiable problems with the estimation of $\mu$ from

253    temporally sampled sequences, and not all these problems will necessarily be solved with more

254    data.

255

256

257    **Conclusions**

258

259    After responding to the comment Ho *et al.* (2015), we find that our general criticisms of both (i) the

260    hypothesis of time-dependent molecular evolution, and (ii) methods to estimate $\mu$ from temporally

261    sampled DNA sequences, are further reinforced. As we have previously pointed out (Emerson &

262    Hickerson 2015), much of the perceived support for the time-dependent molecular evolution

263    hypothesis comes from overestimates of $\mu$ that are derived from  phylogenetic analyses of

264    temporally calibrated aDNA using the Bayesian program BEAST. Such estimates of $\mu$ have been

265    argued to be evidence against calibration error as a sufficient explanation for patterns of TDMR (Ho

266    *et al,* 2011). In this article we clearly identify a positive bias in the estimation of $\mu$ from tip-dated

267    gene trees with BEAST that appears to be associated with the interaction between effective

268    population size and enforcing the age of DNA sequences when reconstructing the topologies of the

269    gene genealogies. Together with previously raised concerns (Debruyne & Poinar 2009; Emerson

270    2007; Emerson & Hickerson 2015; Navascués & Emerson 2009; Ramakrishnan & Hadly 2009) it is

271    now clear that published estimates of $\mu$ using aDNA data should be considered unreliable,

272    particularly if it cannot be shown that analyses underpinning the estimates did not result in

273    topological differences between tip-date constrained and unconstrained trees. As we have pointed

274    out, much of the remaining evidence for patterns of TDMR estimates can be explained without

275    resorting to selection, suggesting no more than a limited temporal contribution of purifying

276    selection to reconcile differences between spontaneous $\mu$ and phylogenetic estimates of $\mu$.

277

278    **Acknowledgements**

279    We thank Sebastián Duchêne for providing access to simulation files from Duchêne *et al.* (2015).

280

281     **References**

282     Anderson CN, Ramakrishnan U, Chan YL, Hadly EA (2005) Serial SimCoal: a population genetics
283             model for data from multiple populations and points in time. *Bioinformatics*, **21**, 1733-1734.

284     Debruyne R, Poinar HN (2009) Time dependency of molecular rates in ancient DNA data sets, a
285             sampling artifact? *Systematic Biology*, **58**, 348-359.

286     Denver DR, Morris K, Lynch M, Thomas WK (2004) High mutation rate and predominance of
287             insertions in the *Caenorhabditis elegans* nuclear genome. *Nature*, **430**, 679-682.

288     Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG (2003) Measurably evolving
289             populations. *Trends in Ecology and Evolution*, **18**, 481-488.

290     Drummond AJ, Suchard MA, Xie D, Rambaut A (2012) Bayesian phylogenetics with BEAUti and
291             the BEAST 1.7. *Molecular Biology and Evolution*, **29**, 1969-1973.

292     Duchêne S, Duchêne D, Holmes EC, Ho SYW (2015) The performance of the date-randomisation
293             test in phylogenetic analyses of time-structured virus data. *Molecular Biology and
294             Evolution*, **32**, 1895-1906.

295     Emerson BC (2007) Alarm Bells for the molecular clock? No support for Ho et al.'s model of time-
296             dependent molecular rate estimates. *Systematic Biology*, **56**, 337-345.

297     Emerson BC, Hickerson MJ (2015) Lack of support for the time-dependent molecular evolution
298             hypothesis. *Molecular Ecology*, **24**, 702-709.

299     Ewing G, Nicholls G, Rodrigo A (2004) Using temporally spaced sequences to simultaneously
300             estimate migration rates, mutation rate and population sizes in measurably evolving
301             populations. *Genetics*, **168**, 2407-2420.

302     Fitch WM, Leiter JME, Li XQ, Palese P (1991) Positive Darwinian evolution in human influenza-a
303             viruses. *Proceedings of the National Academy of Sciences of the United States of America*,
304             **88**, 4270-4274.

305     Hartl DL, Clark AG (2007) *Principles of population genetics* Sinauer Associates, Sunderland, MA.

306     Ho SYW, Duchêne S, Molak M, Shapiro B (2015) Time-dependent estimates of molecular rates:
307             Evidence and causes. *Molecular Ecology*, **this issue**.

308     Ho SYW, Kolokotronis S-O, Allaby RG (2007a) Elevated substitution rates estimated from ancient
309             DNA sequences. *Biology Letters*, **3**, 702-705.

310     Ho SYW, Lanfear R, Phillips MJ*, et al.* (2011) Bayesian estimation of substitution rates from
311             ancient DNA sequences with low information content. *Systematic Biology*, **60**, 366-374.

312     Ho SYW, Shapiro B, Phillips MJ, Cooper A, Drummond A (2007b) Evidence for time dependency
313             of molecular rate estimates. *Systematic Biology*, **56**, 515-522.

314    Keightley PD, Ness RW, Halligan DL, Haddrill PR (2014) Estimation of the spontaneous mutation
315        rate per nucleotide site in a *Drosophila melanogaster* full-sib family. *Genetics*, **196**, 313-
316        320.
317    Keightley PD, Pinharanda A, Ness RW*, et al.* (2015) Estimation of the spontaneous mutation rate in
318        *Heliconious melpomene*. *Molecular Biology and Evolution*, **32**, 239-243.
319    Kuhner MK, Yamato J, Felsenstein, J (1995) Estimating effective population size and mutation rate
320        from sequence data using Metropolis-Hastings sampling. *Genetics*, **140**, 1421-1430.
321    Lambert DM, Ritchie PA, Millar CD*, et al.* (2002) Rates of evolution in ancient DNA from Adelie
322        penguins. *Science*, **295**, 2270-2273.
323    Millar CD, Dodd A, Anderson JM*, et al.* (2008) Mutation and evolutionary rates in Adélie penguins
324        from the Antarctic. *Plos Genetics*, **4**, e1000209.
325    Navascués M, Emerson BC (2009) Elevated substitution rate estimates from ancient DNA: model
326        violation and bias of Bayesian methods. *Molecular Ecology*, **18**, 4390-4397.
327    Ramakrishnan U, Hadly EA (2009) Do complex population histories drive higher estimates of
328        substitution rate in phylogenetic reconstructions? *Molecular Ecology*, **18**, 4341-4343.
329    Ramsden C, Holmes EC, Charleston MA (2009) Hantavirus evolution in relation to its rodent and
330        insectivore hosts: No evidence for codivergence. *Molecular Biology and Evolution*, **26**, 143-
331        153.
332    Scally A, Durbin R (2012) Revising the human mutation rate: implications for understanding
333        human evolution. *Nature Reviews Genetics*, **13**, 745-753.
334    Shapiro B, Drummond AJ, Rambaut A*, et al.* (2004) Rise and fall of the beringian steppe bison.
335        *Science*, **306**, 1561-1565.
336    Shields GF, Wilson AC (1987) Calibration of mitochondrial DNA evolution in geese. *Journal of
337        Molecular Evolution*, **24**, 212-217.
338    Steel MA, Penny P (1993) Distributions of tree comparison metrics - some new results. *Systematic
339        Biology*, **42**, 126-141.
340    Tajima F (1983) Evolutionary relationship of DNA sequences in finite populations. *Genetics*, **105**,
341        437-460.
342    Thomas GWC, Hahn MW (2014) The human mutation rate is increasing, even as it slows.
343        *Molecular Biology and Evolution*, **31**, 253-257.

344

345

346

347 **Author contributions**

348 B.C.E., D.A.S. and M.J.H contributed equally to the preparation of this manuscript. All simulations

349 were conducted by D.A.S.

350

351 **Data Accessibility**

352 All scripts for the simulations conducted within this manuscript and an example BEAST input file

353 are available from https://diegofalvarado-s@bitbucket.org/diegofalvarado-s/tdmra_simulations.git.

354 Bison DNA sequences and their sampling dates can be found within the online supporting

355 information associated with Emerson & Hickerson (2015), doi: 10.1111/mec.13070.

356

357 **Supporting Information**

358 Additional Supporting Information may be found in the online version of this article:

359 **Appendix S1** Maximum clade credibility trees for the *Bison bison* data of Ho *et al.* (2007a) with

360 and without age constraints enforced for the tips.

**Box 1: Tree shape and the overestimation of $\mu$ from tip-dated sequences**

Constraining the tip dates within a phylogeny is expected to change branch lengths, but it might be less clear why topological relationships inferred from identical patterns of sequence variation should change. As can be seen in Appendix S1 (Supporting Information), the maximum clade credibility tree for *Bison bison* (Ho *et al.* 2007b; Shapiro *et al.* 2004) with tip date constraints is topologically very different from an unconstrained tree, with changes involving DNA sequences of older age branching more basally within the tip date-constrained tree, as would be expected if the panmictic effective population size was small. In some cases these rearrangements do not appear to increase the inferred amount of mutational change within the tree, as the change in gene tree topology does not disrupt patterns of shared derived variation, yet in other cases patterns of shared derived variation within the unconstrained tree are disrupted, increasing homoplasy and thus inferring additional mutational change within the tip-dated tree. One obvious outcome of an increase in the inferred number of mutational changes in a tip-date constrained tree is that the estimation of $\mu$ will also increase.

To explore this behavior, we followed a simulation procedure similar to that of Duchêne *et al.* (2015) ― the main difference being the use of an explicit coalescent simulator, BayesSSC (Anderson *et al.* 2005) instead of BEAST (Drummond *et al.* 2012) to generate the input tree topologies given known effective population sizes ($N$) and mutation rates ($\mu$), and a tip date distribution similar to the *B. bison* data (pipeline is available at https://bitbucket.org/diegofalvarado/tdmra_simulations). We have found that trees inferred by BEAST for tip-dated sequences tend to enforce an age-based coalescent pattern on the posterior distribution of gene trees. This pattern would be expected given small effective population sizes, despite true $N$ being 483,827 and 1,451,481 individuals in the simulation models that generated the simulated datasets. One likely culprit is how the the compound demographic parameter ($\theta = 4N\mu$ where $N$ is the effective population size and $\mu$ is the per site per generation per genealogical lineage mutation rate) is decoupled into joint estimates of $N$ and $\mu$ in BEAST. Under a standard panmictic coalescent model, it is only possible to estimate the compound parameter $\theta$ rather than its components ($N$ and $\mu$) unless one of the two parameters are known or assumed (Kuhner et al. 1995). In contrast, the tip-dated panmictic coalescent model employed in BEAST allows decoupling the posterior estimates of $\theta$ into $N$ and $\mu$ using the temporal-mutational information provided from the age-inforced tips of the posterior distribution of gene genealogies. As true $N$ becomes larger, the tip-dated constraints result in inferred gene tree topologies that increasingly depart from the true gene tree topologies (Fig. I). This increasing level of phylogenetic inferential error corresponding with increasing levels of false homoplasy, which in turn corresponds with overestimates of $\mu$ and

395     underestimates of *N*. In other words, underestimates of *N* result in older samples coalescing more

396     basally than younger samples in the inferred topologies, and the consequences of this dynamic

397     appear to be more severe when the true *N* was larger (Fig. I). As true *N* is larger, the magnitude of *N*

398     underestimation and *µ* overestimation becomes more severe with inferred gene tree topologies

399     becoming more age-constrained from the true topologies (Fig. II).

400            Of note is that under these simulations such overestimates of *µ* did not typically pass the

401     date-randomisation test, yet this was less the case under the smaller true *N* (Fig. II). Under a

402     coalescent model with small sized populations, one would expect genealogical coancestry between

403     samples of similar age (i.e., age-based coalescence), and as expected, the simulations reveal that the

404     probability of this is inversely related to population size (Figure III). At the same time, the

405     randomization of tip ages has a stronger impact on rate estimates when disrupting patterns of age-

406     based coalescence in the original tree, and hence, the date-randomization test is more likely passed

407     when the true gene genealogy has a tighter age-coalescent time association (such as under relatively

408     small effective population sizes; Figure III). Accordingly, as can be seen in Fig. IV, the association

409     of coalescence time with sample age is much stronger for the bison data when compared to patterns

410     obtained when simulating under a panmictic coalescent population model. Such a pattern is

411     expected for population structure and/or small *N*. We suggest that even though the bison data was

412     likely generated under scenarios that differed from what we explored in our simulations, the

413     systematic overestimates of *µ* and underestimates of *N* are likely to still be at play with these

414     estimates being biased by the consequences of large effective population sizes, population

415     subdivision and/or local colonisation/extinction. Clearly this is in need of further evaluation with

416     simulations that capture the demographic complexity and the patterns of tip-dates and coalescent

417     times that are observed in real data.

418            Given that topological inconsistencies in BEAST appear to be associated with biasing

419     estimates of both the number and age of DNA mutations together with overestimates of *µ* and

420     underestimates of *N*, we make the following two suggestions. Firstly it would seem relevant to

421     report the agreement between the topologies of tip-date constrained and unconstrained trees when

422     reporting estimates of *µ*. Secondly, we suggest that while previous approaches using coalescent

423     simulation have been useful to demonstrate that, under some conditions, BEAST can successfully

424     estimate *µ* from tip-dated sequences of virus sequences (e.g. Duchêne *et al.* 2015), the complex

425     conditions underlying temporally sampled ancient DNA with respect to sample sizes, effective

426     population sizes, generation times, and subdivision need to be more fully examined to understand

427     when estimates of *µ* from BEAST may be positively biased. Our simulations show that estimates of

428     *µ* from such data can be systematically upwardly biased, and as such a more thorough exploration

429   of the impacts of sample characteristics, historical demographics and analysis settings is needed to

430   better understand the underlying causes of the methodological artifacts we have revealed. Our

431   simulations also suggest that all previous estimates of $\mu$ from temporally sampled DNA sequence

432   data using BEAST need a thorough reexamination before they can be accepted.

433

434   **Figure I.** Comparison of simulated and recovered tree topology for tip-dated sequence data using

435   BEAST (Drummond *et al.* 2012). Note that the tree topology inferred by BEAST (b and d) is

436   markedly different from the tree used to simulate the sequences (a and c) that serve as input to

437   BEAST. This problem is accentuated under comparatively larger population sizes (Robinson-

438   Foulds distance between a and b = 250, between c and d = 260; weighted-path difference (Steel &

439   Penny 1993) between a and b = 0.77, between c and d = 20.09). Tips are coloured based on age to

440   highlight the tendency for age-based coalescent events (i.e. tendency of younger samples to cluster

441   as ingroups to older samples) in BEAST-estimated trees.

442

443   **Figure II.** Estimates of the substitution rate (in log 10 scale) against the width of the calibration

444   window under two different populations sizes: (a) $N = 483,427$; (b) $N = 1,451,481$. The solid

445   horizontal line represent the true simulated rate (mean=1e-8, sd=5%). Symbols represent the mean

446   rate estimate for each simulation, with the error bars showing the 95% credible intervals. We

447   conducted 10 randomizations for the date-randomization test for all data sets. Circles denote rate

448   estimates that failed the test according to both criteria CR1 and CR2 (Duchêne *et al.* 2015), whereas

449   triangles denote those that failed according to CR2 only. Numbers of type I and type II errors are

450   shown for each rate treatment.

451

452   **Figure III.** Clustering of tip ages in BEAST-obtained trees based on simulated samples for (a) $N =$

453   483,427 and (b) $N = 1,451,481$. The ages of pairs of closest related tips is depicted, with original

454   values represented in red, and date-randomised values represented in blue. Note how the difference

455   between date-randomised and original data is smaller when the effective population size is

456   comparatively larger, making it less likely to pass the test proposed by Duchene *et al.* (2015).

457

458   **Figure IV.** Association between tip-age and relative coalescence time. Patristic distance is used as

459   an indicator for the time of coalescence of each sample in the tree. Note the empirical bison dataset

460   (black) (Ho *et al.* 2007b) shows a much tighter association than any of the simulated datasets (small

461   $N = 483,427$ in blue, large $N = 1,451,481$ in red) indicating a strong tendency for samples to

462    coalesce together based on their age in this dataset. Such a pattern is expected under small effective
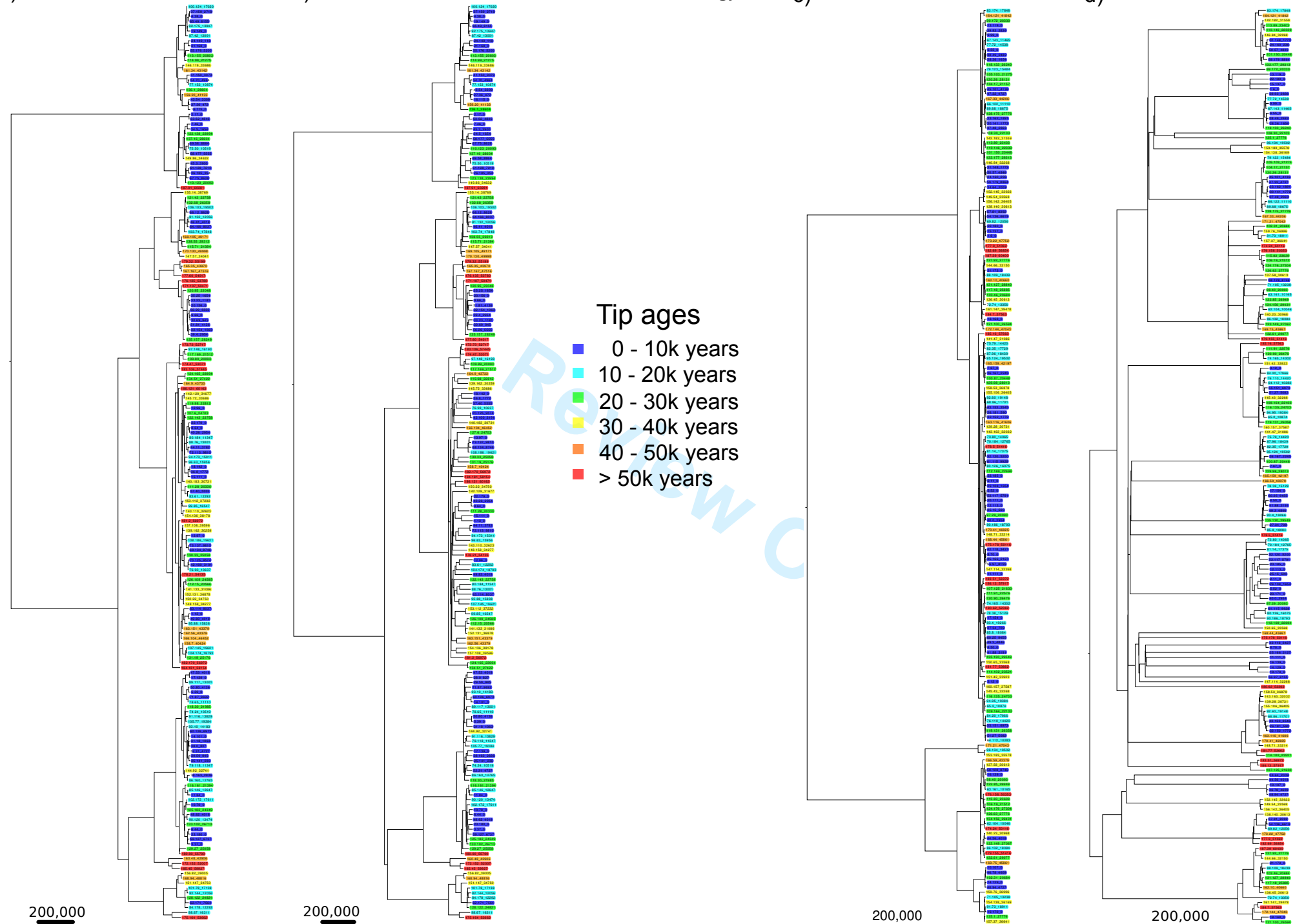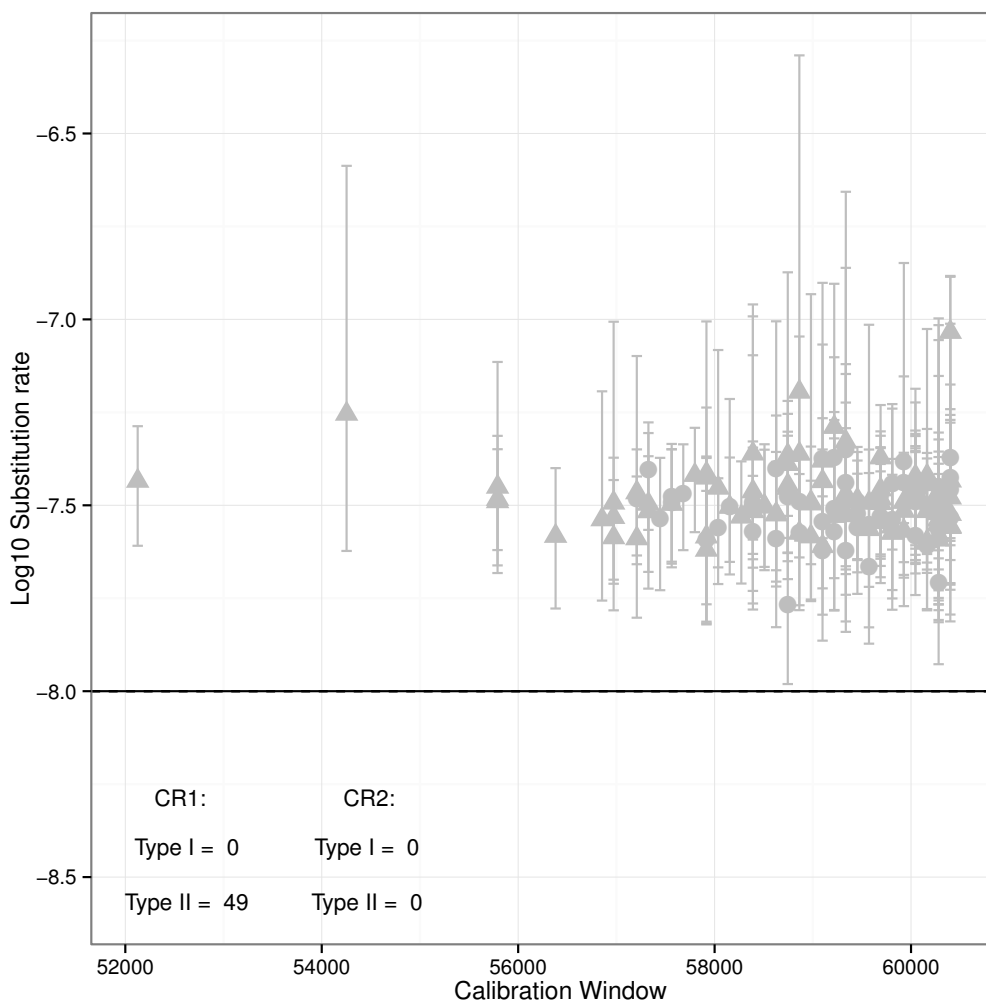
463    population sizes and/or population structure.

Tip ages

- 0 - 10k years
- 10 - 20k years
- 20 - 30k years
- 30 - 40k years
- 40 - 50k years
- > 50k years

Ne = 483,421

Ne = 1,451,481

a)



b)

a)



b)