

1 **A prior-based approach for hypothesis comparison and its utility to**
2 **discern among temporal scenarios of divergence.**

3 **Eugenia Zarza^{1,2*}, Robert B. O’Hara³, Annette Klussmann-Kolb¹ and Markus**
4 **Pfenninger.^{1,4}**

5 ¹ *Senckenberg Biodiversity and Climate Research Centre (BiK-F), Senckenberganlage 25, D-*
6 *60325 Frankfurt am Main, Germany*

7 ² *Colección Nacional de Anfibios y Reptiles, Instituto de Biología, Universidad Nacional*
8 *Autónoma de México*

9 ³ *Department of Mathematical Sciences, Norwegian University of Science and Technology,*
10 *Sentralbygg 2, Gløshaugen*

11 ⁴ *Department of Phylogeny and Systematics, Institute for Ecology, Evolution and Diversity,*
12 *Biosciences, Goethe-University Frankfurt, Max-von-Laue Straße 13, 60438 Frankfurt am*
13 *Main, Germany*

14

15

16 * Corresponding author. E-mail: eugenia.zarza@gmail.com. Current address: Colección

17 Nacional de Anfibios y Reptiles, Departamento de Zoología, Instituto de Biología, Universidad

18 Nacional Autónoma de México, 3er. Circuito Exterior, Ciudad Universitaria, 04510 Ciudad de

19 México, Mexico.

20

21

22

23

24 **Abstract**

25 One of the major problems in evolutionary biology is to elucidate the relationships between
26 historical events and the tempo and mode of lineage divergence. The development of relaxed
27 molecular clock models and the increasing availability of DNA sequences resulted in more
28 accurate estimations of taxa divergence times. However, finding the link between competing
29 historical events and divergence is still challenging. Here we investigate assigning constrained-
30 age priors to nodes of interest in a time-calibrated phylogeny as a means of hypothesis
31 comparison. These priors are equivalent to historic scenarios for lineage origin. The hypothesis
32 that best explains the data can be selected by comparing the likelihood values of the competing
33 hypotheses, modelled with different priors. A simulation approach was taken to evaluate the
34 performance of the prior-based method and to compare it with an unconstrained approach. We
35 explored the effect of DNA sequence length and the temporal placement and span of competing
36 hypotheses (i.e. historic scenarios) on selection of the correct hypothesis and the strength of the
37 inference. Competing hypotheses were compared applying a posterior simulation analogue of the
38 Akaike Information Criterion and Bayes factors (obtained after calculation of the marginal
39 likelihood with three estimators: Harmonic Mean, Stepping Stone and Path Sampling). We
40 illustrate the potential application of the prior-based method on an empirical data set to compare
41 competing geological hypotheses explaining the biogeographic patterns in *Pleurodeles* newts.
42 The correct hypothesis was selected on average 89% times. The best performance was observed
43 with DNA sequence length of 3500-10000 bp. The prior-based method is most reliable when the
44 hypotheses compared are not temporally too close. The strongest inferences were obtained when
45 using the Stepping Stone and Path Sampling estimators. The prior-based approach proved
46 effective in discriminating between competing hypotheses when used on empirical data. The

47 unconstrained analyses performed well but it probably requires additional computational effort.
48 Researchers applying this approach should rely only on inferences with moderate to strong
49 support. The prior-based approach could be applied on biogeographical and phylogeographical
50 studies where robust methods for historical inferences are still lacking.

51

52 **Introduction**

53

54 One of the major problems in evolutionary biology is to elucidate the relationships between
55 historical events and the tempo and mode of lineage divergence and, ultimately, biological
56 diversification. The development of methods to estimate substitution rates with relaxed
57 molecular clock models and the increasing availability of DNA sequences has led to better
58 estimates of species and higher taxa divergence times (Battistuzzi et al., 2010). However, finding
59 the link between historical events, such as past geological and climatic changes, and divergence
60 is still challenging. As phylogeography – and other evolutionary biology disciplines - move away
61 from narrative and traditional null-hypothesis methods towards multiple hypothesis comparison
62 approaches (Johnson & Omland, 2004; Dépraz et al., 2008; Bloomquist, Lemey & Suchard,
63 2010; Carstens et al., 2013), it is necessary to investigate if a hypothesis comparative framework
64 can also be applied at deeper evolutionary times.

65 Hypothesis comparison offers a means to draw inferences from a set of multiple
66 competing hypotheses and to estimate the degree of confidence that can be placed on each of
67 them (Dépraz et al., 2008; Johnson & Crandall, 2009). Competing hypotheses should be
68 thoroughly thought through and formulated as a first step in the research process (Anderson,
69 2008). After data collection and analyses, the competing hypotheses can be compared and ranked
70 to select which of them best explains the data. This can be accomplished using the Bayes factor
71 (BF), the ratio of the marginal likelihood of the data from two models, i.e. the posterior
72 probability of one model to that of another model, divided by the ratio of the prior probabilities,
73 thus BF measures the change in support for one model versus another given the data (Jeffreys,
74 1935; Kass & Raftery, 1995; Suchard, Weiss & Sinsheimer, 2001, 2003).

75 Here we propose a hypothesis comparison approach to evaluate the influence of historic
76 events in lineage divergence. Our main aim is to explore if assigning constrained-age priors to
77 nodes of interest in a time-calibrated phylogeny would serve as a means for hypothesis
78 comparison. These priors would be equivalent to scenarios for lineage divergence under certain
79 competing hypotheses. When comparing the likelihood values of such hypotheses, modelled
80 under different priors, we would be able to select the hypothesis that best explains the data and
81 assign a level of confidence to evolutionary inferences. This hypothesis comparison approach
82 has been employed a few times to empirical data with success to discern among competing
83 temporal biogeographical scenarios in crabs (Klaus et al., 2010; Jesse et al., 2011), and land
84 snails (Pfenninger et al., 2010). However, the efficiency, accuracy and range of validity of the
85 approach have as yet not been rigorously tested in a systematic manner.

86

87 METHODS FOR MODEL SELECTION: BAYES FACTORS AND AKAIKE'S 88 INFORMATION CRITERION.

89

90 Bayes factors allow for hypothesis ranking and evaluation of the relative merits of the competing
91 hypotheses (Jeffreys, 1935; Kass & Raftery, 1995; Baele et al., 2012), placing BF at the core of
92 Bayesian theory of hypothesis (Robert & Wraith, 2009). When using BF, the model, or in this
93 case hypothesis, with the greatest marginal likelihood (for simplicity MLL) is generally
94 preferred. The marginal likelihood is a weighted average of the likelihood, where the weights
95 come from the prior (Xie *et al.* 2011). In a phylogenetic context where the parameter space is
96 very large, calculating MLL, requires integrating over all possible solutions and is not
97 analytically feasible.

98 Until recently, importance sampling approaches were used to calculate the harmonic
99 mean estimator (HM) of MLL (Newton & Raftery, 1994), despite the short-comings of the
100 approach being outlined in the original paper. HM only needs simulations from the posterior
101 distributions and can be easily calculated from an MCMC sample. Consequently, it has been
102 widely used in phylogenetics (e.g. MrBayes and implemented in BEAST). However, HM is not
103 stable and can have infinite variance giving unreliable results for model selection (Lartillot &
104 Philippe, 2006; Xie et al., 2011). Recent developments aim to improve the exploration of the
105 relevant model space via guided transitions across a sequence of intermediate distributions
106 connecting their prior and posterior extremes (Cameron & Pettitt, 2013). Among these methods
107 are thermodynamic integration (Lartillot & Philippe, 2006), also known as path sampling (PS;
108 Ogata, 1989; Gelman & Meng, 1998) and the Stepping Stone method (SS; Xie et al., 2011). Both
109 methods have been implemented in BEAST latest version (from version 1.7.0; Drummond et al.,
110 2012; Baele et al., 2012), together with a posterior simulation analogue of the Akaike's
111 information criterion through MCMC (AICM; Raftery et al., 2007; Baele et al., 2012), forming a
112 useful set of tools for model selection in phylogenetics.

113 Here, using a simulation approach, we evaluate the plausibility of using prior information
114 to compare hypotheses on divergence times between lineages. We apply several model selection
115 techniques (AICM, HM, SS and PS) and evaluate their performance for prior-based hypothesis
116 comparison under several conditions. We varied data amount, relative temporal placement, span
117 and absolute tree location of hypotheses (age priors), but kept the evolutionary and relaxed clock
118 models constant. Using a reduced set of simulations, we compared the prior-based approach with
119 a simpler approach to select among competing scenarios. This consists in executing one analysis
120 to compute the proportion of sampled MCMC steps that fall within date intervals compatible

121 with competing historical events or scenarios. In this way, for each scenario, it would be possible
122 to estimate the posterior probability that a given divergence occurred at the same time as the
123 historical event. This would enable comparing several competing hypotheses through BF. This
124 approach does not require applying constraints on the age prior distribution of the nodes of
125 interest and we refer to it as the “unconstrained” analysis (Uc).

126 We illustrate the potential application of the prior-based hypothesis approach on an
127 empirical data set to compare competing geological hypotheses explaining the biogeographic
128 patterns of *Pleurodeles* newts in Iberia and Northern Africa (Zhang et al., 2008).

129

130

131 **Materials and Methods**

132

133 A simulation study was carried out to evaluate the performance of the prior-based hypothesis
134 comparison approach, to investigate what factors could lead to a reliable selection of the correct
135 hypothesis or scenario and to compare it with Uc. Three sets of simulations were generated. The
136 main difference among them is the age of the simulated correct hypothesis: “deep” correct
137 hypothesis (DCH, 4.2-4.7. Ma), “intermediate” correct hypothesis (ICH, 2.7-3.2 Ma) or
138 “shallow” correct hypothesis (SCH, 1.2-1.7 Ma). The general simulation procedure included
139 several stages. The first step was to simulate trees with 25 taxa with BEAST v1.6.1 (Drummond
140 et al., 2006) in which five nodes were age-constrained to the same time interval (e.g. 2.7-3.2 Ma,
141 the onset of the Northern Hemisphere Glaciation NHG). We constrained this high number of
142 nodes to facilitate divergence time estimation and reflect a scenario where many nodes in the tree
143 were affected by a very significant event. Nodes were defined as two-taxa set, with a uniform-

144 age prior reflecting the age of the event. Trees were built following a birth-death tree prior and a
145 normal prior with 15 Ma mean for the root age, without sequence data. To test performance
146 under a variety of tree shapes, trees were sampled at a frequency resulting in up to 100 final trees
147 from which 20 were randomly selected, using the random.org number generator. All the nodes in
148 the trees were resolved and the topologies and position of the nodes in the tree are shown in File
149 S1. Input files to generate the simulated trees are in File S2.

150 The topologies of the selected trees were used to simulate DNA sequence data with Seq-
151 Gen (Rambaut & Grassly, 1997). Five partitioned DNA-sequences datasets of 3500, 10000 and
152 20000 bp were simulated for each topology under the Jukes-Cantor substitution model. To reflect
153 partition rate heterogeneity, we specified a relative rate of evolution for each partition. As
154 required by SeqGen, the relative rates had a mean of 1.0, but without variation in the substitution
155 rate among taxa. Number of partitions per data set is shown in Table S1. Each data set was used
156 to generate input files for BEAST. The age priors for the nodes of interest in these input files
157 reflect the “correct” hypothesis (i.e. DCH, ICH or SCH), which has the same age priors as those
158 used to generate the simulated topology. The sequence data was also used to create input files for
159 BEAST with age priors reflecting the competing hypotheses described in the following sections
160 (supplementary Table S2, Fig. 1 and Fig. 2). Input files reflecting the correct and competing
161 hypotheses had additional time calibrations on one or two nodes and the tree root. These nodes’
162 age prior follows a normal distribution, whose mean corresponds to the age of that node in the
163 initial simulated topology. In a similar way, input files were created for the unconstrained
164 analyses and are included in File S2 to facilitate analysis replication.

165 The analyses were run under an uncorrelated relaxed molecular clock (UMC). Although
166 the sequences were generated without variation in the substitution rate among taxa, it has been

167 suggested that UMC reliably estimates parameters even when the data follows a strict molecular
168 clock, which is indeed a model comprised within the more complex UMC (Drummond et al
169 2006). Thus we do not consider that this will be detrimental for this study. As the substitution
170 rate is unknown for most no-model organisms, we consider that it will be more informative to
171 estimate this parameter from the data.

172

173 EFFECT OF SEQUENCE LENGTH, HYPOTHESIS RELATIVE TEMPORAL POSITION
174 AND HYPOTHESIS TEMPORAL SPAN.

175

176 In this set of simulations, the correct hypothesis age was fixed to the intermediate time depth (i.e.
177 ICH, 2.7-3.2) whereas sequence length and position and temporal span of competing hypotheses
178 varied. Three historical scenarios were compared: 1) nodes split at the time of a
179 geological/climatic event: the time of the NHG at 2.7-3.2 Ma, and is considered as the correct
180 hypothesis ICH; 2) split occurred before the geological/climatic event; 3) split occurred after the
181 geological/climatic event. These scenarios reflect the situation of a researcher who suspects that
182 a climatic/geological event might have led to a node split in a phylogeny, but would like to know
183 how much better (or worse) the hypothesis explains the data in comparison to the other
184 scenarios.

185 To test sequence length effect of on hypothesis selection, we simulated data sets with
186 3500, 10000 and 20000 bp and compared ICH to competing scenarios where nodes split before
187 or after ICH. To explore how temporally close the competing hypotheses and ICH can be to
188 properly distinguish and select ICH, we used a more or less intermediate data set size (i.e. 10000
189 bp) and varied the temporal location of the competing hypotheses one or two intervals before and

190 after ICH, thus competing hypotheses did not overlap. An interval is defined as equal to the
191 temporal span of ICH: 0.5 million years (Myr). Another important factor to consider is the
192 competing hypotheses duration, specifically is it valid to compare hypotheses with different
193 widths of prior age distributions? To answer this question, we simulated competing hypotheses
194 where the age priors of the nodes of interest were two times wider than, equal to, or half as wide
195 as ICH, and were temporally located before and after ICH.

196 One hundred replicate input files were generated for each type of competing hypothesis
197 (temporal or duration variation), following the general procedure above described. Input files
198 were run in BEASTv1.7.1. MCMC length is shown in supplementary Table S1. The comparable
199 competing hypotheses were run with the same number of iterations.

200

201 EFFECT OF ABSOLUTE AGE OF CORRECT HYPOTHESIS

202

203 To test how correct hypothesis absolute age (i.e. temporal depth) affects hypothesis comparison
204 and selection, we followed the general simulation procedure. Tree topologies were simulated
205 where five nodes of interest were constrained with age priors reflecting DCH or SCH. For each
206 situation, the 20 randomly chosen trees were used to generate DNA-sequence data sets of 10000
207 bp. DCH was compared to more recent competing hypotheses: 2.7-3.2 Ma and 1.2-1.7 Ma;
208 whereas SCH was compared to older competing hypothesis: 2.7-3.2 Ma and 4.2-4.7 Ma.

209 Performance with these two variations of correct hypotheses was compared to performance with
210 ICH. We removed runs that did not converge to keep the run length equal among simulations
211 with similar data set size.

212

213 HYPOTHESIS SELECTION

214

215 The marginal likelihood under the different priors was estimated using the HM, PS and SS
216 methods. The natural logs of the Bayes factors were calculated as $\ln(\text{BF})=H_i-H_j$, where H_i and H_j
217 are the log natural of the competing hypothesis MLL, following the method first implemented in
218 Tracer (Suchard, Weiss & Sinsheimer, 2003; Rambaut & Drummond, 2007) to calculate BF
219 based on HM. The strength of evidence was evaluated according to the table provided by Kass
220 and Raftery (1995) but without multiplying by 2 and without rounding up $\ln(\text{BF})$ values). Thus,
221 $\ln(\text{BF}) < 1.10$ means weak support for H_i over H_j , $1.10 < \ln(\text{BF}) < 2.30$ mean moderate support
222 and a $\ln(\text{BF}) > 2.3$ was considered as strong support ($\text{BF} > 10$). Regarding selection with AICM,
223 a $\Delta \text{AICM} > 10$ between the best ranked hypothesis and the other hypotheses suggests that the
224 latter were very unlikely (Burnham & Anderson, 2002). These calculations were performed in
225 BEASTv1.7.1 with the code of Baele *et al.* (2012). It is expected that the “correct” hypothesis
226 will have higher MLL values than the others if our method is effective.

227

228 UNCONSTRAINED ANALYSES

229

230 The frequency of the MCMC steps falling within each of the correct and competing hypotheses
231 time intervals was calculated to estimate the posterior probability of each hypothesis. We
232 calculated the prior probability of a hypothesis as its interval length divided by the total interval
233 length (i.e. the time from its most recent calibrated ancestral node to the present). BFs were
234 calculated as the ratio of posterior odds to prior odds between the correct hypothesis and a
235 particular competing hypothesis. This was obtained for each of the five nodes, for each
236 competing hypothesis only for the treatments comparing against ICH, and with data sets of 3500,

237 10000 and 20000 bp. To make Uc results comparable with the prior-based approach, we
238 estimated inference strength with this scale: $BF < 1$ false positive; $1 < BF < 3.01$ weak; $3.01 < BF <$
239 10 moderate; $BF > 10$ strong. The frequency of strong, moderate and weak BF per node was
240 calculated. We calculated an average frequency of strong, moderate and weak BF for the five
241 nodes, for each treatment.

242

243 EMPIRICAL DATASET ANALYSIS: SALAMANDERS

244

245 In this section we apply the prior-based approach to compare hypotheses on the time of split
246 between two species of newts and the influence of geological and climatic events. Zhang *et al*
247 (2008) proposed a time-calibrated phylogeny of the family Salamandridae inferred from
248 mitochondrial genomes (10755 bp). The data set comprises 41 taxa, including representatives of
249 all recognized genera. The authors calibrated six nodes with fossil records and one using indirect
250 geological evidence. Based on the results of Bayesian and penalized likelihood analyses the
251 authors proposed a robust time-calibrated phylogeny and postulated several biogeographic
252 hypotheses to account for the distribution patterns between taxa in Salamandridae. We re-
253 analysed their data set to compare three previously suggested competing scenarios to explain the
254 phylogeographic patterns observed in one of the clades, the ribbed newts (*Pleurodeles*), currently
255 distributed in Iberia and Northern Africa (Frost, 2011). According to Veith *et al* (2004) and
256 Zhang *et al.* (2008): 1) The split between *P. waltl* and *P. poireti* could be consistent with the
257 Messinian salinity crisis (ca. 5.33 Ma); or 2) The Betic crisis ca. 14 Ma; or 3) the Betic crisis
258 leading to the split between the north-western and south-eastern populations of *P. waltl*, rather
259 than between the two *Pleurodeles* species, which would imply that the two species split around

260 35 Ma.

261 We used the BEAST input file of Zhang *et al.* (2008) keeping the original fossil
262 calibration points but assigning proper priors to all parameters (Baele *et al.* 2013). In three
263 independent analyses, age priors were added to reflect the competing scenarios. Analysis 1
264 included the original calibration points plus a normal age prior for the most recent common
265 ancestor of *Pleurodeles* species (from now on referred to as Node P) with mean 5.33 Ma,
266 reflecting scenario 1. In analysis 2, in addition to the original calibration points, a normal age
267 prior was assigned to Node P with mean 14.0 Ma, reflecting scenario 2. In analysis 3, the
268 original calibration points plus a normal age prior with mean 35 Ma, reflecting scenario 3, were
269 included. To obtain adequate effective sample sizes of the parameters, five independent runs
270 with 100 million MCMC iterations were executed in BEASTv1.7.1. After MCMC execution,
271 samples of the prior and posterior were collected for later estimation of MLL with HM, PS and
272 SS, following suggestions on the BEAST website (beast.bio.ed.ac.uk/Model_selection). Log files
273 of the five independent runs were combined with LogCombiner of the BEAST package after
274 removing 10% of the samples as burnin. The combined log files were used to calculate the
275 AICM and estimate MLL using HM. PS/SS analyses were executed combining the samples of
276 power posteriors collected at the end of each MCMC. Competing scenario MLLs were then
277 calculated to select the one that best explains the data.

278

279 **Results**

280

281 One hundred simulated replicate datasets were analysed per “treatment”. However, convergence
282 of the MCMC runs for the alternative, the correct hypotheses or/and the unconstrained analyses

283 was not always achieved and acceptable effective sample sizes were not obtained. In the prior-
284 based approach, runs that failed to converge and their competing hypotheses (correct or
285 alternative hypotheses)-even if these converged- were not taken into account to calculate the
286 effectiveness of the method. An improvement of up to 5% in the frequency of success was
287 observed when ignoring the runs lacking convergence in comparison to keeping all runs
288 irrespective of convergence achievement. The PS and SS methods produced similar results under
289 all the simulations strategies, thus only one graph is shown.

290 The unconstrained analyses consisted on executing one run to compare the frequency of
291 MCMC steps falling within the intervals of several competing hypothesis. The Uc analyses were
292 run for the same number of MCMC iterations as the prior-based approach. However with data
293 sets of 3500, 10000 and 20000bp, 7%, 39% and 49% of the Uc runs did not reach convergence,
294 respectively; whereas in average 0%, 8.7% and 18% of the respective competing hypothesis runs
295 in the prior-based approach did not converge.

296

297 EFFECT OF SEQUENCE LENGTH.

298

299 Sequence length was increased from 3500 bp up to 20000 bp as shown in supplementary Table
300 S1. With the prior-based approach, all the MCMC runs analysing 3500 bp data sets achieved
301 convergence. Runs of the correct hypothesis and its competing hypotheses reached convergence
302 78% and 59% with 10000 bp and 20000 bp data sets, respectively. Increasing sequence length
303 leads to an increase in the frequency of selecting the correct hypothesis as the best hypothesis
304 with strong support when using AICM and HM (Fig. 3). However an improvement is not seen
305 when calculating MLL with PS/SS with data sets larger than 10000 bp (Fig. 3). False positives

306 frequency decreases with sequence length from 3500 bp to 10000 bp with all methods (HM:
307 from 12.5 to 4.6 %; AICM: from 9.5% to 3.9%; SS/PP from 6.5% to 5.9%). Only with AICM
308 can a reduction in false positives be seen with 20000 bp data sets (3.3%). Nevertheless, strong
309 inferences frequency is always higher when using PS/SS than with HM, and AICM (Fig. 3). Uc
310 shows better performance than HM and AICM with 10000 and 20000 bp data sets, but performs
311 poorly with small data sets.

312

313 EFFECT OF TEMPORAL SPAN OF COMPETING HYPOTHESES.

314

315 Different sizes for the temporal constraint interval of the competing hypotheses were compared.
316 Regarding the prior based approach, convergence was achieved by 78% of the correct hypothesis
317 and its competing hypotheses MCMC runs. With the AICM calculation, the correct hypothesis
318 was selected above 96% of the times, with no strong false positives. HM performs with a similar
319 rate of success, however the correct hypothesis is selected with strong support more often than
320 with AICM with only 0.65% of strong false positives. In both cases a better performance was
321 obtained when the hypotheses span intervals of similar size or when the competing hypothesis
322 has a narrower temporal range. PS/SS select the correct hypothesis strongly more frequently than
323 the other two methods (Fig. 4). The correct hypothesis was strongly supported slightly more
324 often (90%) when the competing hypotheses had narrower intervals than when the competing
325 hypotheses had an interval as wide as the correct hypothesis (87%; Fig 4). Strong false positives
326 were obtained at a frequency between 2.6 to 3.9%. It should be noted that the AICM does not
327 estimate MLL and thus the results are not entirely comparable. Uc performed better than AICM
328 and HM with all interval sizes and was slightly outperformed by PS/SS.

329

330 EFFECT OF RELATIVE TEMPORAL LOCATION OF HYPOTHESIS.

331

332 Temporal location of competing hypotheses was also varied. In the prior-based approach,

333 convergence was achieved in 76% of the correct hypothesis and its competing hypotheses runs.

334 Our simulations suggest that the closer the competing hypothesis is to the correct hypothesis the

335 less likely it will be to rank the correct hypothesis as the best hypothesis (Fig. 5). A trend

336 towards increase in selection accuracy with increase in temporal distance between hypotheses

337 was observed with all methods. BFs calculated with PS/SS select the correct hypothesis with

338 strong support more often than HM when the hypotheses are the furthest apart (92% and 86%

339 respectively). PS/SS produce stronger inferences than HM when the hypotheses are the closest,

340 although the performance is poor (<50%). Selection of the correct hypothesis with AICM with

341 moderate to strong support occurs above 78% of the times when hypotheses are the furthest

342 apart. High frequency (19%) of false positives was observed when applying HM and hypotheses

343 were very close together, but they are reduced when the hypotheses are further apart (1.9%).

344 False positives frequency obtained with AICM is reduced from 18% to 2.6 % when hypotheses

345 are the furthest away. PS/SS produce the highest frequency of false positives when the

346 hypotheses are close (8.3%), but this is reduced when the hypotheses are temporally apart

347 (0.64%). Similarly, with Uc it is difficult to select among closely located hypotheses.

348

349 EFFECT OF ABSOLUTE AGE OF THE CORRECT HYPOTHESIS.

350

351 To investigate the effect of the absolute age of correct hypothesis in the tree, two sets of

352 simulations were carried out. The first simulated SCH and was compared with less recent
353 hypotheses. In this case, 90% of the MCMC runs achieved convergence. All four methods
354 selected SCH as the correct hypothesis 100% of the times with strong support (Fig. 6 A). Strong
355 false positives occurred only in one case when using PS/SS and HM. When the correct
356 hypothesis was ICH, AICM and HM tended to perform better when the competing hypothesis is
357 deeper than ICH, with a higher frequency of strong inferences. PS/SS led to stronger inferences
358 over more recent hypotheses (Fig. 6 B). No strong false positives were obtained except for one
359 case when using PS/SS. Only 15% of the runs reached convergence in simulations with DCH.
360 Among these runs, PS/SS performed better than the other two methods selecting the DCH above
361 93% of the times with strong support, followed by AICM and HM (Fig. 6 C).

362

363 HYPOTHESIS COMPARISON USING EMPIRICAL DATA

364

365 After combining the MCMC outputs, effective sample sizes above 100 were obtained for all
366 parameters. The three independent competing analyses resulted in the same topology obtained by
367 Zhang *et al.* (2008). Evidence is stronger for Scenario 2 when BFs are estimated based on MLL
368 calculated with PS, SS and HM methods (Table 1). However, AICM ranks Scenario 3 as the best
369 hypothesis. The Bayes factors calculated with PS and SS estimates are larger than those obtained
370 with HM. Δ AICM moderately supports Scenario 3 over the other competing hypotheses. The
371 results from PS, SS and HM are in agreement with results previously obtained with molecular
372 and fossil evidence, suggesting that the split between these species of ribbed newts is associated
373 with the Betic crisis (Zhang *et al.*, 2008).

374

375 **Discussion**

376

377 We evaluated the performance of a hypothesis comparison approach that uses prior information
378 to define competing scenarios of lineage divergence, in which divergence is associated with
379 historic events like climate or geological change. After calculation of their marginal likelihood or
380 AICM, it is possible to rank scenarios and select the one that best explains the data. Our
381 simulation study suggests that under reasonable circumstances, this approach could constitute a
382 reliable tool to compare temporal scenarios: the correct hypothesis is ranked as the best
383 hypothesis over 80% of the time under almost all simulation strategies. However, inference
384 strength varies depending on the method employed to calculate BF or if AICM is used. Most of
385 the times HM ranks the correct hypothesis as the best hypothesis but the BFs are so low that it is
386 difficult to place any confidence in the selection. Generally, PS and SS estimates of MLL differ
387 more strongly between competing hypotheses than HM. We observed that these methods could
388 also lead to few false positives with strong or moderate support. This may, in part, be because the
389 data genuinely support the wrong hypothesis by chance (e.g Kuparinen et al., 2007).

390 Discerning between competing hypotheses is particularly challenging when the
391 hypotheses are located close to each other in time. Interestingly, it was consistently difficult to
392 reach convergence when the node of the correct hypothesis was located deeper in the tree
393 (DCH), especially for runs where the alternative hypotheses were the furthest away from the
394 correct hypothesis. The accuracy and strength of ranking the correct hypothesis as the best
395 hypothesis increase slightly with the amount of data with the AICM and HM methods. However
396 contrary to expectations PS/SS showed a decrease in performance with 20000 bp data sets. There
397 are several factors that can influence this behaviour, for example the path sampling chain length

398 between the prior and the posterior, the number of sample steps and other PS/SS parameters that
399 would need to be adjusted to a particular data set size. PS/SS are relatively new methods in
400 phylogenetics and so far there are only a few studies investigating the influence of these
401 parameters, generally dealing with smaller data sets and number of topologies (Lartillot &
402 Philippe, 2006; Xie et al., 2011; Baele et al., 2013). The computational demand to investigate the
403 possible causes of this behaviour is high and at the moment goes beyond the scope of this study.
404 However, further research is needed especially as the genomic area will allow for the analysis of
405 increasingly larger DNA sequence data sets.

406 We did not test how consistent MLL and AICM estimations are among independent
407 MCMC runs. However Beale *et al* (2012) found that PS and SS produce consistent estimates
408 among MCMC runs more often than the other methods. Thus, considering our results in light of
409 previous studies (Lartillot & Philippe, 2006; Xie et al., 2011; Baele et al., 2012), we suggest that
410 applying PS and SS would produce more reliable results than HM and AICM. However,
411 independent of the method of hypothesis comparison used, it is always advisable to rely only on
412 inferences with moderate to strong support.

413 The prior-based approach proved effective in discriminating between competing
414 hypotheses when applied to empirical data (data set by Zhang et al., 2008). The hypotheses
415 compared reflected scenarios well apart in time and relied on a relatively large data set and a
416 robust phylogeny. Researchers applying this approach should meet these conditions because
417 divergence time and tree topology are estimated at the same time with BEAST, thus changes in
418 topology affect divergence times and vice versa (Heled & Drummond, 2011). Furthermore, as
419 recently demonstrated, the effect of the rate priors could also affect the estimation of divergence
420 times and should be investigated in future studies (Reis, Zhu & Yang, 2014).

421 In the simulated phylogenies we used a relatively high ratio of constrained/no-
422 constrained nodes (five nodes per hypothesis comparison, plus up to three additional calibrated
423 nodes out of 24; see input files in S2). It will be necessary to investigate if reducing the number
424 of constrained nodes could lead to a decrease in the strength of inferences, and if an increase will
425 improve the accuracy of the divergence estimation and thus benefit hypothesis selection. We
426 already observed that constraining 7/40 nodes in the empirical data set analyses led to
427 discrepancies among hypothesis selection methods. This additionally suggests that the direct
428 comparison between these simulated and empirical data analyses should be taken with caution.

429 We executed unconstrained analyses that need to be run only one time to compare several
430 hypotheses simultaneously. Most of the times, Uc was slightly outperformed by PS/SS. However
431 the unconstrained MCMC runs reached convergence less often than the prior-based approach
432 runs. Thus, there might not be a computational benefit in running one very long MCMC instead
433 of several shorter parallel runs reflecting competing hypotheses. Another potential problem with
434 just running a single run and counting the visits to each hypothesis, as we did in the
435 unconstrained analyses, is that if the hypotheses are really disjoint, it will be necessary to throw
436 away MCMC iterations for the times outside the hypotheses. If the hypotheses were overlapping
437 it would be necessary to correct for this when estimating a time that could belong to different
438 hypotheses which is an extra challenge.

439 The development of new methods for model selection, and future research on their
440 performance, will add confidence to inferences led by hypothesis comparison. This could have
441 implications for biogeographical and phylogeographical studies where robust methods for
442 historical inferences are still lacking. Depending on the location of the nodes of interest, the
443 approach here evaluated could also be applied in cases where not only divergence between two

444 taxa, but instead a diversification event is suspected. At this scale, it could complement the
445 traditional method of testing the hypothesis of shifts or heterogeneity in diversification rates
446 against the null hypothesis of constant rates through time and among lineages (Pybus & Harvey,
447 2000; Chan & Moore, 2002; Ricklefs, 2007; Moore & Donoghue, 2009; Steeman et al., 2009;
448 Silvestro, Schnitzler & Zizka, 2011). It would also allow testing the association of such shifts
449 with climate or geological change (Hines, 2008; Schuettelpelz & Pryer, 2009).

450

451 **Acknowledgements**

452

453 Thanks to P. Zhang for kindly providing the Salamandridae data and the input file for BEAST.
454 Thanks to M. Forrest for help with scripting and facilitating access to the Frankfurt Cloud
455 (Goethe University-Deutsche Bank), and for his and C. Weiland's assistance to access and use
456 the BiK-F computer cluster. EZ thanks P. Jansson for providing comments.

457

458 **References**

459

460 Anderson DR. 2008. *Model Based Inference in the Life Sciences - A Primer on Evidence*. New
461 York, USA: Springer.

462 Baele G., Lemey P., Bedford T., Rambaut A., Suchard MA., Alekseyenko AV. 2012. Improving
463 the accuracy of demographic and molecular clock model comparison while
464 accommodating phylogenetic uncertainty. *Molecular Biology and Evolution* 29:2157–
465 2167.

466 Baele G., Li WLS., Drummond AJ., Suchard MA., Lemey P. 2013. Accurate model selection of

- 467 relaxed molecular clocks in bayesian phylogenetics. *Molecular biology and evolution*
468 30:239–243.
- 469 Battistuzzi FU., Filipinski A., Hedges SB., Kumar S. 2010. Performance of Relaxed-Clock
470 Methods in Estimating Evolutionary Divergence Times and Their Credibility Intervals.
471 *Molecular Biology and Evolution* 27:1289–1300.
- 472 Bloomquist EW., Lemey P., Suchard MA. 2010. Three roads diverged? Routes to
473 phylogeographic inference. *Trends in Ecology & Evolution* 25:626–632.
- 474 Burnham KP., Anderson DR. 2002. *Model Selection and Multi-Model Inference*. New York:
475 Springer-Verlag.
- 476 Cameron E., Pettitt A. 2013. *Recursive Pathways to Marginal Likelihood Estimation with Prior-*
477 *Sensitivity Analysis*.
- 478 Carstens BC., Brennan RS., Chua V., Duffie CV., Harvey MG., Koch RA., McMahan CD.,
479 Nelson BJ., Newman CE., Satler JD., Seeholzer G., Posbic K., Tank DC., Sullivan J.
480 2013. Model selection as a tool for phylogeographic inference: an example from the
481 willow *Salix melanopsis*. *Molecular Ecology* 22:4014–4028.
- 482 Chan KMA., Moore BR. 2002. Whole-tree methods for detecting differential diversification
483 rates. *Systematic Biology* 51:855–865.
- 484 Dépraz A., Cordellier M., Hausser J., Pfenninger M. 2008. Postglacial recolonization at a snail's
485 pace (*Trochulus villosus*): confronting competing refugia hypotheses using model
486 selection. *Molecular Ecology* 17:2449–2462.
- 487 Drummond AJ., Ho SYW., Phillips MJ., Rambaut A. 2006. Relaxed phylogenetics and dating
488 with confidence. *PLoS Biology* 4:e88.
- 489 Drummond AJ., Suchard MA., Xie D., Rambaut A. 2012. Bayesian phylogenetics with BEAUti

- 490 and the BEAST 1.7. *Molecular Biology and Evolution* 29:1969–1973.
- 491 Frost DR. 2011. Amphibian Species of the World: an Online Reference. Available at
492 <http://research.amnh.org/vz/herpetology/amphibia> (accessed November 28, 2012).
- 493 Gelman A., Meng X-L. 1998. Simulating normalizing constants: from importance sampling to
494 bridge sampling to path sampling. *Statistical Science* 13:163–185.
- 495 Heled J., Drummond AJ. 2011. Calibrated Tree Priors for Relaxed Phylogenetics and Divergence
496 Time Estimation. *Systematic Biology* 61:138–149.
- 497 Hines HM. 2008. Historical biogeography, divergence times, and diversification patterns of
498 Bumble Bees (Hymenoptera: Apidae: *Bombus*). *Systematic Biology* 57:58–75.
- 499 Jeffreys SH. 1935. Some tests of significance, treated by the theory of probability. *Proceedings*
500 *of the Cambridge Philosophy Society* 31:203–222.
- 501 Jesse R., Grudinski M., Klaus S., Streit B., Pfenninger M. 2011. Evolution of freshwater crab
502 diversity in the Aegean region (Crustacea: Brachyura: Potamidae). *Molecular*
503 *phylogenetics and evolution* 59:23–33.
- 504 Johnson JB., Crandall KA. 2009. Expanding the toolbox for phylogeographic analysis.
505 *Molecular Ecology* 18:4137–4139.
- 506 Johnson JB., Omland KS. 2004. Model selection in ecology and evolution. *Trends in Ecology &*
507 *Evolution* 19:101–108.
- 508 Kass RE., Raftery AE. 1995. Bayes Factors. *Journal of the American Statistical Association*
509 90:773–795.
- 510 Klaus S., Schubart CD., Streit B., Pfenninger M. 2010. When Indian crabs were not yet Asian -
511 biogeographic evidence for Eocene proximity of India and Southeast Asia. *BMC*
512 *Evolutionary Biology* 10:287.

- 513 Kuparinen A., Snäll T., Vänskä S., O’Hara RB. 2007. The role of model selection in describing
514 stochastic ecological processes. *Oikos* 116:966–974.
- 515 Lartillot N., Philippe H. 2006. Computing Bayes Factors Using Thermodynamic Integration.
516 *Systematic Biology* 55:195–207.
- 517 Moore BR., Donoghue MJ. 2009. A Bayesian approach for evaluating the impact of historical
518 events on rates of diversification. *Proceedings of the National Academy of Sciences*
519 106:4307–4312.
- 520 Newton MA., Raftery AE. 1994. Approximating Bayesian inference with the weighted
521 likelihood bootstrap. *J.R. Stat. Soc. B.* 56:3–48.
- 522 Ogata Y. 1989. A Monte Carlo method for high dimensional integration. *Numerische*
523 *Mathematik* 55:137–157.
- 524 Pfenninger M., Véla E., Jesse R., Elejalde MA., Liberto F., Magnin F., Martínez-Ortí A. 2010.
525 Temporal speciation pattern in the western Mediterranean genus *Tudorella* P. Fischer,
526 1885 (Gastropoda, Pomatiidae) supports the Tyrrhenian vicariance hypothesis. *Molecular*
527 *Phylogenetics and Evolution* 54:427–436.
- 528 Pybus OG., Harvey PH. 2000. Testing macro-evolutionary models using incomplete molecular
529 phylogenies. *Proceedings of the Royal Society B: Biological Sciences* 267:2267–2272.
- 530 Raftery AE., Newton MA., Satagopan JM., Krivitsky PN. 2007. Estimating the integrated
531 likelihood via posterior simulation using the harmonic mean identity. In: *Bayesian*
532 *Statistics*. 1–45.
- 533 Rambaut A., Drummond AJ. 2007. *Tracer 1.5*. <http://beast.bio.ed.ac.uk/Tracer>.
- 534 Rambaut A., Grassly NC. 1997. Seq-Gen: an application for the Monte Carlo simulation of DNA
535 sequence evolution along phylogenetic trees. *Computer applications in the biosciences* :

- 536 *CABIOS* 13:235–238.
- 537 Reis M Dos., Zhu T., Yang Z. 2014. The impact of the rate prior on bayesian estimation of
538 divergence times with multiple Loci. *Systematic biology* 63:555–565.
- 539 Ricklefs RE. 2007. Estimating diversification rates from phylogenetic information. *Trends in*
540 *Ecology & Evolution* 22:601–610.
- 541 Robert CP., Wraith D. 2009. Computational methods for Bayesian model choice. In: *Bayesian*
542 *inference and maximum entropy methods in Science and Engineering*. 251–262.
- 543 Schuettpelez E., Pryer KM. 2009. Evidence for a Cenozoic radiation of ferns in an angiosperm-
544 dominated canopy. *Proceedings of the National Academy of Sciences* 106:11200–11205.
- 545 Silvestro D., Schnitzler J., Zizka G. 2011. A Bayesian framework to estimate diversification
546 rates and their variation through time and space. *BMC Evolutionary Biology* 11:311.
- 547 Steeman ME., Hebsgaard MB., Fordyce RE., Ho SYW., Rabosky DL., Nielsen R., Rahbek C.,
548 Glenner H., Sørensen MV., Willerslev E. 2009. Radiation of Extant Cetaceans Driven by
549 Restructuring of the Oceans. *Systematic Biology* 58:573 –585.
- 550 Suchard MA., Weiss RE., Sinsheimer JS. 2001. Bayesian selection of continuous-time Markov
551 Chain evolutionary models. *Molecular Biology and Evolution* 18:1001–1013.
- 552 Suchard MA., Weiss RE., Sinsheimer JS. 2003. Testing a molecular clock without an outgroup:
553 Derivations of induced priors on branch-length restrictions in a Bayesian framework.
554 *Systematic Biology* 52:48–54.
- 555 Veith M., Mayer C., Samraoui B., Barroso DD., Bogaerts S. 2004. From Europe to Africa and
556 vice versa: evidence for multiple intercontinental dispersal in ribbed salamanders (Genus
557 *Pleurodeles*). *Journal of Biogeography* 31:159–171.
- 558 Xie W., Lewis PO., Fan Y., Kuo L., Chen M-H. 2011. Improving Marginal Likelihood

559 Estimation for Bayesian Phylogenetic Model Selection. *Systematic Biology* 60:150–160.
560 Zhang P., Papenfuss TJ., Wake MH., Qu L., Wake DB. 2008. Phylogeny and biogeography of
561 the family Salamandridae (Amphibia: Caudata) inferred from complete mitochondrial
562 genomes. *Molecular Phylogenetics and Evolution* 49:586–597.
563

564 **Supplementary Information**

565

566 Table S1. Characteristics of simulation treatments.

567 Table S2. Properties of competing hypotheses (see also Figures 1 and 2).

568 File S1. Tree topologies used to generate sequence data.

569 File S2. xml files used as input in simulations. Available from:

570 <https://drive.google.com/open?id=0B7P6iuJv3fpiczBrQ3FDcFRGc1E>

571

572 Table 1. Comparison between hypothesised scenarios for the time of split between *Pleurodeles*,
573 using Bayes Factors calculated based on HM, PS, SS Marginal Likelihood estimates and Δ
574 AICM. A value 0 indicates the best ranked hypothesis.

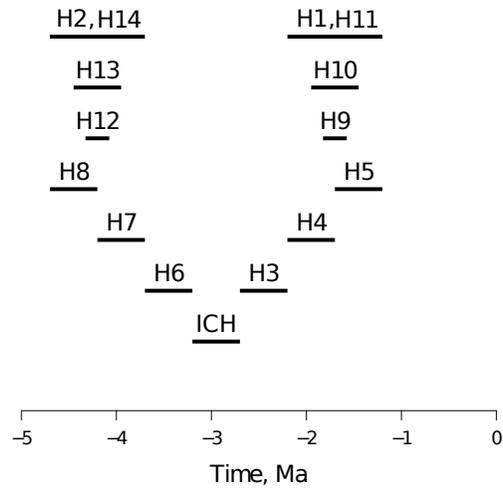
575

576

Scenario	Bayes Factors (PS)	Bayes Factors (SS)	Bayes Factors (HM)	Delta AICM
1	-9.96	-11.06	-0.12	-1.50
2	0	0	0	-2.72
3	-13.8	-15.53	-0.31	0

577

578



579

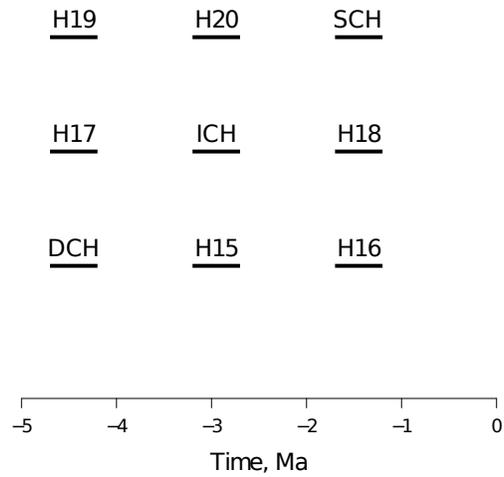
580

581 Figure 1. Competing hypotheses. Lines represent the temporal location and span of competing

582 hypotheses. ICH= correct hypothesis; H1-H14 competing hypotheses (see also Table S2);

583 Ma=million years ago.

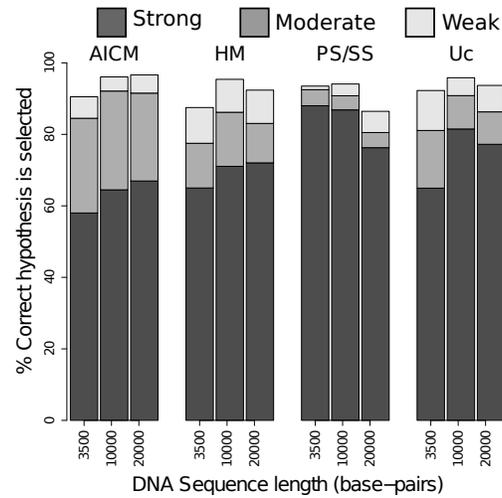
584



585

586 Figure 2. Variations in temporal depth of correct hypotheses. Lines represent the temporal
587 location of the deep (DCH), intermediate (ICH) and shallow (SCH) correct hypotheses, with
588 their respective competing hypotheses shown in the same row.

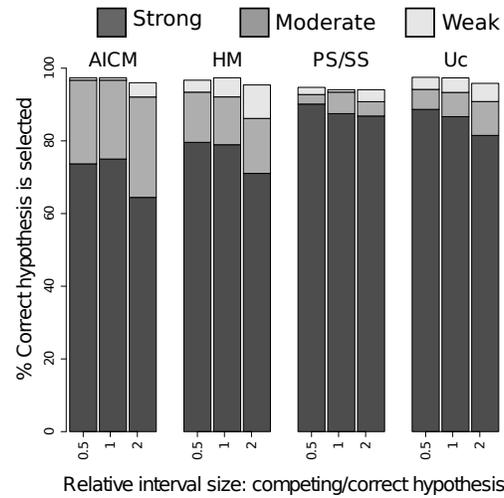
589



590

591 Figure 3. Effect of sequence length on selecting the correct hypothesis. Bars represent the
592 average frequency of ranking ICH as the best hypothesis and strength of inference according to
593 the method employed.

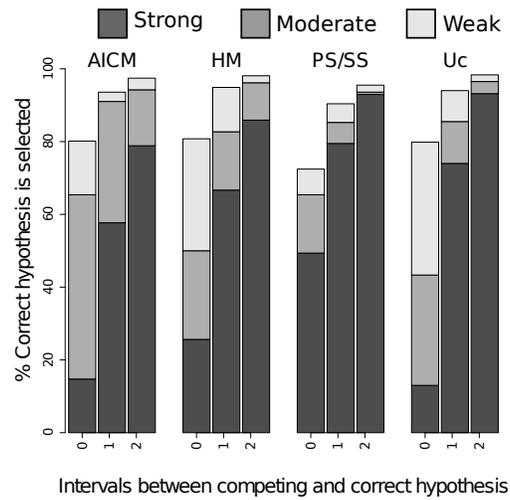
594



595

596 Figure 4. Effect of temporal span (relative interval size) of competing hypotheses on selecting
597 the correct hypothesis. Bars represent the average frequency of ranking ICH as the best
598 hypothesis and strength of inference according to the method employed. Interval=0.5 Million
599 years (Myr).

600

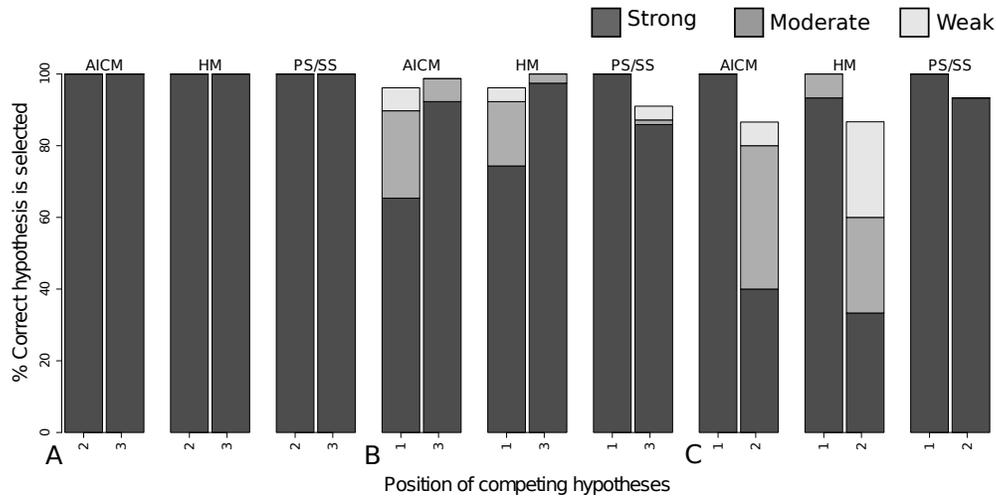


601

602 Figure 5. Effect of temporal location of competing hypothesis on selecting the correct. Bars

603 represent the average frequency of ranking ICH as the best hypothesis and strength of inference

604 according to the method employed Interval=0.5 Myr



605

606 Figure 6. Effect of absolute age (temporal depth) of the correct hypothesis. A) Frequency of
607 selecting the shallow age correct hypothesis (temporal position 1) over deeper competing
608 hypotheses (temporal position 2 and 3). B) Frequency of selecting the intermediate age correct
609 hypothesis (temporal position 2) over a shallower (temporal position 1) and a deeper competing
610 hypothesis (temporal position 3). C) Frequency of selecting the deep age correct hypothesis
611 (temporal position 3) over shallower competing hypotheses (temporal position 1