

1 **Strategies for Partitioning Clock Models in Phylogenomic Dating: Application** 2 **to the Angiosperm Evolutionary Timescale**

3
4 Charles S. P. Foster* and Simon Y. W. Ho
5 School of Life and Environmental Sciences, University of Sydney, Sydney, Australia

6
7 *Corresponding author: E-mail: charles.foster@sydney.edu.au.
8
9

10 **Abstract**

11 Evolutionary timescales can be inferred from molecular sequence data using a Bayesian
12 phylogenetic approach. In these methods, the molecular clock is often calibrated using fossil data.
13 The uncertainty in these fossil calibrations is important because it determines the limiting posterior
14 distribution for divergence-time estimates as the sequence length tends to infinity. Here we
15 investigate how the accuracy and precision of Bayesian divergence-time estimates improve with the
16 increased clock-partitioning of genome-scale data into clock-subsets. We focus on a data set
17 comprising plastome-scale sequences of 52 angiosperm taxa. There was little difference among the
18 Bayesian date estimates whether we chose clock-subsets based on patterns of among-lineage rate
19 heterogeneity or relative rates across genes, or by random assignment. Increasing the degree of
20 clock-partitioning usually led to an improvement in the precision of divergence-time estimates, but
21 this increase was asymptotic to a limit presumably imposed by fossil calibrations. Our clock-
22 partitioning approaches yielded highly precise age estimates for several key nodes in the
23 angiosperm phylogeny. For example, when partitioning the data into 20 clock-subsets based on
24 patterns of among-lineage rate heterogeneity, we inferred crown angiosperms to have arisen 198–
25 178 Ma. This demonstrates that judicious clock-partitioning can improve the precision of molecular
26 dating based on phylogenomic data, but the meaning of this increased precision should be
27 considered critically.

28
29 **Key words:** Angiospermae, molecular dating, phylogenomics, infinite-sites theory, calibration, rate
30 heterogeneity

31 **Introduction**

32 Evolutionary timescales can be estimated from molecular sequence data using phylogenetic
33 methods based on the molecular clock. In practice, most data sets exhibit substantial rate
34 heterogeneity among lineages. These ‘lineage effects’ can be caused by variation in life-history
35 traits, generation time, or exposure to mutagens (Smith and Donoghue 2008; Gaut et al. 2011;
36 Lanfear et al. 2013). Among-lineage rate variation can be taken into account using Bayesian
37 relaxed-clock models, in which the rates can be assumed to be either correlated between
38 neighbouring branches (Thorne et al. 1998; Kishino et al. 2001) or drawn independently from a
39 chosen distribution (Drummond et al. 2006; Rannala and Yang 2007).

40 A number of factors can cause rates to vary across loci in the genome (Wolfe et al. 1987).
41 These ‘gene effects’ can be taken into account by allowing each locus to have a distinct relative
42 rate. Less certain is the best way to deal with interactions between gene effects and lineage effects,
43 which can be caused by differences in selective pressure and other processes (Gaut et al. 2011). In
44 this case, the extent and patterns of among-lineage rate heterogeneity vary across genes or other
45 subsets of the data. This form of rate variation can be captured by assigning separate clock models
46 to different subsets of the data (Ho and Duchêne 2014), a process that we refer to here as clock-
47 partitioning.

48 Appropriate clock-partitioning can improve the precision of Bayesian date estimates (as
49 measured by the associated 95% credibility intervals), but it is rarely done in practice. This is also
50 despite widespread adoption of partitioning schemes for substitution models (Lanfear et al. 2012).
51 The most likely explanation is that the use of clock-partitioning in Bayesian phylogenetics greatly
52 increases the risk of overparameterization, and thus to reduced Markov chain Monte Carlo
53 performance. Overparameterization has been previously addressed in light of the bias-variance
54 trade-off, which is well established in statistical theory (Burnham and Anderson 2003). Compared
55 with a complex, parameter-rich model, a simple model that underfits data is expected to have low
56 accuracy (high bias) but high precision (low variance). Conversely, a parameter-rich model that
57 overfits the data is likely to have higher accuracy, but this comes at the cost of reduced precision.
58 The best model is an intermediate one that simultaneously maximizes accuracy and precision
59 (Wertheim et al. 2010)

60 It is useful to consider the bias-variance trade-off in the context of molecular dating with
61 partitioned clock models. Patterns of among-lineage rate variation are likely to differ across genes
62 (Muse and Gaut 1994), so increasing the number of relaxed clocks will better capture these patterns
63 of rate heterogeneity and should lead to more accurate age estimates (Duchêne and Ho 2014).
64 However, each clock-subset has parameters that need to be estimated, including a distinct set of
65 branch rates. As a consequence, increasing the degree of clock-partitioning should lead to a
66 widening of the posterior distributions of parameters.

67 Contrary to the expectations of the bias-variance trade-off, increasing the degree of clock-
68 partitioning tends improve the precision of Bayesian age estimates (Zhu et al. 2015). One possible
69 explanation for this lies in the treatment of the uncertainty in the estimates of genetic branch
70 lengths. The accuracy and precision of evolutionary rate estimates depend on the accurate inference
71 of branch lengths (in substitutions per site). In the case of molecular dating, branch rates for each
72 clock-subset are combined with node times to give the branch lengths. Therefore, as the number of
73 clock-subsets increases, the node times in the chronogram are estimated from an increasing number
74 of data points, leading to increasing precision. Although branch-length estimation generally
75 improves as the amount of sequence data increases, branch lengths can be estimated with
76 reasonable accuracy even with fairly small amounts of sequence data (Yang and Rannala 2006).
77 This suggests that for a data set of a (large) fixed size, increasing the number of clock-subsets
78 should lead to improved precision in divergence-time estimates until the amount of sequence data in
79 each clock-subset decreases to a critical point.

80 Zhu et al. (2015) explain this phenomenon in their ‘finite sites’ theory, although they use the
81 term ‘loci’ to refer to clock-subsets. Even with sequences of infinite length, there will still be
82 uncertainty in the age estimates, corresponding to the uncertainty in the fossil calibrations (“infinite

83 data limit"; Yang and Rannala 2006; dos Reis and Yang 2013). As the number of clock-subsets (L)
84 increases, the finite-sites theory suggests that the uncertainty in age estimates decreases to the
85 infinite-data limit at the rate of $1/L$ (Zhu et al. 2015). This property has important consequences for
86 analyses of genome-scale data sets, whereby many genes are analysed concurrently. Therefore, it is
87 important that both the finite-sites theory and the bias-variance trade-off are tested comprehensively
88 on a genome-scale data set with clock-partitioning.

89 Persistent uncertainty in molecular date estimates is perhaps best exemplified by studies of
90 the origins of flowering plants (angiosperms) (Foster 2016). The earliest unequivocal angiosperm
91 fossils are tricolpate pollen grains from the Barremian–Aptian boundary, from approximately 125.9
92 million years ago (Ma) (Hughes 1994). Older pollen grains from the Hauterivian provide some
93 evidence of crown-group angiosperms, and are usually accepted as belonging to this group, albeit
94 with less confidence than for the tricolpate pollen grains (Herendeen et al. 2017). Patterns of
95 diversification in the broader fossil record suggest that angiosperms are unlikely to have arisen
96 much earlier than this time (Magallón et al. 2015). The majority of molecular dating analyses tell a
97 vastly different story, with most recent analyses inferring an origin within the Triassic (Foster et al.
98 2017). Additionally, the uncertainty surrounding the age of the angiosperm crown node is large,
99 often spanning an interval of many tens of millions of years, unless strong age constraints are
100 placed on the node. Improving the accuracy and precision of estimates of the age of crown
101 angiosperms thus represents a key goal of molecular dating.

102 In this study, we use a Bayesian phylogenetic approach to investigate the impact of clock-
103 partitioning on the precision of divergence-time estimates. We also investigate whether the criteria
104 used to assign genes to different clocks has an impact on estimation error. To do so, we infer the
105 evolutionary timescale of angiosperms using a plastome-level data set. In analyses with clock-
106 partitioning schemes comprising up to 20 clock-subsets, we allocate genes to clock-subsets based
107 on patterns of among-lineage rate heterogeneity or relative substitution rate, or through random
108 assignment. In all cases, we confirm that increasing the degree of clock-partitioning can lead to vast
109 improvements in the precision of Bayesian date estimates.

110

111 **Materials and Methods**

112 **Data Sets and Clock-Partitioning**

113 We obtained full chloroplast genome sequences for 52 angiosperm taxa and two gymnosperm
114 outgroup taxa from GenBank (supplementary table S1, Supplementary Material online). Each
115 angiosperm taxon was chosen to represent a different order, with our sampling designed to include
116 as many as possible of the 63 angiosperm orders recognized by the Angiosperm Phylogeny Group
117 (2016). We extracted all 79 protein-coding genes from the chloroplast genomes, although some
118 genes were missing from some taxa. We initially translated all genes into amino acid sequences
119 using VirtualRibosome (Wernersson 2006) and aligned them using MAFFT v7.305b (Kato and
120 Standley 2013). We then translated the aligned amino acid sequences back into nucleotide sequence
121 alignments using PAL2NAL (Suyama et al. 2006), made manual adjustments, and filtered out any
122 sites in the alignment at which a gap was present in $\geq 80\%$ of the taxa. Our total core data set
123 consisted of 68,790 nucleotides, of which only 7.54% sites were gaps or missing data (see
124 supplementary file S1, Supplementary Material online).

125 Our primary strategy for clock-partitioning based on patterns of among-lineage rate
126 heterogeneity was to analyse the genes using ClockstaR v2 (Duchêne et al. 2014). ClockstaR takes
127 predefined subsets of the data, along with the estimated gene tree for each subset, and determines
128 the optimal clock-partitioning scheme for the data set. This involves identifying the optimal number
129 of clock-subsets (k), as well as the optimal assignment of the data subsets to each of these clock-
130 subsets. Comparison of clock-partitioning schemes is done by comparing the patterns of among-
131 lineage rate heterogeneity across the gene trees and clustering the gene trees according to the gap
132 statistic (Gap_k) (Tibshirani et al. 2001). Additionally, ClockstaR can determine the optimal clock-

133 partitioning scheme for any value of k . In our case, each of the 79 protein-coding genes was
134 considered as a separate data subset for the ClockstaR analysis.

135 ClockstaR requires all data subsets to share the same tree topology. Since the chloroplast
136 genome does not typically undergo recombination (Birky 1995), all of its genes should share the
137 same topology. Therefore, we first inferred the phylogeny for the concatenated data set using
138 maximum-likelihood analysis in IQ-TREE v1.50a (Nguyen et al. 2015), with node support
139 estimated using 1000 bootstrap replicates with the ultrafast bootstrapping algorithm (Minh et al.
140 2013). We partitioned the data set by codon position using the edge-linked partition model
141 (Chernomor et al. 2016), and implemented the GTR+ Γ_4 model of nucleotide substitution for each
142 subset. The best-scoring tree was very similar to previous estimates of the angiosperm phylogeny
143 based on chloroplast data (Moore et al. 2010; Soltis et al. 2011), and we found strong support for
144 most nodes in the tree (supplementary fig. S1, Supplementary Material online). We used this tree
145 for ClockstaR and optimized the branch lengths for each gene alignment. Finally, we determined
146 the optimal value of k , and then created 12 clock-partitioning schemes using the optimal assignment
147 of genes to clock-subsets for values of k from 1 to 10, 15, and 20 (“ P_{CSTAR} ” schemes). We use the
148 partitioning along medoids (PAM) algorithm, described by Kaufman and Rousseeuw (2009).

149 As a means of comparison with the ClockstaR partitioning schemes, we also chose clock-
150 partitioning schemes based on relative substitution rates across genes (dos Reis et al. 2012). To do
151 so, we focused on a subset of 20 taxa for which sequences of all 79 protein-coding genes were
152 available (supplementary table S1, Supplementary Material online). We then analysed each gene
153 using maximum likelihood in IQ-TREE, in each case partitioning by codon position and
154 implementing the GTR+ Γ_4 model of nucleotide substitution for each codon position. Using the tree
155 lengths as a proxy for the overall substitution rate of each gene, we created 11 partitioning schemes
156 based on relative rates of substitution (“ P_{RATE} ” schemes), in which we assigned genes to clock-
157 subsets for values of k from 2 to 10, 15, and 20.

158 For an additional form of comparison, we generated clock-partitioning schemes with genes
159 randomly allocated to clock-subsets. Genes were randomly sampled without replacement in R
160 v3.3.2 (R Core Team 2016) and assigned to clock-subsets for values of k from 2 to 10, 15, and 20.
161 We repeated this process three times, resulting in a total of 33 clock-partitioning schemes in which
162 genes were randomly assigned to clock-subsets (“ P_{RAND} ” schemes).

163 164 Molecular Dating

165 We inferred the evolutionary timescale using MCMCTREE in PAML v4.8 (Yang 2007) with the
166 GTR+ Γ_4 model of nucleotide substitution. A key requirement of MCMCTREE is a fixed tree
167 topology, so we used the best-scoring tree that we estimated from the total concatenated data set
168 using IQTREE. We primarily analysed our data sets with the UCLN relaxed clock (Drummond et
169 al. 2006; Rannala and Yang 2007), but replicated all analyses to check for any differences under the
170 ACLN relaxed clock (Thorne et al. 1998; Kishino et al. 2001).

171 We estimated the overall substitution rate for each clock-partitioning scheme by running
172 baseml under a strict clock, with a single point calibration at the root. We then used this estimate to
173 select the shape (α) and scale (β) parameters for the gamma-Dirichlet prior on the overall
174 substitution rate across loci in the MCMCTREE analysis according to the formulae $\alpha = (m/s)^2$ and β
175 $= m/s^2$, where m and s are the mean and standard deviation of the substitution rate, respectively. For
176 all analyses, we set the shape and scale parameters for the gamma-Dirichlet prior on rate variation
177 across branches to 1 and 3.3, respectively. The posterior distribution of node ages was estimated
178 with Markov chain Monte Carlo sampling, with samples drawn every 10^3 steps across a total of 10^7
179 steps, after a discarded burn-in of 10^6 steps. We ran all analyses in duplicate to assess convergence,
180 and confirmed sufficient sampling by checking that the effective sample sizes of all parameters
181 were above 200.

182 We repeated the MCMCTREE analysis for all P_{CSTAR} , P_{RATE} , and P_{RAND} schemes. An
183 advantage of MCMCTREE is the option to use approximate likelihood calculation, which is much
184 faster than full likelihood calculation (Thorne et al. 1998; dos Reis and Yang 2011). However, this

185 precludes the calculation of marginal likelihoods using path sampling and similar methods, which
186 require the full likelihood to be computed. Instead, we compared the means and 95% credibility
187 intervals of the posterior estimates of divergence times across our partitioning strategies. We chose
188 to focus on six nodes in the angiosperm phylogeny: the crown groups of all angiosperms,
189 magnoliids, monocots, eudicots, campanulids, and Liliales. The first four of these were chosen
190 because they define major clades in the angiosperm phylogeny. The other two nodes were chosen
191 because they do not have explicit fossil-based calibration priors.

192 193 Fossil Calibrations

194 Calibrations are the most important component of Bayesian molecular dating, with critical impacts
195 on posterior estimates of divergence times. Therefore, we selected a set of 23 calibration priors
196 primarily based on recent studies that carefully considered the phylogenetic affinities of angiosperm
197 fossils (table 1). We also applied two calibration priors to the gymnosperm outgroup. Fossils can
198 strictly only provide a minimum age for the divergence of lineages from their common ancestor, so
199 we chose to implement fossil calibrations primarily as uniform distributions with soft bounds. This
200 approach assigns an equal prior probability for all ages between specified minimum and maximum
201 ages, with a 2.5% probability that the age surpasses each bound (Yang and Rannala 2006).

202 We implemented two maximum age constraints: (i) 350 Ma for the divergence between
203 angiosperms and gymnosperms (the root), a well accepted upper bound for this divergence (Foster
204 et al. 2017); and (ii) 126.7 Ma for the origin of crown eudicots, corresponding to the upper bound of
205 the Barremian–Aptian boundary (reviewed by Massoni et al. 2015a). The latter constraint is widely
206 used and is justified by the complete absence of tricolpate pollen before the latest Barremian, yet
207 some molecular dating results have suggested an earlier origin for eudicots (Smith et al. 2010;
208 Foster et al. 2017; Zeng et al. 2017). Ranunculales, one of the earliest-diverging eudicot orders, has
209 a fossil record dating back to the late Aptian/early Albian. Therefore, implementing the eudicot
210 maximum constraint results in a strong prior being placed on crown-group eudicots appearing
211 between ~126.7–112.6 Ma. As a result, including the eudicot maximum constraint leads to the
212 eudicot crown node being a useful example of a heavily constrained node for downstream
213 comparisons of the uncertainty in posterior age estimates.

214 For comparison, we also performed analyses with our P_{CSTAR} schemes using gamma
215 calibration priors and the UCLN relaxed clock. In this case, the mean of each gamma prior was set
216 to the age of each fossil +10%, with an arbitrary standard deviation of 2 (Table 1). This effectively
217 brackets the age estimates of calibrated nodes within a very narrow interval. In such a calibration
218 scheme, the precision of age estimates is not expected to improve substantially with increased
219 clock-partitioning.

220 221 Results

222 Angiosperm Evolutionary Timescale

223 Our ClockstaR analysis identified the optimal value of k to be 1, suggesting that a single pattern of
224 among-lineage rate heterogeneity is shared across protein-coding genes from the chloroplast
225 genomes. However, despite $k=1$ being optimal, the values of the gap statistic were still higher for all
226 values of $k>5$ (figure 1). Based on our analysis using the optimal clock-partitioning scheme ($k=1$)
227 and the UCLN relaxed clock, we estimated the time to the most recent common ancestor of
228 angiosperms to be 196 Ma (95% credibility interval 237–161 Ma; supplementary fig. S2,
229 Supplementary Material online). We inferred that crown magnoliids first appeared 171–115 Ma,
230 and that crown monocots arose contemporaneously, 167–120 Ma. Crown eudicots were inferred to
231 have arisen 128–124 Ma, with this precise estimate reflecting the strong calibration prior placed
232 upon this node. Finally, our estimates for the time to the most recent common ancestors of
233 campanulids and Liliales were 101–91 Ma and 108–91 Ma, respectively.

234 The true age of crown angiosperms is unknown, so we cannot assess the absolute accuracy
235 of our date estimates. Instead, we consider the consistency of mean age estimates across analyses

Table 1.—The calibration priors used within this study to estimate the angiosperm evolutionary timescale. "CG" and "SG" refer to the crown and stem groups, respectively, of the clade of interest.

Calibration node	Uniform Priors		Gamma Priors		Fossil	Reference
	Min Age Cal (Ma)	Max Age Cal (Ma)	α	β		
CG Alismatales	120.7	350	4332.8	3264.4	<i>Mayoa portugallica</i>	Magallón et al. (2015)
CG Angiospermae	136	350	5245.7	3507.2	Early Cretaceous pollen grains	Magallón et al. (2015)
CG Arecales	83.6	350	1992.0	2167.7	<i>Sabolites carolinensis</i>	Iles et al. (2015)
CG Boraginales	47.8	126.7	806.6	1535.4	<i>Ehretia clausentia</i>	Martinez-Millán 2010
CG Brassicales	89.3	126.7	2530.9	2577.0	<i>Dressiantha bicarpelata</i>	Magallón et al. (2015)
CG Caryophyllales	70.6	126.7	1495.4	1926.2	<i>Coahuilacarpa phytolaccoides</i>	Magallón et al. (2015)
CG Cornales	89.3	126.7	2530.9	2577.0	<i>Tylerianthus crossmanensis</i>	Magallón et al. (2015)
CG Ericales	89.3	126.7	2530.9	2577.0	<i>Pentapetalum trifasciculandricus</i>	Magallón et al. (2015)
CG Fabales	55.8	126.7	897.6	1462.5	<i>Paleosecuridaca curtissi</i>	Magallón et al. (2015)
CG Fagales	96.6	126.7	2689.9	2532.0	<i>Normapolles pollen</i>	Magallón et al. (2015)
CG Gentianales	37.2	126.7	445.7	1086.9	<i>Emmenopterys dilcheri</i>	Magallón et al. (2015)
CG Magnoliales	112.6	350	4197.7	3390.7	<i>Endressinia brasiliana</i>	Massoni et al. (2015)
CG Myrtales	87.5	126.7	2534.2	2632.6	<i>Esgueiria futabensis</i>	Magallón et al. (2015)
CG Oxalidales	100.1	126.7	2918.4	2651.4	<i>Tropidogyne pikei</i>	Chambers et al. (2010)
CG Pandanales	86.3	350	2289.8	2411.3	<i>Mabelia connatifila</i>	Iles et al. (2015)
CG Paracryphiales	79.2	126.7	1926.6	2209.7	<i>Silvianthemum suecicum</i>	Magallón et al. (2015)
CG Ranunculales	112.6	126.7	3867.5	3124.8	<i>Texeiraa lusitanica</i>	Magallón et al. (2015)
CG Saxifragales	89.3	126.7	2530.9	2577.0	<i>Microaltingia apocarpela</i>	Magallón et al. (2015)
CG Zingiberales	72.1	350	1663.3	2096.8	<i>Spirematospermum chandlerae</i>	Iles et al. (2015)
SG Buxales	99.6	126.7	3306.3	3019.9	<i>Spanomera marylandensis</i>	Magallón et al. (2015)
SG Cycadales	268.3	350	21939.8	7434.1	<i>Crossozamia</i>	Nagalingum et al. (2011)
SG gymnosperms	306.8	350	28377.3	8408.2	<i>Cordaixylon iowensis</i>	Clarke et al. (2011)
SG Platanaceae	107.7	126.7	3362.6	2837.3	<i>Sapindopsis variabilis</i>	Magallón et al. (2015)
SG Winteraceae	125	350	4738.5	3419.5	<i>Walkeripollis gabonensis</i>	Massoni et al. (2015)

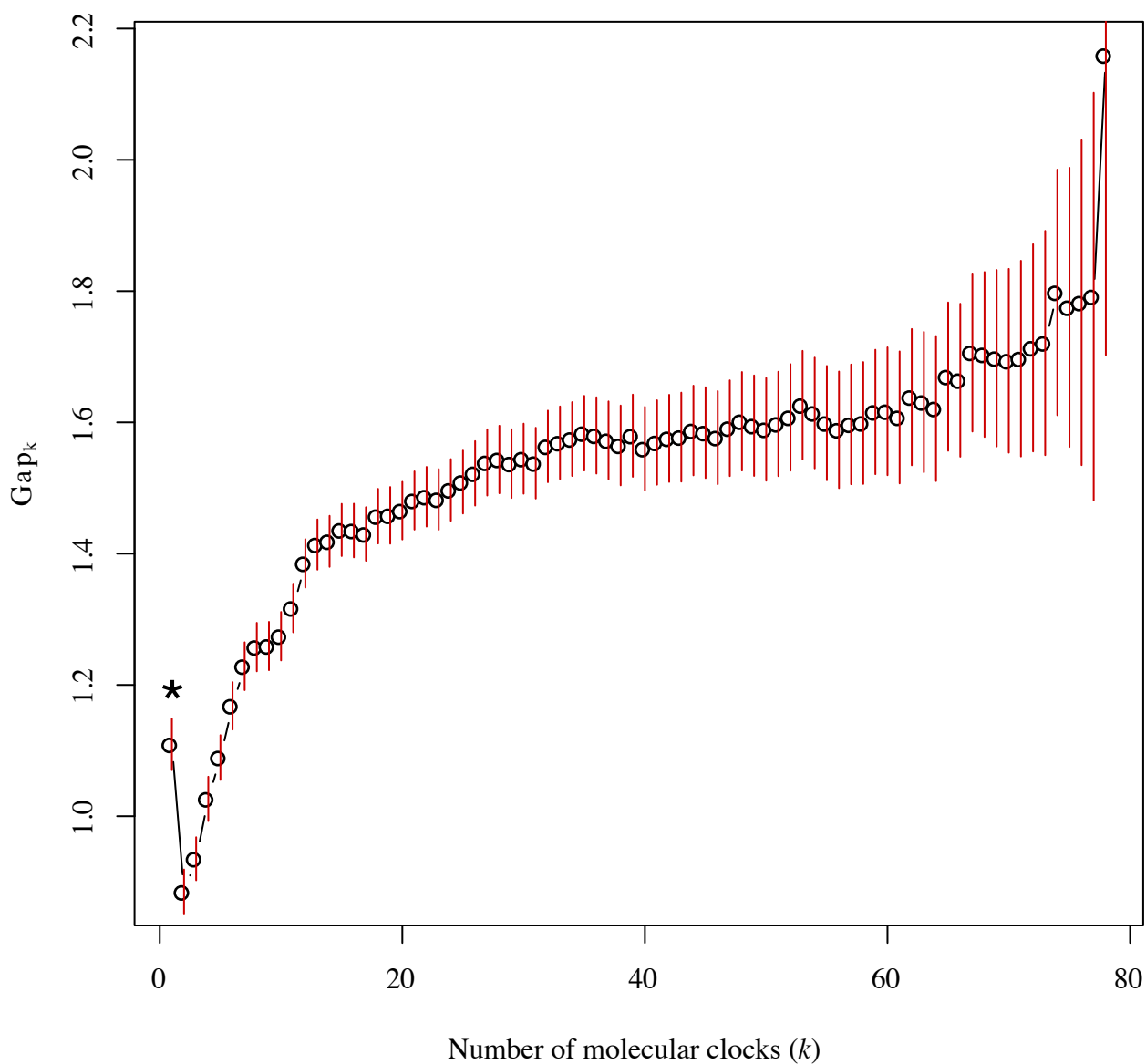


Fig. 1.—Gap statistic values for different numbers of clock-subsets (k) for the plastome-scale angiosperm data set, inferred using partitioning along medoids in ClockstaR. The asterisk indicates the optimal number of clock-subsets.

236 (Hillis 1995). The mean age estimates for all crown angiosperms, magnoliids, and monocots varied
237 slightly across values of k from 1 to 3, but estimates remained stable across all other values of k .
238 Mean age estimates for crown eudicots only varied by approximately 2 myr across all values of k .
239 Mean age estimates for crown Liliales were stable across all clock-partitioning schemes. However,
240 mean estimates for crown campanulids steadily declined by approximately 10–15 myr as the
241 number of loci increased. We observed the same broad trends in accuracy for all nodes of interest
242 when using the ACLN relaxed clock, although mean age estimates were consistently slightly
243 younger than in analyses with the UCLN relaxed clock. In our analyses with the P_{CSTAR} schemes and
244 with gamma calibration priors, mean age estimates for crown angiosperms steadily increased with
245 increasing numbers of clock-subsets, but the mean estimates were stable for all other nodes of
246 interest.

247 248 Precision in Estimates of Divergence Times

249 We focus first on our results when using the UCLN relaxed clock, uniform calibration priors, and
250 with clock-partitioning according to ClockstaR. We report improvements in the precision of node-
251 age estimates by calculating the decrease in 95% CI width, which we standardized by dividing by
252 the posterior mean. The optimal clock-partitioning scheme was inferred to be $k=1$, matching the
253 results of previous analyses (Duchêne et al. 2016). However, increasing the number of clock-
254 subsets generally led to large increases in the precision of node-age estimates. The impact of this is
255 perhaps most striking in the inferred age of crown angiosperms. Increasing the number of clock-
256 subsets from $k=1$ to $k=2$ led to a reduction in statistical fit (figure 1), but also reduced the width of
257 the 95% CI for the inferred age of crown angiosperms from 77 myr to 46 myr (an improvement in
258 precision of 35.4%). Greater clock-partitioning led to further improvement in precision (figure 2).
259 For example, implementing a clock-partitioning scheme with $k=20$ reduced the width of the 95% CI
260 for the inferred age of crown angiosperms to only 20 myr, representing a 73.1% improvement in
261 precision. However, the rate of improvement in precision declined rapidly for increasing numbers
262 of clock-subsets (figure 2).

263 An improvement in precision with the number of clock-subsets can also be observed in the
264 age estimates for both magnoliids and monocots. For example, increasing k from 1 to 20 results in
265 respective increases of 76.1% and 68% in precision in the age estimates for crown magnoliids and
266 crown monocots (figure 2). When considering the nodes corresponding to the crown groups of
267 campanulids and Liliales, a similar trend can be observed, albeit with a less drastic increase in
268 precision. Increasing the number of clock-subsets led to 29.7% and 37.7% increases in precision for
269 the crown groups of campanulids and Liliales, respectively. However, there is a vastly different
270 trend in the age estimate for crown eudicots. In this case, the age estimate for $k=1$ is already precise
271 (95% credibility interval: 128–124 Ma) and increasing the number of clock-subsets actually led to a
272 slight decrease in precision of 0.02%.

273 Compared with the P_{CSTAR} clock-partitioning schemes, very similar trends in precision were
274 observed for both the P_{RATE} scheme (figure 3) and P_{RAND} scheme (figure 4). The only differences
275 were that there was less variation in mean age estimates for smaller values of k compared with the
276 ClockstaR partitioning scheme, and standardized improvements in precision were consistently
277 slightly greater (supplementary table S2, Supplementary Material online). For example, the widths
278 of the 95% CIs, and the mean age estimates, declined monotonically in both classes of clock-
279 partitioning schemes.

280 We observed the same broad trends across all clock-partitioning schemes when using the
281 ACLN relaxed clock. With increasing numbers of clock-subsets, the uncertainty in age estimates
282 rapidly decreased, with the exception of the age estimate for the eudicot crown node. Even with
283 $k=1$, however, the precision of the age estimates was much greater than in the corresponding
284 analysis with the UCLN relaxed clock. For example, when implementing the P_{CSTAR} clock-
285 partitioning schemes, the 95% credibility interval of the age estimate for crown angiosperms
286 spanned 77 myr when using the UCLN relaxed clock, but only 59 myr when using the ACLN
287 relaxed clock. Additionally, age estimates for crown eudicots became less precise as the degree of

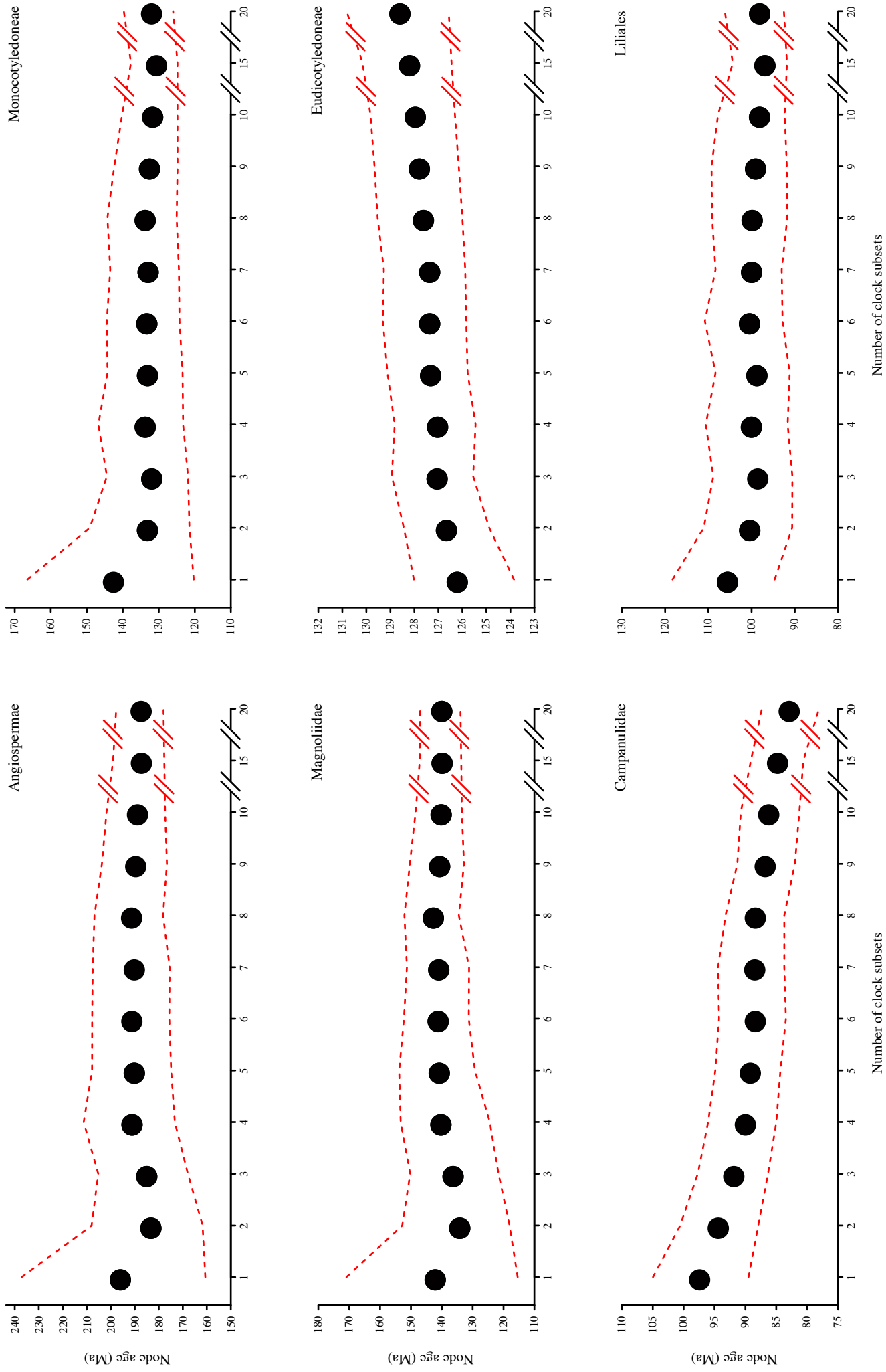


Fig. 2.—Mean age inferences and associated 95% credibility intervals for six nodes in the angiosperm phylogeny with increasing numbers of clock-subsets (k), as inferred using an autocorrelated lognormal relaxed clock, clock-partitioning according to the optimal schemes identified in ClockstAR, and uniform calibration priors.

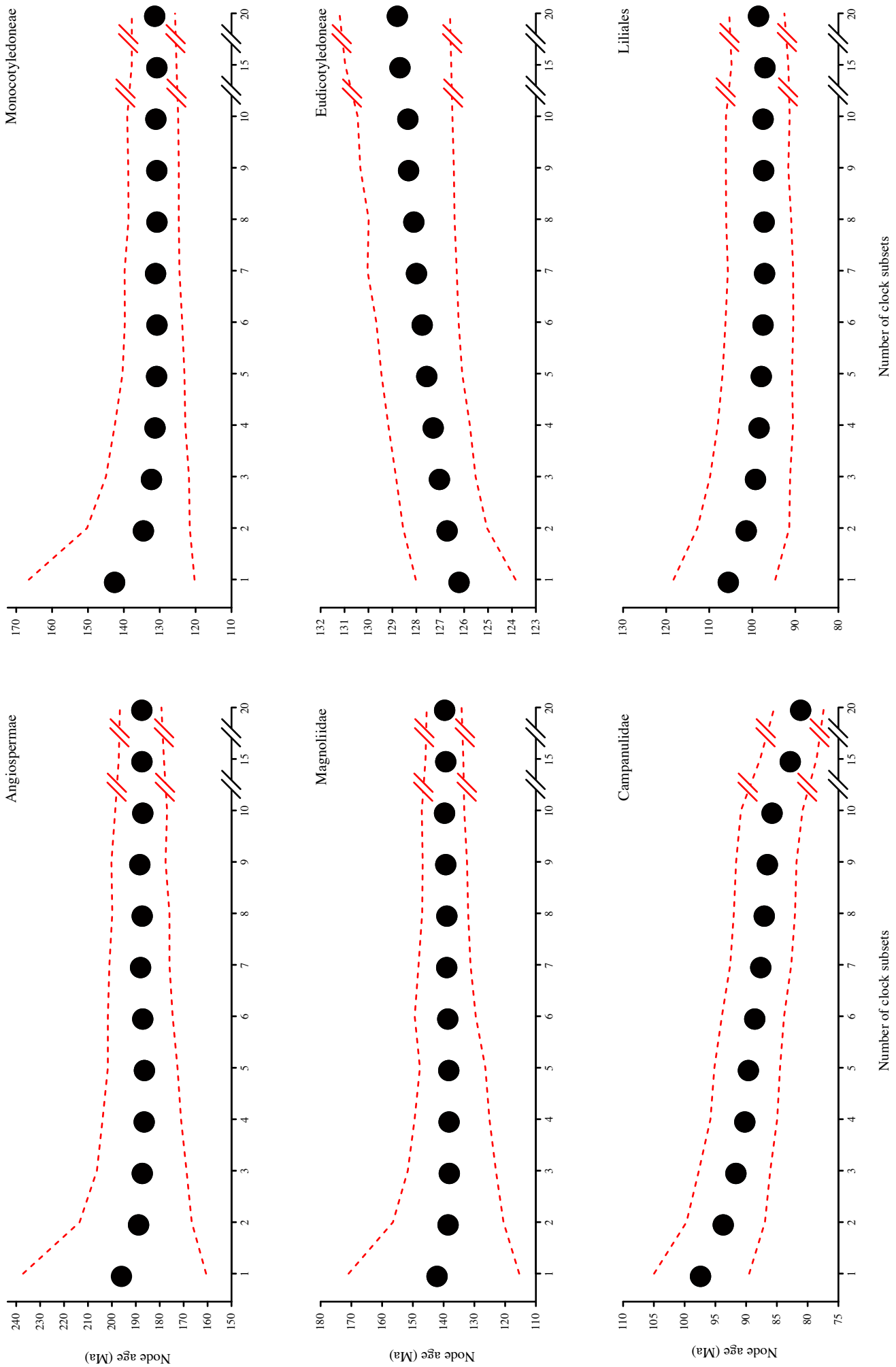


Fig. 3.—Mean posterior age estimates and associated 95% credibility intervals for six nodes in the angiosperm phylogeny with increasing numbers of clock-subsets (k), as inferred using an uncorrelated lognormal relaxed clock, clock-partitioning according to relative rates of substitution, and uniform calibration priors.

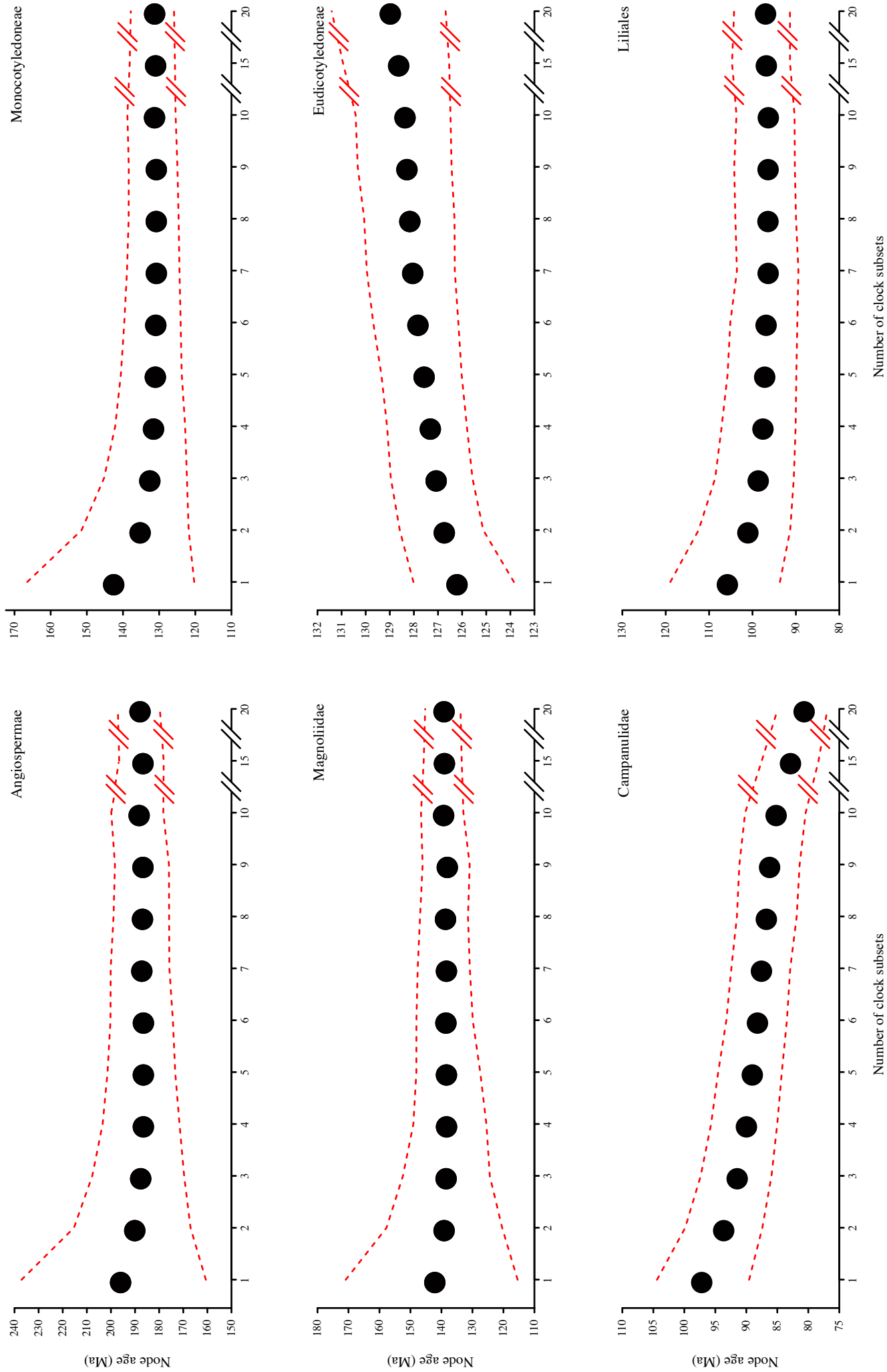


Fig. 4.—Mean posterior age estimates and associated 95% credibility intervals for six nodes in the angiosperm phylogeny with increasing numbers of clock-subsets (k), as inferred using an uncorrelated lognormal relaxed clock, clock-partitioning according to random assignment of genes to clock subsets, and uniform calibration priors. The estimates presented here are the averages of three random assignments of genes to clock-subsets for each value of k .

288 clock-partitioning increased. We observed the same trend for the other nodes of interest across
289 analyses, and the apparent limit to uncertainty appeared to be reached much more rapidly than with
290 the UCLN relaxed clock (supplementary fig. S3–S5, supplementary table S2, Supplementary
291 Material online).

292 When using highly informative gamma calibration priors in our additional analyses of the
293 *PCSTAR* schemes, we found that for the crown groups of angiosperms, monocots, and magnoliids, the
294 increases in precision with greater clock-partitioning were much lower than with uniform
295 calibration priors (supplementary fig. S6 and supplementary table S2, Supplementary Material
296 online). For example, an improvement of only 18.5% occurred in the precision of the age estimate
297 for crown angiosperms. The opposite trend occurred for the crown nodes of eudicots, campanulids
298 and Liliales. When implementing uniform calibration priors, greater clock-partitioning led to either
299 no change or decreases in precision for age estimates of crown-group eudicots, but when using
300 gamma calibration priors the precision improved by 36% with greater clock-partitioning. For
301 crown-group Liliales, increasing k from 1 to 20 led to a 64.3% increase in the precision of age
302 estimates, the greatest improvement of all six key nodes. However, it is worth noting that our age
303 estimates for all six nodes of interest were very precise even when $k=1$. Therefore, in terms of
304 absolute time units, there was generally little improvement in precision with increasing numbers of
305 clock-subsets.

306
307

308 Discussion

309 The primary aim of the present study was not to provide a novel estimate for the angiosperm
310 evolutionary timescale, but it is still useful to consider our results in the context of previous
311 estimates. Our inferred origin for crown-group angiosperms in the late Triassic to early Jurassic is
312 consistent with most modern molecular dating estimates (Bell et al. 2010; Magallón 2010; Clarke et
313 al. 2011; Zeng et al. 2014; Beaulieu et al. 2015; Foster et al. 2017). Similarly, our age estimate for
314 crown magnoliids of 171–115 Ma is very similar to a previous estimate of 179–127 Ma based on
315 the most comprehensive molecular dating analyses of Magnoliidae (Massoni et al. 2015a). Our
316 estimate of 167–120 Ma for the age of crown monocots is compelling, because a recent study of
317 monocots using the fossilized-birth-death model inferred a very similar age of 174–134 Ma (Eguchi
318 and Tamura 2016). Our age estimate for crown eudicots of 128–124 Ma suggests that there was not
319 enough signal within the data to overcome the strong calibration priors placed upon this node.
320 Finally, although our age estimate for the appearance of crown campanulids 101–91 Ma is very
321 similar to those of recent studies (Magallón et al. 2015; Foster et al. 2017), our age estimate of 108–
322 91 Ma for the time to the most recent common ancestor of Liliales was slightly younger than recent
323 estimates.

324 The goal of all molecular dating studies is to estimate the evolutionary timescale with a
325 useful degree of precision and accuracy. We demonstrated that increasing the degree of clock-
326 partitioning leads to increasingly precise age estimates, as predicted by the finite-sites theory (Zhu
327 et al. 2015). Additionally, clock-partitioning schemes based on patterns of among-lineage rate
328 heterogeneity or relative substitution rates did not have any measurable advantage over randomly
329 assigning genes to clock-subsets, at least in terms of the accuracy and precision of the resulting
330 estimates of divergence times. The near-identical patterns of precision across all clock-partitioning
331 schemes stands in contrast with previous suggestions that the assignment of genes to clock-subsets
332 is more important than the number of clock-subsets (Duchêne and Ho 2014).

333 Our results demonstrate that to improve the precision of age estimates, one could simply
334 increase the degree of clock-partitioning by assigning genes to an arbitrarily large number of clock-
335 subsets, until the marginal benefit of increasing the number of clocks is close to zero (Zhu et al.
336 2015). An obvious consequence of this is that one must consider whether such an increase is
337 desirable or biologically meaningful. If there is evidence that a data set conforms to a single pattern
338 of rate variation among lineages, an increase in precision from clock-partitioning is not justifiable
339 because the clock-subsets do not constitute independent realizations of the process of rate variation

340 (Zhu et al. 2015). Our analysis using ClockstaR indicates that within our data set, all genes exhibit
341 the same pattern of rate heterogeneity among lineages, such that they should be analysed using a
342 single clock model. In this case, increasing the degree of clock-partitioning leads to a model that
343 overfits the data, does not appear to accurately predict the data, and is insensitive to the sampled
344 data. Normally this would be expected to occur when a model underfits the data, but the increasing
345 sets of “independent” branch-rate estimates for each clock-subset ensure that estimates of node
346 times remain precise.

347 The uncertainty in posterior divergence times can be divided into three components: (i)
348 uncertainty in branch lengths due to limited sequence length (N); (ii) among-lineage rate variation
349 for each clock-subset, as well as the evolutionary rate variation among clock-subsets; and (iii)
350 uncertainty in fossil calibrations (Zhu et al. 2015). If L is large, then the uncertainty caused by
351 limited sequence length approaches zero at the rate of $1/N$. Additionally, the uncertainty attributable
352 to the second component approaches zero at the rate of $1/L$. As $N \rightarrow \infty$ and $L \rightarrow \infty$, the uncertainty in
353 divergence-time estimates should be wholly attributable to uncertainty in the fossil calibrations
354 (Zhu et al. 2015). For a data set of fixed size, such as our angiosperm data set, increasing L will
355 reduce N , and vice versa. We found that partitioning the data set into increasing numbers of clock-
356 subsets led to improvements in precision, which implies that increasing L has a larger impact on
357 precision than decreasing N has on reducing precision. However, it is likely that for very small
358 values of N , the estimation error in branch lengths will grow rapidly.

359 An important exception to the overall trend was the age inferences for the crown eudicot
360 node. The most common calibration strategy for this node has been to place a maximum bound or a
361 highly informative prior on the age of this node, based on the absence of tricolpate pollen before the
362 Barremian–Aptian boundary (~126 Ma) (Magallón and Castillo 2009; Sauquet et al. 2012; Massoni
363 et al. 2015a; Foster et al. 2017). Additionally, many of the earliest-diverging eudicot lineages have
364 relatively old fossils dating to the late Aptian (~113 Ma). These lines of evidence provide a narrow
365 age bracket for the eudicot crown, often causing age estimates for the eudicot crown node to be
366 necessarily highly precise. As a result, the limit in uncertainty of the fossil calibrations should be
367 reached rapidly. Therefore, the age of the eudicot crown node is useful to evaluate in light of the
368 finite-sites theory. We found that increasing the number of clock-subsets had essentially no effect
369 on the uncertainty in the age estimate of this node. A very similar pattern was observed when using
370 tightly constrained gamma calibration priors, and we expect that the general trend extends to other
371 cases in which calibrated nodes have strongly constrained ages, for example when lognormal or
372 exponential priors are chosen (Smith et al. 2010; Magallón et al. 2015).

373 Our results are especially important for analyses of genome-scale data sets. The size of
374 phylogenomic data sets generally precludes molecular dating with computationally intensive
375 phylogenetic software, such as BEAST (Bouckaert et al. 2014) or MrBayes (Ronquist et al. 2012),
376 unless work-around methods are employed (Ho 2014). For example, some researchers have chosen
377 to analyse each gene or data subset separately and then take the average of the results (Zeng et al.
378 2017). However, this methodology effectively assigns to each gene its own model of nucleotide
379 substitution and its own clock model. Not only does this run the risk of severe
380 overparameterization, but it also raises the question of how the estimates should be combined in a
381 way that takes full account of estimation error. Another method is to apply data filtering to select
382 only a subset of a data set, such as those that are the most clocklike (Jarvis et al. 2014) or the most
383 informative (Tong et al. 2016).

384 In cases where data-filtering approaches are not feasible, less computationally intensive
385 methods can be employed, such as the approximate-likelihood method of MCMCTREE. There are
386 also non-Bayesian alternatives to phylogenomic dating, such as penalized likelihood (Sanderson
387 2002), that have been used to analyse large data sets (Zanne et al. 2014). Additionally, a number of
388 rapid dating methods that can account for among-lineage rate heterogeneity without an explicit
389 statistical model of branch-rate variation have been developed specifically for phylogenomic data
390 sets (Kumar and Hedges 2016). Although these methods appear to have accuracy comparable to
391 that of Bayesian methods, they cannot produce reliable estimates of the uncertainty in the inferred

392 ages (Kumar and Hedges 2016). It is also unclear how well the results of these analyses will
393 conform to the finite-sites theory.

394

395 **Conclusions**

396 In this study, we have demonstrated that the finite-sites theory for molecular dating applies to a
397 typical genome-scale data set from angiosperms, with the exception of nodes that have strong age
398 constraints. In contrast with previous suggestions, the choice of strategy for assigning genes to
399 clocks does not appear to be important. These results imply that the data set can be arbitrarily
400 partitioned into a large number of clock-subsets, up to the point at which there is little marginal
401 benefit in increasing the degree of clock-partitioning. However, we caution that all molecular date
402 estimates should be critically interpreted to determine whether their precision is meaningful or not.
403 To this end, the best approach is to identify the patterns of among-lineage rate heterogeneity in a
404 data set and to apply a clock-partitioning scheme that appropriately captures this variation.

405

406 **Acknowledgements**

407 The authors acknowledge the facilities and the technical assistance of the Sydney Informatics Hub
408 at the University of Sydney and, in particular, access to the high-performance computing facility
409 Artemis. This work was supported by the Research Training Program and the Australian Research
410 Council [grant DP110100383 to S.Y.W.H.].

411

412 **References**

- 413 Angiosperm Phylogeny Group APG. 2016. An update of the Angiosperm Phylogeny Group
414 classification for the orders and families of flowering plants: APG IV. *Bot. J. Linn. Soc.*
415 181: 1–20. doi: 10.1111/boj.12385
- 416 Beaulieu JM, O’Meara B, Crane P, Donoghue MJ. 2015. Heterogeneous rates of molecular
417 evolution and diversification could explain the Triassic age estimate for angiosperms. *Syst.*
418 *Biol.* 64: 869–878.
- 419 Bell CD, Soltis DE, Soltis PS. 2010. The age and diversification of the angiosperms re-revisited.
420 *Am. J. Bot.* 97: 1296–1303.
- 421 Birky CW. 1995. Uniparental inheritance of mitochondrial and chloroplast genes: mechanisms and
422 evolution. *Proc. Natl. Acad. Sci. U.S.A.* 92: 11331–11338.
- 423 Bouckaert R, et al. 2014. BEAST 2: a software platform for Bayesian evolutionary analysis. *PLOS*
424 *Comput. Biol.* 10: e1003537.
- 425 Burnham KP, Anderson DR. 2003. Model selection and multimodel inference: a practical
426 information-theoretic approach. New York: Springer.
- 427 Chambers KL, Poinar Jr G, Buckley R. 2010. *Tropidogyne*, a new genus of Early Cretaceous
428 Eudicots (Angiospermae) from Burmese amber. *Novon* 20: 23–29.
- 429 Chernomor O, von Haeseler A, Minh BQ. 2016. Terrace aware data structure for phylogenomic
430 inference from supermatrices. *Syst. Biol.* 65: 997–1008.
- 431 Clarke JT, Warnock R, Donoghue PCJ. 2011. Establishing a time-scale for plant evolution. *New*
432 *Phytol.* 192: 266–301.
- 433 dos Reis M, et al. 2012. Phylogenomic datasets provide both precision and accuracy in estimating
434 the timescale of placental mammal phylogeny. *Proc. R. Soc. Lond. B Biol. Sci.* 279: 3491–
435 3500.
- 436 dos Reis M, Yang Z. 2011. Approximate likelihood calculation on a phylogeny for Bayesian
437 estimation of divergence times. *Mol. Biol. Evol.* 28: 2161–2172.
- 438 dos Reis M, Yang Z. 2013. The unbearable uncertainty of Bayesian divergence time estimation. *J.*
439 *Syst. Evol.* 51: 30–43.
- 440 Drummond AJ, Ho SYW, Phillips MJ, Rambaut A. 2006. Relaxed phylogenetics and dating with
441 confidence. *PLOS Biol.* 4: e88.

- 442 Duchêne S, Foster CSP, Ho SYW. 2016. Estimating the number and assignment of clock models in
443 analyses of multigene data sets. *Bioinformatics* 32: 1281–1285.
- 444 Duchêne S, Ho SYW. 2014. Using multiple relaxed-clock models to estimate evolutionary
445 timescales from DNA sequence data. *Mol. Phylogenet. Evol.* 77: 65–70.
- 446 Duchêne S, Molak M, Ho SYW. 2014. ClockstaR: choosing the number of relaxed-clock models in
447 molecular phylogenetic analysis. *Bioinformatics* 30: 1017–1019. doi:
448 10.1093/bioinformatics/btt665
- 449 Eguchi S, Tamura MN. 2016. Evolutionary timescale of monocots determined by the fossilized
450 birth-death model using a large number of fossil records. *Evolution* 70: 1136–1144. doi:
451 10.1111/evo.12911
- 452 Foster CSP. 2016. The evolutionary history of flowering plants. *J. Proc. R. Soc. N.S.W.* 149: 65–
453 82.
- 454 Foster CSP, et al. 2017. Evaluating the impact of genomic data and priors on bayesian estimates of
455 the angiosperm evolutionary timescale. *Syst. Biol.* 66: 338–351.
- 456 Gaut B, Yang L, Takuno S, Eguiarte LE. 2011. The patterns and causes of variation in plant
457 nucleotide substitution rates. *Annual Review of Ecology, Evolution, and Systematics* 42:
458 245–266.
- 459 Herendeen PS, Friis EM, Pedersen KR, Crane PR. 2017. Palaeobotanical redux: revisiting the age
460 of the angiosperms. *Nat. Plants* 3: 17015.
- 461 Hillis DM. 1995. Approaches for assessing phylogenetic accuracy. *Syst. Biol.* 44: 3–16.
- 462 Ho SYW. 2014. The changing face of the molecular evolutionary clock. *Trends Ecol. Evol.* 29:
463 496–503.
- 464 Ho SYW, Duchêne S. 2014. Molecular-clock methods for estimating evolutionary rates and
465 timescales. *Mol. Ecol.* 23: 5947–5965.
- 466 Hughes NF. 1994. *The enigma of angiosperm origins*. Cambridge: Cambridge University Press.
- 467 Iles WJD, Smith SY, Gandolfo MA, Graham SW. 2015. Monocot fossils suitable for molecular
468 dating analyses. *Bot. J. Linn. Soc.*
- 469 Jarvis ED, et al. 2014. Whole-genome analyses resolve early branches in the tree of life of modern
470 birds. *Science* 346: 1320–1331. doi: 10.1126/science.1253451
- 471 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:
472 improvements in performance and usability. *Mol. Biol. Evol.* 30: 772–780.
- 473 Kaufman L, Rousseeuw PJ. 2009. *Finding groups in data: an introduction to cluster analysis*.
474 Hoboken, NJ, USA: John Wiley & Sons.
- 475 Kishino H, Thorne JL, Bruno WJ. 2001. Performance of a divergence time estimation method under
476 a probabilistic model of rate evolution. *Mol. Biol. Evol.* 18: 352–361.
- 477 Kumar S, Hedges SB. 2016. Advances in time estimation methods for molecular data. *Mol. Biol.*
478 *Evol.* 33: 863–869.
- 479 Lanfear R, Calcott B, Ho SYW, Guindon S. 2012. PartitionFinder: combined selection of
480 partitioning schemes and substitution models for phylogenetic analyses. *Mol. Biol. Evol.* 29:
481 1695–1701.
- 482 Lanfear R, et al. 2013. Taller plants have lower rates of molecular evolution. *Nat. Commun.* 4:
483 1879.
- 484 Magallón S. 2010. Using fossils to break long branches in molecular dating: a comparison of
485 relaxed clocks applied to the origin of angiosperms. *Syst. Biol.* 59: 384–399.
- 486 Magallón S, Castillo A. 2009. Angiosperm diversification through time. *Am. J. Bot.* 96: 349–365.
- 487 Magallón S, Gómez-Acevedo S, Sánchez-Reyes LL, Hernández-Hernández T. 2015. A
488 metacalibrated time-tree documents the early rise of flowering plant phylogenetic diversity.
489 *New Phytol.* 207: 437–453.
- 490 Martínez-Millán M. 2010. Fossil record and age of the Asteridae. *Bot. Rev.* 76: 83–135.
- 491 Massoni J, Couvreur TLP, Sauquet H. 2015a. Five major shifts of diversification through the long
492 evolutionary history of Magnoliidae (angiosperms). *BMC Evol. Biol.* 15: 49.

- 493 Massoni J, Doyle JA, Sauquet H. 2015b. Fossil calibration of Magnoliidae, an ancient lineage of
494 angiosperms. *Palaeontologia Electronica* 18.1.2FC: 1–25.
- 495 Minh BQ, Nguyen MAT, von Haeseler A. 2013. Ultrafast approximation for phylogenetic
496 bootstrap. *Mol. Biol. Evol.* 30: 1188–1195.
- 497 Moore MJ, Soltis PS, Bell CD, Burleigh JG, Soltis DE. 2010. Phylogenetic analysis of 83 plastid
498 genes further resolves the early diversification of eudicots. *Proc. Natl. Acad. Sci. U.S.A.*
499 107: 4623–4628.
- 500 Muse SV, Gaut BS. 1994. A likelihood approach for comparing synonymous and nonsynonymous
501 nucleotide substitution rates, with application to the chloroplast genome. *Mol. Biol. Evol.*
502 11: 715–724.
- 503 Nagalingum NS, et al. 2011. Recent synchronous radiation of a living fossil. *Science* 334: 796–799.
- 504 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective
505 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32:
506 268–274.
- 507 R Core Team. 2016. R: A language and environment for statistical computing. Vienna, Austria: R
508 Foundation for Statistical Computing.
- 509 Rannala B, Yang Z. 2007. Inferring speciation times under an episodic molecular clock. *Syst. Biol.*
510 56: 453–466.
- 511 Ronquist F, et al. 2012. MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice
512 across a large model space. *Syst. Biol.* 61: 539–542.
- 513 Sanderson MJ. 2002. Estimating absolute rates of molecular evolution and divergence times: a
514 penalized likelihood approach. *Mol. Biol. Evol.* 19: 101–109.
- 515 Sauquet H, et al. 2012. Testing the impact of calibration on molecular divergence times using a
516 fossil-rich group: the case of *Nothofagus* (Fagales). *Syst. Biol.* 61: 289–313.
- 517 Smith SA, Beaulieu JM, Donoghue MJ. 2010. An uncorrelated relaxed-clock analysis suggests an
518 earlier origin for flowering plants. *Proc. Natl. Acad. Sci. U.S.A.* 107: 5897–5902.
- 519 Smith SA, Donoghue MJ. 2008. Rates of molecular evolution are linked to life history in flowering
520 plants. *Science* 322: 86–89.
- 521 Soltis DE, et al. 2011. Angiosperm phylogeny: 17 genes, 640 taxa. *Am. J. Bot.* 98: 704–730.
- 522 Suyama M, Torrents D, Bork P. 2006. PAL2NAL: robust conversion of protein sequence
523 alignments into the corresponding codon alignments. *Nucleic Acids Res.* 34: W609–W612.
- 524 Thorne JL, Kishino H, Painter IS. 1998. Estimating the rate of evolution of the rate of molecular
525 evolution. *Mol. Biol. Evol.* 15: 1647–1657.
- 526 Tibshirani R, Walther G, Hastie T. 2001. Estimating the number of clusters in a data set via the gap
527 statistic. *J. R. Stat. Soc. Ser. B (Statistical Methodol.)* 63: 411–423.
- 528 Tong KJ, Lo N, Ho SYW. 2016. Reconstructing evolutionary timescales using phylogenomics.
529 *Zoological Systematics* 41: 343–351. doi: 10.11865/zs.201640
- 530 Wernersson R. 2006. Virtual Ribosome—a comprehensive DNA translation tool with support for
531 integration of sequence feature annotation. *Nucleic Acids Res.* 34: W385–W388.
- 532 Wertheim JO, Sanderson MJ, Worobey M, Bjork A. 2010. Relaxed molecular clocks, the bias–
533 variance trade-off, and the quality of phylogenetic inference. *Syst. Biol.* 59: 1–8.
- 534 Wolfe KH, Li W-H, Sharp PM. 1987. Rates of nucleotide substitution vary greatly among plant
535 mitochondrial, chloroplast, and nuclear DNAs. *Proc. Natl. Acad. Sci. U.S.A.* 84: 9054–
536 9058.
- 537 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* 24: 1586–
538 1591.
- 539 Yang Z, Rannala B. 2006. Bayesian estimation of species divergence times under a molecular clock
540 using multiple fossil calibrations with soft bounds. *Mol. Biol. Evol.* 23: 212–226.
- 541 Zanne AE, et al. 2014. Three keys to the radiation of angiosperms into freezing environments.
542 *Nature* 506: 89–92.
- 543 Zeng L, et al. 2017. Resolution of deep eudicot phylogeny and their temporal diversification using
544 nuclear genes from transcriptomic and genomic datasets. *New Phytol.* 214: 1338–1354.

- 545 Zeng L, et al. 2014. Resolution of deep angiosperm phylogeny using conserved nuclear genes and
546 estimates of early divergence times. *Nat. Commun.* 5: 4956.
- 547 Zhu T, Dos Reis M, Yang Z. 2015. Characterization of the uncertainty of divergence time
548 estimation under relaxed molecular clock models using multiple loci. *Syst. Biol.* 2015: 267–
549 280.
- 550