

Nature Communications

Article

Pervasive correlation of molecular evolutionary rates in the tree of life

Qiqing Tao^{1,2}, Koichiro Tamura^{3,4}, Fabia Battistuzzi⁵, and Sudhir Kumar^{1,2,*}

¹Institute for Genomics and Evolutionary Medicine, Temple University, Philadelphia, PA19122

²Department of Biology, Temple University, Philadelphia, PA19122

³Department of Biological Sciences, Tokyo Metropolitan University, Tokyo, Japan

⁴Research Center for Genomics and Bioinformatics, Tokyo Metropolitan University, Tokyo, Japan

⁵Department of Biological Sciences, Oakland University, Rochester, MI48309

Correspondence to:

Sudhir Kumar
Temple University
Philadelphia, PA 19122, USA
E-mail: s.kumar@temple.edu

1 **New species arise from pre-existing species and inherit similar genomes and**
2 **environments. This predicts greater similarity of mutation rates and the tempo of**
3 **molecular evolution between direct ancestors and descendants, resulting in**
4 **correlation of evolutionary rates within lineages in the tree of life. Surprisingly,**
5 **molecular sequence data have not confirmed this expectation, possibly because**
6 **available methods lack power to detect correlated rates. Here we present an**
7 **accurate machine learning method used to detect correlation of rates in large**
8 **phylogenies. By applying this method to multigene and genome-scale sequence**
9 **alignments from mammals, birds, insects, metazoans, plants, fungi, and**
10 **prokaryotes, we discover extensive correlation in molecular evolutionary rates**
11 **throughout the tree of life in both DNA and protein sequences. These findings**
12 **suggest concordance between molecular and non-molecular evolutionary patterns**
13 **and will foster unbiased and precise dating of the tree of life.**

14

15 Phylogenomics has revolutionized our understanding of the patterns and timescale of the
16 tree of life^{1,2}. Genome-scale data has revealed that rates of molecular sequence change
17 vary extensively among species³⁻⁵. The causes and consequences of evolutionary rate
18 variation are of fundamental importance in molecular phylogenetics and systematics⁶⁻⁸,
19 not only to inform about the relationship among molecular, biological, and life history
20 traits, but also as a prerequisite for reliable estimation of divergence times among species
21 and genes^{3,5}.

22 Three decades ago, Gillespie⁹ proposed that molecular evolutionary rates within a
23 phylogeny will be correlated due to similarities in genomes, biology and environments
24 between ancestral species and their immediate progeny. This idea led to statistical
25 modelling of the variability of evolutionary rates among branches and formed the basis of
26 the earliest relaxed clock methods for estimating divergence times without assuming a
27 strict molecular clock^{3,5,10-12}. However, the independent branch rate (IBR) model has
28 emerged as a strong alternative to the correlated branch rate (CBR) model. IBR posits
29 that rates vary randomly throughout the tree, such that the evolutionary rate similarity
30 between an ancestor and its descendant is, on average, no more than that between more
31 distantly-related branches in a phylogeny^{5,13}. IBR model is now widely used in estimating

32 divergence times from molecular data for diverse groups of species, including
33 mammals¹³, birds^{14–16}, amphibians¹⁷, plants^{18–24}, and viruses^{13,25,26}. If the IBR model best
34 explains the variability of evolutionary rates, then we must infer a decoupling of molecular
35 and biological evolution, because morphology, behavior, and other life history traits are
36 more similar between closely-related species^{6,27,28} and are correlated with taxonomic or
37 geographic distance^{29,30}.

38 Alternatively, the widespread use of the IBR model^{13–20,22,23,25,26} may be explained
39 by the fact that the currently available statistical tests lack sufficient power to reject the
40 IBR model³¹. This may also explain why some studies report finding extensive branch
41 rate correlation in many datasets (e.g., Lepage et al.³²), but others cannot confirm this
42 using the same tests (e.g., Linder et al.¹⁹). Consequently, many researchers use both
43 CBR and IBR models for the same species groups^{13,23,33–43}, a practice that often
44 generates controversy via widely differing time estimates^{34,37,42,44–46}.

45 Therefore, we need a powerful method to accurately test whether evolutionary
46 rates are correlated among branches. This method should then be applied to molecular
47 datasets representing taxonomic diversity across the tree of life to assess the ubiquity of
48 correlated rates in nature. Here, we introduce a new machine learning approach
49 (CorrTest) with high power to detect correlation between molecular rates. CorrTest is
50 computationally efficient, and its application to a large number of datasets establishes the
51 pervasiveness of rate correlation in the tree of life.

52

53 **RESULTS**

54 **A machine learning approach for detecting rate correlation**

55 Machine learning is widely used to solve problems in many fields, but has not yet been
56 used to address challenges in molecular phylogenetics. We employed a supervised
57 machine learning (McL) framework⁴⁷ to build a predictive model that distinguishes
58 between CBR and IBR models. In our McL approach, the input is a molecular phylogeny
59 with branch lengths (often derived from a multiple sequence alignment), and the output
60 is a classification that corresponds to whether or not the evolutionary rates are correlated
61 (CBR or IBR, respectively). We used a logistic regression to build a predictive model. An
62 overview of our McL approach is presented in **Figure 1**.

63 To build a predictive model, we need measurable properties (features, **Fig. 1g** and
64 **h**) that are derived from the input data. The output is ultimately the assignment of input
65 data as most consistent with either CBR or IBR models. The selection of informative and
66 discriminating features is critical for the success of McL. In CorrTest, we derive relative
67 lineage rates using a given molecular phylogeny with branch lengths⁴⁸ (**Fig. 1e** and **1f**)
68 and use these lineage rates to generate informative features. An evolutionary lineage
69 includes all the branches in the descendant subtree, which is distinct from a branch that
70 only connects an ancestor to one of its immediate descendants. One cannot use branch
71 rates as features, because their computation requires the knowledge of node times in the
72 phylogeny, which cannot be estimated without prior assignment of a branch rate model.

73 The features that we selected for our McL predictive model were the correlation
74 between ancestral and descendant lineage rates (ρ_{ad}), the correlation between the sister
75 lineages (ρ_s), and the decay in ρ_{ad} when one and two parents are skipped (d_1 and d_2 ,
76 respectively). We selected correlation between ancestral and descendant lineage rates
77 (ρ_{ad}) as a feature because our analyses of simulated data showed that ρ_{ad} was much
78 higher for phylogenetic trees in which molecular sequences evolved under CBR model
79 (0.96) than the IBR model (0.54, **Fig. 2a; Supplementary information**). While
80 “independent rate” should imply a lack of correlation, ρ_{ad} is not zero for sequences
81 evolved under the IBR model because the evolutionary rate of an ancestral lineage is
82 necessarily related to the evolutionary rates of its descendant lineages. While ρ_{ad} is
83 greater than zero⁴⁸, this feature shows distinct patterns for both CBR and IBR models and
84 is thus a good candidate feature for McL (**Fig. 2a**). As our second feature, we selected
85 the correlation between the sister lineages (ρ_s), because ρ_s was higher for the CBR model
86 (0.89) than the IBR model (0.00, **Fig. 2b; Supplementary information**). Although our
87 extensive simulations produced some scenarios in which ρ_s was greater than 0.4 for
88 datasets that evolved with the IBR model (because ancestral lineage rates include
89 descendant evolutionary rates) ρ_s remains a highly discriminating feature for McL. Two
90 additional features included in McL measure the decay in ρ_{ad} when one and two parents
91 are skipped (d_1 and d_2), respectively, in ρ_{ad} calculations (**Supplementary information**).
92 We expect that ρ_{ad} will decay slower under CBR than IBR, which was consistent with our
93 observations (**Fig. 2c**).

94 The selected set of candidate features (ρ_s , ρ_{ad} , d_1 , and d_2) can be measured for
95 any phylogeny with branch lengths (e.g., derived from molecular data) and used to train
96 the machine learning classifier (**Fig. 1i**). For this purpose, we need a large set of
97 phylogenies in which branch rates are correlated (CBR = 1, **Fig. 1d**) and phylogenies in
98 which the branch rates are independent (IBR = 0, **Fig. 1c**). By using the four selected
99 features for each phylogeny and the associated numerical output state (0 or 1), we can
100 build a logistic regression that serves as the predictive model (**Fig. 1j**). However, there is
101 a paucity of empirical data for which CBR and IBR rates are firmly established. We
102 therefore trained our McL model on a simulated dataset, a practice that is now widely
103 used in applications when real world training datasets are few in number and often
104 containing high levels of error or uncertainty^{49,50}. We used computer simulations to
105 generate 1,000 phylogenies that evolved with CBR models and 1,000 phylogenies that
106 evolved with IBR models (**Fig. 1a** and **b**). To ensure the general utility of our model for
107 analyses of diverse data, we sampled phylogenies with varying numbers of species,
108 degrees of rate correlation, and degrees of independent rate variation (**Supplementary**
109 **information**). The machine learning process generated a predictive model with an
110 associated correlation score (CorrScore).

111 We evaluated the sensitivity and specificity of our model using standard receiver
112 operating characteristic (ROC) curves, which show the sensitivity of our method to detect
113 rate correlation when it is present (True Positive Rate, TPR) and when it was not present
114 (False Positive Rate, FPR) at different CorrScore thresholds. The ROC curve for McL
115 using all four features was the best, which led to the inclusion of all four features in the
116 predictive model (**Fig. 2d; Supplementary information**). The area under the ROC
117 (AUROC) was 99%, with a 95% TPR (i.e., CBR detection) achieved at the expense of
118 only 5% FPR (**Fig. 2d**, black line). The area under the precision recall (AUPR) curve was
119 also extremely high (0.99; **Fig. 2d** inset), which means that our predictive model detects
120 correlation among branch rates with very high accuracy and precision. We also performed
121 standard cross-validation tests and found that that the predictive models retained high
122 accuracy (>92%, **Fig. 1k** and **Supplementary information**).

123 We developed a conventional statistical test (CorrTest) based on CorrScore (**Fig.**
124 **2e**) that will provide a p-value for researchers to use when deciding whether they should

125 reject a null hypothesis that branch rates within a phylogeny are uncorrelated
126 (independent). A high CorrScore translates into a higher probability that the branch rates
127 are correlated. At a CorrScore greater than 0.5, the Type I error (rejecting the null
128 hypothesis of IBR when it was true) was less than 5%. Type I error of 1% (P-value of
129 0.01) was achieved with a CorrScore greater than 0.83. We found that these CorrScore
130 score thresholds were applicable even when predictive models were developed
131 separately and when the number of sequences in the dataset were small (≤ 100), medium
132 (100 – 200), large (200 – 300), and very large (> 300) (**Supplementary information**).
133 The accuracy obtained using these models (**Fig. S1a-c**) is similar to those presented in
134 **figure 3d - f**. Therefore, we suggest using the general model in CorrTest analysis.

135 **CorrTest performs well in computational tests**

136 We tested the performance of CorrTest on a simulated dataset where the correct rate
137 model is known (**Fig. 1l**). This dataset used 91 angiosperms as a model system for
138 simulating sequence evolution with IBR models (**supplementary information**)⁵¹.
139 CorrTest correctly diagnosed 95% of these datasets to be evolving with independent
140 rates. We also tested a large collection of datasets⁵² generated using diverse evolutionary
141 parameters including both CBR and IBR models (**supplementary information**). CorrTest
142 showed an accuracy greater than 94% in detecting rate autocorrelation for datasets that
143 were simulated with low and high G+C contents (**Fig. 3a**), small and large substitution
144 rate biases (**Fig. 3b**), and different levels of sequence conservation (**Fig. 3c**). As
145 expected, CorrTest performed best on datasets that contain more and longer sequences
146 (**Fig. 3d**). In these analyses, we used the correct tree topology and nucleotide substitution
147 model. We relaxed this requirement and evaluated CorrTest by first inferring a phylogeny
148 using a dataset⁵³ with an oversimplified substitution model⁵⁴. Naturally, many inferred
149 phylogenies contained topological errors, but we found the accuracy of CorrTest to still
150 be high as long as the dataset contained >100 sequences of length $>1,000$ base pairs
151 (**Fig. 3e**). CorrTest performed well even when 20% of the partitions were incorrect in the
152 inferred phylogeny (**Fig. 3f**). Therefore, CorrTest will be most reliable for large datasets,
153 but is relatively robust to errors in phylogenetic inference.

154 **CorrTest versus Bayes factor analysis**

155 We compared the performance of CorrTest with that of the Bayes factor approach.
156 Because the Bayes factor method is computationally demanding, we limited our
157 comparison to 100 datasets containing 100 sequences each (**Supplementary**
158 **information**). We computed Bayes factors (BF) by using the stepping-stone sampling
159 (SS) method (see **Materials and Methods**). BF-SS analysis detected autocorrelation (P
160 < 0.05) for 32% of the datasets that actually evolved with correlated rates (**Fig. 4a**, red
161 curve in the CBR zone). This is because the marginal log-likelihoods under the CBR
162 model for 78% of these datasets were very similar to or lower than the IBR model.
163 Therefore, BF was very conservative in rejecting the null hypothesis (see also ref. ³¹). In
164 contrast, CorrTest correctly detected the CBR model for 88% of the datasets ($P < 0.05$;
165 **Fig. 4b**, red curve in CBR zone). For datasets that evolved with IBR model, BF-SS
166 correctly detected the IBR model for 92% (**Fig. 4a**, blue curves in the IBR zone), whereas
167 CorrTest correctly detected 86% (**Fig. 4b**, blue curve in the IBR zone). Therefore, Bayes
168 Factor analyses generally perform well in correctly classifying phylogenies evolved under
169 IBR, but fail to detect the influence of CBR. The power of CorrTest to correctly infer CBR
170 is responsible for its higher overall accuracy (87%, vs. 62% for BF). Such a difference in
171 accuracy was observed at all levels of statistical significance (**Fig. 4c**). In the future, faster
172 and more advanced BF implementations may allow extensive comparison of traditional
173 Bayesian and CorrTest approaches, as the Bayesian approaches are still evolving⁴³ and
174 currently require extensive computation time. Based on the limited comparisons
175 presented here, we conclude that machine learning enables highly accurate detection of
176 rate correlation in a given phylogeny and presents a computationally feasible alternative
177 to Bayes Factor analyses for large datasets.

178 **Correlation of rates is common in molecular evolution**

179 The high accuracy and fast computational speed of CorrTest enabled us to test the
180 presence of autocorrelation in 16 large datasets from 12 published studies encompassing
181 diverse groups across the tree life. This included nuclear, mitochondrial and plastid DNA,
182 and protein sequences from mammals, birds, insects, metazoans, plants, fungi, and
183 prokaryotes (**Table 1**). CorrTest rejected the IBR model for all datasets ($P < 0.05$). In
184 these analyses, we assumed a time-reversible process for base substitution. However,

185 the violation of this assumption may produce biased results in phylogenetic analysis⁵⁷.
186 We, therefore, applied an unrestricted substitution model for analyzing all the nuclear
187 datasets and confirmed that CorrTest rejected the IBR model in every case ($P < 0.05$).
188 This robustness stems from the fact that the branch lengths estimated under the time-
189 reversible and the unrestricted model show an excellent linear relationship for these data
190 ($r^2 > 0.99$). This is the reason why CorrTest produces reliable results even when an
191 oversimplified model was used in computer simulations (**Fig. 3e** and **f**).

192 These results suggest that the correlation of rates among lineages is the rule,
193 rather than the exception in molecular phylogenies. This pattern contrasts starkly with
194 those reported in many previous studies^{13–24,41}. In fact, all but three datasets^{33,55,56}
195 received very high prediction scores in CorrTest, resulting in extremely significant P -
196 values ($P < 0.001$). The IBR model was also rejected for the other three datasets ($P <$
197 0.05), but their test scores were not as high, likely because they sparsely sample a large
198 phylogenetic space. For example, the metazoan dataset³³ contains sequences primarily
199 from highly divergent species that shared common ancestors hundreds of millions of
200 years ago. In this case, tip lineages in the phylogeny are long and their evolutionary rates
201 are influenced by many un-sampled lineages. Such sampling effects weaken the rate
202 correlation signal. We verified this behavior via analyses of simulated data and found that
203 CorrTest's prediction scores decreased when taxon sampling and density were lowered
204 (**Fig. 5a**). Overall, CorrTest detected rate correlation in all the empirical datasets.

205 Our results establish that the correlated rate model should be the default in
206 molecular clock analysis, and CorrTest can be used to test the independent rate model
207 when sufficient numbers of sequences are available. Use of a correlated rate model is
208 important because model selection has a strong influence on the posterior credible
209 intervals of divergence times⁴⁴. For example, the use of IBR model produces estimates
210 of divergence time of two major groups of grasses that are 66% older⁴⁶ and origin of a
211 major group of mammal (Erinaceidea) to be 30% older³⁵ than estimates under CBR
212 model. In fact, substantial differences between node age estimates under IBR and CBR
213 models have been reported in many studies^{23,34,37,42,44,46}. Thus, the use of an incorrect
214 rate model has a large impact on time estimates, which may not be alleviated by adding

215 calibrations⁴⁴. Knowledge that evolutionary rates are generally correlated within lineages
216 will foster unbiased and precise dating of the tree of life.

217 **Magnitude of the rate correlation in molecular data**

218 CorrScore is influenced by the size of the dataset in addition to the degree of correlation,
219 so it is not a direct measure of the degree of rate correlation (effect size) in a phylogeny.
220 Instead, one should use a Bayesian approach to estimate the degree of rate correlation,
221 for example, under the Kishino et al.'s autocorrelated rate model⁵⁸. In this model, a single
222 parameter (ν) captures the degree of autocorrelation among branches in a phylogenetic
223 tree. A low value of ν indicates high autocorrelation, so, we use the inverse of ν to
224 represent the degree of rate autocorrelation. MCMCTree⁵⁹ analyses of simulated
225 datasets confirmed that the estimated ν is linearly related to the true value (**Fig. 5b**). In
226 empirical data analyses, we find that the inverse of ν is high for all datasets examined,
227 which suggests ubiquitous high rate correlation across the tree of life.

228 Many other interesting patterns emerge from this analysis. First, rate correlation is
229 highly significant not only for mutational rates (= substitution rate at neutral positions),
230 which are expected to be similar in sister species because they inherit cellular machinery
231 from a common ancestor, but also amino acid substitution rates, which are more strongly
232 influenced by natural selection (**Table 1**). For example, synonymous substitution rates in
233 the third codon positions and the four-fold degenerate sites in mammals³⁵, which are
234 largely neutral and are the best reflection of mutation rates⁶⁰, received high CorrScores
235 of 0.99 and 0.98, respectively ($P < 0.001$). Second, our model also detected a strong
236 signal of correlation for amino acid substitution rates in the same proteins (CorrScore =
237 0.99). Bayesian analyses showed that the degree of correlation is high in both cases:
238 inverse of ν was 3.21 in 4-fold degenerate sites and 3.11 in amino acid sequences. Third,
239 mutational and substitution rates in both nuclear and mitochondrial genomes are highly
240 correlated (**Table 1**). These results establish that molecular and non-molecular
241 evolutionary patterns are concordant, because morphological characteristics are also
242 found to be similar between closely-related species^{6,27,28} and correlated with taxonomic
243 or geographic distance^{29,30}.

244 In conclusion, we have successfully addressed an enduring question in
245 evolutionary biology: are the molecular rates of change between species correlated or

246 independent? We have shown that the evolutionary rates of change among closely
247 related species are correlated in diverse species groups. That is, evolutionary rate
248 correlation is likely universal, suggesting concordance between the patterns of
249 evolutionary changes in genomes and higher-level biological attributes. Furthermore,
250 revealing the existence of pervasive correlation in molecular rates throughout the tree of
251 life will improve specification of correct rate models that are essential for molecular clock
252 analyses to provide accurate estimates of evolutionary timing for use in studies of
253 biodiversity, phylogeography, development, and genome evolution.

254

255 **Materials and Methods**

256 **CorrTest analyses.** All CorrTest analyses were conducted using a customized R code
257 (available from <https://github.com/cathyqqtao/CorrTest>). We estimated branch lengths of
258 a tree topology on sequence alignments using maximum likelihood method (or Neighbor-
259 Joining method when we tested the robustness of our model to topological error) in
260 MEGA^{61,62}. Then we used those branch lengths to compute relative lineages rates using
261 RRF^{48,52} and calculated the value of selected features (ρ_s , ρ_{ad} , and two decay measures)
262 to obtain the CorrScore (see detail calculation in **Supplementary information**). We
263 conducted CorrTest on the CorrScore to estimate the *P*-value of rejecting the null
264 hypothesis of independent evolutionary rates. No calibration was needed for CorrTest
265 analyses.

266 **Bayes factor analyses.** We computed the Bayes factor via stepping-stone sampling (BF-
267 SS)⁶³ with $n = 20$ and $a = 5$ using mcmc3r package⁴³. We chose BF-SS because the
268 harmonic mean estimator it has many statistical shortcomings^{32,63,64} and thermodynamic
269 integration^{43,65} is less efficient than BF-SS. Still, BF-SS requires a long computational
270 time, we only finished analyses of 50% of synthetic datasets (**Supplementary**
271 **information**). For each dataset, we computed the log-likelihoods ($\ln K$) of using IBR model
272 and CBR model. The Bayes factor posterior probability for CBR was calculated as shown
273 in dos Reis et al. (2018)⁴³. We used only one calibration point at the root (true age with a
274 narrow uniform distribution) in all the Bayesian analyses, as it is the minimum number of
275 calibrations required by MCMCTree⁵⁹. For other priors, we used diffused distributions of

276 “rgene_gamma = 1 1”, “sigma2_gamma=1 1” and “BDparas = 1 1 0”. In all Bayesian
277 analyses, two independent runs of 5,000,000 generations each were conducted, and
278 results were checked in Tracer⁶⁶ for convergence. ESS values were higher than 200 after
279 removing 10% burn-in samples for each run.

280 ***Analysis of empirical datasets***

281 We used 16 datasets from 12 published studies of eukaryotes and 2 published studies of
282 prokaryotes that cover the major groups in the tree of life (**Table 1**). These were selected
283 because they did not contain too much missing data (<50%) and represented >80
284 sequences. When a phylogeny and branch lengths were available from the original study,
285 we estimated relative rates directly from the branch lengths via the relative rate
286 framework⁴⁸ and computed selected features to conduct CorrTest. Otherwise, maximum
287 likelihood estimates of branch lengths were obtained using the published phylogeny,
288 sequence alignments, and the substitution model specified in the original article^{61,62}.

289 To obtain the autocorrelation parameter (ν), we used MCMCTree⁵⁹ with the same
290 input priors as the original study, but no calibration priors were used in order to avoid
291 undue influence of calibration uncertainty densities on the estimate of autocorrelation
292 parameters. We did, however, provide a root calibration because MCMCTree requires a
293 root calibration. For this purpose, we used the root calibration provided in the original
294 article or selected the median age of the root node in the TimeTree database^{67,68} \pm 50My
295 (soft uniform distribution) as the root calibration, as this does not impact the estimation of
296 ν . Bayesian analyses required long computational times, so we used the original
297 alignments in MCMCTree analyses if alignments were shorter than 20,000 sites. If the
298 alignments were longer than 20,000 sites, we randomly selected 20,000 sites from the
299 original alignments to use in MCMCTree analyses. However, one dataset⁶⁹ contained
300 more than 300 ingroup species, such that even alignments of 20,000 sites required
301 prohibitive amounts of memory. In this case, we randomly selected 2,000 sites from the
302 original alignments to use in MCMCtree analyses (similar results were obtained with a
303 different site subset). Two independent runs of 5,000,000 generations each were
304 conducted, and results were checked in Tracer⁶⁶ for convergence. ESS values were

305 higher than 200 after removing 10% burn-in samples for each run. All empirical datasets
306 are available at <https://github.com/cathyqqtao/CorrTest>.

307 **Code availability statement**

308 The R source code of CorrTest is available at <https://github.com/cathyqqtao/CorrTest>

309 **Data availability statement**

310 All empirical datasets, results, and source code for generating each figure are available
311 at <https://github.com/cathyqqtao/CorrTest>. All simulated datasets will be provided upon
312 request.

313 **Acknowledgements**

314 We thank Xi Hang Cao for assisting on building the machine learning model, and Drs. Bui
315 Quang Minh, Beatriz Mello, Heather Rowe, Ananias Escalante, Maria Pacheco, and S.
316 Blair Hedges for critical comments and editorial suggestions. This research was
317 supported by grants from National Aeronautics and Space Administration (NASA
318 NNX16AJ30G), National Institutes of Health (GM0126567-01; LM012487-02), National
319 Science Foundation (NSF DBI 1356548), and Tokyo Metropolitan University (DB105).
320

321 **References**

- 322 1. Hedges, S. B., Marin, J., Suleski, M., Paymer, M. & Kumar, S. Tree of life reveals clock-like
323 speciation and diversification. *Mol. Biol. Evol.* **32**, 835–845 (2015).
324
- 325 2. Marin, J., Battistuzzi, F. U., Brown, A. C. & Hedges, S. B. The Timetree of Prokaryotes: New
326 Insights into Their Evolution and Speciation. *Mol. Biol. Evol.* **34**, 437–446 (2017).
327
- 328 3. Kumar, S. & Hedges, S. B. Advances in time estimation methods for molecular data. *Mol.*
329 *Biol. Evol.* **33**, 863–869 (2016).
330
- 331 4. Dos Reis, M., Donoghue, P. C. & Yang, Z. Bayesian molecular clock dating of species
332 divergences in the genomics era. *Nat. Rev. Genet.* **17**, 71–80 (2016).
333
- 334 5. Ho, S. Y. & Duchêne, S. Molecular-clock methods for estimating evolutionary rates and
335 timescales. *Mol. Ecol.* **23**, 5947–5965 (2014).
336
- 337 6. Lanfear, R., Welch, J. J. & Bromham, L. Watching the clock: studying variation in rates of
338 molecular evolution between species. *Trends Ecol. Evol.* **25**, 495–503 (2010).
339
- 340 7. Lynch, M. Evolution of the mutation rate. *Trends Genet.* **26**, 345–352 (2010).
341
- 342 8. Kimura, M. *The neutral theory of molecular evolution*. (Cambridge: Cambridge University
343 Press, 1983).
344
- 345 9. Gillespie, J. H. The molecular clock may be an episodic clock. *Proc. Natl. Acad. Sci. U.S.A.*
346 **81**, 8009–8013 (1984).
347
- 348 10. Kumar, S. Molecular clocks: four decades of evolution. *Nat. Rev. Genet.* **6**, 654–662 (2005).
349
- 350 11. Sanderson, M. J. A nonparametric approach to estimating divergence times in the absence
351 of rate constancy. *Mol. Biol. Evol.* **14**, 1218–1231 (1997).
352
- 353 12. Thorne, J. L., Kishino, H. & Painter, I. S. Estimating the rate of evolution of the rate of
354 molecular evolution. *Mol. Biol. Evol.* **15**, 1647–1657 (1998).
355
- 356 13. Drummond, A. J., Ho, S. Y. W., Phillips, M. J. & Rambaut, A. Relaxed phylogenetics and
357 dating with confidence. *PLoS Biol.* **4**, 88–99 (2006).
358
- 359 14. Brown, J. W., Rest, J. S., García-Moreno, J., Sorenson, M. D. & Mindell, D. P. Strong
360 mitochondrial DNA support for a Cretaceous origin of modern avian lineages. *BMC Biol.* **6**, 6
361 (2008).
362
- 363 15. Prum, R. O. *et al.* A comprehensive phylogeny of birds (Aves) using targeted next-
364 generation DNA sequencing. *Nature* **526**, 569–578 (2015).
365
- 366 16. Claramunt, S. & Cracraft, J. A new time tree reveals Earth history’s imprint on the evolution
367 of modern birds. *Sci Adv* **1**, e1501005 (2015).
368
- 369 17. Feng, Y.-J. *et al.* Phylogenomics reveals rapid, simultaneous diversification of three major
370 clades of Gondwanan frogs at the Cretaceous-Paleogene boundary. *Proc. Natl. Acad. Sci.*

- 371 U.S.A. **114**, E5864–E5870 (2017).
372
373 18. Moore, B. R. & Donoghue, M. J. Correlates of diversification in the plant clade Dipsacales:
374 geographic movement and evolutionary innovations. *Am. Nat.* **170** **Suppl 2**, S28–55 (2007).
375
376 19. Linder, M., Britton, T. & Sennblad, B. Evaluation of Bayesian models of substitution rate
377 evolution-parental guidance versus mutual independence. *Syst. Biol.* **60**, 329–342 (2011).
378
379 20. Lu, Y., Ran, J.-H., Guo, D.-M., Yang, Z.-Y. & Wang, X.-Q. Phylogeny and divergence times
380 of gymnosperms inferred from single-copy nuclear genes. *PLoS One* **9**, e107679 (2014).
381
382 21. Barreda, V. D. *et al.* Early evolution of the angiosperm clade Asteraceae in the Cretaceous
383 of Antarctica. *Proc. Natl. Acad. Sci. U.S.A.* **112**, 10989–10994 (2015).
384
385 22. Smith, S. A., Beaulieu, J. M. & Donoghue, M. J. An uncorrelated relaxed-clock analysis
386 suggests an earlier origin for flowering plants. *Proc. Natl. Acad. Sci. U.S.A.* **107**, 5897–5902
387 (2010).
388
389 23. Bell, C. D., Soltis, D. E. & Soltis, P. S. The age and diversification of the angiosperms re-
390 revisited. *Am. J. Bot.* **97**, 1296–1303 (2010).
391
392 24. Barba-Montoya, J., Dos Reis, M., Schneider, H., Donoghue, P. C. J. & Yang, Z.
393 Constraining uncertainty in the timescale of angiosperm evolution and the veracity of a
394 Cretaceous Terrestrial Revolution. *New Phytol.* (2018). doi:10.1111/nph.15011
395
396 25. Metsky, H. C. *et al.* Zika virus evolution and spread in the Americas. *Nature* **546**, 411–415
397 (2017).
398
399 26. Buck, C. B. *et al.* The Ancient Evolutionary History of Polyomaviruses. *PLoS Pathog.* **12**,
400 e1005574 (2016).
401
402 27. Sargis, E. J. & Dagosto, M. *Mammalian evolutionary morphology: a tribute to Frederick S.*
403 *Szalay*. (New York: Springer Netherlands, 2008).
404
405 28. Cox, P. G. & Hautier, L. *Evolution of the Rodents: Volume 5: Advances in Phylogeny,*
406 *Functional Morphology and Development*. (Cambridge: Cambridge University Press, 2015).
407
408 29. Wyles, J. S., Kunkel, J. G. & Wilson, A. C. Birds, behavior, and anatomical evolution. *Proc.*
409 *Natl. Acad. Sci. U.S.A.* **80**, 4394–4397 (1983).
410
411 30. Shao, S. *et al.* Evolution of body morphology and beak shape revealed by a morphometric
412 analysis of 14 Paridae species. *Front. Zool.* **13**, 30 (2016).
413
414 31. Ho, S. Y., Duchêne, S. & Duchêne, D. Simulating and detecting autocorrelation of molecular
415 evolutionary rates among lineages. *Mol. Ecol. Resour.* **15**, 688–696 (2015).
416
417 32. Lepage, T., Bryant, D., Philippe, H. & Lartillot, N. A general comparison of relaxed molecular
418 clock models. *Mol. Biol. Evol.* **24**, 2669–2680 (2007).
419

- 420 33. Erwin, D. H. *et al.* The Cambrian conundrum: early divergence and later ecological success
421 in the early history of animals. *Science* **334**, 1091–1097 (2011).
422
- 423 34. Dos Reis, M. *et al.* Uncertainty in the Timing of Origin of Animals and the Limits of Precision
424 in Molecular Timescales. *Curr. Biol.* **25**, 1–12 (2015).
425
- 426 35. Meredith, R. W. *et al.* Impacts of the Cretaceous Terrestrial Revolution and KPg extinction
427 on mammal diversification. *Science* **334**, 521–524 (2011).
428
- 429 36. Dos Reis, M. *et al.* Phylogenomic datasets provide both precision and accuracy in
430 estimating the timescale of placental mammal phylogeny. *Proc. R. Soc. B* **279**, 3491–3500
431 (2012).
432
- 433 37. Foster, C. S. *et al.* Evaluating the impact of genomic data and priors on Bayesian estimates
434 of the angiosperm evolutionary timescale. *Syst. Biol.* syw086 (2016).
435
- 436 38. Magallón, S., Hilu, K. W. & Quandt, D. Land plant evolutionary timeline: gene effects are
437 secondary to fossil constraints in relaxed clock estimation of age and substitution rates. *Am.*
438 *J. Bot.* **100**, 556–573 (2013).
439
- 440 39. Wikström, N., Savolainen, V. & Chase, M. W. Evolution of the angiosperms: calibrating the
441 family tree. *Proc. R. Soc. B* **268**, 2211–2220 (2001).
442
- 443 40. Hertweck, K. L. *et al.* Phylogenetics, divergence times and diversification from three
444 genomic partitions in monocots. *Bot. J. Linn. Soc.* **178**, 375–393 (2015).
445
- 446 41. Jarvis, E. D. *et al.* Whole-genome analyses resolve early branches in the tree of life of
447 modern birds. *Science* **346**, 1320–1331 (2014).
448
- 449 42. Liu, L. *et al.* Genomic evidence reveals a radiation of placental mammals uninterrupted by
450 the KPg boundary. *Proc. Natl. Acad. Sci. U.S.A.* **114**, E7282–E7290 (2017).
451
- 452 43. Dos Reis, M. *et al.* Using phylogenomic data to explore the effects of relaxed clocks and
453 calibration strategies on divergence time estimation: Primates as a test case. *Syst. Biol.*
454 (2018).
455
- 456 44. Battistuzzi, F. U., Filipowski, A., Hedges, S. B. & Kumar, S. Performance of relaxed-clock
457 methods in estimating evolutionary divergence times and their credibility intervals. *Mol. Biol.*
458 *Evol.* **27**, 1289–1300 (2010).
459
- 460 45. Dos Reis, M., Zhu, T. & Yang, Z. The impact of the rate prior on Bayesian estimation of
461 divergence times with multiple loci. *Syst. Biol.* **64**, 555–565 (2014).
462
- 463 46. Christin, P.-A. *et al.* Molecular dating, evolutionary rates, and the age of the grasses. *Syst.*
464 *Biol.* **63**, 153–165 (2014).
465
- 466 47. Bzdok, D., Krzywinski, M. & Altman, N. Machine learning: supervised methods. *Nat.*
467 *Methods* **15**, 5 (2018).
468
- 469 48. Tamura, K., Tao, Q. & Kumar, S. Theoretical foundation of the RelTime method for
470 estimating divergence times from variable evolutionary rates. *Mol. Biol. Evol.* msy044

- 471 (2018). doi:10.1093/molbev/msy044
472
473 49. Ekbatani, H. K., Pujol, O. & Segui, S. Synthetic Data Generation for Deep Learning in
474 Counting Pedestrians. In *Pattern Recognition Applications and Methods (ICPRAM), 2017*
475 *The International Conference on* 318–323 (2017).
476
477 50. Le, T. A., Baydin, A. G., Zinkov, R. & Wood, F. Using synthetic data to train neural networks
478 is model-based reasoning. In *Neural Networks (IJCNN), 2017 International Joint*
479 *Conference on* 3514–3521 (2017).
480
481 51. Beaulieu, J. M., O'Meara, B. C., Crane, P. & Donoghue, M. J. Heterogeneous rates of
482 molecular evolution and diversification could explain the Triassic age estimate for
483 angiosperms. *Syst. Biol.* **64**, 869–878 (2015).
484
485 52. Tamura, K. *et al.* Estimating divergence times in large molecular phylogenies. *Proc. Natl.*
486 *Acad. Sci. U.S.A.* **109**, 19333–19338 (2012).
487
488 53. Saitou, N. & Nei, M. The neighbor-joining method: a new method for reconstructing
489 phylogenetic trees. *Mol. Biol. Evol.* **4**, 406–425 (1987).
490
491 54. Kimura, M. A simple method for estimating evolutionary rates of base substitutions through
492 comparative studies of nucleotide sequences. *J. Mol. Evol.* **16**, 111–120 (1980).
493
494 55. Calteau, A. *et al.* Phylum-wide comparative genomics unravel the diversity of secondary
495 metabolism in Cyanobacteria. *BMC Genomics* **15**, 977 (2014).
496
497 56. Battistuzzi, F. U. & Hedges, S. B. A major clade of prokaryotes with ancient adaptations to
498 life on land. *Mol. Biol. Evol.* **26**, 335–343 (2009).
499
500 57. Jayaswal, V., Wong, T. K., Robinson, J., Poladian, L. & Jermiin, L. S. Mixture models of
501 nucleotide sequence evolution that account for heterogeneity in the substitution process
502 across sites and across lineages. *Syst. Biol.* **63**, 726–742 (2014).
503
504 58. Kishino, H., Thorne, J. L. & Bruno, W. J. Performance of a divergence time estimation
505 method under a probabilistic model of rate evolution. *Mol. Biol. Evol.* **18**, 352–361 (2001).
506
507 59. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586–
508 1591 (2007).
509
510 60. Kumar, S. & Subramanian, S. Mutation rates in mammalian genomes. *Proc. Natl. Acad. Sci.*
511 *U.S.A.* **99**, 803–808 (2002).
512
513 61. Kumar, S., Stecher, G., Peterson, D. & Tamura, K. MEGA-CC: computing core of molecular
514 evolutionary genetics analysis program for automated and iterative data analysis.
515 *Bioinformatics* **28**, 2685–2686 (2012).
516
517 62. Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular Evolutionary Genetics Analysis
518 version 7.0 for bigger datasets. *Mol. Biol. Evol.* **33**, 1870–1874 (2016).
519

- 520 63. Xie, W., Lewis, P. O., Fan, Y., Kuo, L. & Chen, M.-H. Improving marginal likelihood
521 estimation for Bayesian phylogenetic model selection. *Syst. Biol.* **60**, 150–160 (2011).
522
- 523 64. Baele, G., Li, W. L. S., Drummond, A. J., Suchard, M. A. & Lemey, P. Accurate model
524 selection of relaxed molecular clocks in bayesian phylogenetics. *Mol. Biol. Evol.* **30**, 239–
525 243 (2013).
526
- 527 65. Silvestro, D., Schnitzler, J. & Zizka, G. A Bayesian framework to estimate diversification
528 rates and their variation through time and space. *BMC Evol. Biol.* **11**, 311 (2011).
529
- 530 66. Rambaut, A., Suchard, M., Xie, D. & Drummond, A. Tracer v1.6. (2014). at
531 <<http://beast.bio.ed.ac.uk/Tracer>>
532
- 533 67. Kumar, S., Stecher, G., Suleski, M. & Hedges, S. B. TimeTree: A Resource for Timelines,
534 Timetrees, and Divergence Times. *Mol. Biol. Evol.* **34**, 1812–1819 (2017).
535
- 536 68. Hedges, S. B., Dudley, J. & Kumar, S. TimeTree: a public knowledge-base of divergence
537 times among organisms. *Bioinformatics* **22**, 2971–2972 (2006).
538
- 539 69. Ruhfel, B. R., Gitzendanner, M. A., Soltis, P. S., Soltis, D. E. & Burleigh, J. G. From algae to
540 angiosperms-inferring the phylogeny of green plants (Viridiplantae) from 360 plastid
541 genomes. *BMC Evol. Biol.* **14**, 23 (2014).
542
- 543 70. Misof, B. *et al.* Phylogenomics resolves the timing and pattern of insect evolution. *Science*
544 **346**, 763–767 (2014).
545
- 546 71. Wickett, N. J. *et al.* Phylotranscriptomic analysis of the origin and early diversification of land
547 plants. *Proc. Natl. Acad. Sci. U.S.A.* **111**, E4859–4868 (2014).
548
- 549 72. Shen, X.-X. *et al.* Reconstructing the Backbone of the Saccharomycotina Yeast Phylogeny
550 Using Genome-Scale Data. *G3* **6**, 3927–3939 (2016).
551
552
553

554 **Table 1.** Results from the CorrTest analysis of datasets from a diversity of species.

555

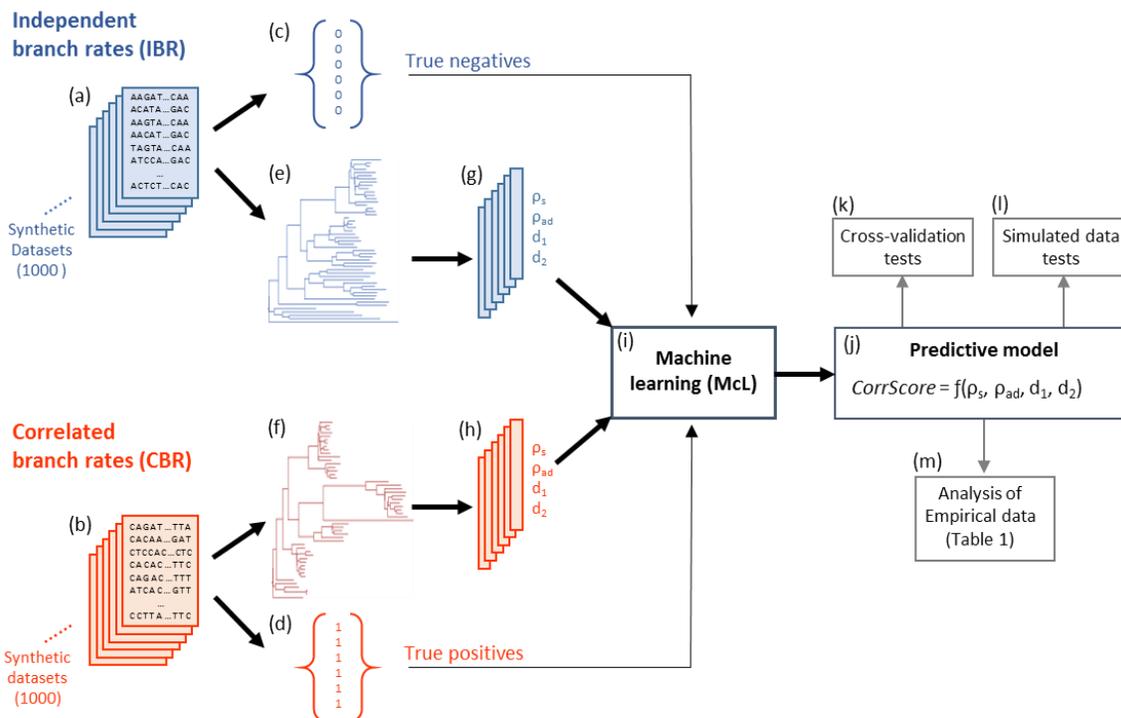
Group	Data type	Taxa number ^a	Sequence length	Substitution model	CorrTest score	P-value	1/v ^b	Reference
Mammals	Nuclear 4-fold degenerate sites	138	1,671	GTR + Γ	0.98	< 0.001	3.21	Meredith et al. (2011) ³⁵
Mammals	Nuclear 3 rd codon	138	11,010	GTR + Γ	0.99	< 0.001	4.42	Meredith et al. (2011) ³⁵
Mammals	Nuclear proteins	138	11,010	JTT + Γ	0.99	< 0.001	3.11	Meredith et al. (2011) ³⁵
Mammals	Mitochondrial DNA	271	7,370	HKY + Γ	0.98	< 0.001	3.77	Dos Reis, et al. (2012) ³⁶
Birds	Nuclear DNA	198	101,781	GTR + Γ	1.00	< 0.001	2.07	Prum et al. (2015) ¹⁵
Birds	Nuclear 3 rd codon	222	1,364	GTR + Γ	1.00	< 0.001	2.11	Claramunt et al. (2015) ¹⁶
Birds	Nuclear 1 st and 2 nd codon	222	2,728	GTR + Γ	1.00	< 0.001	2.53	Claramunt et al. (2015) ¹⁶
Insects	Nuclear proteins	143	220,091	LG + Γ	1.00	< 0.001	8.68	Misof et al. (2015) ⁷⁰
Metazoans	Mitochondrial & nuclear proteins	113	2,049	LG + Γ	0.65	< 0.05	40.0	Erwin et al. (2011) ³³
Plants	Plastid 3 rd codon	335	19,449	GTR + Γ	1.00	< 0.001	2.28	Ruhfel et al. (2014) ⁶⁹
Plants	Plastid proteins	335	19,449	JTT + Γ	1.00	< 0.001	2.46	Ruhfel et al. (2014) ⁶⁹
Plants	Nuclear 1 st and 2 nd codon	99	220,091	GTR + Γ	1.00	< 0.001	5.50	Wickett et al. (2014) ⁷¹
Plants	Chloroplast and nuclear DNA	124	5,992	GTR + Γ	1.00	< 0.001	2.64	Beaulieu et al. (2015) ⁵¹
Fungi	Nuclear proteins	85	609,772	LG + Γ	0.97	< 0.001	3.78	Shen et al. (2016) ⁷²
Prokaryotes	Nuclear proteins	197	6,884	JTT + Γ	0.79	< 0.05	2.54	Battistuzzi et al. (2009) ⁵⁶
Prokaryotes	Nuclear proteins	126	3,145	JTT + Γ	0.83	< 0.05	1.23	Calteau et al. (2014) ⁵⁵

556

557 ^aTaxa number is the number of ingroup taxa only.

558 ^b1/v is the inverse of the autocorrelation parameter that is estimated by MCMCTree with
 559 the autocorrelated rate model in the time unit of 100My.

560



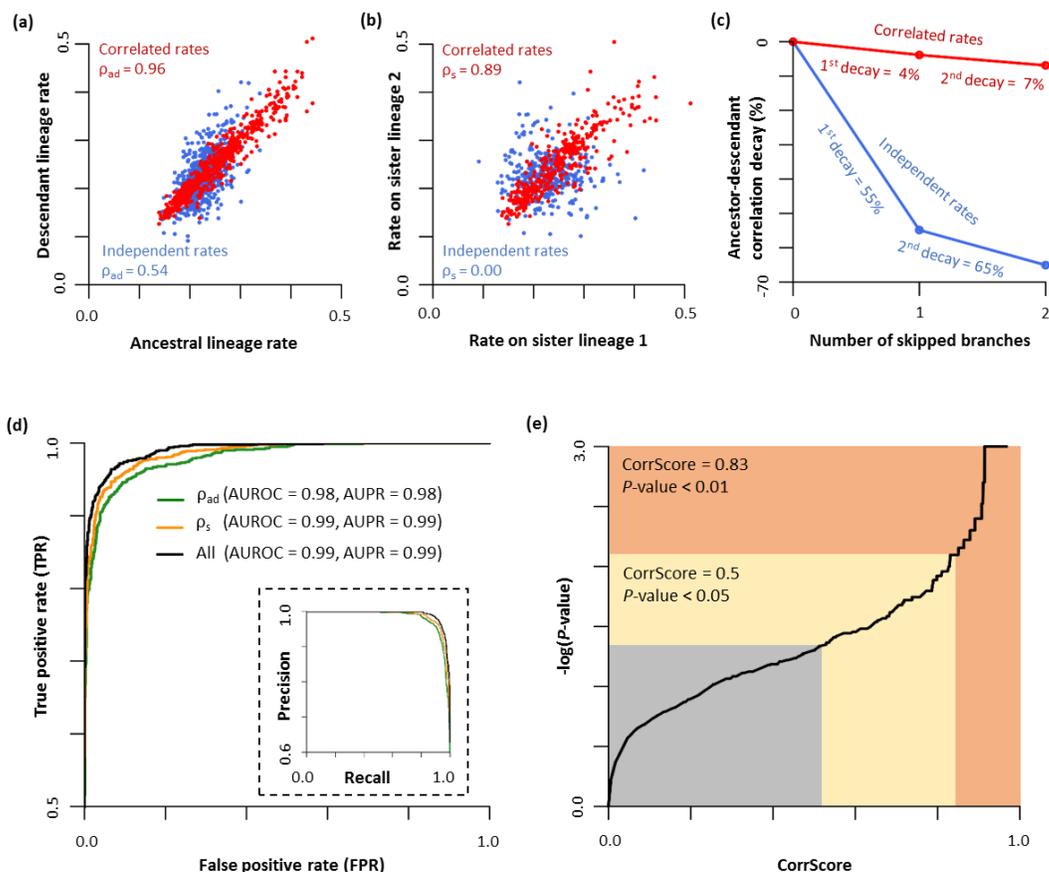
561

562

563 **Figure 1.** A flowchart showing an overview of the machine learning (MCL) approach
 564 applied to develop the predictive model (CorrTest). We generated **(a)** 1,000 synthetic
 565 datasets that were evolved using an IBR model and **(b)** 1,000 synthetic datasets that were
 566 evolved using a CBR model. The numerical label **(c)** for all IBR datasets was 0 and **(d)**
 567 for all CBR datasets was 1. For each dataset, we estimated a molecular phylogeny with
 568 branch lengths **(e and f)** and computed ρ_s , ρ_{ad} , d_1 , and d_2 **(g and h)** that served as features
 569 during the supervised machine learning. **(i)** Supervised machine learning was used to
 570 develop a predictive relationship between the input features and labels. **(j)** The predictive
 571 model produces a CorrScore for an input phylogeny with branch lengths. The predictive
 572 model was **(k)** validated with 10-fold and 2-fold cross-validation tests, **(l)** tested using
 573 external simulated data, and then **(m)** applied to real data to examine the prevalence of
 574 rate correlation in the tree of life.

575

576

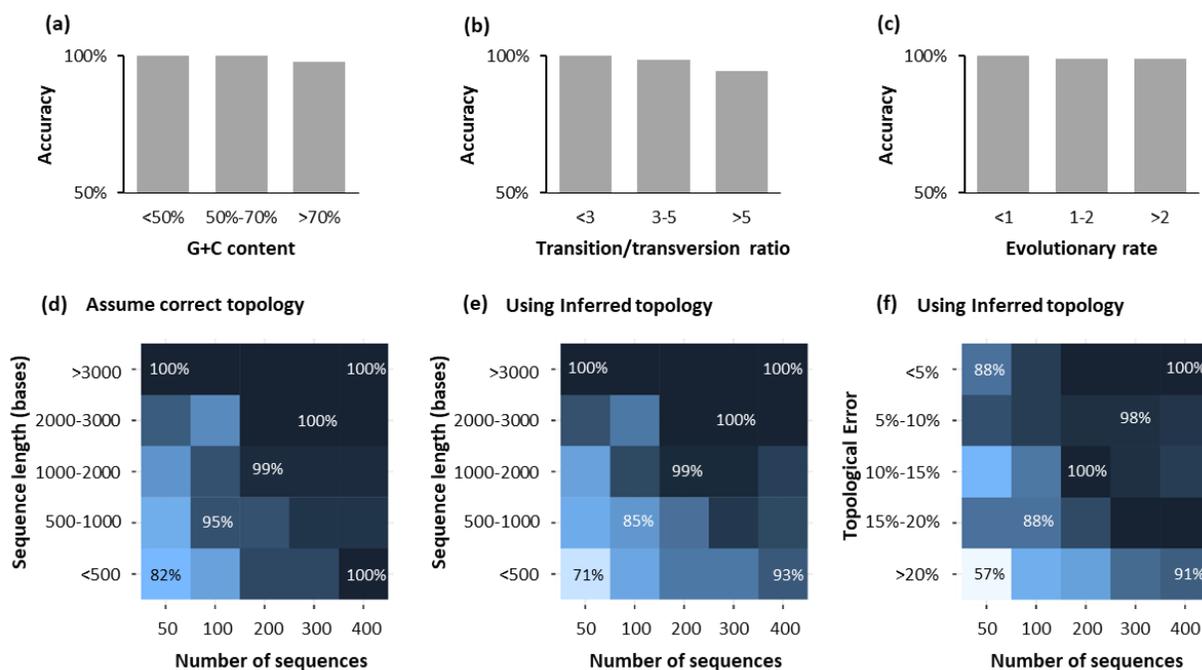


577

578

579 **Figure 2.** The relationship of **(a)** ancestral and direct descendent lineage rates and **(b)**
 580 sister lineage rates when the simulated evolutionary rates were correlated with each other
 581 (red) or varied independently (blue). The correlation coefficients are shown. **(c)** The decay
 582 of correlation between ancestral and descendant lineages when we skip one intervening
 583 branch (1st decay, d_1) and when we skip two intervening branches (2nd decay, d_2). Percent
 584 decay values are shown. **(d)** Receiver Operator Characteristic (ROC) and Precision
 585 Recall (PR) curves (inset) of the CorrTest for detecting branch rate model by using only
 586 ancestor-descendant lineage rates (ρ_{ad} , green), only sister lineage rates (ρ_s , orange), and
 587 all four features (all, black). The area under the curve is provided. **(e)** The relationship
 588 between the CorrScore produced by the machine learning model and the P -value. The
 589 null hypothesis of rate independence can be rejected when the CorrScore is greater than
 590 0.83 at a significant level of $P < 0.01$, or when the CorrScore is greater than 0.5 at $P <$
 591 0.05.

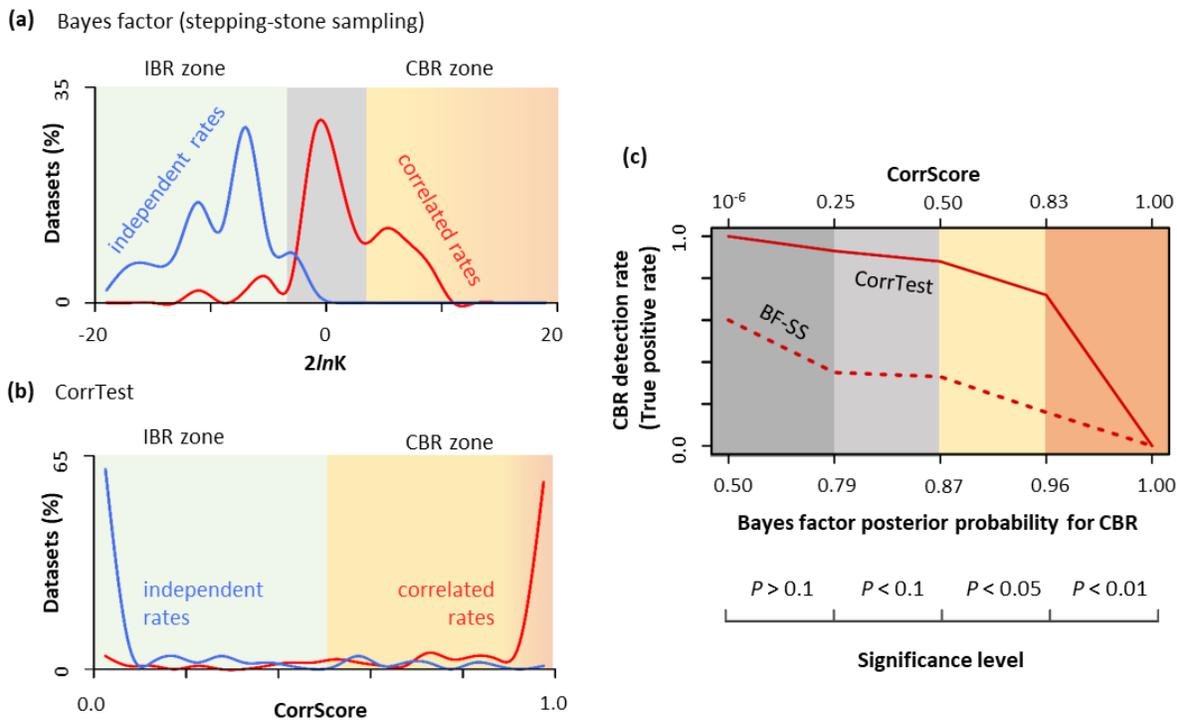
592



593

594 **Figure 3.** The performance of CorrTest in detecting rate correlation in the analysis of
 595 datasets⁵² that were simulated with different (a) G+C contents, (b) transition/transversion
 596 rate ratios, and (c) average molecular evolutionary rates. Darker color indicates higher
 597 accuracy. The evolutionary rates are in the units of 10^{-3} substitutions per site per million
 598 years. (d – f) Patterns of CorrTest accuracy for datasets containing increasing number of
 599 sequences. The accuracy of CorrTest for different sequence length is shown when (d)
 600 the correct topology was assumed and (e) the topology was inferred. (f) The accuracy of
 601 CorrTest for datasets in which the inferred the topology contained small and large number
 602 of topological errors.

603



604

605

606 **Figure 4.** Comparisons of the performance of CorrTest and Bayes Factor analyses. **(a)**

607 Distributions of 2 times the differences of marginal log-likelihood ($2\ln K$) estimated via

608 stepping-stone sampling method for datasets that were simulated with correlated branch

609 rates (CBR, red) and independent branch rates (IBR, blue). CBR is preferred ($P < 0.05$)

610 when $2\ln K$ is greater than 3.841 (CBR zone), and IBR is preferred when $2\ln K$ is less than

611 -3.841 (IBR zone). When $2\ln K$ is between -3.841 and 3.841, the fit of the two rate models

612 is not significantly different (gray shade). **(b)** The distributions of CorrScores in analyses

613 of CBR (red) and IBR (blue) datasets. Rates are predicted to be correlated if the

614 CorrScore is greater than 0.5 ($P < 0.05$, CBR zone) and vary independently if the

615 CorrScore is less than 0.5 (IBR zone). **(c)** The rate of detecting CBR model correctly (True

616 Positive Rate) at different levels of statistical significance in Bayes factor (stepping-stone

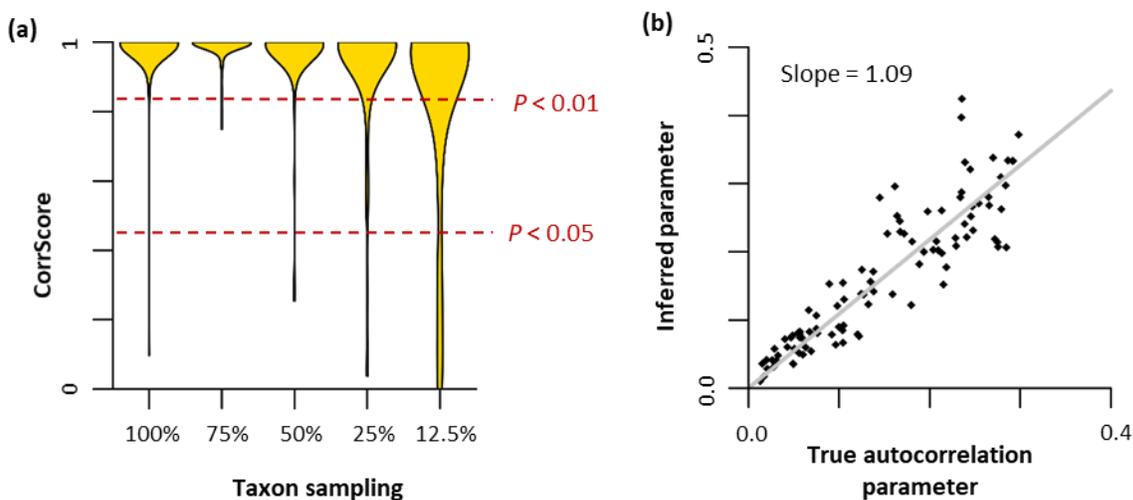
617 sampling) and CorrTest analyses. Posterior probabilities for CBR in BF-SS analysis are

618 derived using the log-likelihood patterns in panel a. CorrTest P -values are derived using

619 the CorrScore pattern in panel b.

620

621
622
623
624



625
626
627
628
629
630
631
632
633

Figure 5. (a) The distribution of CorrScore when data have different taxon sampling densities. The CorrScore decreases when the density of taxon sampling is lower, as there is much less information to discriminate between CBR and IBR. Red, dashed lines mark two statistical significance levels of 5% and 1%. **(b)** The relationship between the inferred autocorrelation parameter from MCMCTree and the true value. The gray line represents the best-fit regression line, which has a slope of 1.09.