1   *Tao et al.   (Supplementary Information)*

2   *Pervasive correlation of molecular evolutionary rates in the tree of life*

3

4   ***Machine learning model building***

5   ***Training data.*** We simulated nucleotide alignments using independent branch rate (IBR)

6   and correlated branch rate (CBR) models using the NELSI package[1]. In IBR, branch-

7   specific rates were drawn from a lognormal distribution with a mean gene rate and a

8   standard deviation (in log-scale) that varied from 0.1 to 0.4, previously used in a study

9   simulating independent rates with different levels of variation[1]. In CBR, branch-specific

10  rates were simulated under an autocorrelated process[2] with an initial rate set as the mean

11  rate derived from an empirical gene and an autocorrelated parameter, $v,$ that was

12  randomly chosen from 0.01 to 0.3, previously used in a study simulating low, moderate

13  and high degrees of autocorrelated rates[1]. We used SeqGen[3] to generate alignments

14  under Hasegawa-Kishino-Yano (HKY) model[4] with 4 discrete gamma categories by using

15  a master phylogeny, consisting of 60-400 ingroup taxa randomly sampled from the bony-

16  vertebrate clade in the Timetree of Life[5]. Mean evolutionary rates, G+C contents,

17  transition/transversion ratios and numbers of sites for simulation were derived from

18  empirical distributions[6]. These 2,000 simulated datasets were used as training data in

19  building the machine learning model.

20  ***Features   acquisition***. Lineage-specific   rate   estimates   ($r_i$'s)   were   obtained   using

21  equations [28] - [31] and [34] - [39] in Tamura et al. (2018)[7]. As mentioned in the main

22  text, a lineage rate is a function of all the branch rates that belong to that lineage. For any

23  given node in the phylogeny, we extracted the relative rates of its ancestral clade ($r_a$) and

24  two   direct   descendant   clades   ($r_1$   and   $r_2$).   Then,   we   calculated   correlation   between

25  ancestral lineage and its direct descendant lineage rate to obtain estimates of ancestor-

26  descendant rate correlation ($\rho_{ad}$). We also calculated correlation between sister lineage

27  rates ($\rho_s$), for which the lineage rates of sister pairs are randomly labeled. The labeling of

28  sister pairs have small impact on $\rho_s$ when the number of sequences in the phylogeny is

29  not too small (>50). However, one can also choose to resample sister pairs for multiple

30    times and use the mean of resampled $\rho_s$ in the CorrTest in order to eliminate any bias

31    that may result from the arbitrary designation of sister rates during the correlation process,

32    which can be a problem when the number of taxa is small. To avoid the assumption of

33    linear correlation between lineages, we used Spearman rank correlation because it can

34    capture both linear and non-linear correlation between two vectors. Two additional

35    features derived from the relative rates in the phylogeny were used in building the

36    machine learning model. We first estimated $\rho_{ad\_skip1}$ as the correlation between rates

37    where the ancestor and descendant were separated by one intervening branch, and

38    $\rho_{ad\_skip2}$ as the correlation between rates where the ancestor and descendant were

39    separated by two intervening branches. This skipping reduces ancestor-descendant

40    correlation, which we then used to derive the decay of correlation values by using

41    equations $(\rho_{ad} - \rho_{ad\_skip1})/\rho_{ad}$ and $(\rho_{ad} - \rho_{ad\_skip2})/\rho_{ad}$. These two features improved the

42    accuracy of our model slightly. In the analyses of empirical datasets, we found that a large

43    amount of missing data (>50%) can result in unreliable estimates of branch lengths and

44    other phylogenetic errors[8–12]. In this case, we recommend computing selected features

45    using only those lineage pairs for which >50% of the positions contain valid data, or

46    remove sequences with a large amount of missing data.

47    ***Predictive model.*** We trained a logistic regression model using the skit-learn module[13],

48    which is a python toolbox for data mining and data analysis using machine learning

49    algorithms, with only $\rho_{ad}$, only $\rho_s$ or all 4 features ($\rho_s$, $\rho_{ad}$, the two decay of correlation

50    features) using 2,000 simulated training datasets (1,000 with CBR model and 1,000 with

51    IBR model). A response value of 1 was given to true positive cases (correlated rates) and

52    0 was assigned to true negative cases (independent rates). Thus, the prediction scores

53    (CorrScore) were between 0 and 1. A high score representing a higher probability that

54    the rates are correlated. Then the global thresholds at 5% and 1% significant levels can

55    be determined. To explore the reliability of the global threshold, we re-trained the model

56    with all 4 features extracted from 4 subsets of training data with ≤ 100 (M100), 100 – 200

57    (M200), 200 – 300 (M300), and > 300 (M400) sequences. A specific threshold was

58    determined for each training subset and then was tested using Tamura *et al.* (2012)'s

59    data[14] with the corresponding size. For example, we used the threshold determined by

60    the model trained with small data (≤ 100 sequences) on the test data that contain less

61    than 100 sequences, and used the threshold determined by the model trained with large

62    data (>300 sequences) on the large test data (400 sequences). We found that the

63    accuracy of using the specific thresholds (**Fig. S1a-c**) is similar to the accuracy when we

64    used a global threshold (**Fig. 3d-f**). This is because the machine learning algorithm has

65    automatically incorporated the impact of the number of sequences when it determined

66    the relationship of four selected features ($\rho_{ad}$, $\rho_s$, and 2 decays).

### *Cross-validation*

68    We performed two cross-validation tests. In 10-fold cross-validation, the predictive model

69    was developed using 90% of the synthetic datasets, and then its performance was tested

70    on the remaining 10% of the datasets. The AUROC was greater than 0.99 and the

71    accuracy was high (>94%). Even in the 2-fold cross-validation, where only half of the

72    datasets were used for training the model and the remaining half were used for testing,

73    the AUROC was still greater than 0.99 with an accuracy greater than 92%. This indicates

74    that the features we used in building the machine learning model are powerful and

75    ensures high accuracy even when the training data are limited.

### *External tests*

77    ***Publicly available data.*** Two previously published simulated dataset were used to

78    evaluate CorrTest's performance. Beaulieu et al.'s data[15] contains 91 ingroup taxa with

79    1,000 base pairs each. For Tamura et al.'s data[14], we present the test results for the data

80    simulated using CBR model (autocorrelated lognormal distribution) and IBR model

81    (independent uniform distribution with 50% rate variation) here. We tested the

82    performance of our model on CBR and IBR data with different GC contents,

83    transition/transversion ratios, and evolutionary rates. We randomly sampled 50, 100, 200,

84    and 300 sequences from the original 400 sequences and conducted CorrTest using the

85    correct, error-prone topology inferred by the Neighbor-joining method[16] with an

86    oversimplified substitution model[17]. We also tested CorrTest's performance on data

87    simulated under an IBR model process with 100% rate variation and found that CorrTest

88    works perfectly (100% accuracy; results not shown).

89    ***Synthetic data***. We conducted another set of simulations using IBR (independent

90    lognormal distribution) and CBR (autocorrelated lognormal distribution)[2] model with 100

91    replicates each using the same strategy as a training data simulation (described above)

92    on a master phylogeny of 100 taxa randomly sampled from the bony-vertebrate clade in

93    the Timetree of Life[5]. These 200 datasets were used to conduct CorrTest and Bayes

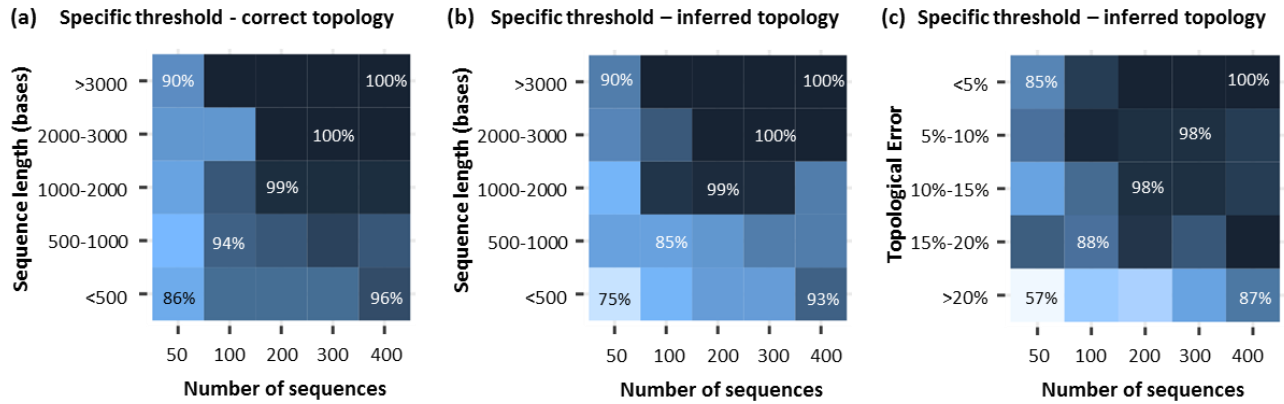94    factor analyses and to obtain the autocorrelation parameter ($v$) in MCMCTree[18].

95

## References

97    1.   Ho, S. Y., Duchêne, S. & Duchêne, D. Simulating and detecting autocorrelation of
98         molecular evolutionary rates among lineages. *Mol. Ecol. Resour.* **15,** 688–696
99         (2015).

101   2.   Kishino, H., Thorne, J. L. & Bruno, W. J. Performance of a divergence time
102        estimation method under a probabilistic model of rate evolution. *Mol. Biol. Evol.* **18,**
103        352–361 (2001).

105   3.   Grassly, N. C., Adachi, J. & Rambaut, A. Seq-Gen: an application for the Monte
106        Carlo simulation of protein sequence evolution along phylogenetic trees. *Comput.*
107        *Appl. Biosci.* **13,** 235–238 (1997).

109   4.   Hasegawa, M., Kishino, H. & Yano, T. Dating of the human-ape splitting by a
110        molecular clock of mitochondrial DNA. *J. Mol. Evol.* **22,** 160–174 (1985).

112   5.   Hedges, S. B. & Kumar, S. *The Timetree of Life*. pp3–18 (New York: Oxford
113        University Press., 2009).

115   6.   Rosenberg, M. S. & Kumar, S. Heterogeneity of nucleotide frequencies among
116        evolutionary lineages and phylogenetic inference. *Mol. Biol. Evol.* **20,** 610–621
117        (2003).

119   7.   Tamura, K., Tao, Q. & Kumar, S. Theoretical foundation of the RelTime method for
120        estimating divergence times from variable evolutionary rates. *Mol. Biol. Evol.*
121        msy044 (2018). doi:10.1093/molbev/msy044

123   8.   Filipski, A., Murillo, O., Freydenzon, A., Tamura, K. & Kumar, S. Prospects for
124        building large timetrees using molecular data with incomplete gene coverage among
125        species. *Mol. Biol. Evol.* **31,** 2542–2550 (2014).

127   9.   Xi, Z., Liu, L. & Davis, C. C. The impact of missing data on species tree estimation.
128        *Mol. Biol. Evol.* **33,** 838–860 (2015).

130   10. Lemmon, A. R., Brown, J. M., Stanger-Hall, K. & Lemmon, E. M. The effect of
131        ambiguous data on phylogenetic estimates obtained by maximum likelihood and
132        Bayesian inference. *Syst. Biol.* **58,** 130–145 (2009).

134  11. Wiens, J. J. & Moen, D. S. Missing data and the accuracy of Bayesian
135      phylogenetics. *Journal of Systematics and Evolution* **46,** 307–314 (2008).
136

137  12. Marin, J. & Hedges, S. B. Undersampling genomes has biased time and rate
138      estimates throughout the tree of life. *Mol. Biol. Evol.* msy103 (2018).
139      doi:doi.org/10.1093/molbev/msy103
140

141  13. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.*
142      **12,** 2825–2830 (2011).
143

144  14. Tamura, K. *et al.* Estimating divergence times in large molecular phylogenies. *Proc.*
145      *Natl. Acad. Sci. U.S.A.* **109,** 19333–19338 (2012).
146

147  15. Beaulieu, J. M., O'Meara, B. C., Crane, P. & Donoghue, M. J. Heterogeneous rates
148      of molecular evolution and diversification could explain the Triassic age estimate for
149      angiosperms. *Syst. Biol.* **64,** 869–878 (2015).
150

151  16. Saitou, N. & Nei, M. The neighbor-joining method: a new method for reconstructing
152      phylogenetic trees. *Mol. Biol. Evol.* **4,** 406–425 (1987).
153

154  17. Kimura, M. A simple method for estimating evolutionary rates of base substitutions
155      through comparative studies of nucleotide sequences. *J. Mol. Evol.* **16,** 111–120
156      (1980).
157

158  18. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24,**
159      1586–1591 (2007).

160



161

162

163 **Figure S1**. Patterns of CorrTest accuracy using M100, M200, M300, and M400 models

164 for the corresponding test datasets[14]. Accuracies are shown for increasing number of

165 sequences. The accuracy of CorrTest for different sequence length is shown when (**a**)

166 the correct topology was assumed and (**b**) the topology was inferred. (**c**) The accuracy of

167 CorrTest for datasets in which the inferred the topology contained small and large number

168 of topological errors.

169

170