Peer Community Journal

Section: Genomics

Research article

Published 2025-05-15

Cite as

Guillaume Louvel and Hugues Roest Crollius (2025) Factors influencing the accuracy and precision in dating single gene trees, Peer Community Journal, 5: e51.

Correspondence

guillaume.louvel@normalesup.org

Peer-review

Peer reviewed and recommended by PCI Genomics, https://doi.org/10.24072/pci. genomics.100292

(cc) BY

This article is licensed under the Creative Commons Attribution 4.0 License.

Factors influencing the accuracy and precision in dating single gene trees

Guillaume Louvel^{⁰,1,2} and Hugues Roest Crollius^{⁰,1}

Volume 5 (2025), article e51

https://doi.org/10.24072/pcjournal.556

Abstract

Molecular dating is the inference of divergence time from genetic sequences. Knowing the time of appearance of a taxon sets the evolutionary context by connecting it with past ecosystems and species. Knowing the divergence times of gene lineages would provide a context to understand adaptation at the genomic level. However, molecular clock inference faces uncertainty due to the variability of the rate of substitution between species, between genes, and between sites within genes. When dating speciations, per-lineage rate variability can be informed by fossil calibrations, and gene-specific rates can be either averaged out or modeled by concatenating multiple genes. By contrast, when dating genespecific events, fossil calibrations only inform about speciation nodes, and concatenation does not apply to divergences other than speciations. This study aims to benchmark the accuracy of molecular dating applied to single gene trees and identify how it is affected by gene tree characteristics. We analyze 5205 alignments of genes from 21 Primates in which no duplication or loss is observed. We also simulated alignments based on characteristics from Primates under a relaxed clock model to analyze the dating accuracy. Divergence times were estimated with the Bayesian program Beast2. From the empirical dataset, we find that the date estimates deviate more from the median age with shorter alignments, high rate heterogeneity between branches, and low average rate, features that underlie the amount of dating information in alignments, hence, statistical power. The smallest deviation is associated with core biological functions such as ATP binding and cellular organization, categories that are expected to be under strong negative selection. We then investigated the accuracy of dating with simulated alignments, by controlling the three above parameters separately. It confirmed the factors of precision, but also revealed biases when branch rates are highly heterogeneous. This suggests that in the case of the relaxed uncorrelated molecular clock, biases arise from the tree prior when calibrations are lacking and rate heterogeneity is high. Our study finally reports the scale of the gene tree features that influence the dating consistency with median ages, so that comparisons can be made with other genes and taxa. To tackle the molecular dating of events only observed in single gene trees, like deep coalescence, horizontal gene transfers, and gene duplications, future models should overcome the lack of power due to limited information from single genes.

¹École normale supérieure, PSL Research University, CNRS, Inserm, Institut de Biologie de l'École normale supérieure (IBENS), F-75005 Paris, France, ²Centre for Anthropobiology and Genomics of Toulouse, CNRS UMR5288/Université de Toulouse, Toulouse, France



Peer Community Journal is a member of the Centre Mersenne for Open Scientific Publishing http://www.centre-mersenne.org/

e-ISSN 2804-3871



Introduction

The phylogenetic tree is the most prevalent concept for representing evolutionary relationships, but a species tree and a gene tree do not necessarily tell the same story: while nodes in a species tree represent speciations, meaning the divergence of sister species from their parent, nodes in a gene tree represent the divergence of genes which do not result directly from speciation, but from the coalescence of alleles or from locus duplication and horizontal transfer. Therefore, reconstructing a species tree generally requires sequences assumed to be orthologous and congruent with the species history; subsequent inferences such as dating taxon divergences usually lump together many genes to increase statistical power, for example by concatenating orthologs. However, the specificities of gene evolution are not only a nuisance that hinders species diversification. In this regard, focusing on single genes instead of concatenations can reveal gene-specific evolutionary processes. In particular, genes have heterogeneous substitution dynamics due to functional specificities or location in the genome, or as a consequence of events of duplication and transfer. Estimating the rate and timing of sequence evolution is the subject of molecular clock dating and we evaluate here its accuracy and precision when applied to single genes.

The molecular clock is a set of methods to date the divergence between DNA or protein sequences by linking absolute time to the number of observed substitutions (Zuckerkandl & Pauling, 1962, 1965; Margoliash, 1963; Doolittle & Blombäck, 1964). The model considers a sequence that evolves by substitution over a given lapse of time *t*. It states an approximately proportional relationship between time *t* and the number of substitutions *N*, by modelling *N* as a Poisson distributed random variable of expectation $r \times t$, where *r* is the rate of substitutions per unit of time. Thus an average of $r \times t$ substitutions occur in the duration *t*. In the simplest model the rate is assumed to be constant over the entire phylogenetic tree (Langley & Fitch, 1974). The number of substitutions can be obtained from the branch lengths inferred with classic phylogenetic tree inference programs (PhyML, Iqtree, RaxML, Phylobayes, MrBayes, PAUP, etc), but to infer the rate, times must be known, and vice-versa. Fossils typically provide independent dating and are used to assign ages to some of the speciation nodes, allowing the rate to be derived and all non-calibrated nodes to be dated.

However, while rate constancy is a convenient simplification, it has been observed that rates of substitution per unit of time vary across taxa (Wu & Li, 1985; Britten, 1986; Pagel et al., 2006). In particular, when considering a single site in an alignment, the across branch rate variation is called heterotachy (literally "different rates"; Philippe et al., 2003). The causes of variations are diverse and may include population size, impact of selection at the molecular level, generation time, efficiency of the DNA repair system, metabolism, etc (Gillespie, 1991). For some of these reasons, rates also vary across the genome (Wolfe et al., 1989). At the scale of a single sequence, the heterogeneity of the rate across branches does not necessarily follow the same pattern between sites, meaning that different sites accelerate or decelerate in an independent manner. This results in a complex pattern of variation across sites and genes and across branches. Some authors (Kumar & Hedges, 1998) have argued that by selecting many genes these variations are attenuated, their mean rate is approximately constant and the molecular clock can still be used, especially in recent clades. Elsewhere it is generally accepted that such simple models can generate overconfident but wrong results (Revell et al., 2005; Wertheim et al., 2010; Duchêne et al., 2015; Lozano-Fernandez et al., 2017).

Initially, tests of constancy such as the relative rate tests (e.g. in Sarich & Wilson, 1966; Fitch & Langley, 1976; Gu & Li, 1992; Tajima, 1993) or the likelihood-ratio test of the clock (Felsenstein, 1981) were developed to assess these rate variations. Nowadays, most models accommodate a non-constant molecular clock providing that enough calibrations are specified. How to relax the clock constancy is not settled: autocorrelated rates assume a progressive change between connected branches (Thorne et al., 1998; Kishino et al., 2001) while uncorrelated clock rates just follow an overall distribution across branches (Drummond et al., 2006). In the uncorrelated case, the

branch-wise rates are most commonly modeled as a lognormal, gamma or exponential distribution, and the branch lengths in units of substitution are usually obtained by multiplying these rates by the absolute time difference. As a consequence, the variance of the branch length scales quadratically with the time difference. Alternatively, the "white noise" model is totally uncorrelated at all times (i.e. within branches) and has the interesting property that the variance of the branch length scales only linearly with the time difference (Lepage et al., 2007). Different empirical studies support either the uncorrelated rates model (Drummond et al., 2006; Rannala & Yang, 2007; Linder et al., 2011; Heath et al., 2012) or the autocorrelated clock rates (Thorne et al., 1998; Lepage et al., 2007; Smith et al., 2018; dos Reis et al., 2018), while a mixture of the two models over different timescales has also been proposed (Lartillot et al., 2016; Bletsa et al., 2019). However, it is uncertain that autocorrelation can actually be detected (Linder et al., 2011; Ho et al., 2015; Tao et al., 2019) and the modalities of the rate relaxation are likely specific to each taxonomic group, timescale studied and calibration scheme. Algorithmic implementations include widely used software such as Beast 2 (Bouckaert et al., 2019) and MCMCTree (Rannala & Yang, 2007) which exploit the Bayesian framework through Markov Chain Monte Carlo (MCMC) sampling, while some faster alternatives take advantage of least squares strategies, like LSD (To et al., 2016) and treedater (Volz & Frost, 2017). Accelerating probabilistic algorithm has also been made possible by using approximate likelihoods (Thorne et al., 1998; Guindon et al., 2010; dos Reis & Yang, 2011).

In contrast to the early ages of molecular biology in the nineteen-sixties where analyses could only be performed on a single protein at a time such as haemoglobins, cytochrome c or fibrinopeptides (Margoliash, 1963; Doolittle & Blombäck, 1964; Zuckerkandl & Pauling, 1965), modern highthroughput sequencing gives access to hundreds or thousands of genes per genome, providing a more general insight. The modern procedure to date species divergences therefore uses many loci whose sequences are concatenated. This concatenation approach aims to increase the amount of information, thus the precision (Bromham et al., 2000; Lanfear et al., 2010; Smith et al., 2018; Bletsa et al., 2019). Even in a favorable situation with many genes, the accuracy of molecular dating has been questioned, because of flawed procedures such as reusing inferred dates as secondary calibrations without propagating uncertainty (Graur & Martin, 2004; Schenk, 2016), but also because of inherent limitations of the problem (Pulquério & Nichols, 2007; Burbrink & Pyron, 2008; dos Reis & Yang, 2013; Zhu et al., 2015; Kumar & Hedges, 2016; Warnock et al., 2017); the most critical limitations are the difficulty of characterizing the type of clock relaxation and the uncertainty in calibration points themselves (reviewed in Ho, 2014; dos Reis et al., 2015; Mello & Schrago, 2024). Assessing the adequacy of different models of rate variation has been investigated by a number of simulation-based studies (Aris-Brosou & Yang, 2002; Ho et al., 2005; Rannala & Yang, 2007; Battistuzzi et al., 2010; dos Reis et al., 2014; Duchêne et al., 2015) but with a focus on multiple-loci datasets and speciation dates.

If accurate dating requires sufficient sequence information as well as calibrations, this means that inferring dates in single gene or transposable element evolutionary histories is even more difficult than with concatenations of hundreds of genes. Indeed, on such short sequences, low statistical power amplifies the following known limitations of dating by molecular clock: methodological errors in tree topologies and substitution rates, rate variation of multiple biological origins and inherent stochasticity in the substitution process, all of which create a disproportionate amount of noise.

Here we assess molecular clock dating from single gene trees taken individually, and what characteristics of a gene tree are related to precision and accuracy. We first evaluate a dataset of primate genes under a cross-validation procedure where we date speciation events in each tree and compare each individual tree estimate with the median age over all trees. Using the median age as a point of reference produces a measure of the precision of dating. In these primate gene trees we find characteristics correlated with the deviation from the median, mainly the length of the sequence alignment, the heterogeneity of the rate between branches and the mean rate of the tree. In addition to measuring the deviation from the median we also compare with reference ages from TimeTree, and observe a bias towards younger ages. Because several explanations for such bias cannot be disentangled from the empirical dataset, we simulated gene trees to also measure

how the accuracy of dating vary depending on alignment length, rate heterogeneity between branches and mean rate of substitution.

Results

Single gene estimates show a high dispersion and a bias towards younger ages

We perform a benchmark on dating speciation nodes from single gene trees with an empirical dataset of trees that include genes from 21 primate species obtained from Ensembl version 93. We focus on the Simiiformes clade as ingroup, which contains all Primates but the Lemuriformes that we use as outgroup. Speciation ages for all nodes have been estimated many times independently based on fossil-calibrated dating on concatenation of genes, and resulting consensus ages can be obtained from TimeTree (Kumar et al., 2017). Because our aim is to replicate the uncertainty in dating nodes that lack fossil calibrations, we do not set any calibration except for the Similformes ancestor at 43.2 My (C.I [41.0, 45.7]). We then quantify the uncertainty relative to the surrounding interval of 43.2 My. The choice of this root calibration and its associated uncertainty is arbitrary because all trees are then compared by this yardstick. Likewise, the choice of Primates for the source trees is arbitrary; the specific selection of species does not matter, what matters is that we collect natural replicates of the same tree. We selected 5205 gene trees that do not display any duplication or loss in this tree of 21 primate species, so that each gene tree shares exactly the same topology. These were dated with Beast 2, providing 5205 age estimates at each internal node (Figure 1). Because the gene trees were independently subjected to dating, and the dated output trees were all scaled to 43.2 My, their differences in mean rate (i.e. the "gene effect", Ho 2014) are not considered in this part. The impact of genes specific features is investigated in the subsequent part.

By comparison with the reference age estimates from TimeTree, we obtain clearly younger estimates on all nodes, sometimes even outside of the confidence interval from the reference (median, Figure 1). The highest shift of 10 My affects Cercopithecidae. This suggests that either our estimates are biased, or less plausibly that the reference ages themselves are inaccurate. In our dating procedure, several simplifications may bias ages. To start with, calibrating only one node is unlikely to lead to accurate variable rates, but this is precisely the purpose of our analysis, since we study gene trees for which nodes lack calibrations. With one calibration but variable branch rate, rate and time are unknown because they are conflated into the branch length. In fact, molecular clock "dating" is as much about estimating rates using calibrations than estimating dates from rates, this shift in focus being due to the high variability of molecular rates. Without many calibrations, variation in rates cannot be faithfully inferred, and in turn dates remain uncertain. Accordingly, our empirical dataset is expected to display heterogeneous branch rates of different origins. First there are species-specific rate variations, in which all genes experience the same trend in a specific lineage (also called lineage effects; Gillespie 1989; Muse and Gaut 1997; Smith and Eyre-Walker 2003). For example, the average branch length of Cercopithecidae is higher than their Hominoidea sister (supp. info. S1). This is consistent with the known generation times in these clades: based on Pacifici et al. (2013) the averages are 4035 days for Cercopithecidae versus 6132 for Hominoidea. In addition to species trends, each gene tree may experience particular variations of branch rates, therefore producing dispersed ages when compared (called residual effects, i.e. variation remaining after gene and lineage effects; Gillespie 1991). Finally, independent across-branch rate variation also likely occurs within genes, between sites. However, these gene and site heterotachies should just cause dispersion but not loss of accuracy, as we expect that the distortions on branch lengths should compensate themselves on average (after correction for lineage specific rates). Another simplification of our inference is to consider instantaneous segregation of genes at speciations. In reality, speciation takes up to several million years and segregated genes can correspond to older allelic divergence. This phenomenon of deep coalescence should cause speciations to appear older than they are, and not younger as is the case in our results. Conversely, shallow coalescence caused by introgression between recently diverged species would produce younger age, but we deem implausible that introgression would be so pervasive

among genes so as to strongly bias ages. In summary, variation in lineage specific rates is the most plausible explanation for our bias towards younger ages, unless the inference obtained from Beast is not statistically consistent.



Figure 1 - Distribution of speciation ages estimated on primate gene trees independently (5205 trees). Each histogram is rescaled to display the same height, but their total count is each equal to 5205. Dating was performed with Beast under the nucleotide substitution model HKY and a relaxed clock rate (see Methods).

Aside from comparisons with the reference ages, we quantified the dispersion of the estimates at each node with the following metrics: the mean absolute deviation from the median (MAD) of the node age and the 95% inter-quantile range (IQR95); the averages of these metrics over the 12 speciations yield 3 My and 16 My, respectively. These dispersions are larger for deeper nodes, such as Catarrhini with a MAD of 5.6 My and an IQR95 of 32.4 My. Divided by the Simiiformes calibration of 43.2 My, they represent 13 % and 75 % of this interval, respectively. In other words, the confidence interval for the age of Catarrhini is barely shorter than the age of Simiiformes.

In the following, we break down the causes leading to the dispersion of our estimates by measuring gene features and searching for a relationship with the deviation from the median ages.

The dispersion of age estimates is associated with low statistical signal and branch rate variation

We computed 56 features from each gene tree related to alignment quality and substitution rate estimations (supp. info. S2). We regressed them against the absolute deviation from the median age averaged across all speciation nodes, chosen to describe the dating imprecision in a single gene tree. In order to select a restricted set of independent variables we used a Lasso regression with a regularization parameter 'alpha' of 0.02, which retained 10 parameters further fitted by Ordinary Least Squares (adjusted $R^2 = 0.288$, Figure 2). Variables were normalized so that coefficients can be compared and they were centered. Apparent outliers, such as trees with extreme rates, were discarded from the regression (see Methods), which was therefore based on 5170 trees.





By order of coefficient size, the length of the sequence alignment comes first (coefficient 0.35), associating short alignments with high dispersion. It is followed by the rate heterogeneity across branches which is positively associated with the dispersion (0.30). Next comes the mean rate of substitution (-0.14). We hypothesize two distinct causes for these correlations: first, the strength of the statistical signal is linked to the alignment length and to the mean rate, because both influence the number of observable substitutions, thus the amount of information available to infer ages; also because of rescaling all trees to the same height of 43.2 My, trees with the lower mean rate are expected to display a higher dispersion. Second, some model parameters are in practice difficult to infer individually. In particular, a variability in branch rates cannot be estimated when there is only one calibration, meaning that a branch specific rate could take arbitrary values as long as the product of rate by time equals the branch length in substitutions. Consequently, it would take values dictated by the tree prior (the Birth-Death branch process) that cannot possibly recover the exact Similformes speciation dates (dos Reis & Yang, 2013). Here the rate heterogeneity is the standard deviation of the rate across branches, as inferred by Beast. This indicates that while Beast manages to detect a high-rate heterogeneity in a given gene tree, the ages it infers in such a tree are guite distinct from the median ages, and probably far from their true age. The fourth largest association is a measure of the incongruence between the species tree and the gene tree reestimated using IQtree. This means that lack of support for the species topology is associated with high deviation from the median ages. This could be attributable to either actual incongruence being masked by the reconciliation step in our dataset, to low signal from short sequences or to such departure from a strict clock that the true topology cannot be recovered.

The next five coefficients have p-values below 0.05, and correspond to characteristics of the substitution process or the alignment, for which we provide the following interpretations. The kappa parameter (for codon position 3), or ratio of transitions over transversions, is negatively correlated. It may relate to the degree of saturation: saturation brings this ratio closer to one, and it should increase dating uncertainty, although we think that very few gene trees show saturation in our Primates dataset. Next is the proportion of synonymous substitutions at equilibrium, negatively correlated. This value that we measured with 'codeml' from PAML depends on the alignment and represents the amount of substitutions that would be synonymous if the ratio dN/dS was equal to one. This result indicates that gene sequences that offer more opportunity for synonymous changes favor precise dating. For the standard deviation of the GC-content between sequences of the alignment, the correlation is positive. This could have two causes that are not mutually exclusive, one methodological, one biological: in the substitution model, sequence composition is

assumed homogeneous. Only few phylogenetic models can actually model composition changes over the tree (Foster, 2004) while this is known to highly impede phylogenetic inferences. Under a purely biological explanation so even if we fitted the appropriate model, genes subject to composition shifts could be evolving heterogeneously in other aspects, in particular with heterogeneous rates. The next variable is the standard deviation of root-to-tip path lengths, an approximation of the branch rate heterogeneity that can easily be computed on any phylogram. It is, as expected negatively correlated but weakly, probably because our more precise measure of across-branch rate heterogeneity already captures most of the effect. In Smith et al. (2018) it is used as a criterion for selecting trees which are suitable for dating, and we can expect a better predictive power by using a more refined measure of rate heterogeneity instead; we also performed the regression with only the root-to-tip variance as the sole measure of rate heterogeneity, and it has a lower coefficient than the more refined measure (supp. info. S3). Last to be significant, the shape of the site-wise gamma distribution of rates is a parameter in which a lower value means a more skewed gamma distribution. Thus, as indicated by the negative coefficient (Figure 2), rates that are heterogeneous site-wise (not only branch-wise) appear to also increase dating uncertainty.

To situate these features in their real scale, we summarized them in the 10% of trees with highest or lowest predicted dispersion (Table 1). Additionally, we searched for gene functions overrepresented in these two extreme subsets compared to the full set of 5170 genes, based on human Gene Ontology (GO) annotations (Table 2, supp. info. S4). For the most dispersed trees no enriched function was detected. On the other hand, well dated trees (low dispersion) appear to have constrained functions. Negative selection on this set is further evidenced by 67 "Human Phenotype Ontology" overrepresented terms related to various developmental abnormalities (nervous system, eye, ear morphology, movement, see supp. info. S4). These functions are a subset of the functions retrieved when using the 10% longest genes (supp. info. S5), indicating that the length is the main driver of the functions found in the set of genes that have the lowest dating deviation from the median.

In this regression of gene tree characteristics, the explained variable (deviation from the median) is a measure of dispersion, not of accuracy. However, we can only measure accuracy if we have access to the "true" age. Besides, the distribution of ages (Figure 1 suggests a bias in our estimates that empirical data alone is powerless to explain. Furthermore, empirical data is the product of several confounding factors, where for example an apparent rate heterogeneity might instead be caused by undetected incongruence in the gene tree. We therefore performed a simulation to confirm the features that impact dating accuracy.

Subset of gene fam- ilies	Observed dispersion (My)	Alignment length (nucl)	Rate heterogeneity (10-4 subst/nucl/My)	Mean rate (10-3 subst/nucl/My)
10% highest pre- dicted dispersion	3.45 ± 1.6	1008 ± 506	11.3 ± 9.76	2.16 ± 1.25
10% lowest predicted dispersion	1.41 ± 0.576	5255 ± 3247	4.69 ± 3.94	2.75 ± 1.13

 Table 1 - Characteristics of gene families whose dispersion is predicted to be the lowest/highest (mean ± standard deviation)

Subset of gene fam- ilies	GO category	GO term IDs	GO term names
10% highest pre- dicted dispersion	_	_	_
10% lowest predicted dispersion	Molecular Function	GO:0005524 GO:0032559 GO:0030554	ATP binding Adenyl ribonucleotide binding Adenyl nucleotide binding
	Biological Process	GO:0071840 GO:0016043 GO:0007010 GO:0006996	Cellular component organization or biogenesis Cellular component organization Cytoskeleton organization Organelle organization
	Cellular Compartment	GO:0005604 GO:0043228 GO:0043232	Basement membrane Non-membrane-bounded organelle intracellular non-membrane-bounded organelle

Table 2 - Overrepresented Gene Ontology terms in gene families with lowest/highest predicted dispersion, using g:Profiler with default parameters: significance threshold 0.05 and multiple testing correction g:SCS.

Simulating gene tree variability impact on dating

Simulations allow to measure two facets of the dating quality: first the precision or how finegrained the estimates are; second and more importantly, the accuracy or how close the estimates are from the true value on average. We focused on the top 3 variables obtained from the regression: alignment length, across-branch rate heterogeneity and mean rate of substitution. Our three variables of interest were set to representative values spanning their observed range in the real trees (supp. info. S6). For each set of parameters, we first simulated branch lengths on the fixed primate tree according to a relaxed clock model, then simulated codon sequence evolution along this tree (see Methods- *Simulating Alignments*). Dates were then reconstructed from simulated data with the same method as for the empirical dataset. We expected the three variables to reproduce the same effect on dispersion as with the above empirical dataset. However, we expected the accuracy to be unaffected because the model we use to simulate alignments and trees is very similar to the model used to reconstruct dates: Beast dating was thus performed without discrepancy between its inference model and the process that generated the data, and in this condition, we expect it to be statistically consistent, i.e. converge to the true value as the amount of data increases.

We ran 500 simulations for each parametrization, by fixing two parameters and varying the third, around a central point defined by length=3000 nucleotides, mean rate=0.004 substitutions/codon/My and branch rate heterogeneity=1/4 (defined as the standard deviation divided by the mean of the branch rates, abbreviated σ/μ).

The dispersion, as seen by the interquantile range covering 95% of the data (IQR95) (Figure 3) varies with each variable in the same direction as observed in the regression of primate gene trees: it increases for shorter alignments (Figure 3b), higher rate heterogeneity Figure 3c) and

lower evolutionary rate (Figure 3d). In almost all speciation nodes and parameter values, we recover an unbiased median age, falling accurately on the age from the underlying species tree. However, we find shifts for the highest rate heterogeneity between branches (Figure 3d, $\sigma/\mu = 1$). Cebidae appears younger while Catarrhini appears older than in reality, an effect that we find when sampling from the prior (supp. info. S7). This shows that in presence of very high across-branch rate variation and uninformative calibrations, the prior on the time tree (Birth-Death) strongly influences the ages.



Figure 3 - Speciation dates on simulated alignments, under a relaxed clock model and a fixed species tree of Similformes, under various parameter values. Central dots represent the median date, and segments represent the range [0.025, 0.975] of the distribution. a) Dated species tree from TimeTree; b) the number of sites in the alignments, in number of nucleotides; c) the diffusion parameter, defining the degree of relaxation of the clock rate, such that the ratio σ/μ (standard deviation of the rates over the mean rate) equals one twentieth, one quarter, and one; d) The mean rate of substitution, in number of substitutions per site per million years. When not specified, length is 3000, rate is 40e-4 and σ/μ is 1⁄4.

Discussion

Factors of consistency between ages inferred from separate gene trees

From an empirical dataset of 5205 gene trees in 21 Primates, we identified characteristics related to the deviation of age estimates from the median age. More statistical power, provided by higher rates or longer alignments, lowers the dating deviation while an increase in rate variation between branches increases the deviation. This confirms previous observations that the number of sites and substitutions is a limiting factor with regard to statistical power (Bromham et al., 2000; Lanfear et al., 2010; Smith et al., 2018; Bromham, 2019; Bletsa et al., 2019), and also that the clock rate does not appear to be constant across branches. In our analysis the deviation also correlates with characteristics of the substitution process (kappa, GC content), indicating either that the dating model is too simple, or that multiple kinds of heterogeneities (composition, rates) co-occur in some trees. The 5205 selected gene trees only contain strict orthologs and as such represent a biased fraction of the 24,614 gene trees extracted from Simiiformes: this selection likely represents a more reliable subset for dating, because the absence of gene duplication and loss may be associated with a more stable evolution in terms of constancy of rate and selective pressure. On the other hand, even this selection of trees can contain incongruent trees not identified as such, because of the reconciliation method that forces gene trees into the species topology without a model for incomplete lineage sorting (ILS) or introgression. This is suggested by the fourth detected factor in our regression being the discordance between the recomputed gene trees and the species tree. Despite having kept a few polytomic species nodes to lift some of the constraints on the topology, it is likely that it still distorts branch lengths (Mendes & Hahn, 2016; Carruthers et al., 2022).

In the regression based on empirical gene trees, the adjusted R^2 is low (0.29), indicating a low fraction of explained variance. To investigate the missing variance, we should incorporate additional biological characteristics of the sequences, but as we show in simulations, stochastic properties of the mutational and evolutionary processes already account for a large uncertainty. Furthermore, in our empirical dataset, ILS, introgression or reciprocal paralog loss may also obfuscate the relationship between molecular change and time, but we have not used models able to identify these processes. Among other simplifications we used, our dating method employs a nucleotide substitution model which cannot distinguish neutral from non-neutral substitutions. According to the nearly neutral theory the majority of amino-acid changes is slightly deleterious which causes molecular divergence patterns to be more clock-like in absolute time, whereas strictly neutral substitutions should show a generation time effect (Ho, 2014). A codon model would be able to separate synonymous from non-synonymous substitutions and test this hypothesis, but being more complex (using 61 character frequencies instead of 4) it also requires more data, which as we have seen is the limiting factor at the scale of single gene trees.

Limiting the bias when dating single gene trees

In the empirical estimates, we observed a shift towards ages younger than the TimeTree reference. Then we simulated sequences without species-specific rates and that bias was not reproduced. This is unsurprising because we used the same model for inference as for simulating the data, and thus were expecting statistical consistency. However, in the simulated case with very high-rate heterogeneity, the ages are biased by the time tree prior (Birth-Death model). Other sources of bias have been demonstrated, in particular incongruence with the species tree as caused by incomplete lineage transfer (ILS) (Mendes & Hahn, 2016; Carruthers et al., 2022). In the latter extensive empirical and simulation study, it was shown that the length of branches descending from incongruent nodes (such as terminal branches) was overestimated, while the length of branches predating incongruent nodes was underestimated. This produces older ages than the actual speciation times, and the authors propose to mitigate this by considering congruent branches only. In the real primate trees, we observe younger ages instead so it is likely that the bias is instead driven by species-specific rates, a factor that was absent from our simulations. Other studies on empirical data have found biases caused by heterotachy (Wertheim et al., 2012), especially if the rate change is punctuated (Dornburg et al., 2012). In presence of substantial lineage rate variation, the relaxed clock model requires calibrations on internal nodes to properly infer branch rates and times (Duchêne et al., 2014). Here, the choice of a single calibration was made precisely with the aim of measuring the dating accuracy on uncalibrated nodes, such as those occurring in gene trees if we consider events other than speciations. Some studies use the relative ordering of events from different taxa (Lutzoni et al., 2018) or from different gene families (Pittis & Gabaldón, 2016; Vosseberg et al., 2021), but across-branch rate heterogeneity is a major hurdle to doing so and conditions for applicability need to be rigorously verified, in particular the absence of systematic bias (Susko et al., 2021).

How to reduce uncertainty

Stochastic uncertainty is the amount of information that is lost assuming that the model is correct, as opposed to inductive uncertainty due to inadequate modelling (Holland, 2013). Inadequate modelling generally produces systematic errors. Our simulations show what level of stochastic uncertainty to expect in gene trees with realistic rate characteristics, and they confirm the drivers of uncertainty identified in empirical gene trees. We highlight the requirement for a high number of sites and substitutions, which is critical for statistical power. This problem can probably only be circumvented by incorporating information beyond single genes, such as genomic context and linked model parameters between multiple gene trees, as recent developments suggest (Duchêne et al., 2020). However fossil calibrations are quantitatively as important as sequence data alone to accurate dating (dos Reis et al., 2015). These solutions are unfortunately harder to design for gene trees with duplications or transfers because they do not share the same set of branches. Some pragmatic approaches that have been proposed to handle gene specific heterogeneity are the removal of loci with apparently non constant rate prior to the analysis (Jarvis et al., 2014), a "geneshopping" approach (Smith et al., 2018), but discarding data might be unsatisfactory for genome scale analyses that look for the most general picture. Our results may be extended to a similar gene filtering approach as it pinpoints gene features that are correlated with the dating uncertainty, although its predictive power would need to be improved. Beyond identifying outliers, it would be interesting to understand why they are so, in terms of function, selection pressure or genomic context. In this high-throughput era, homologous sequences from a wide variety of taxa are accessible. A few genes have received considerable attention because of very distinctive evolutionary dynamics, such as PRDM9 which is generally evolving under positive selection and is a "speciation gene" in mammals (Oliver et al., 2009) or MHC immune genes that display elevated polymorphism in populations (Piertney & Oliver, 2006). Since our evaluation used only one calibration, it is a worst case setting that could occur in gene trees with many duplications or transfers, events for which fossil calibration is less informative (but see Davín et al., 2018, who use horizontal transfers as relative time constraints). Key adaptations or transitions may result from these gene specific events, in particular duplications (Aguileta et al., 2006; Vosseberg et al., 2021), horizontal and endosymbiotic transfers of genes (Ochman et al., 2000; Koonin, 2016), or movements of transposable elements (Boissinot et al., 2000; Ovchinnikov et al., 2002; Khan et al., 2006).

Regarding *inductive uncertainty*, richer models avoid overconfident but erroneous results. Notably, there is room for refining the models of sequence evolution, for example by considering indels evolution, or the domain composition of proteins, or their 3D structure and the consequence on residue coevolution. Also, when the inference of tree topology, substitutions and clock rates is performed jointly by an integrative method such as Beast, it guarantees that uncertainty is properly cumulated at all steps. However, in our experience it appears challenging to run a complex parametric inference such as Beast at the scale of comparative genomics. Such integrative methods might not be numerically tractable because they generally require MCMC which implies prohibitive running times and in addition, MCMC requires a thorough and experimented human inspection to validate the inference output. In the face of this, numerous projects aiming at sequencing broad segments of biodiversity on Earth are emerging, under the auspices of the Earth Biogenome Project (EBP; Lewin et al., 2018, 2022). Extremely large datasets, composed of tens of thousands of genes, will become common place. Therefore, the development of faster non-parametric algorithms is still relevant, alongside improvements in the computational footprint of probabilistic algorithms (Mello & Schrago, 2024).

Research on the molecular clock enters an exciting time, with huge amounts of data and increasingly sophisticated methods to dissect the hidden mechanisms of sequence evolution. How the rate of molecular evolution varies still holds mysteries, both at the scale of lineages and at the scale of genomes. Future developments answering this question will be methodologically challenging but shall shed light on important evolutionary processes.

Methods

Source and number of trees

The species tree and gene trees were obtained from Ensembl Compara 93 (July 2018; Zerbino et al., 2018) Metazoa dataset: 99 species descending from Opisthokonta (the last common ancestor of fungi and metazoan) and 23,904 reconciled gene trees. Low quality genomes, aberrant gene branch lengths, split genes were removed (supp. info. S8, Species and gene trees preprocessing). Nodes from the species tree were assigned an age from TimeTree (data retrieved Jan. 2019; Kumar et al., 2017). The corresponding sequences of the coding domain of the longest transcribed isoforms (fasta format) were downloaded as fasta multiple sequence alignments via the Ensembl Perl API. Focusing on the 21 primate species (18 Simiiformes and 3 Lemuriformes), we extracted 24,614 subtrees at their Simiiformes node plus two of the outgroup sequences with shortest branch lengths. From those we selected the 5235 gene trees without duplications and losses within Simiiformes, so that they share a single tree topology.

Multiple Sequence Alignment building and cleaning

We built codon alignments: protein sequences were aligned using FSA version 1.15.9 (Fast Statistical Alignment, Bradley et al., 2009) with default parameters, then back-translated to the corresponding nucleotide (codon) alignment. HmmCleaner (version 0.180750; Di Franco et al., 2019) finds segments of sequence that appear inconsistent with the other aligned sequences, interpreting them as sequencing errors. We apply it and replace the detected segments by gaps. It was applied to the 5235 alignments that do not have duplications or losses, on the amino-acid data which were then back translated to codons. Out of these, 30 alignments caused HmmCleaner to fail, resulting in the 5205 subtrees under study here.

Bayesian node dating with Beast 2

Using the alignments as input for Beast 2.6.3 (Bouckaert et al., 2019) we ran the following model: a HKY nucleotide substitution model with two site partitions corresponding to the codon positions {1,2} and 3. The "Birth-Death model" is set as the tree prior (Gernhard, 2008), and a relaxed uncorrelated lognormal clock rate is fitted on branches, with mean and standard deviation being estimated. All 12 Similformes clades (see species tree in Figure 1) as well as Primates and Lemuriformes were constrained to be monophyletic (note that Macaca, Papionini and Platyrrhini are polytomic). Only Primates and Similformes were calibrated precisely, the others having uninformative priors. The calibrations were specified as gamma distributions manually tuned to match TimeTree point estimate and 95% interval: offset 70.8 My, shape 4.6 and scale 0.656 for Primates, and offset 40.9 My, shape 4.0 and scale 0.575 for Similformes (therefore the mean prior age for Primates is 70.8 + 4.6×0.656 = 73.8 My, and for Similformes it is 40.9 + 4.0×0.575 = 43.2 My). BEASTGen (v1.0.1) was used to automatically generate the parameter file for each alignment from a common template. The MCMC was run in one chain of 20,000,000 iterations, after a pre-burnin of 1,000,000. The 390 trees which had an ESS (effective sample size) under 200 for any variable were resumed and extended with 20,000,000 additional iterations, leading to 13 trees with one ESS under 200 that we discarded. The resulting mean dates were annotated on the species tree using TreeAnnotator from Beast 2.

Global age dispersion from all primate trees

Given that all speciation nodes in the primate species tree are replicated in our set of 5205 gene trees, we measured the dispersion of estimated ages for each speciation. We used two metrics of dispersion, the MAD and the IQR95. The mean absolute deviation from the median (MAD) is a robust estimate of dispersion: it is less influenced by skewed distributions than when using the mean as center, and less influenced by outliers than the standard deviation. It is also still sensitive

enough compared to median deviation from the median or interquartile range. The IQR95, or interquantile range covering 95% of the data (between percentiles 2.5 and 97.5) is more robust but less sensitive.

Regression

Statistical analyses on the output dates were performed with Python 3.8 and additional scientific computing packages including Numpy, Pandas, Statsmodels, Scikit-learn, Biopython and Ete3 (supp. info. S9).

Quantifying the average deviation of ages for a single gene tree

As the dependent variable to regress, we used the average deviation of ages in each gene tree: for one gene tree, the absolute deviation from the median speciation age is obtained at each node, then the deviations of the 12 nodes are averaged.

Retrieving mean rate and rate heterogeneity per tree

The clock model that we fit in Beast is unlinked between codon positions {1,2} and position {3}, meaning that the rate parameters are inferred separately for each of these site partitions. Beast outputs different types of rates: First, it outputs the sampled parameters of the uncorrelated log-normal clock model (its mean and standard deviation). Additionally, it computes a posteriori the mean and variance of the branch rate. The latter estimation differs in that it is not a parameter of the model, but a statistic that is computed at the end of each iteration on the proposed tree. We monitor both estimates because they yield quite different values, although being correlated (supp. info. S10). For each sampled tree, it then produces the mean and variance. We use the latter rate statistics as variables in the regression. Since there are as many values as sampled trees of the MCMC chain, we take the median over the chain as final value.

As we estimated separately the mean rates m1,2 and m3 for the corresponding codon positions, we finally combined them to obtain the total rate of substitutions per codon: $m = 2 \times m_{1,2} + m_3$

This rate by codon equals three times the average rate by nucleotide, but we chose the codon metric for consistency with the simulation parameters based on a codon model in INDELible.

Similarly, as measure of across-branch rate heterogeneity, we summed the rate standard deviations based on the returned variances v1,2 and v3 from Beast:

$$s = 2 \times \sqrt{v_{1,2}} + \sqrt{v_3}$$

Collecting features of the gene families

We measured 56 characteristics of the gene families as explanatory variables in the regression. They include (full list in supp. info. S2):

• Tree features:

- mean bootstrap values of nodes, from Ensembl;

- the root-to-tip deviation, or standard deviation of the path lengths from the root to the leaves, from the Ensembl trees.

- whether the gene tree was a posteriori edited to fit the species tree topology (See supp. info. S8, Species and gene trees preprocessing).

• Alignment features:

- global statistics such as alignment length and proportion of pairs of sequences that do not share sites (indicative of "split genes");

- statistics over sequences like the mean frequency of gaps, the mean percentage of ambiguous nucleotides, the mean GC-content, the mean CpG content;

- statistics over sites, like the column entropy and the parsimony score;

- proportions of sequences cleaned by HmmCleaner;

 Substitution process features estimated with the program 'codeml' from PAML 4.9e (Yang, 2007): - proportion of synonymous substitutions at equilibrium, based on measured codon frequencies and the fitted substitution matrix (using names from 'codeml' output, it is S/(N+S));

- mean and standard deviation of ω (dN/dS) under the free-ratio model;
- Substitution process features estimated during the Beast dating inference, separately for codon positions {1,2} and {3}:
 - the ratio of transitions over transversions (kappa);
 - the gamma shape for site variable rates;
 - the mean clock rate;
 - the standard deviation of the clock rate;
- MCMC related statistics of the Beast runs:
 - the number of iterations the chain was run (either 20 or 40 million);
 - whether any parameter had an ESS below 200.

Transforming regressed features

We chose the appropriate transformation of features in order to stay close to the assumptions of a linear modelling framework with ordinary least squares, that is explanatory variables with low skew. For this, a semi-automated procedure was set up, where the transform that minimizes the skew was selected from the predefined set made of 'no transform', 'logarithm base 10', and 'square root', with additional increment or sign modification if needed by the function domain of validity. Some variables were binary encoded in cases where a large proportion of values was constant, or where distributions were clearly bimodal. Transforms for each variable can be found in supp. info. S11. Finally, all variables were normalized (divided by their standard deviation) and centered before the regressions.

Reducing multicollinearity

Multicollinearity being known to impede ordinary least squares fitting, we reduce it prior to the regression with two strategies:

- First, we obtained covariances from a Factor Analysis (analogous to Principal Component Analysis, but accounting for discrete and ordinal data). From the results, we removed some features from clusters of heavily correlated features. We also "decorrelated" pairs of features by dividing one by the other or by computing the residues of the simple regression between the two linked variables. For instance, we decorrelated standard deviation metrics when they correlate with the mean of the same variable (supp. info. S11).
- Afterwards we checked the multicollinearity condition number, which is the square root of the highest eigenvalue divided by the smallest eigenvalue of X^TX (X being the design matrix). Since the multicollinearity condition number was already less than the usual cutoff of 20, it was not needed to remove additional features.

Removing trees with outlier values

To avoid a misleading impact of outliers on the linear regression, trees with excessive rates computed by Beast were excluded. The cutoffs were chosen based on the histograms, which display very long tails. As mentioned above, the 13 trees with insufficient ESS (less than 200) were discarded as well. Finally, alignments including pairs of sequences which do not share any common site, i.e. unaligned pairs, were also excluded, as this can lead to meaningless branch lengths. In total these filters removed 35 outlier gene trees (supp. info. S12).

Final feature selection with Lasso and OLS refitting

Lasso (Tibshirani, 1996) is a regression fitting algorithm which simultaneously performs feature selection, thus dealing with highly dimensional design matrices. We therefore apply it in a first pass to discard coefficients with an absolute value inferior to 0.01, using a penalty value (parameter alpha) of 0.02. However, Lasso coefficients are biased (by design), so that p-values and *R*² cannot be readily computed. For this reason, we subsequently refit by ordinary least squares (OLS) using

the selected features, with a covariance set to MacKinnon and White's (1985) heteroscedasticity robust type 1. Finally, the p-values were subjected to a Bonferroni correction by multiplying them by the number of features used in the Lasso step. Statsmodels 0.13.2 implementations were used. Detailed regression output statistics are in supp. info. S13.

Gene functional annotation overrepresentation

From the 5170 gene trees without duplication or loss retained in the regression, human genes were used as the reference set of functions. Overrepresentation of the 517 human genes from the most dispersed and least dispersed trees was done with the online tool g:Profiler version e111_eg58_p18_30541362 (Raudvere et al., 2019) using default parameters (significance threshold 0.05 and multiple testing correction by g:SCS).

As specified in table 1, the set of most (or least) dispersed trees is obtained from the fitted regression line. It is therefore a predicted dispersion, distinct from the observed dispersion displayed in column 1.

Simulating alignments

Alignments were generated by simulating the evolution of a sequence along the above tree of Primates, using INDELible v1.03 (Fletcher & Yang, 2009). The three tested parameters were the mean rate, the rate heterogeneity across branches and the alignment length, with 500 simulation replicates for each set of parameter values. These parameters were chosen to represent the extent of variation in the real data of 5205 trees (as measured to be used in the regression, see below and supp. info. S6). To generate the variable branch rates, we applied an independent log-normal relaxed clock simulation with the 'simclock' R package; the tested mean rates were 5×10^{-4} , 40×10^{-4} and 80×10^{-4} substitutions.site⁻¹.My⁻¹. To obtain the desired rate heterogeneities, we tuned the "diffusion" parameter s^2 so that the standard deviation of the branch rates is comparable with empirical estimates, as per the relation between standard deviation σ , mean μ and diffusion s^2 of a log-normal distribution:

$$s^2 = \log\left(1 + \frac{\sigma^2}{\mu^2}\right)$$

The chosen diffusion values (s^2) of 0.6931, 0.0606 and 0.0025 correspond to a standard deviation of the branch rate equal to 1, ¹/₄ and 1/20 of the mean rate (σ/μ), respectively. For the alignment length, we set the root sequence lengths to 300, 3000 or 30,000 nucleotides in INDELible. The evolutionary model includes insertions and deletions, which were parametrized to occur at a rate of 8.4×10^{-7} and 16.8×10^{-7} My⁻¹ respectively (corresponding to 1/14 and 2/14 of the median substitution rate 1.177 codon^{-1} .My⁻¹, a ratio taken from Fletcher and Yang 2009), with lengths following a Lavalette distribution with parameters (2, 300). The sequences evolve according to a codon model with *kappa* = 4 (transitions/transversions) based on estimates from the Primates gene trees (supp. info. S14), *omega* = 0.175 (dN/dS) and codon frequencies taken from the concatenated alignments of Similformes genes. No across-site rate heterogeneity was modeled. The random seed was set to 9342.

We then dated Primates speciations from these simulated sequences with Beast 2, as for the real Primates dataset above, except that the tree prior is "Calibrated Yule" for panels b and d of Figure 3. This tree prior was updated to Birth-Death in panel c because it was showing a bias towards younger ages. The number of trees with any ESS under 200 is given in supp. info. S15.

Data, scripts, code, and supplementary information availability

The Python statistical analysis, the source data and the intermediate data are archived at Zenodo.org (https://doi.org/10.5281/zenodo.14000603; Louvel, 2024). The core phylogenetics/bioinformatics library is also available at https://github.com/DyogenIBENS/Phylorgs (version 0.1.0). Supplementary information is available online (https://doi.org/10.5281/zenodo.15203103; Louvel & Roest Crollius, 2024).

Acknowledgment

We wish to thank Alexandra Louis for the management of data and software resources used in this work, and Pierre Vincens for computing support.

Preprint version 7 of this article has been peer-reviewed and recommended by Peer Community In Genomics (https://doi.org/10.24072/pci.genomics.100292; Hoffmann, 2024)

Funding

This work was supported by a grant from Fondation pour la Recherche Médicale to G.L. (FRM FDT201904008392).

Conflict of interest disclosure

The authors declare that they comply with the PCI rule of having no financial conflicts of interest in relation to the content of the article.

References

- Aguileta G, Bielawski JP, Yang Z (2006) Evolutionary rate variation among vertebrate β globin genes: Implications for dating gene family duplication events. *Gene*, **380**, 21–29. https://doi.org/10.1016/j.gene.2006.04.019
- Aris-Brosou S, Yang Z (2002) Effects of Models of Rate Evolution on Estimation of Divergence Dates with Special Reference to the Metazoan 18S Ribosomal RNA Phylogeny. *Systematic Biology*, **51**, 703–714. https://doi.org/10.1080/10635150290102375
- Battistuzzi FU, Filipski A, Hedges SB, Kumar S (2010) Performance of Relaxed-Clock Methods in Estimating Evolutionary Divergence Times and Their Credibility Intervals. *Molecular Biology and Evolution*, **27**, 1289–1300. https://doi.org/10.1093/molbev/msq014
- Bletsa M, Suchard MA, Ji X, Gryseels S, Vrancken B, Baele G, Worobey M, Lemey P (2019)
 Divergence dating using mixed effects clock modelling: An application to HIV-1. *Virus Evolution*, **5**. https://doi.org/10.1093/ve/vez036
- Boissinot S, Chevret P, Furano AV (2000) L1 (LINE-1) Retrotransposon Evolution and Amplification in Recent Human History. *Molecular Biology and Evolution*, **17**, 915–928. https://doi.org/10.1093/oxfordjournals.molbev.a026372
- Bouckaert R, Vaughan TG, Barido-Sottani J, Duchêne S, Fourment M, Gavryushkina A, Heled J, Jones G, Kühnert D, Maio ND, Matschiner M, Mendes FK, Müller NF, Ogilvie HA, Plessis L du, Popinga A, Rambaut A, Rasmussen D, Siveroni I, Suchard MA, Wu C-H, Xie D, Zhang C, Stadler T, Drummond AJ (2019) BEAST 2.5: An advanced software platform for Bayesian evolutionary analysis. *PLOS Computational Biology*, **15**, e1006650. https://doi.org/10.1371/journal.pcbi.1006650
- Bradley RK, Roberts A, Smoot M, Juvekar S, Do J, Dewey C, Holmes I, Pachter L (2009) Fast Statistical Alignment. *PLOS Computational Biology*, **5**, e1000392. https://doi.org/10.1371/journal.pcbi.1000392
- Britten RJ (1986) Rates of DNA sequence evolution differ between taxonomic groups. *Science*, **231**, 1393–1398. https://doi.org/10.1126/science.3082006
- Bromham L (2019) Six Impossible Things before Breakfast: Assumptions, Models, and Belief in Molecular Dating. *Trends in Ecology & Evolution*, **34**, 474–486. https://doi.org/10.1016/j.tree.2019.01.017
- Bromham L, Penny D, Rambaut A, Hendy MD (2000) The Power of Relative Rates Tests Depends on the Data. *Journal of Molecular Evolution*, **50**, 296–301. https://doi.org/10.1007/s002399910034

- Burbrink FT, Pyron RA (2008) The Taming of the Skew: Estimating Proper Confidence Intervals for Divergence Dates. *Systematic Biology*, **57**, 317–328. https://doi.org/10.1080/10635150802040605
- Carruthers T, Sun M, Baker WJ, Smith SA, De Vos JM, Eiserhardt WL (2022) The Implications of Incongruence between Gene Tree and Species Tree Topologies for Divergence Time Estimation (M Alfaro, Ed,). *Systematic Biology*, **71**, 1124–1146. https://doi.org/10.1093/sysbio/syac012
- Davín AA, Tannier E, Williams TA, Boussau B, Daubin V, Szöllősi GJ (2018) Gene transfers can date the tree of life. *Nature Ecology & Evolution*, **2**, 904–909. https://doi.org/10.1038/s41559-018-0525-3
- Di Franco A, Poujol R, Baurain D, Philippe H (2019) Evaluating the usefulness of alignment filtering methods to reduce the impact of errors on evolutionary inferences. *BMC Evolutionary Biology*, **19**, 21. https://doi.org/10.1186/s12862-019-1350-2
- Doolittle RF, Blombäck B (1964) Amino-Acid Sequence Investigations of Fibrinopeptides from Various Mammals: Evolutionary Implications. *Nature*, **202**, 147–152. https://doi.org/10.1038/202147a0
- Dornburg A, Brandley MC, McGowen MR, Near TJ (2012) Relaxed Clocks and Inferences of Heterogeneous Patterns of Nucleotide Substitution and Divergence Time Estimates across Whales and Dolphins (Mammalia: Cetacea). *Molecular Biology and Evolution*, **29**, 721–736. https://doi.org/10.1093/molbev/msr228
- Drummond AJ, Ho SYW, Phillips MJ, Rambaut A (2006) Relaxed phylogenetics and dating with confidence. *PLoS biology*, **4**, e88. https://doi.org/10.1371/journal.pbio.0040088
- Duchêne DA, Duchêne S, Holmes EC, Ho SYW (2015) Evaluating the Adequacy of Molecular Clock Models Using Posterior Predictive Simulations. *Molecular Biology and Evolution*, **32**, 2986–2995. https://doi.org/10.1093/molbev/msv154
- Duchêne S, Lanfear R, Ho SYW (2014) The impact of calibration and clock-model choice on molecular estimates of divergence times. *Molecular Phylogenetics and Evolution*, **78**, 277–289. https://doi.org/10.1016/j.ympev.2014.05.032
- Duchêne DA, Tong KJ, Foster CSP, Duchêne S, Lanfear R, Ho SYW (2020) Linking Branch Lengths Across Sets of Loci Provides the Highest Statistical Support for Phylogenetic Inference. *Molecular Biology and Evolution*, **37**, 1202–1210. https://doi.org/10.1093/molbev/msz291
- Felsenstein J (1981) Evolutionary trees from DNA sequences: A maximum likelihood approach. *Journal of Molecular Evolution*, **17**, 368–376. https://doi.org/10.1007/BF01734359
- Fitch WM, Langley CH (1976) Evolutionary Rates in Proteins: Neutral Mutations and the Molecular Clock. In: *Molecular Anthropology: Genes and Proteins in the Evolutionary Ascent of the Primates* (eds Goodman M, Tashian RE, Tashian JH), pp. 197–219. Springer US, New York, NY. https://doi.org/10.1007/978-1-4615-8783-5_10
- Fletcher W, Yang Z (2009) INDELible: A Flexible Simulator of Biological Sequence Evolution. Molecular Biology and Evolution, 26, 1879–1888. https://doi.org/10.1093/molbev/msp098
- Foster PG (2004) Modeling compositional heterogeneity. *Systematic Biology*, **53**, 485–495. https://doi.org/10.1080/10635150490445779
- Gernhard T (2008) The conditioned reconstructed process. *Journal of Theoretical Biology*, **253**, 769–778. https://doi.org/10.1016/j.jtbi.2008.04.005
- Gillespie JH (1989) Lineage effects and the index of dispersion of molecular evolution. *Molecular Biology and Evolution*, **6**, 636–647. https://doi.org/10.1093/oxfordjournals.molbev.a040576
- Gillespie JH (1991) The Causes of Molecular Evolution. Oxford University Press. https://doi.org/10.1093/oso/9780195068832.001.0001
- Graur D, Martin W (2004) Reading the entrails of chickens: molecular timescales of evolution and the illusion of precision. *Trends in Genetics*, **20**, 80–86. https://doi.org/10.1016/j.tig.2003.12.003
- Gu X, Li W-H (1992) Higher rates of amino acid substitution in rodents than in humans. *Molecular Phylogenetics and Evolution*, **1**, 211–214. https://doi.org/10.1016/1055-7903(92)90017-B

- Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O (2010) New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance of PhyML 3.0. Systematic Biology, **59**, 307–321. https://doi.org/10.1093/sysbio/syq010
- Heath TA, Holder MT, Huelsenbeck JP (2012) A Dirichlet Process Prior for Estimating Lineage-Specific Substitution Rates. *Molecular Biology and Evolution*, **29**, 939–955. https://doi.org/10.1093/molbev/msr255
- Ho SYW (2014) The changing face of the molecular evolutionary clock. *Trends in Ecology & Evolution*, **29**, 496–503. https://doi.org/10.1016/j.tree.2014.07.004
- Ho SYW, Duchêne S, Duchêne DA (2015) Simulating and detecting autocorrelation of molecular evolutionary rates among lineages. *Molecular Ecology Resources*, **15**, 688–696. https://doi.org/10.1111/1755-0998.12320
- Ho SYW, Phillips MJ, Drummond AJ, Cooper A (2005) Accuracy of rate estimation using relaxedclock models with a critical focus on the early metazoan radiation. *Molecular Biology and Evolution*, **22**, 1355–1363. https://doi.org/10.1093/molbev/msi125
- Hoffmann F (2024) Dating single gene trees in the age of phylogenomics. *Peer Community in Genomics*, **1**, 100292. https://doi.org/10.24072/pci.genomics.100292
- Holland BR (2013) The Rise of Statistical Phylogenetics. *Australian & New Zealand Journal of Statistics*, **55**, 205–220. https://doi.org/10.1111/anzs.12035
- Jarvis ED, Mirarab S, Aberer AJ, Li B, Houde P, Li C, Ho SYW, Faircloth BC, Nabholz B, Howard JT, Suh A, Weber CC, Fonseca RR da, Li J, Zhang F, Li H, Zhou L, Narula N, Liu L, Ganapathy G, Boussau B, Bayzid MS, Zavidovych V, Subramanian S, Gabaldón T, Capella-Gutiérrez S, Huerta-Cepas J, Rekepalli B, Munch K, Schierup M, Lindow B, Warren WC, Ray D, Green RE, Bruford MW, Zhan X, Dixon A, Li S, Li N, Huang Y, Derryberry EP, Bertelsen MF, Sheldon FH, Brumfield RT, Mello CV, Lovell PV, Wirthlin M, Schneider MPC, Prosdocimi F, Samaniego JA, Velazquez AMV, Alfaro-Núñez A, Campos PF, Petersen B, Sicheritz-Ponten T, Pas A, Bailey T, Scofield P, Bunce M, Lambert DM, Zhou Q, Perelman P, Driskell AC, Shapiro B, Xiong Z, Zeng Y, Liu S, Li Z, Liu B, Wu K, Xiao J, Yinqi X, Zheng Q, Zhang Y, Yang H, Wang J, Smeds L, Rheindt FE, Braun M, Fjeldsa J, Orlando L, Barker FK, Jønsson KA, Johnson W, Koepfli K-P, O'Brien S, Haussler D, Ryder OA, Rahbek C, Willerslev E, Graves GR, Glenn TC, McCormack J, Burt D, Ellegren H, Alström P, Edwards SV, Stamatakis A, Mindell DP, Cracraft J, Braun EL, Warnow T, Jun W, Gilbert MTP, Zhang G (2014) Whole-genome analyses resolve early branches in the tree of life of modern birds. *Science*, **346**, 1320–1331. https://doi.org/10.1126/science.1253451
- Khan H, Smit A, Boissinot S (2006) Molecular evolution and tempo of amplification of human LINE-1 retrotransposons since the origin of primates. *Genome Research*, **16**, 78–87. https://doi.org/10.1101/gr.4001406
- Kishino H, Thorne JL, Bruno WJ (2001) Performance of a Divergence Time Estimation Method under a Probabilistic Model of Rate Evolution. *Molecular Biology and Evolution*, **18**, 352–361. https://doi.org/10.1093/oxfordjournals.molbev.a003811
- Koonin EV (2016) Viruses and mobile elements as drivers of evolutionary transitions. *Philosophical Transactions of the Royal Society B: Biological Sciences*, **371**, 20150442. https://doi.org/10.1098/rstb.2015.0442
- Kumar S, Hedges SB (1998) A molecular timescale for vertebrate evolution. *Nature*, **392**, 917–920. https://doi.org/10.1038/31927
- Kumar S, Hedges SB (2016) Advances in Time Estimation Methods for Molecular Data. *Molecular Biology and Evolution*, **33**, 863–869. https://doi.org/10.1093/molbev/msw026
- Kumar S, Stecher G, Suleski M, Hedges SB (2017) TimeTree: A Resource for Timelines, Timetrees, and Divergence Times. *Molecular Biology and Evolution*, **34**, 1812–1819. https://doi.org/10.1093/molbev/msx116
- Lanfear R, Welch JJ, Bromham L (2010) Watching the clock: Studying variation in rates of molecular evolution between species. *Trends in Ecology & Evolution*, **25**, 495–503. https://doi.org/10.1016/j.tree.2010.06.007

- Langley CH, Fitch WM (1974) An examination of the constancy of the rate of molecular evolution. *Journal of Molecular Evolution*, **3**, 161–177. https://doi.org/10.1007/BF01797451
- Lartillot N, Phillips MJ, Ronquist F (2016) A mixed relaxed clock model. *Philosophical Transactions* of the Royal Society of London. Series B, Biological Sciences, **371**. https://doi.org/10.1098/rstb.2015.0132
- Lepage T, Bryant D, Philippe H, Lartillot N (2007) A general comparison of relaxed molecular clock models. *Molecular Biology and Evolution*, **24**, 2669–2680. https://doi.org/10.1093/molbev/msm193
- Lewin HA, Richards S, Lieberman Aiden E, Allende ML, Archibald JM, Bálint M, Barker KB, Baumgartner B, Belov K, Bertorelle G, Blaxter ML, Cai J, Caperello ND, Carlson K, Castilla-Rubio JC, Chaw S-M, Chen L, Childers AK, Coddington JA, Conde DA, Corominas M, Crandall KA, Crawford AJ, DiPalma F, Durbin R, Ebenezer TE, Edwards SV, Fedrigo O, Flicek P, Formenti G, Gibbs RA, Gilbert MTP, Goldstein MM, Graves JM, Greely HT, Grigoriev IV, Hackett KJ, Hall N, Haussler D, Helgen KM, Hogg CJ, Isobe S, Jakobsen KS, Janke A, Jarvis ED, Johnson WE, Jones SJM, Karlsson EK, Kersey PJ, Kim J-H, Kress WJ, Kuraku S, Lawniczak MKN, Leebens-Mack JH, Li X, Lindblad-Toh K, Liu X, Lopez JV, Marques-Bonet T, Mazard S, Mazet JAK, Mazzoni CJ, Myers EW, O'Neill RJ, Paez S, Park H, Robinson GE, Roquet C, Ryder OA, Sabir JSM, Shaffer HB, Shank TM, Sherkow JS, Soltis PS, Tang B, Tedersoo L, Uliano-Silva M, Wang K, Wei X, Wetzer R, Wilson JL, Xu X, Yang H, Yoder AD, Zhang G (2022) The Earth BioGenome Project 2020: Starting the clock. *Proceedings of the National Academy of Sciences*, 119, e2115635118. https://doi.org/10.1073/pnas.2115635118
- Lewin HA, Robinson GE, Kress WJ, Baker WJ, Coddington J, Crandall KA, Durbin R, Edwards SV, Forest F, Gilbert MTP, Goldstein MM, Grigoriev IV, Hackett KJ, Haussler D, Jarvis ED, Johnson WE, Patrinos A, Richards S, Castilla-Rubio JC, van Sluys M-A, Soltis PS, Xu X, Yang H, Zhang G (2018) Earth BioGenome Project: Sequencing life for the future of life. *Proceedings* of the National Academy of Sciences, **115**, 4325–4333. https://doi.org/10.1073/pnas.1720115115
- Linder M, Britton T, Sennblad B (2011) Evaluation of Bayesian models of substitution rate evolution—parental guidance versus mutual independence. *Systematic Biology*, **60**, 329–342. https://doi.org/10.1093/sysbio/syr009
- Louvel G (2024) Code and data for molecular dating benchmark based on real and simulated Primates gene trees. Version 4. Zenodo. https://doi.org/10.5281/zenodo.14000603
- Louvel G, Roest Crollius H (2024) Supplementary information for "Factors influencing the accuracy and precision in dating single gene trees." Version 7. Zenodo. https://doi.org/10.5281/ze-nodo.15203103
- Lozano-Fernandez J, dos Reis M, Donoghue PCJ, Pisani D (2017) RelTime Rates Collapse to a Strict Clock When Estimating the Timeline of Animal Diversification. *Genome Biology and Evolution*, **9**, 1320–1328. https://doi.org/10.1093/gbe/evx079
- Lutzoni F, Nowak MD, Alfaro ME, Reeb V, Miadlikowska J, Krug M, Arnold AE, Lewis LA, Swofford DL, Hibbett D, Hilu K, James TY, Quandt D, Magallón S (2018) Contemporaneous radiations of fungi and plants linked to symbiosis. *Nature Communications*, **9**, 5451. https://doi.org/10.1038/s41467-018-07849-9
- MacKinnon JG, White H (1985) Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties. *Journal of Econometrics*, **29**, 305–325. https://doi.org/10.1016/0304-4076(85)90158-7
- Margoliash E (1963) Primary Structure and Evolution of Cytochrome C. *Proceedings of the National Academy of Sciences of the United States of America*, **50**, 672–679. https://doi.org/10.1073/pnas.50.4.672
- Mello B, Schrago CG (2024) Modeling Substitution Rate Evolution across Lineages and Relaxing the Molecular Clock. *Genome Biology and Evolution*, **16**, evae199. https://doi.org/10.1093/gbe/evae199
- Mendes FK, Hahn MW (2016) Gene Tree Discordance Causes Apparent Substitution Rate Variation. *Systematic Biology*, **65**, 711–721. https://doi.org/10.1093/sysbio/syw018

- Muse SV, Gaut BS (1997) Comparing Patterns of Nucleotide Substitution Rates Among Chloroplast Loci Using the Relative Ratio Test. *Genetics*, **146**, 393–399. <u>https://doi.org/10.1093/genetics/146.1.393</u>
- Ochman H, Lawrence JG, Groisman EA (2000) Lateral gene transfer and the nature of bacterial innovation. *Nature*, **405**, 299–304. https://doi.org/10.1038/35012500
- Oliver PL, Goodstadt L, Bayes JJ, Birtle Z, Roach KC, Phadnis N, Beatson SA, Lunter G, Malik HS, Ponting CP (2009) Accelerated Evolution of the Prdm9 Speciation Gene across Diverse Metazoan Taxa. PLOS Genetics, 5, e1000753. https://doi.org/10.1371/journal.pgen.1000753
- Ovchinnikov I, Rubin A, Swergold GD (2002) Tracing the LINEs of human evolution. *Proceedings* of the National Academy of Sciences, **99**, 10522–10527. https://doi.org/10.1073/pnas.152346799
- Pacifici M, Santini L, Di Marco M, Baisero D, Francucci L, Grottolo Marasini G, Visconti P, Rondinini C (2013) Generation length for mammals. *Nature Conservation*, **5**, 89–94. https://doi.org/10.3897/natureconservation.5.5734
- Pagel M, Venditti C, Meade A (2006) Large Punctuational Contribution of Speciation to Evolutionary Divergence at the Molecular Level. *Science*, **314**, 119–121. https://doi.org/10.1126/science.1129647
- Philippe H, Casane D, Gribaldo S, Lopez P, Meunier J (2003) Heterotachy and functional shift in protein evolution. *IUBMB life*, **55**, 257–265. https://doi.org/10.1080/1521654031000123330
- Piertney SB, Oliver MK (2006) The evolutionary ecology of the major histocompatibility complex. *Heredity*, **96**, 7–21. https://doi.org/10.1038/sj.hdy.6800724
- Pittis AA, Gabaldón T (2016) Late acquisition of mitochondria by a host with chimaeric prokaryotic ancestry. *Nature*, **531**, 101–104. https://doi.org/10.1038/nature16941
- Pulquério MJF, Nichols RA (2007) Dates from the molecular clock: how wrong can we be? *Trends in Ecology & Evolution*, **22**, 180–184. https://doi.org/10.1016/j.tree.2006.11.013
- Rannala B, Yang Z (2007) Inferring Speciation Times under an Episodic Molecular Clock. Systematic Biology, 56, 453–466. https://doi.org/10.1080/10635150701420643
- Raudvere U, Kolberg L, Kuzmin I, Arak T, Adler P, Peterson H, Vilo J (2019) g:Profiler: a web server for functional enrichment analysis and conversions of gene lists (2019 update). *Nucleic Acids Research*, **47**, W191–W198. https://doi.org/10.1093/nar/gkz369
- dos Reis M, Gunnell GF, Barba-Montoya J, Wilkins A, Yang Z, Yoder AD (2018) Using Phylogenomic Data to Explore the Effects of Relaxed Clocks and Calibration Strategies on Divergence Time Estimation: Primates as a Test Case. Systematic Biology, **67**, 594–615. https://doi.org/10.1093/sysbio/syy001
- dos Reis M, Thawornwattana Y, Angelis K, Telford MJ, Donoghue PCJ, Yang Z (2015) Uncertainty in the Timing of Origin of Animals and the Limits of Precision in Molecular Timescales. *Current Biology*, **25**, 2939–2950. https://doi.org/10.1016/j.cub.2015.09.066
- dos Reis M, Yang Z (2011) Approximate Likelihood Calculation on a Phylogeny for Bayesian Estimation of Divergence Times. *Molecular Biology and Evolution*, **28**, 2161–2172. https://doi.org/10.1093/molbev/msr045
- dos Reis M, Yang Z (2013) The unbearable uncertainty of Bayesian divergence time estimation. Journal of Systematics and Evolution, **51**, 30–43. https://doi.org/10.1111/j.1759-6831.2012.00236.x
- dos Reis M, Zhu T, Yang Z (2014) The Impact of the Rate Prior on Bayesian Estimation of Divergence Times with Multiple Loci. Systematic Biology, 63, 555–565. https://doi.org/10.1093/sysbio/syu020
- Revell LJ, Harmon LJ, Glor RE (2005) Under-parameterized Model of Sequence Evolution Leads to Bias in the Estimation of Diversification Rates from Molecular Phylogenies (P Linder, Ed,). *Systematic Biology*, **54**, 973–983. https://doi.org/10.1080/10635150500354647
- Sarich VM, Wilson AC (1966) Quantitative Immunochemistry and the Evolution of Primate Albumins: Micro-Complement Fixation. Science, 154, 1563–1566. https://doi.org/10.1126/science.154.3756.1563

- Schenk JJ (2016) Consequences of Secondary Calibrations on Divergence Time Estimates. *PloS One*, **11**, e0148228. https://doi.org/10.1371/journal.pone.0148228
- Smith SA, Brown JW, Walker JF (2018) So many genes, so little time: A practical approach to divergence-time estimation in the genomic era (H Escriva, Ed,). *PLOS ONE*, **13**, e0197433. https://doi.org/10.1371/journal.pone.0197433
- Smith NGC, Eyre-Walker A (2003) Partitioning the Variation in Mammalian Substitution Rates. *Molecular Biology and Evolution*, **20**, 10–17. https://doi.org/10.1093/oxfordjournals.molbev.a004231
- Susko E, Steel M, Roger AJ (2021) Conditions under which distributions of edge length ratios on phylogenetic trees can be used to order evolutionary events. *Journal of Theoretical Biology*, 526, 110788. https://doi.org/10.1016/j.jtbi.2021.110788
- Tajima F (1993) Simple methods for testing the molecular evolutionary clock hypothesis. *Genetics*, **135**, 599–607. https://doi.org/10.1093/genetics/135.2.599
- Tao Q, Tamura K, Mello B, Kumar S (2019) Reliable Confidence Intervals for RelTime Estimates of Evolutionary Divergence Times. *Molecular Biology and Evolution*, **37**, 280–290. https://doi.org/10.1093/molbev/msz236
- Thorne JL, Kishino H, Painter IS (1998) Estimating the rate of evolution of the rate of molecular evolution. *Molecular Biology and Evolution*, **15**, 1647–1657. https://doi.org/10.1093/oxfordjournals.molbev.a025892
- Tibshirani R (1996) Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **58**, 267–288. https://doi.org/10.1111/j.2517-6161.1996.tb02080.x
- To T-H, Jung M, Lycett S, Gascuel O (2016) Fast Dating Using Least-Squares Criteria and Algorithms. *Systematic Biology*, **65**, 82–97. https://doi.org/10.1093/sysbio/syv068
- Volz EM, Frost SDW (2017) Scalable relaxed clock phylogenetic dating. *Virus Evolution*, **3**. https://doi.org/10.1093/ve/vex025
- Vosseberg J, van Hooff JJE, Marcet-Houben M, van Vlimmeren A, van Wijk LM, Gabaldón T, Snel B (2021) Timing the origin of eukaryotic cellular complexity with ancient duplications. *Nature Ecology & Evolution*, 5, 92–100. https://doi.org/10.1038/s41559-020-01320-z
- Warnock RCM, Yang Z, Donoghue PCJ (2017) Testing the molecular clock using mechanistic models of fossil preservation and molecular evolution. *Proceedings of the Royal Society B: Biological Sciences*, **284**, 20170227. https://doi.org/10.1098/rspb.2017.0227
- Wertheim JO, Fourment M, Kosakovsky Pond SL (2012) Inconsistencies in Estimating the Age of HIV-1 Subtypes Due to Heterotachy. *Molecular Biology and Evolution*, **29**, 451–456. https://doi.org/10.1093/molbev/msr266
- Wertheim JO, Sanderson MJ, Worobey M, Bjork A (2010) Relaxed Molecular Clocks, the Bias– Variance Trade-off, and the Quality of Phylogenetic Inference. *Systematic Biology*, **59**, 1–8. https://doi.org/10.1093/sysbio/syp072
- Wolfe KH, Sharp PM, Li W-H (1989) Mutation rates differ among regions of the mammalian genome. *Nature*, **337**, 283–285. https://doi.org/10.1038/337283a0
- Wu CI, Li WH (1985) Evidence for higher rates of nucleotide substitution in rodents than in man. *Proceedings of the National Academy of Sciences*, **82**, 1741–1745. https://doi.org/10.1073/pnas.82.6.1741
- Yang Z (2007) PAML 4: phylogenetic analysis by maximum likelihood. *Molecular biology and evolution*, **24**, 1586–91. https://doi.org/10.1093/molbev/msm088
- Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, Billis K, Cummins C, Gall A, Girón CG, Gil L, Gordon L, Haggerty L, Haskell E, Hourlier T, Izuogu OG, Janacek SH, Juettemann T, To JK, Laird MR, Lavidas I, Liu Z, Loveland JE, Maurel T, McLaren W, Moore B, Mudge J, Murphy DN, Newman V, Nuhn M, Ogeh D, Ong CK, Parker A, Patricio M, Riat HS, Schuilenburg H, Sheppard D, Sparrow H, Taylor K, Thormann A, Vullo A, Walts B, Zadissa A, Frankish A, Hunt SE, Kostadima M, Langridge N, Martin FJ, Muffato M, Perry E, Ruffier M, Staines DM, Trevanion SJ, Aken BL, Cunningham F, Yates A, Flicek P (2018) Ensembl 2018. *Nucleic Acids Research*, **46**, D754–D761. https://doi.org/10.1093/nar/gkx1098

- Zhu T, dos Reis M, Yang Z (2015) Characterization of the Uncertainty of Divergence Time Estimation under Relaxed Molecular Clock Models Using Multiple Loci. Systematic Biology, 64, 267– 280. https://doi.org/10.1093/sysbio/syu109
- Zuckerkandl E, Pauling L (1962) Molecular disease, evolution, and genetic heterogeneity. *Horizons in biochemistry*, 189–225.
- Zuckerkandl E, Pauling L (1965) Molecules as documents of evolutionary history. *Journal of Theoretical Biology*, **8**, 357–366. https://doi.org/10.1016/0022-5193(65)90083-4