

# Substitution Rate Analysis and Molecular Evolution

Lindell Bromham

► **To cite this version:**

Lindell Bromham. Substitution Rate Analysis and Molecular Evolution. Scornavacca, Celine; Delsuc, Frédéric; Galtier, Nicolas. Phylogenetics in the Genomic Era, No commercial publisher | Authors open access book, pp.4.4:1–4.4:21, 2020. hal-02536329

**HAL Id: hal-02536329**

**<https://hal.archives-ouvertes.fr/hal-02536329>**

Submitted on 10 Apr 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# Chapter 4.4 Substitution Rate Analysis and Molecular Evolution

**Lindell Bromham**

Macroevolution & Macroecology, Division of Ecology & Evolution  
Research School of Biology Australian National University  
Canberra, ACT, 0200 Australia  
lndell.bromham@anu.edu.au

---

## Abstract

The study of the tempo and mode of molecular evolution has played a key role in evolutionary biology, both as a stimulant for theoretical enrichment and as the foundation of useful analytical tools. When protein and DNA sequences were first produced, the surprising constancy of rates of change brought molecular evolution into conflict with mainstream evolutionary biology, but also stimulated the formation of new theoretical understanding of the processes of genetic change, including the recognition of the role of neutral mutations and genetic drift in genomic evolution. As more data were collected, it became clear that there were systematic differences in the substitution rate between species, which prompted further elaboration of ideas such as the generation time effect and the nearly neutral theory. Comparing substitution rates between species continues to provide a window on fundamental evolutionary processes. However, investigating patterns of substitution rates requires attention to potential complicating factors such as the phylogenetic non-independence of rates estimates and the time-dependence of measurement error. This chapter compares different analytical approaches to study the tempo and mode of molecular evolution, and considers the way a richer biological understanding of the causes of variation in substitution rate might inform our attempts to use molecular data to uncover evolutionary history.

**How to cite:** Lindell Bromham (2020). Substitution Rate Analysis and Molecular Evolution. In Scornavacca, C., Delsuc, F., and Galtier, N., editors, *Phylogenetics in the Genomic Era*, chapter No. 4.4, pp. 4.4:1–4.4:21. No commercial publisher | Authors open access book. The book is freely available at <https://hal.inria.fr/PGE>.

## 1 Substitution rates and the shape of evolutionary theory

Evolutionary genetics was founded on the patterns of inheritance of phenotypically measurable differences, and their change in frequency in populations over time. Rates of change were measured in terms of shifts in the mean trait values over time (e.g. Haldane, 1949). Mutation rates were estimated from careful detection of visible differences in members of wild populations or through laboratory crosses (e.g. Dobzhansky and Wright, 1941). While many of the leaders of the neo-Darwinian synthesis were keen to incorporate molecular data into their view of evolution, they expected it to join the party on their terms, adhering to the hard-won principle that natural selection was the composer of the molecular message, and that the genotype was servant to the phenotype (Simpson, 1964). Change in the genes and proteins, it was assumed, would reflect the changes wrought on the phenotype by selection, and would, therefore, match the phenotype in tempo and mode of evolution, varying over time as organisms responded to change in environment and selective regime (Aronson, 2002; Dietrich, 1994; Stoltzfus, 2017). Some evolutionary biologists even objected to the very notion of “molecular evolution”, on the grounds that evolution as a process of phenotypic



© Lindell Bromham.

Licensed under Creative Commons License CC-BY-NC-ND 4.0.

*Phylogenetics in the genomic era.*

Editors: Céline Scornavacca, Frédéric Delsuc and Nicolas Galtier; chapter No. 4.4; pp. 4.4:1–4.4:21

A book completely handled by researchers.



No publisher has been paid.

#### 4.4:2 Substitution Rate Analysis and Molecular Evolution

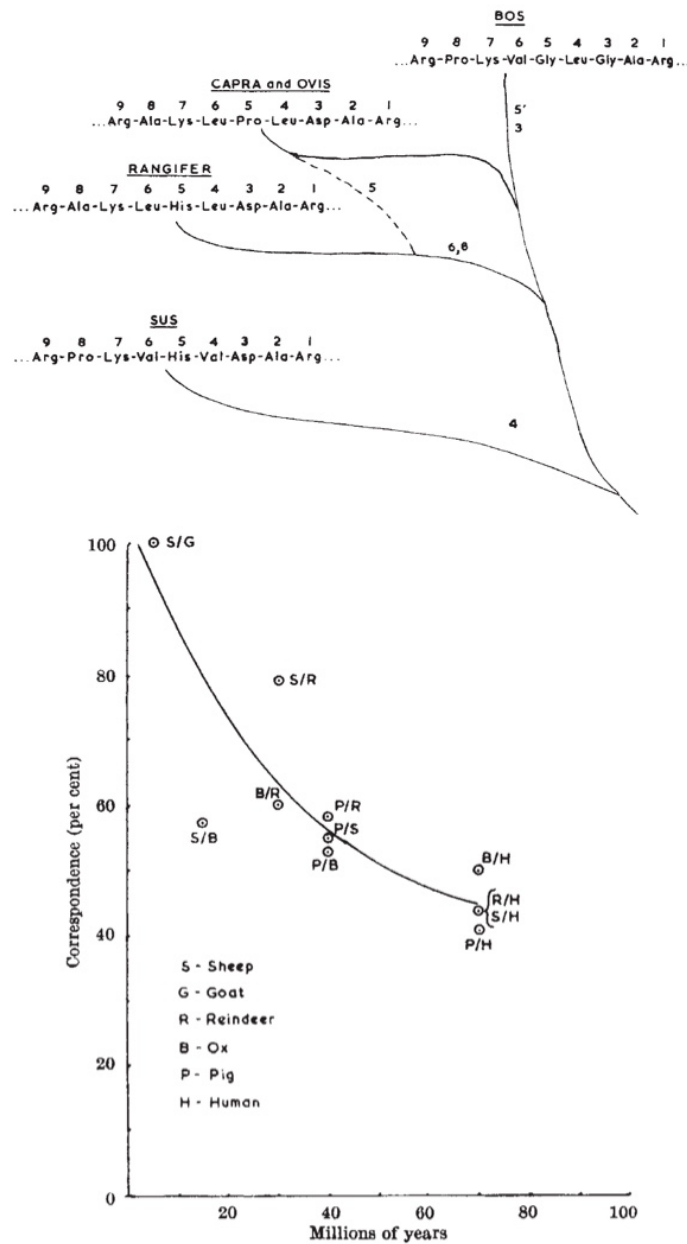
change was reflected in molecular variation, not driven by molecular change itself (Anfinsen, 1965).

Many of the pioneers of the field of molecular evolution emphasized an essentially Darwinian approach to understanding evolution, with change at the molecular level affecting short term processes of individual development as well as connecting to patterns change at longer timescales: “A general description of the evolutionary process is applicable to all levels of complexity, including the chemical level” (Dayhoff and Eck, 1969). Yet, unlike the neo-Darwinian synthesis, the new molecular view did not endow explanatory privilege on the individual level of biological organisation (e.g. Zuckerkandl and Pauling, 1965). Furthermore, it was recognized that change at the molecular level might follow different patterns to phenotypic change. In particular, the potential for neutral evolution was recognized from the beginning of the development of molecular evolution, as researchers acknowledged that change in some proteins, or parts of proteins, may be less impacted by selection than others (e.g. Anfinsen, 1959; Buettner-Janusch and Hill, 1965a). Yet these early workers were not able to make direct connections between the variation generation by mutation and the processes of evolutionary change by substitution within populations, contributing to the divergence of lineages.

In the 1960s, three new molecular techniques finally allowed scientists to peer beneath the phenotypic skin to the genotypic variation within, allowing comparison of genetic variants both within populations and between lineages. And what they saw sent shockwaves through biology. These early studies of molecular rates –from DNA hybridization, protein electrophoresis and amino acid sequences– revealed that the genome was moving out of step with the phenotype. Protein electrophoresis allowed, for the first time, some semblance of random sampling of genetic diversity within populations, measuring variability of many different proteins chosen more or less arbitrarily, revealing that a surprisingly large proportion of loci varied between individuals (Harris, 1966; Hubby and Lewontin, 1966; Lewontin and Hubby, 1966). The amount of variation at the molecular level was far higher than had been predicted from theoretical and empirical studies of rate of change at the phenotypic level (Charlesworth et al., 2016). Furthermore, DNA hybridization experiments, which used the disassociation rates of DNA from different species to indicate overall genome similarity between lineages, showed that the genome evolved continuously, and even faster than proteins. These experiments also suggested that a substantial part of the genome was made up not of unique gene sequences, each with a specific function determined by its sequence, but of vast numbers of repeats of the same short sequences (Britten and Kohne, 1968). The connection of this “repetitious DNA”, if any, to the phenotype was unknown.

But the most controversial observation to come out of these early days of molecular evolutionary biology arose from the comparison of protein sequences across species. As sequences accumulated, it became possible not only to reconstruct the history of change of these molecules over evolutionary time, but also to estimate rates of change (Figure 1). Molecular change seemed to accumulate at a relatively steady rate (Doolittle and Blomback, 1964; Zuckerkandl and Pauling, 1965). This observation of constant rates was immediately put to practical use. Given an average rate of change based on fossil evidence, genetic distance between species - estimated from protein sequence comparisons, immunological distance or DNA hybridization - could be used to infer the age of their last common ancestor (Doolittle and Blomback, 1964; Zuckerkandl and Pauling, 1965; Margoliash, 1963; Wilson and Sarich, 1969).

The surprising observation of that amino acid sequences seemed to change at a roughly constant rates led to fundamental theory change, because it was used to support an argu-



■ **Figure 1** The advent of protein sequencing led to the first analyses of substitution rates, kicking off the controversies about molecular dating analyses that continue to this day. Reprinted by permission from Springer Nature “Amino-Acid Sequence Investigations of Fibrinopeptides from Various Mammals: Evolutionary Implications” Russell F. Doolittle, Birger Blomback *Nature* 1964 202(4928):147-152. Rightslink licence: 4410641268537.

ment that a substantial fraction of changes at the molecular level were neutral, and therefore not influenced by selection but only by random sampling (Kimura, 1968; King and Jukes, 1969). The possibility of neutral mutations had been recognised since the beginnings of evolutionary biology (Darwin, 1859), but had been largely rejected by whole-organism biologists (e.g. Simpson, 1964). The evolution of characters by drift was generally regarded

#### 4.4:4 Substitution Rate Analysis and Molecular Evolution

as of little practical impact in evolution (e.g. Fisher and Ford, 1950), or at very least as generally unproven (e.g. Cain, 1951). But those who were moving into the wild uncharted territories of comparative protein sequence analysis recognized that some changes to amino acid sequences might have no significant functional impact on the resulting protein, and might make no contribution to phenotype (Jukes 1966; Buettner-Janusch and Hill 1965b; Chapter 4.2 [Robinson-Rechavi 2020]). Such changes would not be under the influence of natural selection.

Constant rates of protein change formed one of the pillars on which the neutral theory was built (Kimura, 1968). If many mutations have little or no effect on relative fitness, then they will not be governed by selection. Their fate will be determined by chance events. Since each neutral mutation has an equal chance of drifting to fixation, their overall rate of substitution is governed by the rate at which they are generated. So Kimura (1968) proposed that the neutral substitution rate should be determined only by the neutral mutation rate.

Ironically, given the key role the molecular clock played in launching neutral theory, neutrality is neither necessary nor sufficient to explain constant rates. In fact, the apparently clock-like nature of molecular change had been debated in terms of selection for many years (e.g. Simpson, 1964), and many people working in the field were content to consider both selective and neutral explanations for constancy of rates (Zuckerandl and Pauling, 1965). A steady rate of change could occur under selection if mutation regularly supplies variants of slight selective advantage which then undergo substitution by selection, accumulating at a roughly constant rate when considered over long time periods. Conversely, neutral evolution need not lead to constant rates. The core conclusion of the neutral theory, that the neutral substitution rate is determined by the mutation rate, leads directly to the prediction that rates of genome evolution will vary with differences in the mutation rate. It is also important to note that early molecular clock studies did not assume that the rate of change was invariant, but that any variation was random, and that the long term average rate did not differ substantially between different lineages (Margoliash, 1963). But, nonetheless, these examples show how important consideration of substitution rates has been in the debate about the causes of genomic evolution, both in the early days and continuing to the present day (e.g. Fay and Wu, 2001; Gossmann et al., 2012; Kern and Hahn, 2018; Lynch et al., 2016; Nei et al., 2010; Zhang and Yang, 2015).

In fact, it soon became apparent that rates of molecular evolution showed far more complex patterns. DNA hybridization studies revealed different rates of genomic change in different species, consistent with the prediction that species with faster generation times would generate more mutations per unit time (Laird et al., 1969; Ohta, 1972). The perceived lack of a generation time effect in protein sequence change was interpreted as a result of the interaction of several influences on rates of molecular evolution, both at the level of the mutation rate (smaller species have faster generations so generate more copy errors per year) and the substitution rate (smaller species have larger populations which have less fixation of nearly neutral changes, Ohta 1972, 1973). We now recognise a tangle of different forces that influence both mutation rate and substitution rate, which all come together to shape rates of molecular evolution, at both the DNA and protein level (Bromham, 2011).

Even in the phylogenomic era of ginormous databases, it is worth taking the time to read the earliest papers on the analysis of substitution rates, back when the challenge was to derive big theoretical conclusions from very small amounts of data (Lewontin, 1974). The foundations of the field of molecular evolution were built at a time there were few available protein sequences, each one of which had been painstakingly acquired by skilful and persistent lab work. As a consequence, a feature of this early work is the degree of

biochemical knowledge and attention to detail. Each residue that differed between species was interrogated in terms of structure and function, reactivity and charge, and interpreted in light of the principle that natural selection operates on the working properties of a three-dimensional molecule not a linear sequence of amino acids or nucleotides (Dickerson, 1971).

As the number of protein sequences grew, the first comparative databases were established. Notably, Margaret Dayhoff laid the foundations for modern phylogenomics, by bringing together biochemistry, database construction, computational tools and evolutionary principles. Her “Atlas of Protein Sequence and Structure” (Dayhoff, 1965) was the forerunner of the giant electronic databases such as GenBank. Not surprisingly, given the effort taken to generate the data, some scientists were a little possessive of their data, so Dayhoff and her collaborators had to persuade people to contribute their hard won sequences<sup>1</sup> (Strassman, 2012). Dayhoff also pioneered bioinformatic analysis, using computational models to examine patterns of molecular evolution (Eck and Dayhoff, 1966), constructing the first phylogeny generated through computational analysis of molecular sequences, using empirically derived frequencies to calibrate transition probabilities (Dayhoff and Eck, 1966). This work formalised the view of the sequence as a document of evolutionary history (Zuckerkanndl and Pauling, 1965).

We now have so much sequence data that we are awash with information. As sequencing vast amounts of DNA becomes routine, the emphasis has shifted to large-scale computation. In only a few decades, the major challenge in molecular evolutionary biology has shifted from the problem of generating sequences and deriving evolutionary history and processes from limited data, to the problems of analysing and making sense of too much data. And so the emphasis has shifted from biochemistry to computing. As a result, we have stepped away from the sequence as representing a real molecule and are more inclined to view the sequence as a string of information. But to read the traces of evolutionary history and mechanism from the comparison of DNA, RNA or protein sequences, we need to know something of the processes that generated those traces. To do so, we need to appreciate that the sequences we analyse are a simplified representation of intricate biomolecular devices operating within living organisms, subject to a complex interacting web of biological processes and evolutionary forces. We need to remind ourselves that the string is the representation, not the reality.

## 2 Comparing substitution rates

Studying substitution rate is much trickier than it first appears. It would seem to be straightforward to compare sequences to come up with an estimate of the number of changes that have happened over evolutionary time from the branch lengths of a molecular phylogeny. But branch lengths reflect the amount of genetic change that has occurred, the rate at which change occurs, and the time period elapsed. None of these things is easy to measure, and often two or more of the quantities are imperfectly known, making the solution to the problem non-identifiable. If we only know only one out of the three qualities – genetic distance, time and rate – there is an infinite set of possible branch length solutions for any observed sequence data (Bromham, 2019).

For many messy problems in biology, we expect the more data we get, the more ability we

---

<sup>1</sup> As an aside, even as the gene databases expanded and went online in 1990s, many lab-based scientists who generated sequence data were somewhat reluctant to share their DNA sequences with “data parasites” who specialised in comparative analysis of sequences that other people had produced.

#### 4.4:6 Substitution Rate Analysis and Molecular Evolution

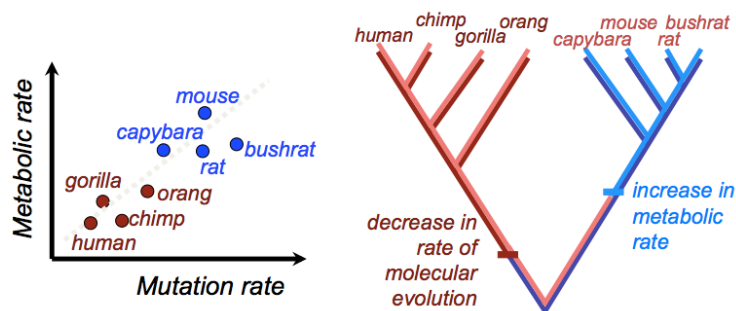
will have to detect signal over noise. But this is not necessarily the case with characterizing substitution rates (Bromham et al., 2017; dos Reis and Yang, 2013; Zhu et al., 2015). In fact, as the amount of data increases and the uncertainty on parameter estimates decreases, it may result in increasing confidence in the wrong answer. For example, an artefact such as long branch attraction, which can cause lineages with a rapid substitution rate to cluster together on the tree, will not necessarily be overcome by using phylogenomic datasets consisting of thousands of genes (e.g. Boussau et al., 2014; Lin et al., 2014). Given that there are many features of organismal biology that will affect the whole genome, we do not necessarily expect rate variation to simply contribute noise to substitution rate estimates, but also systematic bias. If each gene is subject to the same bias, increasing the number of loci can increase precision, but may not increase accuracy, potentially converging on the incorrect estimate (Kubatko and Degnan 2007; Kumar et al. 2012; Philippe et al. 2011; Chapter 2.1 [Simion et al. 2020]). An additional challenge with phylogenomic datasets is the possibility that loci sampled from across the genome may contain different historical narratives. Incongruence between loci may influence estimates of substitution rate, a problem that is likely to increase as more loci are analysed (Mendes and Hahn 2016; Chapters 3.3 and 3.4 [Rannala et al. 2020; Bryant and Hahn 2020]).

The problem is compounded by the evolutionary lability of rates. Substitution rates are shaped by species life history, and therefore can vary between even closely related species. To cite a few examples, rates of molecular evolution vary between mammal species according to their size, generation time, fecundity and longevity (Welch et al., 2008); closely related rockfish species can have different substitution rates if they differ in longevity (Hua et al., 2015); taller plants have slower rates of substitution (Lanfear et al., 2013); flight loss in insects leads to increased substitution rates (Mitterboeck and Adamowicz, 2013); and parasitic plants have faster rates of molecular evolution than their free-living relatives (Bromham et al., 2013). Given the large number of factors that can influence substitution rates, many of which can vary between close relatives, we expect the rate of molecular evolution to evolve as species evolve (Bromham, 2011).

Currently, there are two common approaches to dealing with evolving rates of molecular evolution when estimating substitution rates along a phylogeny for a set of sequences. One is to draw a rate for each branch independently from a convenient distribution, and choose the set of branch rates that maximizes the fit to the data, given a particular model and assumptions (generally referred to as an uncorrelated model, e.g. Drummond et al. 2006). The other is to fit an evolutionary model of rate-change to the data, allowing rates to step up and down at phylogenetic nodes or change continuously along the branches of the phylogeny (an autocorrelated rates model, e.g. Thorne et al. (1998)). All of these models are stochastic in nature and biologically arbitrary (Bromham et al., 2017). They allow rates to vary but are not informed by any special understanding of why or how they do so. There is nothing wrong with this, as long as these stochastic models can reliably capture real patterns of rate variation. But a problem arises when different rate models suggest different solutions, and we have little or no a priori information to help us decide which solution is correct (e.g. Duchêne et al., 2014; Foster et al., 2016; Lepage et al., 2007). There is some evidence that our ability to accurately infer branch length rates (i.e. distance, rates and times) using these stochastic models declines as the level of rate variation across lineages increases (Duchêne et al., 2017). In any case, the substantial variation in rate estimates generated using different methods, models and assumptions tells us that we are not yet able to precisely infer rates with the tools currently available to us. It would be helpful to have a means of studying rate variation independently of variable-rate molecular dating methods (“relaxed clocks”),

so that we can use the knowledge of patterns gained to test the validity of the relaxed clock solutions.

Estimates made independently of the relaxed clock methods may provide something of a reality-check for the phylogenetic rate estimates. Genomic analysis can provide a means of making direct estimates of rates of genome change across generations, for example by tracking genome sequence change in microbes from lab assays (e.g. Bradwell et al., 2013), from serially sampled viruses (e.g. Duffy et al., 2008), dated ancient DNA sequences (e.g. Tong et al., 2018), or in well-studied pedigrees (e.g. Thomas et al., 2018). This direct approach to rate estimation is useful for setting empirically determined bounds on likely mutation rate values, and has been used to seek correlates of variation in rate of molecular evolution (e.g. Thomas et al., 2018). But it has its limitations. Firstly, it is applicable to only a small subset of taxa, though advances in sequencing will put pedigree analysis within reach for an increasing range of species. Secondly, mutation rates estimated in the lab or from pedigrees sometimes seem to have little direct correspondence to the values estimated from phylogenetic studies (Moorjani et al., 2016; Obbard et al., 2012), which suggests that per-generation mutation rates do not necessarily reflect long term substitution rates, even for supposedly neutral substitutions (Ho et al., 2011). Thirdly, it is important to recognise that the rates estimated from related species are likely to be more similar to each other than to randomly chosen species, due to the heritability of factors that influence mutation rate evolution (Lanfear et al., 2010). This complicates the search for consistent patterns in rate variation, because rates from different lineages cannot be treated as independent observations in a statistical analysis. So if rate estimates from each species are plotted against some other feature, such as body size or average temperature, it is not appropriate to conduct a statistical test of the association between rates and traits without correcting for covariation due to relatedness, as treating the observed rates as independent observations does not satisfy the assumption of any general statistical test such as correlation analysis (Figure 2).



■ **Figure 2** Why independent contrasts are necessary for the study of correlates of substitution rates. A toy example showing that if rates change along phylogenies, they can appear to be correlated with species traits that also vary between clades. In this case, because primates have undergone a slowdown in rates, rates will be correlated with anything that differs consistently between primates and rodents – for example having nails instead of claws. Reproduced from *Trends in Ecology and Evolution* 25, 2010 R. Lanfear, J. J. Welch, L. Bromham “Watching the Clock: studying variation in rate of molecular evolution between species” pages 495-503 with permission from Elsevier.

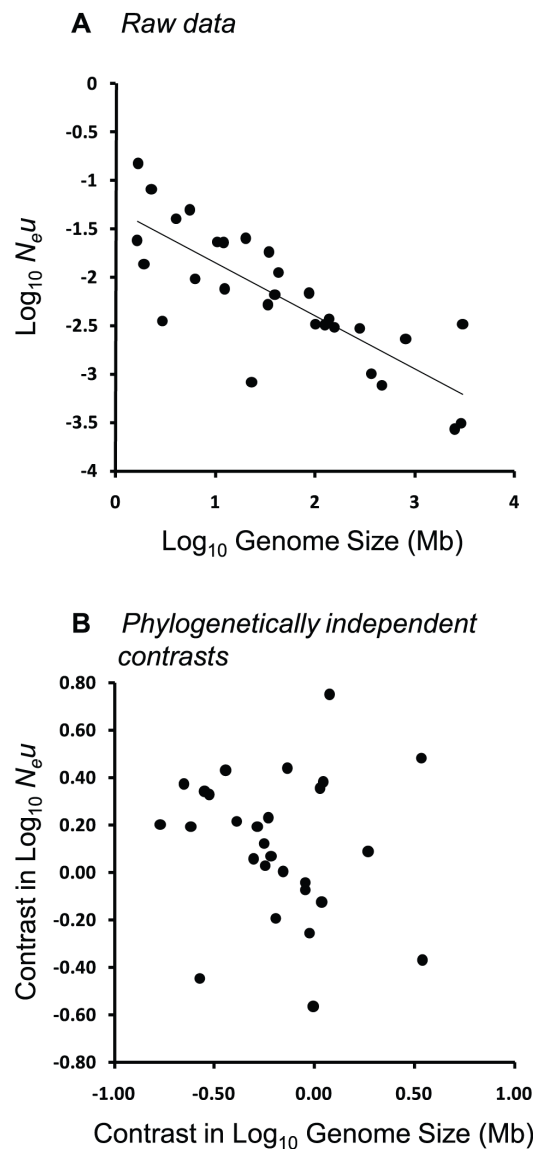


## 2.1 Phylogenetic non-independence of substitution rates

The problem of the non-independence of substitution rates due to shared descent is pertinent when we use substitution rates to answer questions about the driving forces in evolution. For example, the relationship between genome size and rate of genome change has been used to support a hypothesis that genetic drift is a major factor shaping genome evolution in organisms with small effective population sizes (Lynch and Conery, 2003; Lynch, 2010). One of the pieces of evidence provided in support of this hypothesis was a linear relationship between genome size and  $N_e\mu$ , a composite parameter (effective population size and mutation rate) estimated from genetic variability within a species at “silent sites”. But the species-specific estimates compared in such analyses cannot be considered statistically independent observations of the influence of genome size on molecular evolution, because genome size shows phylogenetic signal in at least some groups, meaning that close relatives are more likely to have similar values than they are to randomly chosen species (e.g. Grotkopp et al., 2004; Sessegolo et al., 2016; Waltari and Edwards, 2002). Since mutation rate is influenced by species traits, it too should show phylogenetic inertia. Because species traits that could influence mutation rates, such as population size and genome size, will be more similar between relatives, this generates the potential for spurious correlations between genome size, species traits, and mutation rates (Bromham et al., 2015). A reanalysis using Phylogenetic Least Squares (PGLS) regression indicated that the significant association between effective population size, substitution rates and genome size disappears under correction for phylogenetic nonindependence (Figure 3). This does not invalidate the hypothesis, but suggests that more evidence is needed to give it empirical support, at least as far as cross-species comparisons are concerned.

If a reanalysis using methods that control for phylogenetic relatedness fails to confirm the original study, it may be tempting to conclude that the loss of significance is due to the reduction in number of datapoints reducing statistical power. A better interpretation is that the original study erroneously inflated statistical power by including data points that are effectively replicates of each other, a statistical problem that has long been recognized in evolutionary studies (including by the guy who first formulated the concept of statistical correlation analyses (Galton, 1889)). Using family-averages or including taxonomic levels as a factor in the analysis does not remove the problem of phylogenetic non-independence, because lineages within a family will still show hierarchical structuring according to relatedness, as will between-family contrasts (Bromham et al., 2018).

While there are a number of established methods for dealing with phylogenetic non-independence, care must be taken in applying standard phylogenetic comparative methods to substitution rates. Methods such as PGLS make strong assumptions about the nature of the qualities being analysed and about the way those qualities evolve over time. Specifically, they require inference of states on the internal branches, which cannot be directly measured. Typically, this involves describing the internal nodes as having values consistent with their production from a common ancestral value via a random walk along the connecting branches of the phylogenetic tree. Brownian motion is a handy way to describe random walks in character states along evolutionary trees (Lartillot and Poujol, 2011). However, there are cases of traits where this mode of inference will not provide an accurate inference of ancestral states, for example where the rate of change of traits has varied among lineages, or where there has been a directional trend in values over time (Finarelli and Flynn, 2006; Oakley and Cunningham, 2000). This may be particularly problematic for adaptive radiations such as placental mammal orders, where average body size has increased in most lineages since their last common ancestor (Bromham, 2003; Lartillot and Delsuc, 2012; Phillips, 2015). Since



■ **Figure 3** Controlling for phylogenetic non-independence can influence the statistical support for hypotheses about the drivers of substitution rate variation. (A) Plots of genome size against  $\log_{10}N_e u$  (a composite parameter representing the effective population size and mutation rate, estimated from site variability within each species) have been used to support a causal link between genome size and rate of molecular evolution. (B) The relationship is less distinctly linear when relatedness between taxa is taken into account, and is now not statistically significant to  $p < 0.05$ . Reproduced under Creative Commons Attribution license (CC BY 4.0) from Whitney and Garland (2010).

substitution rates are correlated with body size in mammals, we would not expect change in rate of molecular evolution to follow a Brownian motion model for the placental mammal radiation.

PGLS and related methods are often applied to determining the patterns of molecular evolution as if substitution rates were just like other species traits, such as genome size or body mass, which can be represented as a continuous variable measured with some degree of

#### 4.4:10 Substitution Rate Analysis and Molecular Evolution

error. But substitution rates are something rather different. They are ultimately based on counts of changes that accrue by a stochastic process over time. This gives substitution rate measures a number of important properties that distinguish them from most other species traits, like body size or metabolic rate or niche (Welch and Waxman, 2008). One such property is that past fluctuations in rate can leave a signature in contemporary rate estimates. When we infer average ancestral states for a species trait, like body size, we typically use only the current state at the tips to derive the likely ancestral value. But, because rates reflect substitutions accruing over time, they do not really represent instantaneous measures of a trait value occurring at the tips, coincident with other species trait measurements. They represent a history of accumulation of substitutions, occurring over a protracted period of time. Because of this, a transient increase in substitution rate at some point in the past may have peppered the genome with substitutions that contribute to the estimate of substitution rate assigned to species at the tips, even after the rate has returned to the average value (Lanfear et al., 2010). For example, changes in population size over time could influence on mutation fixation rates, which might then be unrepresentative of species trait values at the tips. This is why substitution rate estimates should not be treated as if they were instantaneous measures of a species rate of genome change.

There is another property of substitution rates that sets them apart from other species traits. The accuracy of most measures of species traits is not dependent on measurements made on other species: the value of metabolic rate for a mouse is independent of whether a rat, a guinea pig or a monkey is also included in the analysis. But the estimation of the substitution rate for the mouse does depend on which other taxa are included in the analysis, as rate estimates are influenced by both the number of species included in a phylogenetic analysis and their relationship to each other. As we add in more species, we have more chance of breaking up long branches with subtending nodes, and this gives more purchase for uncovering past changes now obscured by multiple hits. More species, more nodes, more substitutions, faster rates. While the node density effect is particularly heinous for parsimony analysis, it also applies to estimates of branch length in likelihood and Bayesian methods as well. The practical upshot is that taxon sampling not only influences molecular dating analyses, but can also affect estimates of substitution rates made from phylogenies (Duchêne et al., 2015; Hugall and Lee, 2007; Linder et al., 2005; Phillips, 2015).

One approach to these problems is to simultaneously solve both rates and trait evolution for a phylogeny, then look for evidence of correlation between traits and rates over the whole tree. Whole tree analyses are increasingly being used in substitution rate analyses (e.g. Lourenco et al., 2012; Qiu et al., 2014; Santos, 2012; Wollenberg et al., 2011; Wong, 2014). While some whole-tree methods take the phylogenetic topology and branch lengths as fixed (e.g. Pagel et al., 2006), new methods jointly model rate changes and trait evolution, then assess covariation between trait and rate estimates (e.g. Lartillot and Poujol, 2011). Any use of internal edges of a phylogeny relies on being able to accurately infer past states using only the information at the tips of the tree, and this in turn relies on being able to adequately model evolutionary trajectories (typically using something like a Brownian motion model or Ornstein-Uhlenbeck process). Inference of rates changes along internal edges also requires that relative or absolute dates of divergence are known for all nodes in the phylogeny. Few phylogenies have independent dates for every node (e.g. fossil or biogeographic calibrations), so node heights must be either fixed from a molecular dating analysis (which, of course, relies on making prior assumptions about the way rates evolve over the tree), or co-estimated along with rates and traits (see Chapter 5.1 [Pett and Heath 2020]).

The flipside of using life history traits to understand evolving rates of molecular evolution is to use patterns of molecular evolution to reconstruct ancestral life histories (Wu et al., 2017). For example, models that reconstruct both life history and rates throughout the mammalian tree, using only sequences and species data derived at the tips, have led to the unexpected prediction of an ancestral placental mammal that was larger than the earliest known placental fossils, at least ten times larger than current median value for living mammals species (Jones et al., 2009), slow to mature and with a lifespan over a decade (Lartillot and Delsuc, 2012; Nabholz et al., 2013; Romiguier et al., 2013). The mammals may be a particularly challenging case study for these methods, for although the effect of life history on rates is well-studied for mammals, there is also strong directional trends in life history evolution in most placental mammalian orders (Figuier et al., 2017). For groups with a reliable fossil record, it may be possible to get more traction on evolving rates by using fossils not only to provide a prior distribution on node times, but also a prior distribution on life history traits at ancestral nodes. This would allow the estimation of ancestral rates on a phylogeny to break away from purely stochastic models and be, to some extent, ground-truthed by what we know about the biology of substitution rate variation.

## 2.2 Sister pairs analysis

The method of analysis that is most robust to the problems of comparative analysis of substitution rates is also the simplest<sup>2</sup>. If you compare the differences between two sequences that were originally copied from the same ancestral sequence, then any difference between them must have accrued since their last common ancestor. And if you have information that allows you to guess the position of that ancestor on the path of genetic change that separates them, such as an outgroup or ancestral lineage identified on a phylogeny, then you can compare the relative numbers of substitutions that have accumulated in each lineage since they split. A sister pairs approach does not produce absolute rates of change. But it does produce phylogenetically independent observations of differences in substitution rates that can be profitably used to search for correlates of rates of molecular evolution. More particularly, you can design a test where each sister pair differs in some particular trait of interest, such as life history, niche or behaviour, and you can ask whether the lineages with the greater value of the trait tend to have faster or slower rates than their sisters (Lanfear et al., 2010).

The sister pairs method has a number of advantages. Unlike PGLS and PIC, a sister pairs approach does not require a fully resolved dated phylogeny, because any information on relatedness (e.g. taxonomic information) can be used to choose non-overlapping pairs (pairs that are each others' closest relatives, with respect to any members of any other pairs in your analysis [Bromham et al. 2018]). No calibrations are required, because rates are anchored by the last common ancestor, so each member of the pair has had the same amount of time to accumulate changes. Sister pairs analyses make minimal assumptions about the model of evolution that produced the data (so, for example, they should work even when traits violate a Brownian motion model of change). Choosing a single species to represent each sister lineage removes the possibility of node density, but also forgoes the increased precision of rate estimates that comes from denser taxon sampling. Having a balanced number of taxa per sister clade should improve rate estimates, but cannot guarantee to avoid node density

---

<sup>2</sup> Which, ironically, is something of a disadvantage, as it can be hard to publish simple analyses when more complex methods are available – a kind of reverse Ockham's razor.

#### 4.4:12 Substitution Rate Analysis and Molecular Evolution

entirely if the distribution of speciation events or rate changes is uneven (Bromham et al., 2015; Lanfear et al., 2010). Similarly, choosing only a single locus will avoid artefacts due to gene tree discordance (Mendes and Hahn, 2016), but at the expense of including fewer informative sites.

Sister pairs analyses will solve some of the special problems of comparative analysis of rates, but not all of them. One pervasive challenge is that error in substitution rates is time-dependent, so that accurate inference of rates is tricky at both the “shallow end” and “deep end” of the evolutionary scale (e.g. van Tuinen and Torres, 2015). Systematic patterns of error in rates over time can impact on the assumptions of standard statistical tests, making correlation analyses unreliable. Accurate estimate of substitution rates from recently diverged sequences is tricky as the variance around such estimates is large due to the stochastic accrual of sequence changes. It may be tempting to dismiss poor estimates of rate due to few observable substitutions as inconsequential noise that should be overwhelmed by more robust rate estimates. But for a comparative analysis this need not necessarily be true. Welch and Waxman (2008) show how including poorly informative contrasts at the shallow end of divergence can reduce the power of comparative tests, and they recommend using simple diagnostic tests to remove these troublesome contrasts from analyses. While deleting data points can lead to a deep sense of loss, associated with the feeling that one is “throwing away data”, it is preferable to being misled due to the inclusion of poor quality datapoints in an analysis, and it could lead to an ability to detect a pattern that was previously marred by the shallow datapoints (Welch and Waxman, 2008).

However, the Welch & Waxman test requires some estimate of comparison depth so that variance can be plotted against time for all contrasts. For most phylogenies, time depth comes from molecular dates, which introduces a worrying circularity for the study of the correlates of substitution rate variation. An alternative approach does not require divergence dates yet allows inclusion of shallow contrasts, by modelling the accumulation of substitutions as a Poisson process (Hua et al., 2015). The power of such comparative tests depends not only on the amount of data, but also the absolute substitution rate and also the rate of change in related species characteristics. Increasing the number of loci analysed in phylogenomic studies will help to determine the substitution rate, particularly for shallow contrasts. But for most studies, adding more independent comparisons will bring the greatest benefit in increasing the ability to detect meaningful patterns in the evolution and determination of substitution rates.

### 2.3 Phylogenomic data and substitution rate analysis

Phylogenomic data may help at the shallow end if including more sequence data provides a larger sample of substitutions. But it will not necessarily help at the deeper end, if too many changes have overwritten past changes. Multiple hits cause irreversible erasure of historical signal: when a site in a sequence changes more than once, the previous nucleotide states are overwritten. Overwritten history cannot be recovered, no matter how many saturated sites you look at (Bromham, 2019). Instead, we rely on models of the substitution process to guess how many changes we might no longer be able to observe, based on the pattern of those that we can see. Phylogenomic datasets may allow you a greater choice of markers to identify sequences or sites that are evolving slowly enough to avoid saturation at deep time depths, but this advantage might be lost if all loci are analysed together without discrimination. Of course, neither the deep end nor the shallow end are defined by absolute time, but by the combination of rate, time and number of observed changes (shaped by both the number of observed sites free to vary and the ability to estimate unseen changes using an evolutionary

model).

Thus far studies of correlates of substitution rates have been limited in their use of phylogenomic data. But there are many possible advantages of using a larger sample of genomic loci (Wilson Sayres et al., 2011). Multi-locus datasets provide the potential to decompose rates into gene specific and lineage-specific components (Rasmussen and Kellis, 2007). Large, genome-wide datasets may help estimate rates for shallower comparisons, allowing more meaningful comparisons between sister species. However, more loci do not necessarily provide more power to detect significant patterns in rates. For example, a phylogenomic study of rates in herbaceous and woody plants identified 5 independent comparisons between sister lineages (Yang et al., 2015). The large number of loci may provide a more comprehensive sample of sites to characterise rates across the genome, such that the rate difference for each comparison has greater confidence, but the power of the test to detect a correlation between growth habit and rates is determined by the number of independent comparisons (equivalent to the sample size in an experiment or observational study). To provide convincing test of a link between woodiness and rates, more sister comparisons would be needed, regardless of the amount of sequence data available.

### **3** Substitution rates shape our view of evolutionary history

The analysis of substitution rates has played an important part of developing and testing hypotheses of the drivers of molecular evolution, and the connection between change at the genotypic and phenotypic levels. But, curiously, the study of patterns of substitution rates has thus far had relatively little impact on one of the fields where you would expect it to play a most important role. Modern molecular dating methods depend entirely on an ability to infer patterns of substitution rates over the tree, but currently the models they use are almost entirely biologically arbitrary. Very few molecular dating studies use any empirically-derived information about the way substitution rates evolve. That does not matter if our current models are up to the job. But the range of answers it is possible to get from molecular dating analyses, and the difference between published studies using the same sequence data but different methods, models and prior assumptions, suggests that we still have some way to go before we can trust molecular date estimates.

Placental mammals provide an interesting case study, for two reasons. Firstly, rates of molecular evolution have been intensively studied in mammals, and clear patterns have emerged that substitution rates are significantly associated with body size and other aspects of life history (Bromham et al., 1996; Galtier et al., 2009; Welch et al., 2008). Secondly, molecular dates for the mammalian radiation are as old as the concept of the molecular clock itself (e.g. Doolittle and Blomback, 1964; Margoliash, 1963; Zuckerkandl and Pauling, 1965), and have, for much of that history, been controversially out of step with the story told from fossil evidence alone (e.g. Bininda-Emonds et al., 2007; Hasegawa et al., 2003; Sarich and Wilson, 1967; Murphy et al., 2001). While newer molecular dating studies also tend to put the diversification of placentals in the Cretaceous, the gap between fossil and molecular dates is perceived to be shrinking (e.g. dos Reis et al., 2016; Goswami, 2012; Phillips, 2015; Ronquist et al., 2016).

This looks like a progress: more sophisticated methods and bigger datasets give us an answer that fits more snugly with both the paleontological record (fossil evidence of modern placental orders confined to post-Cretaceous) and our understanding of mammalian molecular evolution (smaller species have faster rates). It has become “a dating success story” (Goswami, 2012). But there is reason to pause for thought. The new molecular dates are

driven by two features of the new Bayesian molecular dating methods: variable-rate models and prior distributions on node height based on fossil evidence. If fossil calibrations are enforced as providing strong bounds on maximum ages, then the solution must infer very fast substitution rates on the early branches of the tree, in order to fit the sequence data to the fossil dates (O’Leary et al., 2013). If the bounds on ages are relaxed, to allow a distribution of possible ages informed by fossil data, then this allows lower rate estimates and older dates (dos Reis et al., 2014). Comparison of the prior and posterior distributions on node heights suggests that the calibrating information is strongly informative, and that estimated nodes rarely fall outside the joint prior, which may be shaped by the prior distributions on calibrations, rates, and tree shape (dos Reis et al., 2012).

It has also been suggested that the molecular dates for the placental radiation are systematically biased by uneven sampling of living mammal species, because larger-bodied contemporary species overestimate rates for the presumably smaller-bodied ancestral lineages (Phillips, 2015). A related size-biased effect has been proposed for molecular dates for the radiation of modern birds (Berv and Field, 2018). The case of the placental mammals illustrates how decisions made regarding data inclusion, calibration and other aspects of analysis can lead to substantial differences in the estimates of substitution rates and dates of divergence (e.g. dos Reis et al., 2014; Gatesy and Springer, 2017; Phillips, 2015; Springer et al., 2018; Wu et al., 2017). Despite growing confidence in molecular dating methods, there is still plenty of disagreement on molecular dates for the placental mammal radiation. So even in this case study, where we have the best understanding of the determinants of substitution rate evolution of any taxonomic group, we still have quite a long way to go before we can be sure that our molecular date estimates are not just telling us what we wanted to hear.

What has all this got to do with phylogenomics? These are systemic problems in our analysis that will not necessarily be solved by adding more data. We cannot have faith that our molecular dates will be better the more loci we include. But phylogenomic datasets give us a fantastically useful tool for understanding the way rates evolve, across the genome, over time and between lineages. The hope is that more we know about the way the historic record is written in the genome, the better we will get at reading it.

## Acknowledgements

Thanks to Matt Hahn, Rob Lanfear and Xia Hua, specifically for their helpful comments on this chapter, but more generally for many wonderfully interesting and stimulating conversations on this and other topics.

## References

- Anfinsen, C. B. (1959). *The molecular basis of evolution*. John Wiley and Son, New York.
- Anfinsen, C. B. (1965). Evolution of proteins I: Chairman’s remarks. In Bryson, V. and Vogel, H. J., editors, *Evolving genes and proteins*, page 95. Academic Press, New York.
- Aronson, J. D. (2002). Molecules and monkeys: George Gaylord Simpson and the challenge of molecular evolution. *History and Philosophy of the Life Sciences*, 24(3-4):441–465.
- Berv, J. S. and Field, D. J. (2018). Genomic signature of an avian lilliput effect across the K-Pg extinction. *Systematic Biology*, 67(1):1–13.
- Bininda-Emonds, O., Cardillo, M., Jones, K., MacPhee, R., Beck, R., Grenyer, R., Price, S.,

- Vos, R., Gittleman, J., and Purvis, A. (2007). The delayed rise of present-day mammals. *Nature*, 446:507–512.
- Boussau, B., Walton, Z., Delgado, J. A., Collantes, F., Beani, L., Stewart, I. J., Cameron, S. A., Whitfield, J. B., Johnston, J. S., Holland, P. W. H., Bachtrog, D., Kathirithamby, J., and Huelsenbeck, J. P. (2014). Strepsiptera, phylogenomics and the long branch attraction problem. *PLOS ONE*, 9(10):e107709.
- Bradwell, K., Combe, M., Domingo-Calap, P., and Sanjuán, R. (2013). Correlation between mutation rate and genome size in riboviruses: mutation rate of bacteriophage Q $\beta$ . *Genetics*, 195(1):243–251.
- Britten, R. J. and Kohne, D. E. (1968). Repeated sequences in DNA: Hundreds of thousands of copies of DNA sequences have been incorporated into the genomes of higher organisms. *Science*, 161(3841):529–540.
- Bromham, L. (2003). Molecular clocks and explosive radiations. *Journal of Molecular Evolution*, 57(1):S13–S20.
- Bromham, L. (2011). The genome as a life-history character: Why rate of molecular evolution varies between mammal species. *Philosophical Transactions of the Royal Society of LondonB: Biological Sciences*, 366(1577):2503–2513.
- Bromham, L. (2019). Six impossible things before breakfast: Assumptions, models, and belief in molecular dating. *Trends in Ecology and Evolution*, 34(5):474–486.
- Bromham, L., Cowman, P. F., and Lanfear, R. (2013). Parasitic plants have increased rates of molecular evolution across all three genomes. *BMC Evolutionary Biology*, 13:126.
- Bromham, L., Duchêne, S., Hua, X., Ritchie, A., Duchêne, D., and Ho, S. (2017). Bayesian molecular dating: Opening up the black box. *Biological Reviews*, 93(2):1165–1191.
- Bromham, L., Hua, X., Cardillo, M., Schneemann, H., and Greenhill, S. J. (2018). Parasites and politics: why cross-cultural studies must control for relatedness, proximity and covariation. *Royal Society Open Science*, 5(8):181100.
- Bromham, L., Hua, X., Lanfear, R., and Cowman, P. (2015). Exploring the relationships between mutation rates, life history, genome size, environment and species richness in flowering plants. *American Naturalist*, 185(4):507–524.
- Bromham, L., Rambaut, A., and Harvey, P. H. (1996). Determinants of rate variation in mammalian DNA sequence evolution. *Journal of Molecular Evolution*, 43(6):610–621.
- Bryant, D. and Hahn, M. W. (2020). The concatenation question. In Scornavacca, C., Delsuc, F., and Galtier, N., editors, *Phylogenetics in the Genomic Era*, chapter 3.4, pages 3.4:1–3.4:23. No commercial publisher | Authors open access book.
- Buettner-Janusch, J. and Hill, R. L. (1965a). Evolution of haemoglobins in primates. In Bryson, V. and Vogel, H. J., editors, *Evolving genes and proteins*, pages 167–181. Academic Press, New York.
- Buettner-Janusch, J. and Hill, R. L. (1965b). Molecules and monkeys. *Science*, 147(3660):836–842.
- Cain, A. J. (1951). Non-adaptive or neutral characters in evolution. *Nature*, 168:1049.
- Charlesworth, B., Charlesworth, D., Coyne, J. A., and Langley, C. H. (2016). Hubby and Lewontin on protein variation in natural populations: When molecular genetics came to the rescue of population genetics. *Genetics*, 203(4):1497–1503.
- Darwin, C. (1859). *The origin of species by means of natural selection: or the preservation of favoured races in the struggle for life*. John Murray, London, first edition.
- Dayhoff, M. O. (1965). *Atlas of protein sequence and structure*, volume 1. National Biomedical Research Foundation, Washington.



#### 4.4:16 REFERENCES

- Dayhoff, M. O. and Eck, R. V. (1966). *Atlas of protein sequence and structure*. Biomedical Research Foundation, Washington.
- Dayhoff, M. O. and Eck, R. V. (1969). Inferences from protein sequence studies. In Dayhoff, M. O., editor, *Atlas of protein sequence and structure*, volume 4. Biomedical Research Foundation, Washington.
- Dickerson, R. E. (1971). The structure of cytochrome c and rates of molecular evolution. *Journal of Molecular Evolution*, 1(1):26–45.
- Dietrich, M. (1994). The origins of the neutral theory of molecular evolution. *Journal of the History of Biology*, 27(1):21–59.
- Dobzhansky, T. and Wright, S. (1941). Genetics of natural populations. V. Relations between mutation rate and accumulation of lethals in populations of *Drosophila pseudoobscura*. *Genetics*, 26(1):23.
- Doolittle, R. F. and Blomback, B. (1964). Amino-acid sequence investigations of fibrinopeptides from various mammals: Evolutionary implications. *Nature*, 202(4928):147–152.
- dos Reis, M., Donoghue, P. C., and Yang, Z. (2014). Neither phylogenomic nor palaeontological data support a Palaeogene origin of placental mammals. *Biology Letters*, 10(1):20131003.
- dos Reis, M., Donoghue, P. C., and Yang, Z. (2016). Bayesian molecular clock dating of species divergences in the genomics era. *Nature Reviews Genetics*, 17(2):71–80.
- dos Reis, M., Inoue, J., Hasegawa, M., Asher, R. J., Donoghue, P. C., and Yang, Z. (2012). Phylogenomic datasets provide both precision and accuracy in estimating the timescale of placental mammal phylogeny. *Proceedings of the Royal Society of London B: Biological Sciences*, 279(1742):3491–3500.
- dos Reis, M. and Yang, Z. (2013). The unbearable uncertainty of Bayesian divergence time estimation. *Journal of Systematics and Evolution*, 51(1):30–43.
- Drummond, A. J., Ho, S. Y. W., Phillips, M. J., and Rambaut, A. (2006). Relaxed phylogenetics and dating with confidence. *PLOS Biology*, 4(5):e88.
- Duchêne, D. A., Duchêne, S., and Ho, S. Y. W. (2015). Tree imbalance causes a bias in phylogenetic estimation of evolutionary timescales using heterochronous sequences. *Molecular Ecology Resources*, 15(4):785–794.
- Duchêne, D. A., Hua, X., and Bromham, L. (2017). Phylogenetic estimates of diversification rate are affected by molecular rate variation. *Journal of Evolutionary Biology*, 30(10):1884–1897.
- Duchêne, S., Lanfear, R., and Ho, S. Y. W. (2014). The impact of calibration and clock-model choice on molecular estimates of divergence times. *Molecular Phylogenetics and Evolution*, 78:277–289.
- Duffy, S., Shackelton, L. A., and Holmes, E. C. (2008). Rates of evolutionary change in viruses: patterns and determinants. *Nature Reviews Genetics*, 9(4):267–276.
- Eck, R. V. and Dayhoff, M. O. (1966). Evolution of the structure of ferredoxin based on living relics of primitive amino acid sequences. *Science*, 152(3720):363–366.
- Fay, J. C. and Wu, C.-I. (2001). The neutral theory in the genomic era. *Current Opinion in Genetics and Development*, 11(6):642–646.
- Figuet, E., Ballenghien, M., Lartillot, N., and Galtier, N. (2017). Reconstruction of body mass evolution in the Cetartiodactyla and mammals using phylogenomic data. *bioRxiv*, page 139147.

- Finarelli, J. A. and Flynn, J. J. (2006). Ancestral state reconstruction of body size in the Caniformia (Carnivora, Mammalia): The effects of incorporating data from the fossil record. *Systematic Biology*, 55(2):301–313.
- Fisher, R. A. and Ford, E. B. (1950). The "Sewall Wright" effect. *Heredity*, 4:117–19.
- Foster, C. S. P., Sauquet, H., Van der Merwe, M., McPherson, H., Rossetto, M., and Ho, S. Y. W. (2016). Evaluating the impact of genomic data and priors on Bayesian estimates of the angiosperm evolutionary timescale. *Systematic Biology*, 66(3):338–351.
- Galtier, N., Blier, P. U., and Nabholz, B. (2009). Inverse relationship between longevity and evolutionary rate of mitochondrial proteins in mammals and birds. *Mitochondrion*, 9(1):51–57.
- Galton, F. (1889). Comment on 'On a method of investigating the development of institutions; applied to laws of marriage and descent' by E. B. Tylor. *The Journal of the Anthropological Institute of Great Britain and Ireland*, 18:245–272.
- Gatesy, J. and Springer, M. S. (2017). Phylogenomic red flags: Homology errors and zombie lineages in the evolutionary diversification of placental mammals. *Proceedings of the National Academy of Sciences*, page 201715318.
- Gossmann, T. I., Keightley, P. D., and Eyre-Walker, A. (2012). The effect of variation in the effective population size on the rate of adaptive molecular evolution in eukaryotes. *Genome Biology and Evolution*, 4(5):658–667.
- Goswami, A. (2012). A dating success story: genomes and fossils converge on placental mammal origins. *EvoDevo*, 3(1):18.
- Grotkopp, E., Rejmánek, M., Sanderson, M. J., and Rost, T. L. (2004). Evolution of genome size in pines (*Pinus*) and its life-history correlates: supertree analyses. *Evolution*, 58(8):1705–1729.
- Haldane, J. B. S. (1949). Suggestions as to quantitative measurement of rates of evolution. *Evolution*, 3(1):51–56.
- Harris, H. (1966). C. genetics of man enzyme polymorphisms in man. *Proceedings of the Royal Society of London B: Biological Sciences*, 164(995):298–310.
- Hasegawa, M., Thorne, J. L., and Kishino, H. (2003). Time scale of eutherian evolution estimated without assuming a constant rate of molecular evolution. *Genes and Genetic Systems*, 78(4):267–283.
- Ho, S. Y. W., Lanfear, R., Bromham, L., Phillips, M. J., Soubrier, J., Rodrigo, A. G., and Cooper, A. (2011). Time-dependent rates of molecular evolution. *Molecular Ecology*, 20(15):3087–3101.
- Hua, X., Cowman, P., Warren, D., and Bromham, L. (2015). Longevity is linked to mitochondrial mutation rates in rockfish: a test using Poisson regression. *Molecular Biology and Evolution*, 32(10):2633–2645.
- Hubby, J. L. and Lewontin, R. C. (1966). A molecular approach to the study of genic heterozygosity in natural populations. I. The number of alleles at different loci in *Drosophila pseudoobscura*. *Genetics*, 54(2):577–594.
- Hugall, A. F. and Lee, M. S. Y. (2007). The likelihood node density effect and consequences for evolutionary studies of molecular rates. *Evolution*, 61(10):2293–2307.
- Jones, K. E., Bielby, J., Cardillo, M., Fritz, S. A., O'Dell, J., Orme, C. D. L., Safi, K., Sechrest, W., Boakes, E. H., and Carbone, C. (2009). PanTHERIA: a species-level database of life history, ecology, and geography of extant and recently extinct mammals. *Ecology*, 90(9):2648–2648.
- Jukes, T. H. (1966). *Molecules and evolution*. Columbia University Press, New York.

- Kern, A. D. and Hahn, M. W. (2018). The neutral theory in light of natural selection. *Molecular Biology and Evolution*, 35(6):1366–1371.
- Kimura, M. (1968). Evolutionary rate at the molecular level. *Nature*, 217(5129):624–626.
- King, J. L. and Jukes, T. H. (1969). Non-Darwinian evolution. *Science*, 164(3881):788–798.
- Kubatko, L. S. and Degnan, J. H. (2007). Inconsistency of phylogenetic estimates from concatenated data under coalescence. *Systematic Biology*, 56(1):17–24.
- Kumar, S., Filipski, A. J., Battistuzzi, F. U., Kosakovsky P., S. L., and Tamura, K. (2012). Statistics and truth in phylogenomics. *Molecular Biology and Evolution*, 29(2):457–472.
- Laird, C. D., McConaughy, B. L., and McCarthy, B. J. (1969). Rate of fixation of nucleotide substitutions in evolution. *Nature*, 224:149 – 154.
- Lanfear, R., Ho, S. Y. W., Davies, T. J., Moles, A. T. Aarssen, L., Swenson, N. G., Warman, L., Zanne, A. E., and Allen, A. P. (2013). Taller plants have lower rates of molecular evolution: the rate of mitosis hypothesis. *Nature Communications*, 4(1):1879.
- Lanfear, R., Welch, J. J., and Bromham, L. (2010). Watching the clock: Studying variation in rates of molecular evolution. *Trends in Ecology and Evolution*, 25(9):495–503.
- Lartillot, N. and Delsuc, F. (2012). Joint reconstruction of divergence times and life-history evolution in placental mammals using a phylogenetic covariance model. *Evolution*, 66(6):1773–1787.
- Lartillot, N. and Poujol, R. (2011). A phylogenetic model for investigating correlated evolution of substitution rates and continuous phenotypic characters. *Molecular Biology and Evolution*, 28(1):729–744.
- Lepage, T., Bryant, D., Philippe, H., and Lartillot, N. (2007). A general comparison of relaxed molecular clock models. *Molecular Biology and Evolution*, 24(12):2669–80.
- Lewontin, R. C. (1974). *The genetic basis of evolutionary change*. Columbia University Press, New York.
- Lewontin, R. C. and Hubby, J. L. (1966). A molecular approach to the study of genic heterozygosity in natural populations. II. Amount of variation and degree of heterozygosity in *Drosophila pseudoobscura*. *Genetics*, 54(2):595–609.
- Lin, J., Chen, G., Gu, L., Shen, Y., Zheng, M., Zheng, W., Hu, X., Zhang, X., Qiu, Y., Liu, X., and Jiang, C. (2014). Phylogenetic affinity of tree shrews to glires is attributed to fast evolution rate. *Molecular Phylogenetics and Evolution*, 71:193–200.
- Linder, H. P., Hardy, C. R., and Rutschmann, F. (2005). Taxon sampling effects in molecular clock dating: an example from the African Restionaceae. *Molecular Phylogenetics and Evolution*, 35(3):569–582.
- Lourenco, J. M., Glemin, S., Chiari, Y., and Galtier, N. (2012). The determinants of the molecular substitution process in turtles. *Journal of Evolutionary Biology*, 26:38–50.
- Lynch, M. (2010). Evolution of the mutation rate. *Trends in Genetics*, 26(8):345–352.
- Lynch, M., Ackerman, M. S., Gout, J.-F., Long, H., Sung, W., Thomas, W. K., and Foster, P. L. (2016). Genetic drift, selection and the evolution of the mutation rate. *Nature Reviews Genetics*, 17:704.
- Lynch, M. and Conery, J. S. (2003). The origins of genome complexity. *Science*, 302(5649):1401–1404.
- Margoliash, E. (1963). Primary structure and the evolution of cytochrome c. *Proceedings of the National Academy of Sciences*, 50(4):672–679.
- Mendes, F. K. and Hahn, M. W. (2016). Gene tree discordance causes apparent substitution rate variation. *Systematic Biology*, 65(4):711–721.

- Mitterboeck, T. F. and Adamowicz, S. J. (2013). Flight loss linked to faster molecular evolution in insects. *Proceedings of the Royal Society of London B: Biological Sciences*, 280(1767):20131128.
- Moorjani, P., Gao, Z., and Przeworski, M. (2016). Human germline mutation and the erratic evolutionary clock. *PLOS Biology*, 14(10):e2000744.
- Murphy, W. J., Eizirik, E., O'Brien, S. J., Madsen, O., Scally, M., Douady, C. J., Teeling, E., Ryder, O. A., Stanhope, M. J., and de Jong, W. W. (2001). Resolution of the early placental mammal radiation using Bayesian phylogenetics. *Science*, 294(5550):2348–2351.
- Nabholz, B., Uwimana, N., and Lartillot, N. (2013). Reconstructing the phylogenetic history of long-term effective population size and life-history traits using patterns of amino acid replacement in mitochondrial genomes of mammals and birds. *Genome Biology and Evolution*, 5(7):1273–1290.
- Nei, M., Suzuki, Y., and Nozawa, M. (2010). The neutral theory of molecular evolution in the genomic era. *Annual Review of Genomics and Human Genetics*, 11:265–289.
- Oakley, T. H. and Cunningham, C. W. (2000). Independent contrasts succeed where ancestor reconstruction fails in a known bacteriophage phylogeny. *Evolution*, 54(2):397–405.
- Obbard, D. J., Maclennan, J., Kim, K.-W., Rambaut, A., O'Grady, P. M., and Jiggins, F. M. (2012). Estimating divergence dates and substitution rates in the *Drosophila* phylogeny. *Molecular Biology and Evolution*, page mss150.
- Ohta, T. (1972). Evolutionary rate of cistrons and DNA divergence. *Journal of Molecular Evolution*, 1(2):150–157.
- Ohta, T. (1973). Slightly deleterious mutant substitutions in evolution. *Nature*, 246(5428):96–98.
- O'Leary, M. A., Bloch, J. I., Flynn, J. J., Gaudin, T. J., Giallombardo, A., Giannini, N. P., Goldberg, S. L., Kraatz, B. P., Luo, Z.-X., Meng, J., Ni, X., Novacek, M. J., Perini, F. A., Randall, Z. S., Rougier, G. W., Sargis, E. J., Silcox, M. T., Simmons, N. B., Spaulding, M., Velazco, P. M., Weksler, M., Wible, J. R., and Cirranello, A. L. (2013). The placental mammal ancestor and the post-K-Pg radiation of placentals. *Science*, 339(6120):662–667.
- Pagel, M., Venditti, C., and Meade, A. (2006). Large punctuational contribution of speciation to evolutionary divergence at the molecular level. *Science*, 314(5796):119–121.
- Pett, W. and Heath, T. A. (2020). Inferring the timescale of phylogenetic trees from fossil data. In Scornavacca, C., Delsuc, F., and Galtier, N., editors, *Phylogenetics in the Genomic Era*, chapter 5.1, pages 5.1:1–5.1:18. No commercial publisher | Authors open access book.
- Philippe, H., Brinkmann, H., Lavrov, D. V., Littlewood, D. T. J., Manuel, M., Wörheide, G., and Baurain, D. (2011). Resolving difficult phylogenetic questions: why more sequences are not enough. *PLOS Biology*, 9(3):e1000602.
- Phillips, M. J. (2015). Geomolecular dating and the origin of placental mammals. *Systematic Biology*, 65(3):546–557.
- Qiu, F., Kitchen, A., Burleigh, J. G., and Miyamoto, M. M. (2014). Scombroid fishes provide novel insights into the trait/rate associations of molecular evolution. *Journal of Molecular Evolution*, 78(6):338–348.
- Rannala, B., Edwards, S. V., Leaché, A., and Yang, Z. (2020). The multi-species coalescent model and species tree inference. In Scornavacca, C., Delsuc, F., and Galtier, N., editors, *Phylogenetics in the Genomic Era*, chapter 3.3, pages 3.3:1–3.3:21. No commercial publisher | Authors open access book.

#### 4.4:20 REFERENCES

- Rasmussen, M. D. and Kellis, M. (2007). Accurate gene-tree reconstruction by learning gene- and species-specific substitution rates across multiple complete genomes. *Genome Research*, 17(12):1932–1942.
- Robinson-Rechavi, M. (2020). Molecular evolution and gene function. In Scornavacca, C., Delsuc, F., and Galtier, N., editors, *Phylogenetics in the Genomic Era*, chapter 4.2, pages 4.2:1–4.2:20. No commercial publisher | Authors open access book.
- Romiguier, J., Ranwez, V., Douzery, E. J. P., and Galtier, N. (2013). Genomic evidence for large, long-lived ancestors to placental mammals. *Molecular Biology and Evolution*, 30(1):5–13.
- Ronquist, F., Lartillot, N., and Phillips, M. J. (2016). Closing the gap between rocks and clocks using total-evidence dating. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 371(1699):20150136.
- Santos, J. C. (2012). Fast molecular evolution associated with high active metabolic rates in poison frogs. *Molecular Biology and Evolution*, 29(8):2001–2018.
- Sarich, V. M. and Wilson, A. C. (1967). Immunological time scale for hominid evolution. *Science*, 158(805):1200.
- Sessegolo, C., Burlet, N., and Haudry, A. (2016). Strong phylogenetic inertia on genome size and transposable element content among 26 species of flies. *Biology Letters*, 12(8):20160407.
- Simion, P., Delsuc, F., and Philippe, H. (2020). To what extent current limits of phylogenomics can be overcome? In Scornavacca, C., Delsuc, F., and Galtier, N., editors, *Phylogenetics in the Genomic Era*, chapter 2.1, pages 2.1:1–2.1:34. No commercial publisher | Authors open access book.
- Simpson, G. G. (1964). Organisms and molecules in evolution. *Science*, 146(3651):1535–1538.
- Springer, M. S., Murphy, W. J., and Roca, A. L. (2018). Appropriate fossil calibrations and tree constraints uphold the Mesozoic divergence of solenodons from other extant mammals. *Molecular Phylogenetics and Evolution*, 121:158–165.
- Stoltzfus, A. (2017). Why we don't want another 'synthesis'. *Biology Direct*, 12(1):23.
- Strassman, B. J. (2012). Dayhoff, M. O. In *Encyclopedia of Life Sciences*. Wiley.
- Thomas, G. W. C., Wang, R. J., Puri, A., Harris, R. A., Raveendran, M., Hughes, D., Murali, S., Williams, L., Doddapaneni, H., and Muzny, D. (2018). Reproductive longevity predicts mutation rates in primates. *Current Biology*, 28(19).
- Thorne, J. L., Kishino, H., and Painter, I. S. (1998). Estimating the rate of evolution of the rate of molecular evolution. *Molecular Biology and Evolution*, 15(12):1647–1657.
- Tong, K. J., Duchêne, D. A., Duchêne, S., Geoghegan, J. L., and Ho, S. Y. W. (2018). A comparison of methods for estimating substitution rates from ancient DNA sequence data. *BMC Evolutionary Biology*, 18(1):70.
- van Tuinen, M. and Torres, C. R. (2015). Potential for bias and low precision in molecular divergence time estimation of the canopy of life: an example from aquatic bird families. *Frontiers in Genetics*, 6:203.
- Waltari, E. and Edwards, S. V. (2002). Evolutionary dynamics of intron size, genome size, and physiological correlates in archosaurs. *The American Naturalist*, 160(5):539–552.
- Welch, J. J., Bininda-Emonds, O. R. P., and Bromham, L. (2008). Correlates of substitution rate variation in mammalian protein-coding sequences. *BMC Evolutionary Biology*, 8:53.
- Welch, J. J. and Waxman, D. (2008). Calculating independent contrasts for the comparative study of substitution rates. *Journal of Theoretical Biology*, 251(4):667–678.

- Whitney, K. D. and Garland, T., J. (2010). Did genetic drift drive increases in genome complexity? *PLOS Genetics*, 6(8):e1001080.
- Wilson, A. C. and Sarich, V. M. (1969). A molecular timescale for human evolution. *Proceedings of the National Academy of Sciences*, 63(4):1088–1093.
- Wilson Sayres, M. A., Venditti, C., Pagel, M., and Makova, K. D. (2011). Do variations in substitution rates and male mutation bias correlate with life-history traits? a study of 32 mammalian genomes. *Evolution*, 65(10):2800–2815.
- Wollenberg, K. C., Vieites, D. R., Glaw, F., and Vences, M. (2011). Speciation in little: the role of range and body size in the diversification of Malagasy mantellid frogs. *BMC Evolutionary Biology*, 11(1):217.
- Wong, A. (2014). Covariance between testes size and substitution rates in primates. *Molecular Biology and Evolution*, 31(6):1432–1436.
- Wu, J., Yonezawa, T., and Kishino, H. (2017). Rates of molecular evolution suggest natural history of life history traits and a post-K-Pg nocturnal bottleneck of placentals. *Current Biology*, 27(19):3025–3033.e5.
- Yang, Y., Moore, M. J., Brockington, S. F., Soltis, D. E., Wong, G. K.-S., Carpenter, E. J., Zhang, Y., Chen, L., Yan, Z., Xie, Y., Sage, R. F., Covshoff, S., Hibberd, J. M., Nelson, M. N., and Smith, S. A. (2015). Dissecting molecular evolution in the highly diverse plant clade Caryophyllales using transcriptome sequencing. *Molecular Biology and Evolution*, 32(8):2001–2014.
- Zhang, J. and Yang, J.-R. (2015). Determinants of the rate of protein sequence evolution. *Nature Reviews Genetics*, 16:409–420.
- Zhu, T., Dos Reis, M., and Yang, Z. (2015). Characterization of the uncertainty of divergence time estimation under relaxed molecular clock models using multiple loci. *Systematic Biology*, 64(2):267–280.
- Zuckerkandl, E. and Pauling, L. (1965). Evolutionary divergence and convergence in proteins. *Evolving Genes and Proteins*, 97:97–166.