

The phylodynamic threshold of measurably evolving populations

Ariane Weber^{1,*}, Julia Kende^{2,3}, Camila Duitama González^{3,4}, Sanni Översti^{1,‡} and Sebastian Duchene^{3,4,5,‡,*}.

¹ Max Planck Institute of Geoanthropology, Jena, Germany.

² Bioinformatics and Biostatistics Hub, Institut Pasteur, Paris, France.

³ Université Paris Cité, Paris, France.

⁴ ED-ID unit, Dept of Computational Biology, Institut Pasteur, Paris, France.

⁵ Peter Doherty Institute for Infection and Immunity, Dept of Microbiology and Immunology, University of Melbourne, Melbourne, Australia.

*email: weber@gea.mpg.de, sduchene@pasteur.fr

‡ Equal contribution to the supervision of this work.

Abstract The molecular clock is a fundamental tool for understanding the time and pace of evolution, requiring calibration information alongside molecular data. Sampling times are often used for calibration since some organisms accumulate enough mutations over the course of their sampling period. This practice ties together two key concepts: measurably evolving populations and the phylodynamic threshold. Our current understanding suggests that populations meeting these criteria are suitable for molecular clock calibration via sampling times. However, the definitions and implications of these concepts remain unclear. Using Hepatitis B virus-like simulations and analyses of empirical data, this study shows that determining whether a population is measurably evolving or has reached the phylodynamic threshold does not only depend on the data, but also on model assumptions and sampling strategies. In Bayesian applications, a lack of temporal signal due to a narrow sampling window results in a prior that is overly informative relative to the data, such that a prior that is potentially misleading typically requires a wider sampling window than one that is reasonable. In our analyses we demonstrate that assessing prior sensitivity is more important than the outcome of tests of temporal signal. Our results offer guidelines to improve molecular clock inferences and highlight limitations in molecular sequence sampling procedures.

Keywords: Measurably evolving population, phylodynamic threshold, molecular clock, Bayesian phylogenetics, microbial evolution.

1 Introduction

Molecular sequence data have become nearly ubiquitous for studying the evolution of modern and ancient organisms. A fundamental concept in molecular evolution is the ‘molecular clock’, which posits that substitutions accumulate roughly constantly over time (Zuckerkandl and Pauling, 1965). An underlying assumption of the classic molecular clock is that selective constraints are negligible for most sites and over time. The development of molecular clock models as statistical processes relaxes this and other assumptions by allowing for rate variation among branches (and sometimes sites see Ho 2014) in phylogenetic trees (reviewed by Guindon (2020), Ho and Duchêne (2014)).

Molecular clock models necessarily involve two key quantities, the evolutionary timescale and the ‘evolutionary rate’, with the latter representing the combination of mutations and substitutions that accrue over time. However, evolutionary times and rates are unidentifiable (Dos Reis and Yang (2013), as reviewed by Bromham et al. 2018, Guindon 2020), and therefore cannot be jointly estimated using genetic sequence data alone. To make inferences from genetic sequences, all molecular clock methods require prior assumption about evolutionary times or rates, known as a ‘molecular clock calibration’. Three main calibrations exist: First, the age of the most recent common ancestor between two samples can be constrained to a given time point or interval (‘internal node calibration’). Second, a known estimate of the evolutionary rate can be incorporated (e.g. as a prior in Bayesian frameworks). Third, in cases where molecular sequences are sampled at different points in time (heterochronous sampling), the tips of the phylogeny can be anchored to these time points (‘tip calibration’; reviewed by Rieux and Balloux (2016)). The choice of calibration depends on the information available and its reliability (Duchêne et al., 2014, Warnock et al., 2012). For instance, it would be remiss to ignore evidence about when two lineages shared a common ancestor if the fossil record is compelling (Gavryushkina et al., 2017, Ronquist et al., 2016). Crucially, multiple sources of calibration information can be provided for the molecular clock.

1.1 Measurably evolving populations

Rapidly evolving organisms, notably many viruses and bacteria, have been found to accrue an appreciable number of mutations over the sampling timescale. Influenza viruses, for example, have evolutionary rates of around 6×10^{-3} subs/site/year (substitutions per genomic site per year) (Ghafari et al., 2022, Sanjuán, 2012). Assuming a genome size of 13,500 nucleotides, one would expect to observe one mutation every 4 to 5 days ($\frac{365 \text{ days/year}}{13,500 \text{ sites} \times 6 \times 10^{-3} \text{ subs/site/year}} \approx 4.5 \text{ days/subs}$). If genome samples are collected over the course of a few weeks, the sampling times themselves can be used to calibrate the molecular clock and tip calibration is therefore warranted. Data sets for which tip calibration is feasible are considered to have been sampled from a ‘measurably evolving population’ (Drummond et al., 2003b) and to have ‘temporal signal’.

Measurably evolving populations are typically characterised either by a sampling period that is long relative to the evolutionary rate, a sufficiently big data set (long molecular sequences or many samples), or both. Traditionally, such characteristics were mainly found in rapidly evolving organisms, typically RNA viruses. Nowadays, advances in sequencing technologies have dramatically expanded the range of organisms from which data sets can be considered to have been sampled from a measurably evolving population. Namely, ancient DNA techniques have effectively expanded the genome sampling window for many organisms (Duchene et al., 2020b, Spyrou et al., 2019a), and whole genome sequencing has meant that data sets of ‘slowly’ evolving microbes often carry sufficient information for calibrating the molecular clock (Biek et al., 2015) even when the sampling period covers only a few decades (Menardo et al., 2019).

1.2 The phylodynamic threshold

Genomic data sets collected during the early stages of an outbreak, for example, often pose two problems: low genetic diversity and a narrow sampling window. Both can lead to highly uncertain estimates of evolutionary rate and time of origin. The point at which an organism has accumulated sufficient genetic changes since its emergence to allow for informative tip calibration is referred to as the ‘phylodynamic threshold’ (Duchene

et al., 2020a). At a minimum, tip calibration requires that one mutation has occurred over the sampling period for the method to be informative. For a given organism, the minimum sampling period can be calculated as the inverse of the product of genome size and the evolutionary rate (i.e. $\frac{1}{\text{genome size (sites)} \times \text{evol. rate (subs/site/year)}} =$ years to observe one mutation). We refer to this amount of time as the expected phylodynamic threshold.

The terms phylodynamic threshold and measurably evolving population are different, albeit related, concepts. A population is measurably evolving if the samples available are sufficiently informative as to allow for tip calibration. In contrast, the phylodynamic threshold is the amount of time over which we would need to draw samples after their emergence for them to behave as from a measurably evolving population. For a recently evolving pathogen the phylodynamic threshold would simply correspond to the time until it can be considered a measurably evolving population, under the condition that the data have been collected constantly over time. In contrast, an organism that emerged further in the past may have accumulated considerable genetic diversity over time, effectively reaching its phylodynamic threshold. However, if samples are drawn from a very short time window they may fail to capture a representative amount of such genetic diversity.

1.3 Tests of temporal signal

Our ability to extract information from a tip calibration framework can be assessed through tests of temporal signal. The importance of performing such tests arises from the observation that a lack of temporal signal is associated with unreliable evolutionary rate estimates (Duchêne et al., 2015, Rieux and Balloux, 2016), although the presence and direction of a potential bias remain poorly understood. However, it is important to note that a lack of temporal signal does not necessarily preclude estimating evolutionary rates and timescales because alternative sources of calibration, such as prior estimates of evolutionary rates or constraints on internal node ages, can still be used to inform analyses.

In principle, frameworks developed to test for temporal signal do not differentiate between recently emerging organisms (fig. 1a) and those with narrow sampling windows (fig. 1d), both of which may lack temporal signal. As most of these tests involve fitting a phylogenetic model to the data, they implicitly assume that the model adequately captures the evolutionary process and thus their performance also highly depends on model fit. Recent research, for example, suggests that the choice of tree prior and molecular clock model significantly impacts the sensitivity and specificity of temporal signal tests (Tay et al., 2024). Thus, temporal signal is not solely a property of the data but also depends on the choice of model.

Various methods exist for assessing temporal signal. The root-to-tip regression (Buonagurio et al., 1986, Drummond et al., 2003a, Gojobori et al., 1990) fits a regression to the distance from the root to the tips in a phylogenetic tree against sampling time. High R^2 values of the regression suggest that phylogenetic distance can be sensibly modelled as a linear function of time and can thus be used as an indication of informative tip-calibration. Date-randomisation tests (Duchêne et al., 2015, Duchene et al., 2018, Ramsden et al., 2009, Trovão et al., 2015) compare evolutionary rate estimates using correct sampling times against those from permutations. Bayesian Evaluation of Temporal Signal (BETS; Duchene et al. (2020c)) evaluates whether a model with sampling times performs better than a model that assumes isochronous sampling using Bayes factors. Each method comes with a set of limitations and strengths, such that tests of temporal signal should rather be used in combination than being mutually exclusive (Duchene et al., 2020c, Rieux and Balloux,

116 2016).

117 1.4 Concepts of measurably evolving populations, the phylodynamic threshold, and 118 temporal signal in practice

119 In fig. 1, we present four simple example cases to illustrate the relationships among the concepts of mea-
120 surably evolving populations, the phylodynamic threshold, and temporal signal. The first example depicts
121 an organism that has emerged recently and therefore has not yet reached its phylodynamic threshold (with
122 a phylogenetic time tree shown in panel (a)). Due to its recent origin, there has not been enough time for
123 the accumulation of a sufficient number of substitutions (represented in the phylogram in panel (b)), such
124 that it is not possible to establish a statistical relationship between molecular evolution (i.e., substitutions)
125 and time (as shown in panel (c)). A real-world example of such a case comes from the early phase of the
126 SARS-CoV-2 outbreak: initial efforts to estimate the evolutionary rate and time of origin had substantial
127 uncertainty due to a narrow sampling window and low genetic diversity (Boni et al., 2020). In Duchene et al.
128 (2020a), Bayesian phylodynamic analyses were conducted on genome data as the outbreak unfolded. The
129 number of available genomes and the width of the sampling window increased over time and ranged from
130 22 genomes sampled over 31 days to 122 genomes sampled over 63 days. Although early estimates of the
131 evolutionary rate and time of origin were highly uncertain, they quickly converged to stable values as more
132 data became available (Ghafari et al., 2022).

133 The second example in fig. 1 illustrates a case in which an organism has evolved over a long period,
134 but the available sequence data have been collected within a very narrow timeframe, insufficient to treat the
135 dataset as a measurably evolving population (time tree in panel (d) and phylogram in panel (e)). This results
136 in no temporal signal, as demonstrated by the lack of correlation in the root-to-tip regression in panel (f).
137 The causative agent of tuberculosis, the bacterium *Mycobacterium tuberculosis*, was commonly considered to
138 evolve too slowly for calibrating the molecular clock using samples collected over a few years Duchene et al.
139 (2016). A range of studies have shown, however, that for *M. tuberculosis* a genome sampling window of a few
140 decades might be sufficient for reliable clock calibration (Eldholm et al., 2015, Kühnert et al., 2018, Menardo
141 et al., 2019, Merker et al., 2022).

142 The third example in fig. 1 describes a data set that may involve a wide sampling window of time and
143 for which samples have been drawn from a population that has attained its phylodynamic threshold, but
144 with substantial rate variation among lineages – i.e. overdispersed molecular clock –, resulting in a lack of
145 temporal signal (panels (g) – (i)). This pattern appears to be the case in *Yersinia pestis*, the bacterium
146 that causes the plague, for which some localised outbreaks display obvious temporal signal, but its long-term
147 evolution has pervasive evolutionary rate variation (Andrades Valtueña et al., 2022, Eaton et al., 2023).

148 In the final example in fig. 1, a hypothetical organism has attained its phylodynamic threshold, has been
149 sampled for sufficiently long time, and evolutionary rate variation among lineages is low. These conditions
150 together produce a clear relationship between molecular evolution and time, thus providing unequivocal
151 temporal signal (panels (j) – (l)). The long term evolution of *Vibrio cholerae*, the causative agent of cholera,
152 and H3N2 influenza virus are exemplar microbes whose molecular evolution has been fairly constant across
153 long periods of time (Devault et al., 2014, Rambaut et al., 2016).

154 In summary, the concepts of measurably evolving population, phylodynamic threshold and temporal signal

155 describe the information that can be drawn from a sampled population about its evolutionary timescale.
156 Because populations that are not measurably evolving have been observed to yield biased estimates, they
157 remain important to consider (Gharbi et al., 2024). In Bayesian inference, such biases can be the result
158 of complex interactions between prior distributions or model settings that do not align with the true data
159 generating model, as these drive the inference in the absence of informative data. Traditionally, potential
160 biases due to prior interactions (Tay et al., 2024) or model misspecification (Möller et al., 2018) have been
161 found through simulations studies, while data analyses often involve little validations of the results (Mendes
162 et al., 2025). Here, we illustrate through a range of examples the degree to which differing levels of temporal
163 signal in a data set can interact with prior settings and model assumptions, both on simulated and empirical
164 data.

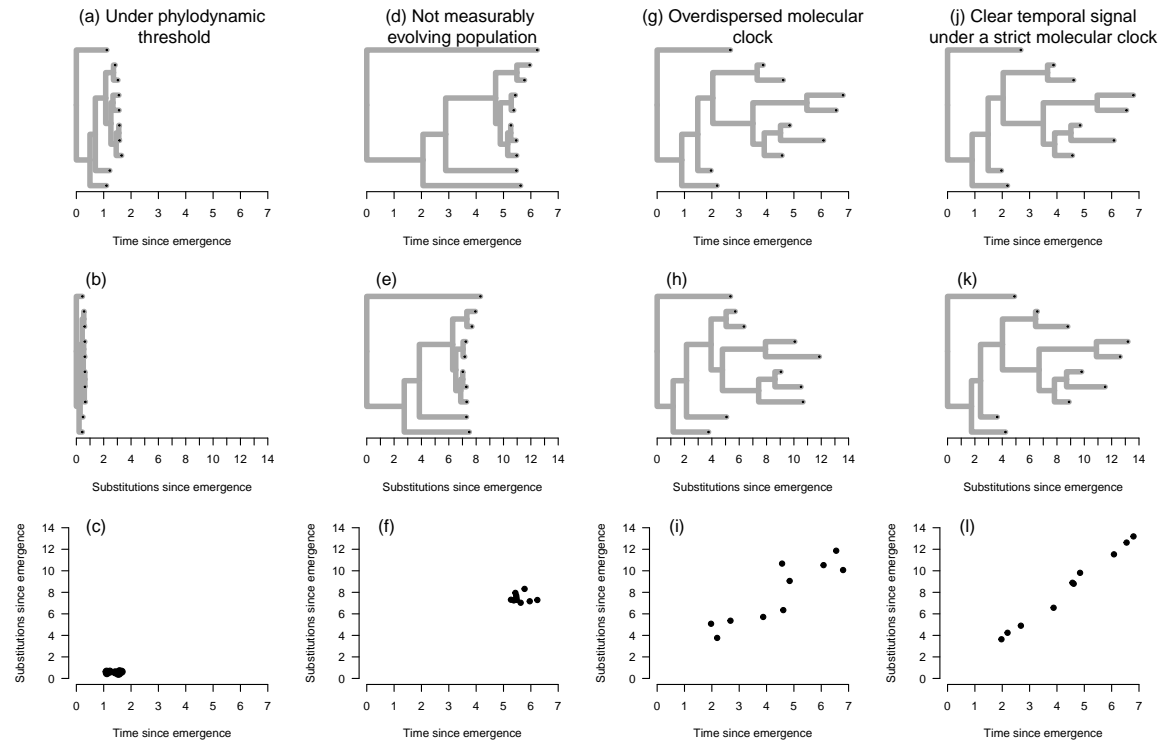


Figure 1: Examples of situations where temporal signal may or may not be detected. An organism that has not attained its phylodynamic threshold has a recent time of emergence (with a phylogenetic time tree shown in (a)) because it has not had sufficient time to accrue an appreciable number of substitutions (phylogenetic tree with branch lengths in subs/site, i.e. a ‘phylogram’, shown in (b)), such that it is not possible to establish a statistical relationship between molecular evolution (substitutions) and time (shown in (c)). Sequence data from an organism that has evolved for a substantial amount of time may have been sampled over a very narrow window of time that is not sufficient to treat it as a measurably evolving population (time tree in (d) and phylogram in (e)), which results in no temporal signal (root-to-tip regression in (f)). A data set may involve a wide sampling window of time and from a population that has attained its phylodynamic threshold, but an overdispersed molecular clock (substantial rate variation among lineages; panels (g) - (i)) may result in a lack of temporal signal. In (j) through (l) we show the situation where an organism has attained its phylodynamic threshold, it has been sampled for sufficiently long, and where evolutionary rate variation among lineages is negligible, as to produce a clear relationship between molecular evolution and time, and thus unequivocal temporal signal.

2 Results

We sought to pinpoint the impact of sampling strategies on molecular clock estimates. We focused our attention on two major problems for emerging microbes and studies involving ancient DNA. First, we conducted simulations varying the sampling window of a population that had attained its phylodynamic threshold. In the second simulation scenario, we subsampled a population over time to vary the number of ancient samples, leading to a temporal sampling bias. Finally, we illustrate these results in an empirical data set of Hepatitis B virus (HBV) that includes a large number of ancient samples (Kocher et al., 2021). This virus has been the subject of intense research due to its close association with human populations and complex evolutionary dynamics (Kahila Bar-Gal et al., 2012, Paraskevis et al., 2013, Ross et al., 2018).

2.1 Sampling windows relative to the phylodynamic threshold

We simulated sequence data that resembled the evolution of HBV, a double-stranded DNA (dsDNA) virus that has evolved in humans at least for around ten thousand years (Kocher et al., 2021). Our synthetic data had a genome length of 3,200 nucleotides and an evolutionary rate of 1.5×10^{-5} subs/site/year (Kocher et al., 2021, Mühlemann et al., 2018) with a moderate amount of rate variation among lineages (see Materials and methods). Under these conditions we expect to observe one mutation every 20 years ($\frac{1}{3,200 \text{ sites} \times 1.5 \times 10^{-5} \text{ subs/site/year}} \approx 20 \text{ years/subs}$). This number is important for the design of our simulation experiments: 20 years is the expected phylodynamic threshold, as introduced above, and typically serves as a good reference point from when on to expect temporal signal. We analysed the data under a Bayesian phylogenetic framework and considered whether the posterior contained the true value used to generate the data, known as coverage, and the width of the posterior, known as precision (a precise estimate has a narrow posterior distribution).

We conceived a simulation process under which the evolutionary timescale had an expectation of ten thousand years and with a sampling window of 0, 10, 20, 200, or 2,000 years. A sampling window spanning 0 years results in ultrametric trees with the sampling times providing no calibration information. In contrast, a sampling window of 10 years is half of the expected phylodynamic threshold and is likely to have weak temporal signal (see fig. 1(d)-(f)). Sampling windows of 20 years (the expected phylodynamic threshold) or wider are more likely to behave as measurably evolving populations with increasingly strong temporal signal (see fig. 1(j)-(l)). Our synthetic data sets were analysed under Bayesian phylogenetic framework, as implemented in the BEAST 2 platform (Bouckaert et al., 2019).

To investigate the impact of the prior we considered several configurations for the prior on the mean evolutionary rate. In our analyses the molecular clock model is a uncorrelated relaxed molecular clock model with an underlying lognormal distribution, with mean M . For this parameter we set nine possible prior Gamma distributions, for which the prior mean could be the value used to generate the data, or one order of magnitude higher or lower. We also included three degrees of uncertainty in this prior (see Fig 2 and table 5). In this respect, a prior with low uncertainty, and a mean that is much higher or lower than the true should result in more bias than one that has higher uncertainty or is centred on the true value.

Our simulations for which the prior on M was centred on the true value had very high coverage. At least 94 out of 100 simulation replicates across sampling windows included the true value of M within their

95% credible interval (CI) (table 1 and fig 2). Coverage, however, was associated with the degree to which the prior was biased and the sampling window width. For example, when the prior was highly precise (95% CI/mean=1.0) but biased downwards, even the simulations with a sampling window of 100× the phylodynamic threshold (i.e. 2,000 years before present) still had low coverage (only 1 out of 100 simulation replicates included the true value in the 95% credible interval, table 1).

Table 1: Coverage of the mean evolutionary rate, M . We consider coverage as whether a posterior distribution contains the value used to generate the data within its 95% credible interval (CI) and we show the number of simulation replicates out of 100 for which we found coverage. Data were simulated under trees with five possible sampling window depths relative to the phylodynamic threshold, where $Sw=\times D$ is for simulations with a sampling window of D times the phylodynamic threshold (20 years in our data; see fig 2). Each column is a sampling window depth, and the rows denote the configuration of the prior on the mean rate, M , for which we show the prior mean and the width of the 95% CI divided by the mean (larger values imply higher uncertainty). The true value of M is 1.5×10^{-5} subs/site/year. The horizontal lines separate prior configurations with different prior means.

Prior configuration (mean, 95%CI width)	Sw=0×	Sw=0.5×	Sw=1×	Sw=10×	Sw=100×
1.5×10^{-5} , 3.04	100	100	100	96	95
1.5×10^{-5} , 6.33	100	100	99	95	94
1.5×10^{-5} , 1.00	100	100	100	94	94
1.5×10^{-4} , 3.04	100	97	97	93	93
1.5×10^{-4} , 6.33	100	100	99	94	92
1.5×10^{-4} , 1.00	0	0	0	77	91
1.5×10^{-6} , 3.04	0	0	0	68	91
1.5×10^{-6} , 6.33	100	88	85	95	94
1.5×10^{-6} , 1.00	0	0	0	0	1

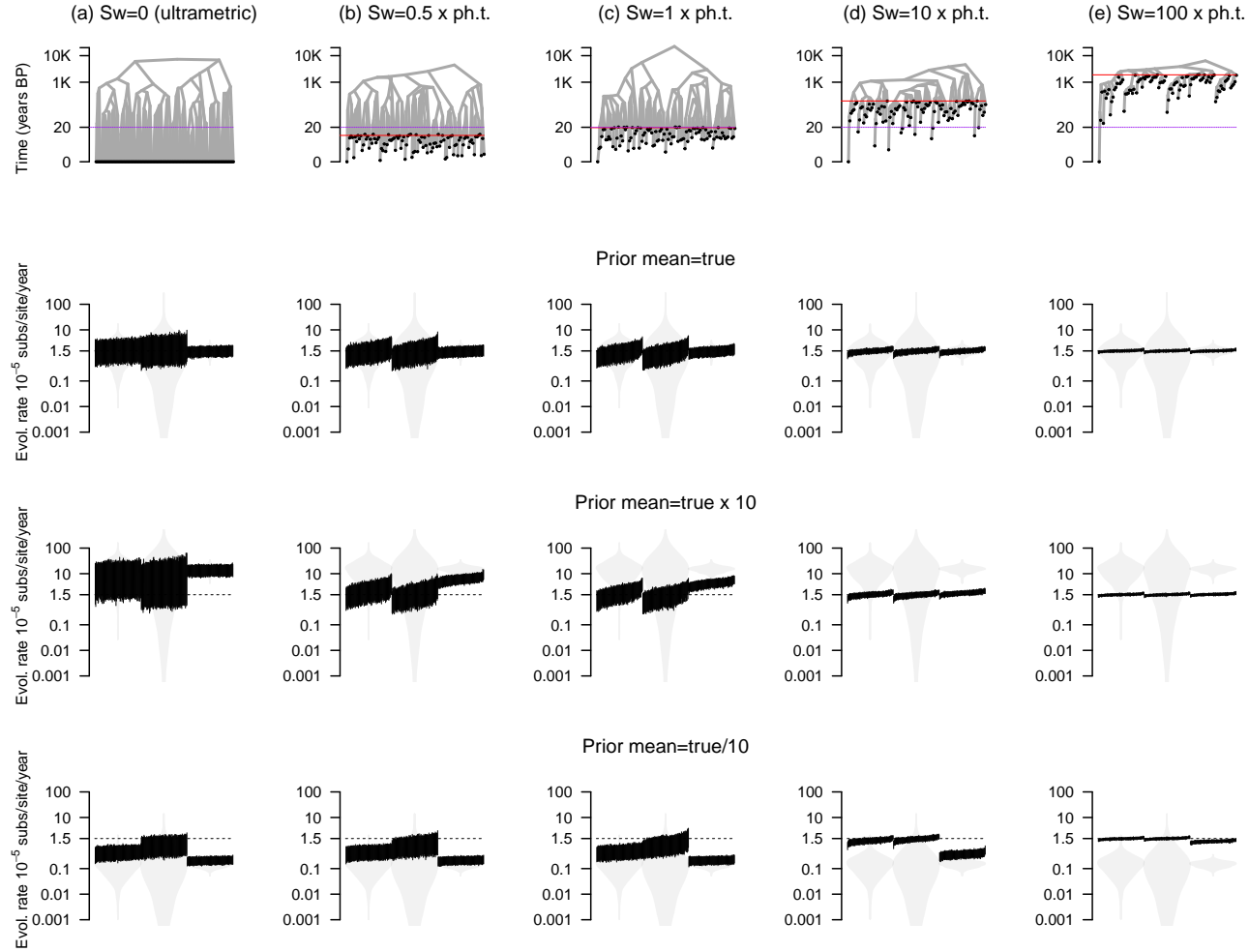


Figure 2: Estimates of evolutionary rates, M for simulations of varying sampling window widths. Each column corresponds to a simulation setting: (a) is for ultrametric trees where all samples are collected at the same point in time (sampling window, $sw=0$), (b) is for the situation where the sampling window is 10 years (half the expected phylodynamic threshold; sampling window, $sw=0.5 \times ph.t$), (c) is where the sampling window is exactly the expected phylodynamic threshold of 20 years ($sw=1 \times ph.t$). Scenarios (d) and (e) denote sampling windows that are 10 and 100 \times the expected phylodynamic threshold ($sw=10 \times ph.t$ and $sw=100 \times ph.t$, respectively). The rows denote example phylogenetic trees and prior configurations where the mean is set to the correct value (first row), an order of magnitude higher (second row), or an order of magnitude lower, last row). The prior is shown with the grey violins and each black bar is the 95% credible interval of the posterior. The dashed line in each case denotes the correct value.

208 The uncertainty in estimates of evolutionary rates was associated with the width of the sampling win-
 209 dow, but also with the uncertainty in the prior. In the situation where the sampling window was $1\times$ the
 210 phylodynamic threshold or less we consistently found that uncertainty in the prior was commensurate with
 211 uncertainty in the posterior. Estimates using a prior uncertainty of 1.00 (the width of the 95% credible
 212 interval is the same as the mean) are narrower than those using a prior of uncertainty of 3.04 or 6.33 (see
 213 table 2 and fig 2).

214 When the sampling window was $10\times$ the phylodynamic threshold or more we found a more complicated
 215 picture. When the prior had a downward bias (mean= 1.5×10^{-6} subs/site/year) and a narrow credible interval
 216 of 1.00, the posterior estimate was wider than when this prior was less uncertain (table 2 and fig 2). This
 217 finding can be explained because a very wide sampling window provides a large amount of information for
 218 inferring the evolutionary rate, which can yield high uncertainty if the prior stands in conflict. Importantly,
 219 the widest sampling window in our experiments, of $100\times$ the phylodynamic threshold produced consistently
 220 high precision, although it is important to note that when the prior is highly biased downward (mean of
 221 1.5×10^{-6} subs/site/year and uncertainty of 1.00) coverage is very low, with only one simulation replicate
 222 containing the true value used to generate the data within the 95% credible interval.

Table 2: Average uncertainty of the mean evolutionary rate, M across 100 simulation replicates in each case. We quantify uncertainty as the posterior width of the 95% credible interval (CI) divided by the mean. The rows and columns match those from table 1, with the posterior for each simulation replicate shown in fig 2.

Prior configuration (mean, 95%CI width)	Sw=0 \times	Sw=0.5 \times	Sw=1 \times	Sw=10 \times	Sw=100 \times
1.5×10^{-5} , 3.04	2.24	1.87	1.55	0.48	0.21
1.5×10^{-5} , 6.33	3.06	2.01	1.72	0.48	0.22
1.5×10^{-5} , 1.00	0.96	0.94	0.91	0.45	0.21
1.5×10^{-4} , 3.04	1.36	1.83	1.50	0.48	0.22
1.5×10^{-4} , 6.33	5.39	2.11	1.70	0.49	0.21
1.5×10^{-4} , 1.00	1.04	0.87	0.80	0.43	0.21
1.5×10^{-6} , 3.04	1.36	1.37	1.38	0.52	0.22
1.5×10^{-6} , 6.33	1.99	1.90	1.72	0.48	0.22
1.5×10^{-6} , 1.00	0.75	0.75	0.76	0.75	0.28

223 To understand the directionality of posterior evolutionary rate estimates relative to the value used to
 224 generate the data, we quantified the amount of bias, as difference between the true evolutionary rate and the
 225 posterior mean divided by the true value (i.e. $\frac{\text{true value} - \text{posterior mean}}{\text{true value}} = \frac{1.5\times 10^{-5} - \text{posterior mean}}{1.5\times 10^{-5}}$). As expected,
 226 when the prior was centred in the true value, we observed no or minimum bias, with values ranging from
 227 an average of 0.00 and 8.31 (a posterior estimate that was on average 8 times higher than the truth). The
 228 most marked average bias was found for our prior with mean 1.5×10^{-4} subs/site/year and an uncertainty
 229 of 1.00, which could be as high as 8.31 (table 3). Increasingly wide sampling windows had lower amounts of
 230 bias. A sampling window of $100\times$ the phylodynamic threshold had a maximum average bias of -0.22 for the
 231 prior with downward bias and low uncertainty, with other prior configurations resulting in a bias of at most
 232 -0.22.

233 Overall, these simulations demonstrate that increasingly wide sampling windows result in evolutionary

Table 3: Average bias of the mean evolutionary rate, M across 100 simulation replicates in each case. Bias is calculated as the true evolutionary rate, 1.5×10^{-5} subs/site/year - posterior mean, divided by the true value.

Prior configuration (mean, 95%CI width)	Sw=0×	Sw=0.5×	Sw=1×	Sw=10×	Sw=100×
1.5×10^{-5} , 3.04	0.16	0.06	0.01	0.00	0.01
1.5×10^{-5} , 6.33	0.28	0.03	-0.03	0.00	0.01
1.5×10^{-5} , 1.00	0.00	-0.01	-0.01	0.00	0.01
1.5×10^{-4} , 3.04	4.93	0.67	0.30	0.02	0.01
1.5×10^{-4} , 6.33	3.37	0.27	0.08	0.00	0.01
1.5×10^{-4} , 1.00	8.31	3.22	1.84	0.17	0.04
1.5×10^{-6} , 3.04	-0.71	-0.70	-0.69	-0.16	-0.02
1.5×10^{-6} , 6.33	-0.45	-0.43	-0.37	-0.04	0.00
1.5×10^{-6} , 1.00	-0.86	-0.86	-0.86	-0.75	-0.22

rate estimates that are more accurate, precise, and less biased, than those from data sets with narrow sampling windows. Although increasingly wide sampling windows are more robust to prior misspecification, we emphasise the importance of choosing the prior for this parameter carefully.

For all trees that were non ultrametric (i.e. those with a sampling window width of at least 0.5× the phylodynamic threshold) a downward bias in the evolutionary rate prior appears to be more detrimental than one with an upward bias. In our simulations with sampling windows of width 100× the phylodynamic threshold a highly biased prior (mean= 1.5×10^{-6} subs/site/year and uncertainty of 1.00) still had very low coverage. In contrast, a prior with a similar upward bias (mean= 1.5×10^{-4} and uncertainty of 1.00) produced much higher coverage for sampling windows from 10× the phylodynamic threshold (table 1).

These results indicate that priors with high uncertainty should be advised for practical studies. In our simulations posterior estimates using a prior uncertainty of 6.33 seemed to produce a good trade-off between uncertainty and accuracy. Such a prior means that the 95% credible interval spans just over six orders of magnitude. For a sampling window of 1× the phylodynamic threshold, the posterior distribution had an average uncertainty of around 1.7, which may be sufficient for biological interpretation of estimates of evolutionary rates and timescales.

2.2 Hierarchical priors and the phylodynamic threshold

Here we need to explain the hyper prior, that it is very upwardly biased, but htat it was easy to overcome, even with a sampling window of 0.5× the phylodynamic threshold, with a biasof less than 0.09 (less than 10%)

Table 4: Coverage, average uncertainty and average bias of the mean evolutionary rate, M using a hierarchical prior on M . Statistics are reported over 100 simulation replicates under each of the five sampling window widths (Sw) relative to the phylodynamic threshold (of 20 years).

Statistic	Sw=0×	Sw=0.5×	Sw=1×	Sw=10×	Sw=100×
Coverage	100	100	99	96	94
Uncertainty	57.69	2.44	1.86	0.49	0.22
Bias	1.01	0.09	0.01	0.00	0.01

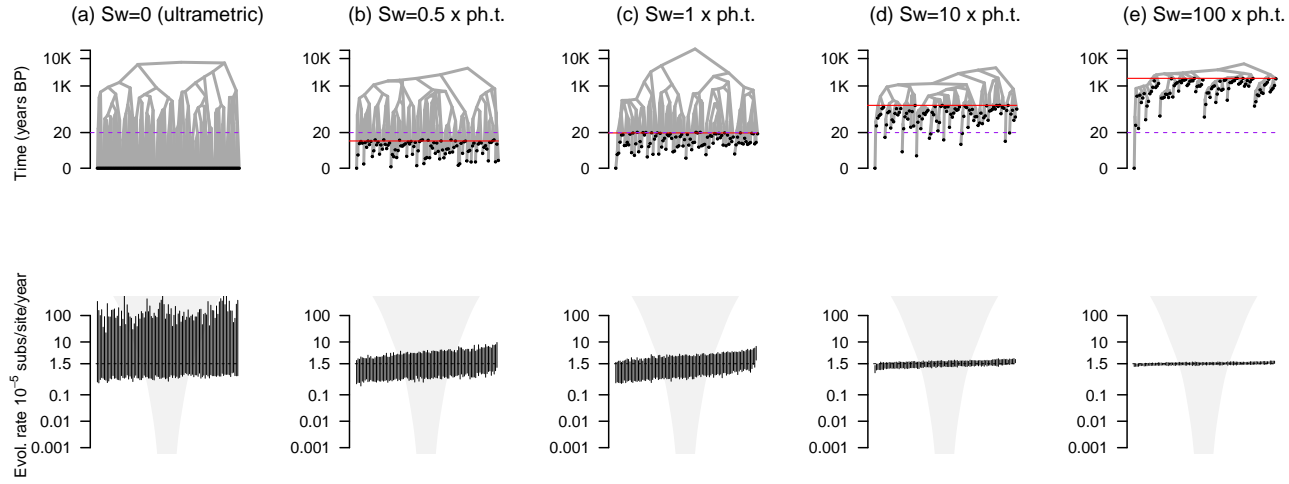


Figure 3: Estimates of evolutionary rates, M for simulations of varying sampling window widths. Each column corresponds to a simulation setting: (a) is for ultrametric trees where all samples are collected at the same point in time (sampling window, $sw=0$), (b) is for the situation where the sampling window is 10 years (half the expected phylodynamic threshold; sampling window, $sw=0.5 \times \text{ph.t.}$), (c) is where the sampling window is exactly the expected phylodynamic threshold of 20 years ($sw=1 \times \text{ph.t.}$). Scenarios (d) and (e) denote sampling windows that are 10 and $100 \times$ the expected phylodynamic threshold ($sw=10 \times \text{ph.t.}$ and $sw=100 \times \text{ph.t.}$, respectively). The first row denotes example phylogenetic trees, while the second row is the corresponding posterior estimates. The prior is shown with the grey violins (a hierarchical prior configuration) and each black bar is the 95% credible interval of the posterior. The dashed line in each case denotes the correct value.

Takeaway points for discussion:

- That a downward rate prior is likely more problematic because it has no lower bound, whereas the sampling window effectively brackets the rate. The maximum rate is the total genetic divergence divided by the sampling window width. The exception is when the sampling window is very narrow (even 0), where we can observe higher bias when the evol. rate prior is upwardly biased, but this is because the bracketing does not occur.
- That the phylodynamic threshold does not guarantee correct estimates, because we also depend on the number of samples. It just means that we are guaranteed to have some genetic diversity in the sampling window, but note that its coverage, uncertainty and bias, are often comparable to what we obtain with $0.5 \times \text{ph.t.}$
- Also note that in Bayesian analyses, the sampling times are never the only source of calibration. The prior on the evol. rate is an obvious one, but the tree prior is also very informative. In our case, we assume a constant population size coalescent, where the θ parameter governs the root height (explain here), so it is a prior on time, and therefore on the rate given the sequence divergence ($M \propto \frac{1}{\theta}$, $M = \frac{D}{T_i}$). The CTMC-rate reference, commonly used on the rate in *beastx* is thought to be uninformative, but it effectively favours low values.
- That using priors that are vague seems like a good idea. Here the widest priors, with uncertainty of 6.33 were not overly informative to overcome the signal from the data.

produced posterior estimates that included the correct value used to generate the data (i.e. 100% coverage; fig. 2). Increasingly wide sampling windows improved the precision of the estimates, while still including the correct value. This result is unsurprising, given our configuration of the prior. Here, the tree prior is a constant-size coalescent for which the prior on the population size (known as θ) is an exponential distribution with mean of 5,000, which matches the value used to generate the data. Similarly, the evolutionary rate had a prior in the form of a Gamma distribution with *shape* = 1.5 and *rate* = 10^5 , whose mean is *shape*/*rate* = 1.5×10^{-5} and thus also matches the ‘true’ value. Although these priors are centred on the correct values, they are vague, and it is important to note that in all cases, the posterior distributions of the evolutionary rate and tree height were narrower than their priors, meaning that even in the absence of sampling times the sequence data themselves provide some information about these two parameters.

We reanalysed these data with deliberately ‘misleading’ priors on the population size and the evolutionary rate. The prior on the population size was an exponential distribution with mean of 50,000, whereas the prior on the evolutionary rate was *Gamma*(*shape* = 1.5, *rate* = 10^6) (mean = 1.5×10^{-6}). Under this configuration the mean prior mass corresponds to trees that are one order of magnitude older than the truth and evolutionary rates that are an order of magnitude slower. The objective of this experiment is to determine whether the sampling window is sufficiently informative to overcome such misleading prior information.

The posterior distribution was largely contained within the prior, resulting in low coverage for the evolutionary rate for sampling windows of 0, 10, and 20 years (0% coverage; fig. 3). A sampling window of 200 years was necessary to obtain 92% coverage, while a sampling window of 2,000 years had 100% coverage

and even higher precision (fig. 3). These results demonstrate that a misleading prior that places a very low probability on the true value, requires a sampling window that is potentially much wider than the expected phylodynamic threshold.

Contrary to the expectation that low temporal signal necessarily results in an underestimation of the evolutionary rate and an overestimation of the tree height (Duchêne et al., 2015), we find that a lack of temporal signal due to narrow sampling windows may simply lend more influence to the prior. To confirm this point we conducted an additional set of simulations where the mean evolutionary rate prior was $\text{Gamma}(\text{shape} = 1.5, \text{rate} = 10^4)$ (mean = 1.5×10^{-4} , and 95% range from 1.1×10^{-5} to 4.7×10^{-4}), and thus should lead to an overestimation of this parameter. As expected, we also found that increasing the width of the sampling window resulted in a prior that was less influential on the posterior and with the latter converging to the value used to generate the data (see Supplementary material). Compared to the previous setting with incorrectly lower prior, a narrower sampling window already resulted in good coverage. This is unsurprising, as higher rate values are less likely when only few mutations are observed in a relatively long sampling window, while lower rate values cannot be excluded.

2.3 Temporal sampling bias

We investigated the impact of temporal sampling bias on the precision and accuracy of molecular clock estimates. For this purpose we simulated data with the same genomic characteristics as HBV and where genome sampling was conducted over five periods of time uniformly distributed between the present and the root of the trees (fig. 4(a)). The fully sampled trees contained 500 genome samples, with 100 for each of the five sampling times. Such stratified sampling is expected in ancient DNA studies, for example when a set of samples are drawn from archaeological sites (e.g. Spyrou et al. (2019b)). We sampled the complete data sets by randomly selecting 20 samples from each strata, which we refer to as ‘time-uniform’ sampling, and by sampling with a probability that is inversely proportional to the age of the strata, referred to as ‘time-biased’. The time-uniform and time-biased sampling strategies both contain 100 samples ($1/5^{th}$ of the complete data), but the time-biased only includes a small number of ancient samples.

The coverage of the evolutionary rate estimate was comparable across simulation treatments, at 88% for the complete data set, 83% for the time-uniform, and 89% for the time-biased (fig. 4(b)). The somewhat higher coverage for the estimates from the time-biased analyses is likely because this sampling treatment has the lowest precision in the posterior, rather than an improvement in both accuracy and precision.

We also calculated a measure of bias for both sampling strategies by counting the number of simulations for which the posterior mean with either sampling strategy was higher than that with the complete data. We found that 50% of the estimates under time-uniform sampling had higher means than the complete data, while the same was true for 45% of those with time-biased sampling (fig. 5(a)). Although these numbers do not indicate substantial bias, such as a systematic over- or underestimation, we do note that time-biased sampling tends to produce lower mean evolutionary rate estimates than those obtained from the complete data or time-uniform sampling.

The most striking result of the temporal sampling strategies was in the precision of the posterior. Both sampling treatments resulted in posterior distributions that were wider than with the complete data, which is to be expected because they are effectively smaller data sets with less information. However, the time-biased

330 sampling data sets almost invariably have posterior distributions that were less precise than those from the
331 time-uniform sampling (Fig. 5(b)), implying that the distribution of samples, and not just the number, is
332 important for estimation precision.

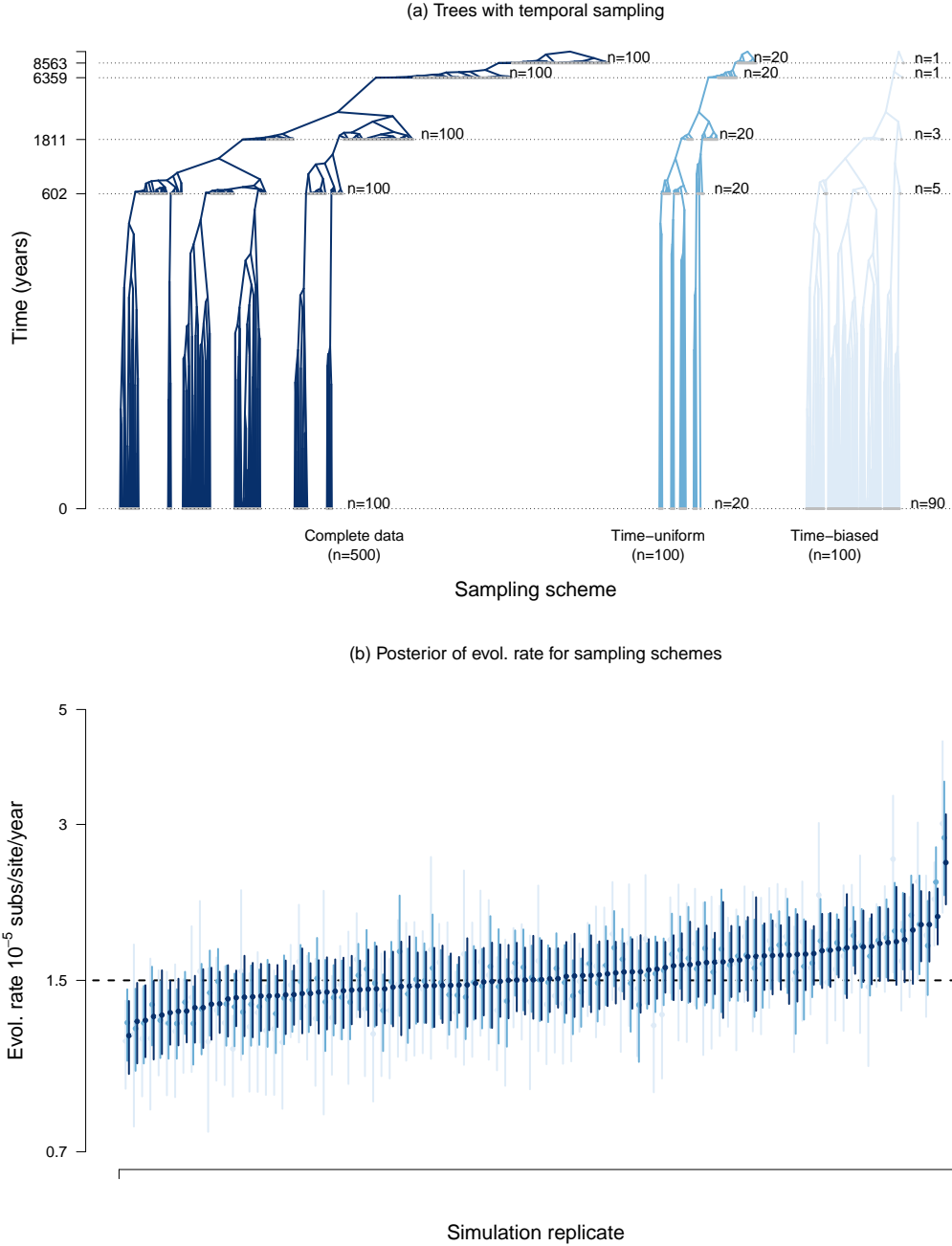


Figure 4: Analyses under sampling treatments over time. In (a) we show an example of the trees for a simulation replicate, with branch lengths and time in \log_{10} scale. The complete data set consists of 500 genome samples, collected in five points in time, with an equal number of samples per time point ($n=100$). The first sampling strategy is unbiased, where 20 samples are drawn from each time point, and is known here as ‘time-uniform’. The ‘time-biased’ regime is where sampling intensity decreases over time. Note that the total number of samples in the time-uniform and time-biased treatments is identical. In (b) we show the posterior estimates of the evolutionary rate for each treatment. Each simulation replicate is represented by three error bars: dark blue for the complete data, and lighter shades of blue for the estimates from the time-uniform and time-biased sampling treatments. The width of the error bars denotes the 95% quantile range and the dots are the mean value. The dashed line shows the true value used to generate the data.

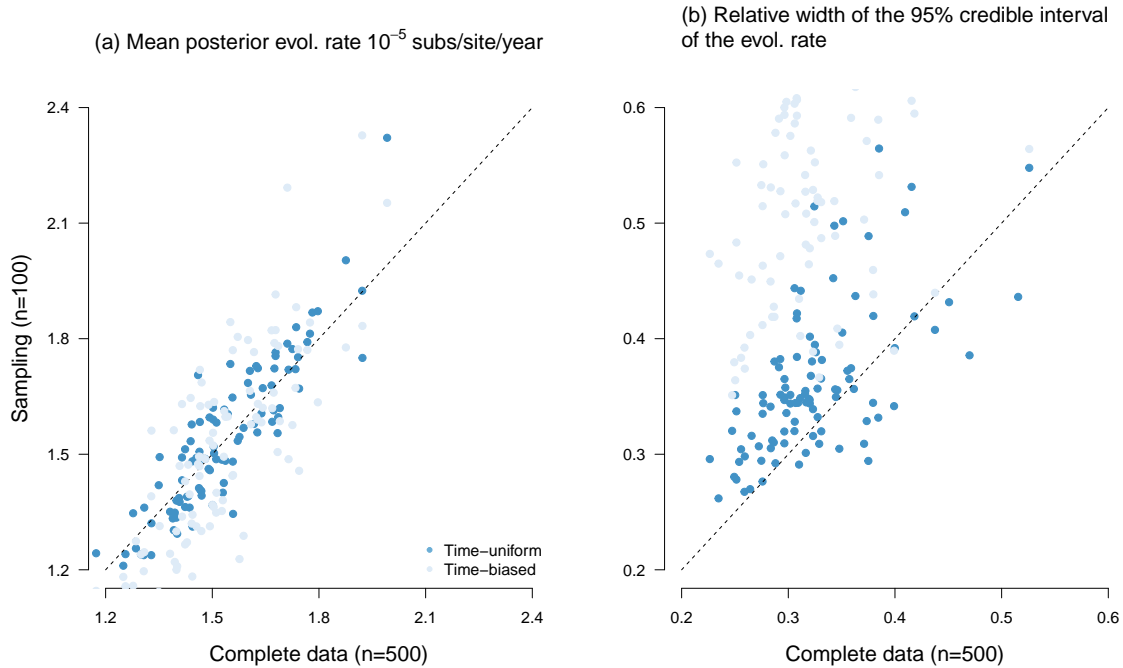


Figure 5: Comparison of posterior evolutionary rate estimates between complete data (x-axis) and two sampling treatments (y-axis): time-uniform (dark blue) and time-biased (light blue). Each dot is a simulation replicate. In (a) we show the mean posterior evolutionary rate estimate. Points that fall along the $y=x$ line (dashed line) represent identical mean posterior for the sampling treatment and the complete data, while those above or below represent higher or lower estimates, respectively, relative to the complete data. In (b) we show the width of the credible interval (a measure of precision or uncertainty), calculated as the upper minus the lower 95% credible interval range divided by the mean value. Values that fall along the $y=x$ line denote those for which the complete data and either sampling strategies are equally precise, while those above and below the $y=x$ line are more or less precise, respectively.

2.4 Empirical analyses of Hepatitis B virus (HBV) ancient and modern genomes

To explore the impact of the width of the sampling window and the temporal sampling bias on the estimates of evolutionary rates and times, we performed analyses of a HBV data set that includes modern and ancient genomes, from Kocher et al. (2021). The complete data set consisted of 232 genomes of length 3,344 nucleotides and with a sampling window of 10,535 years. HBV is an ancient pathogen that has likely codiverged with human populations for thousands of years (Locarnini et al., 2021, Mühlemann et al., 2018, Paraskevis et al., 2013, Zehender et al., 2014), and thus its phylodynamic threshold has been reached while it has not been empirically established if it can be considered to be a measurably evolving population, as is the case for recent outbreaks, like SARS-CoV-2 (Duchene et al., 2020a).

For our first set of analyses we varied the width of the sampling window. We drew 100 genomes with different sampling window widths: 0 (only modern samples), up to 500, 1,000, or 5,000 years before present. Increasing the sampling window resulted in estimates of the evolutionary rate that were more precise and closer to the estimate from the complete data set (fig. 6). Here we find that the evolutionary rate is estimated to be higher for shorter sampling windows, with correspondingly older estimates for the tree height (see

Supplementary material). This pattern can be due to one or a combination of other factors influencing the inference, for example the vagaries of evolutionary rate variation in this virus, particularly time-dependency (Vrancken et al., 2017). Similarly, population structure that is unaccounted for has been shown to produce an overestimation of the evolutionary rate, because under the tree prior samples that are genetically linked are expected to have been sampled at the same point in time (Möller et al., 2018).

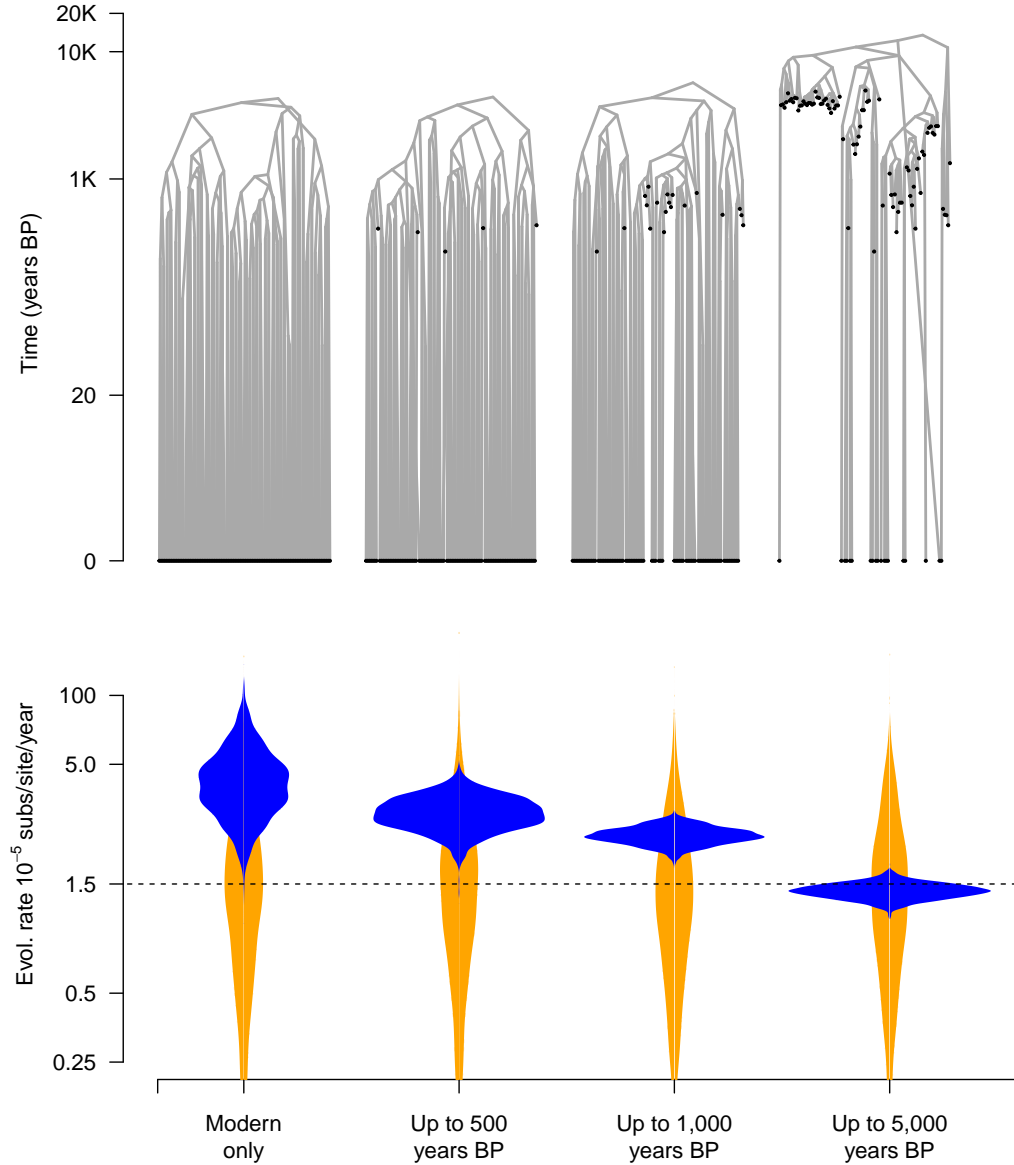


Figure 6: Results from empirical analyses of Hepatitis B virus (HBV) ancient DNA data. The phylogenetic trees correspond to highest clade credibility trees from three analyses where the data were subsampled to increase the width of the sampling window progressively. First, we only consider modern samples, then those up to 500, 1,000, and 5,000 years before present. In all cases the data sets consist of 100 genome samples. The violin plots show the posterior distribution of the evolutionary rate in blue and its corresponding prior in orange. The dashed line shows the mean evolutionary rate estimate from the complete data set.

For our second set of empirical analyses, we subsampled the data to produce a range of temporally biased sampling scenarios. We varied the proportion of modern samples, from 95% to 10%. Here the sampling window is constant because we always retain the most ancient samples. Each data set consisted of 100 genomes, such that they only differ in the distribution of modern and ancient samples (fig. 7).

The impact of temporal sampling bias was less clear than in our simulations above (figs 4 and 5). The data set with 95% modern genomes had the highest uncertainty in the evolutionary rate estimate, but uncertainty did not decrease monotonically with the proportion of modern genomes. Moreover, the posterior estimate of the evolutionary rate from the data set with 50% modern genomes deviated the most from the that obtained using the full data set. Critically, when we analysed these data using the ‘misleading’ prior setting we found that increasing the number of ancient samples resulted in a less influential prior (see Supplementary material). These results demonstrate that the exact impact of temporal sampling bias may be difficult to predict in practice, but that generally increasing the number of ancient samples results in data sets that are more informative, in that the difference between the prior and posterior is more pronounced than when only a few ancient samples are available.

3 Discussion

The concepts of the phylodynamic threshold, measurably evolving populations, and temporal signal are helpful for our understanding of rapidly evolving organisms or data sets with ancient DNA. Our analyses help us disentangle the definition of these concepts and their practical implications.

The phylodynamic threshold and measurable evolution are not discrete bounds. Increasing the sampling window generally improves precision and accuracy, but there is no clear cut-off for when the estimates become accurate and objectively ‘precise’.

Notably, the prior on the phylogenetic tree and the evolutionary rate are particularly influential for estimating evolutionary rates and timescales (for a detailed investigation see Tay et al. (2024)). In our simulations with a reasonable prior the posterior always included the correct value, but when we set a misleading prior a sampling window of 10 times the time expected to accrue one mutation (i.e. the so called expected phylodynamic threshold) was necessary to obtain a posterior that included the correct value. As a consequence, the phylodynamic threshold and measurable evolution depend on the extent to which the data inform the posterior, which is ultimately a measure of the relative contribution of the prior and the data (via the likelihood function).

Ideally, the prior and posterior of key parameters, particularly the height of the root node should overlap, while the posterior should be narrower than the prior (i.e. more informative), meaning that the data and the prior are not in conflict (for recent work on quantifying prior-data conflict see: Nott et al. 2020). In this vein, assessing the adequacy of the model and prior via predictive checks can be illuminating (McElreath, 2018), particularly in situations where the joint prior is poorly understood (see Baele and Lemey 2014, Wang and Yang 2014 for discussions about the prior in Bayesian phylogenetics). Recent years have seen the development of a range of methods for assessing phylogenetic model adequacy (Brown and Thomson, 2018, Duchêne et al., 2018, Duchene et al., 2019, McElreath, 2018), for instance one can simulate phylogenetic trees under the posterior estimates to verify whether the height of the root node and the topology could have been

390 generated by the model in question.

391 Measurably evolving populations are those for which the phylodynamic threshold has been attained and
392 the sampling window is *sufficiently* wide. The criteria for determining the phylodynamic threshold and
393 whether a population is measurably evolving are the same, and are typically assessed via temporal signal.
394 Statistical tests for this purpose quantify the strength of the association between sampling times and genetic
395 distance (Duchêne et al., 2015, Featherstone et al., 2024, Murray et al., 2016, Rambaut et al., 2016, Rieux
396 and Balloux, 2016). That is, the degree to which the sampling times on their own constitute an informative
397 molecular clock calibration. We contend that assessing prior sensitivity is more important than the outcome
398 of tests of temporal signal for obtaining reliable molecular clock estimates. In fact, a poor choice of prior
399 can mislead tests of temporal signal (Tay et al., 2024). If data are drawn from a sampling window that
400 spans the expected phylodynamic threshold, the presence of temporal signal is likely supported by most
401 tests, suggesting accurate estimates. Yet, if the prior used for the estimation is misleading and informative,
402 it might actually obscure the ‘correct’ signal from the data. In contrast, if the data are drawn from a narrow
403 sampling window but the prior is reasonable then the estimates may be still be reliable, despite a lack of
404 temporal signal. It also has to be noted that an increasing sampling effort does not necessarily lead to
405 increasingly correct inferences, because misspecification not only in prior distributions of hyper-parameters,
406 but also in the underlying model, can introduce biase (Ferretti et al., 2024, Möller et al., 2018).

407 An obvious concern about molecular clock calibrations using sequence sampling times is sampling bias.
408 We find that temporal sampling bias, where data are overwhelmingly collected at a particular period of
409 time does not have a substantial impact in estimation accuracy on a simple coalescent model, but that
410 increasing the number of ancient sequences can improve precision. An other form of sampling bias is when
411 genetic diversity is not uniformly sampled or the underlying population is structured. Previous work has
412 demonstrated that in such cases, the evolutionary rate and the age of the root node tend to be overestimated
413 (Möller et al., 2018), a problem that diminishes when sequence data are are increasingly informative or by
414 using a tree prior that explicitly models population structure (e.g. Kühnert et al. 2016, Müller et al. 2017).

415 Our study has a few limitations that have been partly addressed elsewhere. The number of sequences
416 is fixed in most of our experiments, but it is well known that increasing the number of sequences generally
417 means that data are more informative and thus the estimates are more precise (see Möller et al. (2018) and
418 our fig. 4), and therefore it is likely that the width of the sampling window needed to obtain reliable estimates
419 also depends on the number of sequences. Moreover, our simulation experiments involved a low degree of
420 evolutionary rate variation among lineages. In this respect, it is expected that the width of the sampling
421 window scales with the amount of dispersion in the molecular clock. In addition, our simulations are based
422 on a simple population dynamic model, the constant coalescent. The impact of the sampling scheme and
423 width is likely to be more complex for models with more parameters. We also assume the correct model of
424 evolution and population growth in all simulation-based inferences. With the empirical analysis, we, however,
425 highlight how conclusions drawn from these do not directly extend to real-world data. Rather, the isolated
426 effects found therein describe only one of many elements impacting the inference from real data. Further
427 scrutiny of these factors is warranted, but the main implication is that the necessary sampling time window
428 is combination of the data set, the organism, and the model at hand.

429 Overall, our study elucidates some of the fundamental intricacies of molecular clock calibration strategies.

430 We urge researchers to carefully question their model and its underlying prior assumptions, not only via tests
 431 of temporal signal, but also through careful choice of the prior, an understanding of the information content
 432 in the data, and the implications of model misspecification.

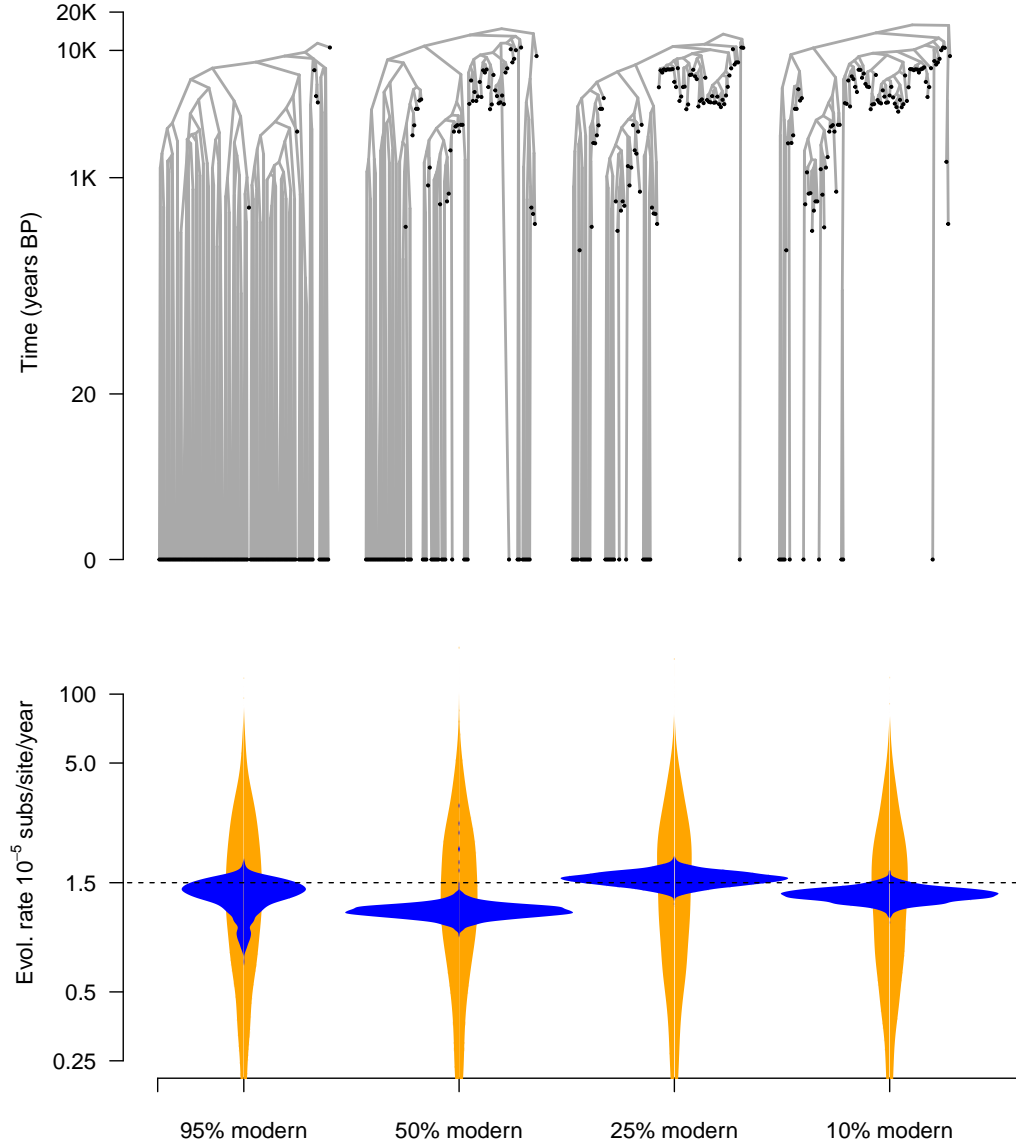


Figure 7: Results from empirical analyses of Hepatitis B virus (HBV) ancient DNA data. The phylogenetic trees correspond to highest clade credibility trees from three analyses where the data were subsampled to include an increasing number of ancient samples. First, we consider a data set for which the samples are 95% modern and the remaining 5% being the most ancient. Then, we reduce the number of modern samples to 50%, 25% and 10%, and the rest being ancient. Note that the sampling window is constant because we always retain the most ancient samples. In all cases the data sets consist of 100 genome samples. The violin plots show the posterior distribution of the evolutionary rate in blue and its corresponding prior in orange. The dashed line shows the mean evolutionary rate estimate from the complete data set.

4 Materials and methods

4.1 Simulations

4.1.1 Data generation

We simulated phylogenetic trees and the evolution of nucleotide sequences to assess the impact of varying the sampling window and temporal sampling bias. We parameterised our simulations to resemble an HBV population evolving over 10,000 years before present, as described by Kocher et al. (2021).

We generated phylogenetic trees under a coalescent process in which the population size has been constant over time using the package ReMaster (Vaughan, 2024), part of the BEAST 2 (version 2.7) software. We set the population size to 5,000, which results in trees with an average age of about 10,000 time units (years) (see Arbisser et al. 2018). The number of samples (i.e. tips) drawn from the trees and their ages was defined in ReMaster, according to the simulation scenario described below.

The simulation of sequence data requires trees with branch lengths in units of subs/site, instead of time. For this purpose, we multiplied the branch lengths of the simulated time tree with the rate of evolution, a lognormally distributed random variable with parameters $\mu = \log(1.5 \times 10^{-5}) - \frac{0.25^2}{2}$ and $\sigma = 0.25$. This procedure equates to simulating an uncorrelated relaxed molecular clock model with an underlying lognormal distribution (Drummond et al., 2006) with mean, M of 1.5×10^{-5} subs/site/year and a standard deviation of 0.25 subs/site/year (following that $\mu = \log(M) - \frac{\sigma^2}{2}$). Because we multiply the branch lengths in units of years by a variable in subs/site/year, the resulting trees have branch lengths in units of subs/site, formally known as phylograms, in contrast to chronograms where the branch lengths correspond to time. We obtained sequence alignments using the R package phangorn (v2.8.1) (Schliep, 2011), according to a HKY+ Γ_4 substitution model, with parameters $\kappa = 2$, $\alpha = 4$, and equal base frequencies. The alignments consisted of 3,200 nucleotides to match the average genome size of HBV.

We considered the expected phylodynamic threshold of our data to be about 20 years. For our simulations where we varied the sampling window, we set the ages of 100 tips to be sampled at present (all have an age of 0), or to be drawn from a uniform distribution between 0 and 10 (1/2 of the expected phylodynamic threshold), 0 and 20, 0 and 200, or 0 and 2,000.

To investigate the impact of temporal sampling bias we initially simulated trees with 500 tips with sampling times distributed in 5 time points, with 100 tips per time point. The distribution of sampling times followed an exponential distribution with mean of 4,000, such that sampling times were concentrated towards the present. We simulated these complete trees in ReMaster. We followed the procedure above to simulate a molecular clock model and sequence alignments.

We conducted two sampling schemes for the trees with 500 tips: the ‘time-uniform’ scheme consisted of drawing 20 samples from each time point, whereas the ‘time-biased’ scheme included of mostly modern samples (90 from the present, and 5, 3, 1, and 1 from the remaining time points). For each simulation scenario we generated 100 replicates.

4.1.2 Analysis of simulated data

We analysed all data sets in BEAST 2 under a model that matched that used to generate the data: the HKY+ Γ_4 substitution model, a uncorrelated relaxed molecular clock model with an underlying lognormal distribution, and a constant size coalescent tree prior. We used the default prior configuration in the program, except for the mean evolutionary rate (M) and the population size of the coalescent (θ). For the population size we assumed $\theta \sim \text{Exponential}(\text{mean} = 5,000)$, which is centred in the value used to generate the data. The expected height of the root node is roughly 10,000 years (expected time to coalescent = $2 \times \theta$ for an ultrametric tree, see Nordborg (2019)).

For the mean evolutionary rate we considered a range of priors with different degrees of information content (uncertainty) and for which the mean was either the value used to generate the data, or one order of magnitude higher or lower, as shown in table 5. We also included three degrees of uncertainty in such prior distributions, where the 95% quantile width was equal to the mean, or around three or six times as wide. In all cases we set the sampling times for calibration.

In the case of ultrametric trees, sampling times are set to the present, such that all calibration information is provided by the tree prior. Concretely, let the tree length (sum of all branch lengths b_i in units of time) be T_l , the tree height T_h , branch rates r_i (the vector of branch rates $\vec{R} \in \{r_1, \dots, r_{\text{num. branches}}\}$ (see Douglas et al. 2021), and D is the total genetic distance in the tree (sum of branch lengths in subs/site).

In Bayesian frameworks branch lengths in genetic distance are the product of rates and times (Douglas et al., 2021, Drummond et al., 2006), such that $D = \sum r_i b_i$. The average evolutionary rate (an estimate of parameter M), $\bar{r} = \frac{\sum r_i b_i}{\sum b_i} = \frac{D}{T_l}$. However, T_l is a function of θ through $\mathbb{E}[T_l] = 2\theta \times \log(n)$, where n is the number of sampled lineages (Arbissier et al., 2018, Tavaré et al., 1997). Thus, the prior on θ pertains to the tree height (through $\mathbb{E}[T_h] = 2 \times \theta(1 - \frac{1}{n}) \approx 2 \times \theta$ for large n), the tree length, and therefor the evolutionary rate.

Table 5: Prior configuration for the mean evolutionary rate, M of the lognormal distribution of branch rates. Note that the mean of the *Gamma* distribution here is *shape/rate* and that the true value used to generate the data is 1.5×10^{-5} subs/site/year.

Mean value of M	Prior configuration	95% quantile width / mean M
1.5×10^{-5}	$\text{Gamma}(\text{shape} = 1.5, \text{rate} = 1 \times 10^5)$	3.04
1.5×10^{-5}	$\text{Gamma}(\text{shape} = 0.3, \text{rate} = 2 \times 10^4)$	6.33
1.5×10^{-5}	$\text{Gamma}(\text{shape} = 15, \text{rate} = 1 \times 10^6)$	1.00
1.5×10^{-4}	$\text{Gamma}(\text{shape} = 1.5, \text{rate} = 1 \times 10^4)$	3.04
1.5×10^{-4}	$\text{Gamma}(\text{shape} = 0.3, \text{rate} = 2 \times 10^3)$	6.33
1.5×10^{-4}	$\text{Gamma}(\text{shape} = 15, \text{rate} = 1 \times 10^5)$	1.00
1.5×10^{-6}	$\text{Gamma}(\text{shape} = 1.5, \text{rate} = 1 \times 10^6)$	3.04
1.5×10^{-6}	$\text{Gamma}(\text{shape} = 0.3, \text{rate} = 2 \times 10^5)$	6.33
1.5×10^{-6}	$\text{Gamma}(\text{shape} = 15, \text{rate} = 1 \times 10^7)$	1.00

We used an additional configuration for the prior M using a hierarchical structure as follows:

$$M \sim \text{Gamma}(\text{shape}, \text{rate})$$

$$shape \sim LogNormal(1, 5)$$

$$rate \sim LogNormal(1, 5)$$

To set up this model, we simply treat the *shape* and *rate* as parameters that are sampled in the model.

We used Markov chain Monte Carlo (MCMC) to sample the posterior distribution for all analyses. We set the chain length to 10^8 steps, sampling every 5×10^4 steps. We deemed sufficient sampling by verifying that the effective sample size was at least 200, by using the R package CODA (version 0.19) (Plummer et al., 2006). When this criterion was not met we extended the chain length to 5×10^8 steps.

To visualise the prior on M in our model with a hierarchical prior we drew MCMC samples from the marginal prior of M . That is, the prior integrating over the hyperprior distributions and all other parameters. We obtained such samples by setting the option `sampleFromPrior="true"` in the input xml files in BEAST 2, which conducts the MCMC while ignoring the phylogenetic likelihood.

4.1.3 HBV empirical data

We selected a complete genome data set of HBV published by Kocher et al. (2021). The complete alignment included 232 genomes of length 3,344 nucleotides, with 1,807 variable sites, and 1,498 site patterns. The sampling times ranged from the present to 10,535 years before present. To investigate the impact of varying the sampling window and on temporal sampling bias we subsampled the data as described in our Results section. We analysed each data set using the same model and prior settings as in our simulations, including the use of the reasonable and misleading prior configuration.

5 Data availability

Computer code, analysis files, and data sets in this study are available at:

https://github.com/sebastianduchene/phylo_threshold_code_data

6 Competing interests

None.

7 Acknowledgments

Pending.

8 Funding

This work received funding from the Inception program (Investissement d'Avenir grant ANR-16-CONV-0005 awarded to SD) and the Australian National Health and Medical Research Council (2017284 awarded to SD),

References

- A. Andrades Valtueña, G. U. Neumann, M. A. Spyrou, L. Musralina, F. Aron, A. Beisenov, A. B. Belinskiy, K. I. Bos, A. Buzhilova, M. Conrad, et al. Stone age yersinia pestis genomes shed light on the early evolution, diversity, and ecology of plague. *Proceedings of the National Academy of Sciences*, 119(17): e2116722119, 2022.
- I. M. Arbisser, E. M. Jewett, and N. A. Rosenberg. On the joint distribution of tree height and tree length under the coalescent. *Theoretical population biology*, 122:46–56, 2018.
- G. Baele and P. Lemey. Bayesian model selection in phylogenetics and genealogy-based population genetics. In M. Chen, K. L. and L. PO, editors, *Bayesian phylogenetics, methods, algorithms, and applications*, chapter 4, pages 59–93. CPC Press, Boca Raton (Florida), 2014.
- R. Biek, O. G. Pybus, J. O. Lloyd-Smith, and X. Didelot. Measurably evolving pathogens in the genomic era. *Trends in ecology & evolution*, 30(6):306–313, 2015.
- M. F. Boni, P. Lemey, X. Jiang, T. T.-Y. Lam, B. W. Perry, T. A. Castoe, A. Rambaut, and D. L. Robertson. Evolutionary origins of the sars-cov-2 sarbecovirus lineage responsible for the covid-19 pandemic. *Nature microbiology*, 5(11):1408–1417, 2020.
- R. Bouckaert, T. G. Vaughan, J. Barido-Sottani, S. Duchêne, M. Fourment, A. Gavryushkina, J. Heled, G. Jones, D. Kühnert, N. De Maio, et al. Beast 2.5: An advanced software platform for bayesian evolutionary analysis. *PLoS computational biology*, 15(4):e1006650, 2019.
- L. Bromham, S. Duchêne, X. Hua, A. M. Ritchie, D. A. Duchêne, and S. Y. Ho. Bayesian molecular dating: opening up the black box. *Biological Reviews*, 93(2):1165–1191, 2018.
- J. M. Brown and R. C. Thomson. Evaluating model performance in evolutionary biology. *Annual Review of Ecology, Evolution, and Systematics*, 49(1):95–114, 2018.
- D. A. Buonagurio, S. Nakada, J. D. Parvin, M. Krystal, P. Palese, and W. M. Fitch. Evolution of human influenza a viruses over 50 years: rapid, uniform rate of change in ns gene. *Science*, 232(4753):980–982, 1986.
- A. M. Devault, G. B. Golding, N. Waglechner, J. M. Enk, M. Kuch, J. H. Tien, M. Shi, D. N. Fisman, A. N. Dhody, S. Forrest, et al. Second-pandemic strain of vibrio cholerae from the philadelphia cholera outbreak of 1849. *New England Journal of Medicine*, 370(4):334–340, 2014.
- M. Dos Reis and Z. Yang. The unbearable uncertainty of bayesian divergence time estimation. *Journal of Systematics and Evolution*, 51(1):30–43, 2013.
- J. Douglas, R. Zhang, and R. Bouckaert. Adaptive dating and fast proposals: Revisiting the phylogenetic relaxed clock model. *PLoS computational biology*, 17(2):e1008322, 2021.
- A. Drummond, O. G. Pybus, and A. Rambaut. Inference of viral evolutionary rates from molecular sequences. *Adv Parasitol*, 54:331–358, 2003a.

551 A. J. Drummond, O. G. Pybus, A. Rambaut, R. Forsberg, and A. G. Rodrigo. Measurably evolving popu-
552 lations. *Trends in ecology & evolution*, 18(9):481–488, 2003b.

553 A. J. Drummond, S. Y. W. Ho, M. J. Phillips, and A. Rambaut. Relaxed phylogenetics and dating with
554 confidence. *PLoS Biology*, 4(5):e88, 2006.

555 D. A. Duchêne, S. Duchêne, and S. Y. Ho. Phylomad: efficient assessment of phylogenomic model adequacy.
556 *Bioinformatics*, 34(13):2300–2301, 2018.

557 S. Duchêne, R. Lanfear, and S. Y. Ho. The impact of calibration and clock-model choice on molecular
558 estimates of divergence times. *Molecular phylogenetics and evolution*, 78:277–289, 2014.

559 S. Duchêne, D. Duchêne, E. C. Holmes, and S. Y. Ho. The performance of the date-randomization test in
560 phylogenetic analyses of time-structured virus data. *Molecular Biology and Evolution*, 32(7):1895–1906,
561 2015.

562 S. Duchene, K. E. Holt, F.-X. Weill, S. Le Hello, J. Hawkey, D. J. Edwards, M. Fourment, and E. C. Holmes.
563 Genome-scale rates of evolutionary change in bacteria. *Microbial genomics*, 2(11):e000094, 2016.

564 S. Duchene, D. A. Duchene, J. L. Geoghegan, Z. A. Dyson, J. Hawkey, and K. E. Holt. Inferring demographic
565 parameters in bacterial genomic data using bayesian and hybrid phylogenetic methods. *BMC evolutionary*
566 *biology*, 18:1–11, 2018.

567 S. Duchene, R. Bouckaert, D. A. Duchene, T. Stadler, and A. J. Drummond. Phylodynamic model adequacy
568 using posterior predictive simulations. *Systematic biology*, 68(2):358–364, 2019.

569 S. Duchene, L. Featherstone, M. Haritopoulou-Sinanidou, A. Rambaut, P. Lemey, and G. Baele. Temporal
570 signal and the phylodynamic threshold of sars-cov-2. *Virus evolution*, 6(2):veaa061, 2020a.

571 S. Duchene, S. Y. Ho, A. G. Carmichael, E. C. Holmes, and H. Poinar. The recovery, interpretation and use
572 of ancient pathogen genomes. *Current Biology*, 30(19):R1215–R1231, 2020b.

573 S. Duchene, P. Lemey, T. Stadler, S. Y. Ho, D. A. Duchene, V. Dhanasekaran, and G. Baele. Bayesian
574 evaluation of temporal signal in measurably evolving populations. *Molecular Biology and Evolution*, 37
575 (11):3363–3379, 2020c.

576 K. Eaton, L. Featherstone, S. Duchene, A. G. Carmichael, N. Varlik, G. B. Golding, E. C. Holmes, and
577 H. N. Poinar. Plagued by a cryptic clock: insight and issues from the global phylogeny of yersinia pestis.
578 *Communications Biology*, 6(1):23, 2023.

579 V. Eldholm, J. Monteserin, A. Rieux, B. Lopez, B. Sobkowiak, V. Ritacco, and F. Balloux. Four decades of
580 transmission of a multidrug-resistant mycobacterium tuberculosis outbreak strain. *Nature communications*,
581 6(1):7119, 2015.

582 L. A. Featherstone, A. Rambaut, S. Duchene, and W. Wirth. Clockor2: Inferring global and local strict
583 molecular clocks using root-to-tip regression. *Systematic biology*, 73(3):623–628, 2024.

584 L. Ferretti, T. Golubchik, F. Di Lauro, M. Ghafari, J. Villabona-Arenas, K. E. Atkins, C. Fraser, and M. Hall.
585 Biased estimates of phylogenetic branch lengths resulting from the discretised gamma model of site rate
586 heterogeneity. *bioRxiv*, pages 2024–08, 2024.

587 A. Gavryushkina, T. A. Heath, D. T. Ksepka, T. Stadler, D. Welch, and A. J. Drummond. Bayesian total-
588 evidence dating reveals the recent crown radiation of penguins. *Systematic biology*, 66(1):57–73, 2017.

589 M. Ghafari, L. du Plessis, J. Raghvani, S. Bhatt, B. Xu, O. G. Pybus, and A. Katzourakis. Purifying selec-
590 tion determines the short-term time dependency of evolutionary rates in sars-cov-2 and ph1n1 influenza.
591 *Molecular Biology and Evolution*, 39(2):msac009, 2022.

592 N. Gharbi, E. Rousseau, and T. Wirth. Clock rates and bayesian evaluation of temporal signal. In *Phyloge-*
593 *nomics*, pages 153–175. Elsevier, 2024.

594 T. Gojobori, E. N. Moriyama, and M. Kimura. Molecular clock of viral evolution, and the neutral theory.
595 *Proceedings of the National Academy of Sciences*, 87(24):10015–10018, 1990.

596 S. Guindon. Rates and rocks: strengths and weaknesses of molecular dating methods. *Frontiers in Genetics*,
597 11:526, 2020.

598 S. Y. Ho. The changing face of the molecular evolutionary clock. *Trends in Ecology & Evolution*, 29(9):
599 496–503, 2014.

600 S. Y. Ho and S. Duchêne. Molecular-clock methods for estimating evolutionary rates and timescales. *Molecular*
601 *Ecology*, 23(24):5947–5965, 2014.

602 G. Kahila Bar-Gal, M. J. Kim, A. Klein, D. H. Shin, C. S. Oh, J. W. Kim, T.-H. Kim, S. B. Kim, P. R.
603 Grant, O. Pappo, et al. Tracing hepatitis b virus to the 16th century in a korean mummy. *Hepatology*, 56
604 (5):1671–1680, 2012.

605 A. Kocher, L. Papac, R. Barquera, F. M. Key, M. A. Spyrou, R. Hübner, A. B. Rohrlach, F. Aron, R. Stahl,
606 A. Wissgott, et al. Ten millennia of hepatitis b virus evolution. *Science*, 374(6564):182–188, 2021.

607 D. Kühnert, T. Stadler, T. G. Vaughan, and A. J. Drummond. Phylodynamics with migration: a computa-
608 tional framework to quantify population structure from genomic data. *Molecular biology and evolution*, 33
609 (8):2102–2116, 2016.

610 D. Kühnert, M. Coscolla, D. Brites, D. Stucki, J. Metcalfe, L. Fenner, S. Gagneux, and T. Stadler. Tuber-
611 culosis outbreak investigation using phylodynamic analysis. *Epidemics*, 25:47–53, 2018.

612 S. A. Locarnini, M. Littlejohn, and L. K. Yuen. Origins and evolution of the primate hepatitis b virus.
613 *Frontiers in Microbiology*, 12:653684, 2021.

614 R. McElreath. *Statistical rethinking: A Bayesian course with examples in R and Stan*. Chapman and
615 Hall/CRC, 2018.

616 F. Menardo, S. Duchêne, D. Brites, and S. Gagneux. The molecular clock of mycobacterium tuberculosis.
617 *PLoS pathogens*, 15(9):e1008067, 2019.

618 F. K. Mendes, R. Bouckaert, L. M. Carvalho, and A. J. Drummond. How to validate a bayesian evolutionary
619 model. *Systematic Biology*, 74(1):158–175, 2025.

620 M. Merker, J.-P. Rasigade, M. Barbier, H. Cox, S. Feuerriegel, T. A. Kohl, E. Shitikov, K. Klaos, C. Gaudin,
621 R. Antoine, et al. Transcontinental spread and evolution of mycobacterium tuberculosis w148 euro-
622 pean/russian clade toward extensively drug resistant tuberculosis. *Nature Communications*, 13(1):5105,
623 2022.

624 S. Möller, L. du Plessis, and T. Stadler. Impact of the tree prior on estimating clock rates during epidemic
625 outbreaks. *Proceedings of the National Academy of Sciences*, 115(16):4200–4205, 2018.

626 B. Mühlemann, T. C. Jones, P. d. B. Damgaard, M. E. Allentoft, I. Shevnina, A. Logvin, E. Usmanova, I. P.
627 Panyushkina, B. Boldgiv, T. Bazartseren, et al. Ancient hepatitis b viruses from the bronze age to the
628 medieval period. *Nature*, 557(7705):418–423, 2018.

629 N. F. Müller, D. A. Rasmussen, and T. Stadler. The structured coalescent and its approximations. *Molecular*
630 *biology and evolution*, 34(11):2970–2981, 2017.

631 G. G. Murray, F. Wang, E. M. Harrison, G. K. Paterson, A. E. Mather, S. R. Harris, M. A. Holmes,
632 A. Rambaut, and J. J. Welch. The effect of genetic structure on molecular dating and tests for temporal
633 signal. *Methods in Ecology and Evolution*, 7(1):80–89, 2016.

634 M. Nordborg. Coalescent theory. *Handbook of Statistical Genomics: Two Volume Set*, pages 145–30, 2019.

635 D. J. Nott, X. Wang, M. Evans, and B.-G. Englert. Checking for prior-data conflict using prior-to-posterior
636 divergences. *Statistical Science*, 35(2):234–253, 2020.

637 D. Paraskevis, G. Magiorkinis, E. Magiorkinis, S. Y. Ho, R. Belshaw, J.-P. Allain, and A. Hatzakis. Dating
638 the origin and dispersal of hepatitis b virus infection in humans and primates. *Hepatology*, 57(3):908–916,
639 2013.

640 M. Plummer, N. Best, K. Cowles, K. Vines, et al. Coda: convergence diagnosis and output analysis for
641 mcmc. *R news*, 6(1):7–11, 2006.

642 A. Rambaut, T. T. Lam, L. Max Carvalho, and O. G. Pybus. Exploring the temporal structure of hete-
643 rochronous sequences using tempest (formerly path-o-gen). *Virus evolution*, 2(1):vew007, 2016.

644 C. Ramsden, E. C. Holmes, and M. A. Charleston. Hantavirus evolution in relation to its rodent and
645 insectivore hosts: no evidence for codivergence. *Molecular biology and evolution*, 26(1):143–153, 2009.

646 A. Rieux and F. Balloux. Inferences from tip-calibrated phylogenies: a review and a practical guide. *Molecular*
647 *ecology*, 25(9):1911–1924, 2016.

648 F. Ronquist, N. Lartillot, and M. J. Phillips. Closing the gap between rocks and clocks using total-evidence
649 dating. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 371(1699):20150136, 2016.

650 Z. P. Ross, J. Klunk, G. Fornaciari, V. Giuffra, S. Duchêne, A. T. Duggan, D. Poinar, M. W. Douglas, J.-S.
651 Eden, E. C. Holmes, et al. The paradox of hbv evolution as revealed from a 16th century mummy. *PLoS*
652 *pathogens*, 14(1):e1006750, 2018.

653 R. Sanjuán. From molecular genetics to phylodynamics: evolutionary relevance of mutation rates across
654 viruses. *PLoS pathogens*, 8(5):e1002685, 2012.

655 K. P. Schliep. phangorn: phylogenetic analysis in r. *Bioinformatics*, 27(4):592–593, 2011.

656 M. A. Spyrou, K. I. Bos, A. Herbig, and J. Krause. Ancient pathogen genomics as an emerging tool for
657 infectious disease research. *Nature Reviews Genetics*, 20(6):323–340, 2019a.

658 M. A. Spyrou, M. Keller, R. I. Tukhbatova, C. L. Scheib, E. A. Nelson, A. Andrades Valtueña, G. U.
659 Neumann, D. Walker, A. Alterauge, N. Carty, et al. Phylogeography of the second plague pandemic
660 revealed through analysis of historical yersinia pestis genomes. *Nature communications*, 10(1):4470, 2019b.

661 S. Tavaré, D. J. Balding, R. C. Griffiths, and P. Donnelly. Inferring coalescence times from dna sequence
662 data. *Genetics*, 145(2):505–518, 1997.

663 J. H. Tay, A. Kocher, and S. Duchene. Assessing the effect of model specification and prior sensitivity on
664 bayesian tests of temporal signal. *PLoS Computational Biology*, 20(11):e1012371, 2024.

665 N. S. Trovão, G. Baele, B. Vrancken, F. Bielejec, M. A. Suchard, D. Fargette, and P. Lemey. Host ecology
666 determines the dispersal patterns of a plant virus. *Virus evolution*, 1(1):vev016, 2015.

667 T. G. Vaughan. Remaster: improved phylodynamic simulation for beast 2.7. *Bioinformatics*, 40(1):btae015,
668 2024.

669 B. Vrancken, M. A. Suchard, and P. Lemey. Accurate quantification of within-and between-host hbv evolu-
670 tionary rates requires explicit transmission chain modelling. *Virus Evolution*, 3(2):vex028, 2017.

671 Y. Wang and Z. Yang. Priors in bayesian phylogenetics. *Bayesian phylogenetics: methods, algorithms, and*
672 *applications*, pages 5–24, 2014.

673 R. C. Warnock, Z. Yang, and P. C. Donoghue. Exploring uncertainty in the calibration of the molecular
674 clock. *Biology letters*, 8(1):156–159, 2012.

675 G. Zehender, E. Ebranati, E. Gabanelli, C. Sorrentino, A. L. Presti, E. Tanzi, M. Ciccozzi, and M. Galli.
676 Enigmatic origin of hepatitis b virus: an ancient travelling companion or a recent encounter? *World*
677 *Journal of Gastroenterology: WJG*, 20(24):7622, 2014.

678 E. Zuckerkandl and L. Pauling. Evolutionary divergence and convergence in proteins. In *Evolving genes and*
679 *proteins*, pages 97–166. Elsevier, 1965.