

1  
2

### 3 **Bayesian Evaluation of Temporal Signal in Measurably Evolving** 4 **Populations**

5  
6  
7  
8  
9

Sebastian Duchene<sup>\*,1</sup>, Philippe Lemey<sup>2</sup>, Tanja Stadler<sup>3</sup>, Simon YW Ho<sup>4</sup>, David A Duchene<sup>5</sup>,  
Vijaykrishna Dhanasekaran<sup>6</sup>, Guy Baele<sup>2</sup>

10 <sup>1</sup>Department of Microbiology and Immunology, Peter Doherty Institute for Infection and Immunity,  
11 University of Melbourne, Melbourne, VIC, Australia.

12 <sup>2</sup>Department of Microbiology, Immunology and Transplantation, Rega Institute, KU Leuven,  
13 Leuven, Belgium.

14 <sup>3</sup>Department of Biosystems Science and Engineering, ETH Zürich; and Swiss Institute of  
15 Bioinformatics; Basel, Switzerland

16 <sup>4</sup>School of Life and Environmental Sciences, University of Sydney, Sydney, NSW, Australia

17 <sup>5</sup>Research School of Biology, Australian National University, Canberra, ACT, Australia.

18 <sup>6</sup>Department of Microbiology, Biomedicine Discovery Institute, Monash University,  
19 Melbourne, VIC, Australia.

20  
21

22  
23

\*Corresponding author: E-mail: [sebastian.duchene@unimelb.edu.au](mailto:sebastian.duchene@unimelb.edu.au)

24  
25

#### **Abstract (250 words max)**

26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

Phylogenetic methods can use the sampling times of molecular sequence data to calibrate the molecular clock, enabling the estimation of substitution rates and time scales for rapidly evolving pathogens and data sets containing ancient DNA samples. A key aspect of such calibrations is whether a sufficient amount of molecular evolution has occurred over the sampling time window, that is, whether the data can be treated as being from a measurably evolving population. Here we investigate the performance of a fully Bayesian evaluation of temporal signal (BETS) in molecular sequence data. The method involves comparing the fit of two models: a model in which the data are accompanied by the actual (heterochronous) sampling times, and a model in which the samples are constrained to be contemporaneous (isochronous). We conduct extensive simulations under a range of conditions to demonstrate that BETS accurately classifies data sets according to whether they contain temporal signal or not, even when there is substantial among-lineage rate variation. We explore the behaviour of this classification in analyses of five data sets: modern samples of *A/H1N1 influenza virus*, the bacterium *Bordetella pertussis*, and coronaviruses from mammalian hosts, and ancient DNA data sets of *Hepatitis B virus* and of mitochondrial genomes of dog species. Our results indicate that BETS is an effective alternative to other measures of temporal signal. In particular, this method has the key advantage of allowing a coherent assessment of the entire model, including the molecular clock and tree prior which are essential aspects of Bayesian phylodynamic analyses.

**Key words:** Bayesian phylogenetics, ancient DNA, measurably evolving population, marginal likelihood, molecular clock, temporal signal.

## 48 Introduction

49 The molecular clock has become a ubiquitous tool for studying evolutionary processes in rapidly  
50 evolving organisms and in data sets that include ancient DNA. In its simplest form, the molecular  
51 clock posits that evolutionary change occurs at a predictable rate over time (Zuckerandl and  
52 Pauling 1965). The molecular clock can be calibrated to estimate divergence times by using  
53 sampling time information, the timing of known divergence events, or a previous estimate of the  
54 substitution rate (Hipsley and Müller 2014). For example, Korber et al. (2000) used sampling times  
55 to calibrate the molecular clock and to infer the time of origin of HIV group 1. Their approach  
56 consisted of estimating a phylogenetic tree and conducting a regression of the distance from the  
57 root to each of the tips as a function of sequence sampling time. In this method, the slope of the  
58 regression is an estimate of the substitution rate in substitutions per site per unit of time, the  
59 intercept with the time axis is the age of the root node, and the coefficient of determination ( $R^2$ ) is  
60 the degree to which the data exhibit clocklike behaviour (Rambaut et al. 2016). Despite the  
61 practicality of root-to-tip regression, its use as a statistical tool for molecular dating has several  
62 well-known limitations. In particular, data points are not independent because they have shared  
63 ancestry (i.e., internal branches are traversed multiple times) and a strict clocklike behaviour is  
64 assumed by necessity.

65  
66 The past few decades have seen a surge in molecular clock models that explicitly use phylogenetic  
67 information. Bayesian methods have gained substantial popularity, largely due to the wide array of  
68 complex models that can be implemented and the fact that independent information, including  
69 calibrations, can be specified via prior distributions (Nascimento et al. 2017). Of particular  
70 importance is the availability of molecular clock models that relax the assumption of strict clock  
71 behaviour by explicitly modelling rate variation among lineages (reviewed by Ho and Duchene  
72 (2014) and by Bromham et al. (2018)).

73  
74 Regardless of the methodology used to analyse time-stamped sequence data, a sufficient amount  
75 of molecular evolution must have occurred over the sampling time window to allow reliable  
76 estimates of substitution rates and timescales. In such cases, the population can be considered to  
77 be 'measurably evolving' (Drummond et al. 2003). The degree of 'temporal information' in sequence  
78 data is determined by the sequence length, the substitution rate, and the range of available  
79 sampling times. Some viruses evolve at a rate of around  $5 \times 10^{-3}$  subs/site/year (Duchene et al. 2014),  
80 such that samples collected over a few weeks can be sufficient to calibrate the molecular clock. In  
81 more slowly evolving organisms, such as mammals, a sampling window of tens of thousands of  
82 years might be necessary; this can be achieved by including ancient DNA sequences (Drummond et  
83 al. 2003; Biek et al. 2015).

84  
85 Testing for temporal signal is an important step for verifying that the molecular clock can be  
86 calibrated using the sampling times (Rieux and Balloux 2016). For this purpose a date-  
87 randomization test has been proposed that compares actual substitution rate estimates to those  
88 obtained by repeatedly permuting the sequence sampling times (Ramsden et al. 2009). A data set is  
89 considered to have strong temporal signal if the rate estimated using the correct sampling times  
90 does not overlap with those of the permutation replicates (Duchene et al. 2015, 2018; Murray et al.  
91 2015). An implementation of this test is also available that performs the permutation during a single  
92 Bayesian inference (Trovão et al. 2015). The interpretation of the date-randomization test is  
93 essentially frequentist in nature, which leads to an inconsistent mixture of statistical frameworks  
94 when Bayesian phylogenetic methods are used. Moreover, the procedure is not applicable in cases  
95 with small numbers of sampling times, owing to the limited number of possible permutations  
96 (Duchene et al. 2015).

97

98 We propose a full Bayesian model test, which we refer to as BETS (Bayesian Evaluation of Temporal  
99 Signal), to assess temporal signal based on previous analyses by Baele et al. (2012). The approach  
100 involves quantifying statistical support for two competing models: a model in which the data are  
101 accompanied by the actual sampling times (i.e., the data are treated as heterochronous) and a  
102 model in which the sampling times are contemporaneous (i.e., the data are treated as isochronous).  
103 Therefore, the sampling times are treated as part of the model and the test can be understood as a  
104 test of ultrametricity of the phylogenetic tree. If incorporating sampling times improves the  
105 statistical fit, then their use for clock calibration is warranted. The crux of BETS, as with Bayesian  
106 model selection, is that it requires calculating the marginal likelihood of the model in question. The  
107 marginal likelihood measures the evidence for a model given the data, and calculating it requires  
108 integration of its likelihood across all parameter values, weighted by the prior (Kass and Raftery  
109 1995).

110  
111 Because the marginal likelihood is a measure of model evidence, the ratio of the marginal  
112 likelihoods of two competing models, known as the Bayes factor, is used to assess support for one  
113 model relative to the other. In the case of applying BETS, let  $M_{\text{het}}$  represent the heterochronous  
114 model,  $M_{\text{iso}}$  the isochronous model, and  $Y$  the sequence data, such that  $P(Y|M_{\text{het}})$  and  $P(Y|M_{\text{iso}})$  are  
115 their respective marginal likelihoods. These models differ in the number of parameters; in  $M_{\text{iso}}$  the  
116 substitution rates and times are nonidentifiable, so the rate is fixed to an arbitrary value, while in  
117  $M_{\text{het}}$  it is a free parameter. Differences in the number of parameters do not need to be taken into  
118 account separately, because the marginal likelihood naturally penalizes excessive parameterization.  
119 Kass and Raftery (1995) gave guidelines to interpreting Bayes factors, where a (log) Bayes factor  
120  $\log(P(Y|M_{\text{het}})) - \log(P(Y|M_{\text{iso}}))$  of at least 5 indicates 'very strong' support for  $M_{\text{het}}$  over  $M_{\text{iso}}$ , a value of  
121 3 indicates 'strong' support, and a value of 1 is considered as positive evidence for  $M_{\text{het}}$  over  $M_{\text{iso}}$ .

122  
123 The importance of model selection in Bayesian phylogenetics has prompted the development of  
124 various techniques to calculate marginal likelihoods (reviewed by Baele et al. (2014) and by Oaks et  
125 al. (2019)). These techniques can be broadly classified into prior-based and/or posterior-based  
126 estimators and path-sampling approaches. Prior- and posterior-based estimators, also known as  
127 importance sampling, include the widely used harmonic-mean estimator (Newton and Raftery  
128 1994) and the AICM and BICM (Bayesian analogues to the Akaike information criterion and the  
129 Bayesian information criterion, respectively) (Raftery et al. 2007). These scores are easy to compute  
130 because they only require samples from the posterior distribution as obtained through Markov  
131 chain Monte Carlo (MCMC) integration. However, the harmonic-mean estimator has been shown to  
132 have unacceptably high variance when the prior is diffuse relative to the posterior, and, together  
133 with the AICM, has shown poor performance in practical settings (Baele et al. 2012, 2013). The BICM  
134 requires a sample size to be specified for each parameter, which is far from trivial for phylogenetic  
135 inference and therefore remains unexplored for such applications.

136  
137 Path-sampling approaches include path sampling (originally introduced in phylogenetics as  
138 'thermodynamic integration') (Lartillot and Philippe 2006), stepping-stone sampling (Xie et al.  
139 2011), and generalized stepping-stone (GSS) sampling (Fan et al. 2011; Baele et al. 2016). These  
140 methods depend on drawing samples using MCMC from a range of power posterior distributions  
141 that represent the path from the posterior to the (working) prior, and therefore require additional  
142 computation. Another numerical technique that was recently introduced to phylogenetics is nested  
143 sampling (NS) (Maturana et al. 2019), which approximates the marginal likelihood by simplifying  
144 the marginal-likelihood function from a multi-dimensional to a one-dimensional integral over the  
145 cumulative distribution function of the marginal likelihood (Skilling 2006). Fourment et al. (2019)  
146 recently compared the accuracy of a range of marginal-likelihood estimation methods and found  
147 GSS to be the most accurate, albeit at increased computational cost. Clearly, the reliability of the  
148 marginal-likelihood estimator is a key consideration for applying BETS.

149

150 We conducted a simulation study to assess the reliability of BETS under a range of conditions that  
151 are typical for data sets of rapidly evolving organisms and of those involving ancient DNA. We also  
152 analysed five empirical data sets to showcase the performance of the test in practice. Our analyses  
153 demonstrate the utility of BETS to provide accurate evaluation of temporal signal across a wide  
154 range of conditions.

155

## 156 **Results**

157

### 158 **Simulations of Measurably Evolving Populations**

159 In our simulations we considered sequence data from heterochronous and isochronous trees.  
160 Heterochronous trees represent a situation where there is sufficient temporal signal, whereas  
161 isochronous trees lack temporal signal altogether. We simulated heterochronous phylogenetic  
162 trees under a stochastic birth-death process with between 90 and 110 tips. To generate isochronous  
163 trees we used similar settings, but we assumed a single sampling time. We then simulated  
164 substitution rates along the trees according to an uncorrelated relaxed clock with an underlying  
165 lognormal distribution with a mean of  $5 \times 10^{-3}$  subs/site/unit time and a standard deviation,  $\sigma$ , of 0,  
166 0.1, 0.5, or 1, where  $\sigma=0$  is equivalent to simulating under a strict clock. We then simulated  
167 sequence evolution using an HKY+ $\Gamma$  substitution model, with parameter values similar to those  
168 estimated for influenza virus (Hedge et al. 2013), to generate alignments of 10,000 nucleotides.

169

170 Our main simulation conditions produced data sets in which about 50% of the sites were variable.  
171 We refer to this simulation scenario as (i) 'high substitution rate and wide sampling window', and we  
172 considered three other simulation scenarios that involved (ii) a lower substitution rate of  $10^{-5}$   
173 subs/site/unit time, (iii) a narrower sampling window, and (iv) both of the last two conditions. We  
174 analysed the sequence data using a strict clock and an uncorrelated relaxed clock with an  
175 underlying lognormal distribution (Drummond et al. 2006). We considered three configurations for  
176 sampling times: birth-death sampling times, which are correct for the heterochronous data but not  
177 for the isochronous data; identical sampling times, which is correct for isochronous data but not for  
178 the heterochronous data; and permuted birth-death sampling times, which are incorrect for both  
179 heterochronous and isochronous data.

180

181 We estimated the log marginal likelihoods of these six combinations of sampling times and clock  
182 models using NS and GSS as implemented in BEAST 2.5 (Bouckaert et al. 2019) and BEAST 1.10  
183 (Suchard et al. 2018), respectively. Our BETS approach ranked the models according to their log  
184 marginal likelihoods and computed log Bayes factors of the best heterochronous model ( $M_{\text{het}}$ )  
185 compared with the best isochronous model ( $M_{\text{iso}}$ ).

186

#### 187 *(i) Simulations with High Substitution Rate and Wide Sampling Window*

188 Both NS and GSS correctly classified data as being heterochronous or isochronous in 10 out of 10  
189 simulations, including in the presence of a high degree of among-lineage rate variation (i.e.,  $\sigma=1$ ;  
190 fig. 1 for heterochronous data and supplementary fig. S1, Supplementary Material online, for  
191 isochronous data). Although both marginal-likelihood estimators detected temporal signal, NS  
192 supported the relaxed clock over the strict clock for three heterochronous data sets simulated  
193 without among-lineage rate variation ( $\sigma=0$ ) and for six data sets simulated with low among-lineage  
194 rate variation ( $\sigma=0.1$ ). In the simulations of isochronous data, NS often favoured the relaxed clock  
195 over the strict clock when there was low among-lineage rate variation ( $\sigma=0.0$  and  $\sigma=0.1$ ), albeit  
196 mostly with log Bayes factors below 5 (supplementary fig. S1, Supplementary Material online). In  
197 contrast, GSS always selected the strict clock under these conditions (fig. 1 and supplementary fig.  
198 S1, Supplementary Material online).

199

200 For the heterochronous data sets, NS and GSS always displayed very strong support for  $M_{het}$  over  
201  $M_{iso}$ , with log Bayes factors of at least 90. For the isochronous data sets, the log Bayes factors for  
202  $M_{iso}$  relative to  $M_{het}$  were overall much lower, but still decisive, ranging from 30 to 50. Furthermore,  
203 log Bayes factors tended to decline with an increasing degree of among-lineage rate variation in the  
204 data. Another important observation is that in the heterochronous data, the relaxed clock was  
205 consistently selected over the strict clock when assuming that the data were isochronous, or when  
206 the sampling times had been permuted (fig. 1 and supplementary fig. S1, Supplementary Material  
207 online). Moreover, the strict clock with permuted sampling times yielded the lowest log marginal  
208 likelihoods for heterochronous data. Both of these patterns are likely to be due to an apparently  
209 higher degree of among-lineage rate variation when sampling times are misspecified.

210

#### 211 *(ii) Simulations with Low Substitution Rate and Wide Sampling Window*

212 Our simulations with a low substitution rate of  $10^{-5}$  subs/site/unit time produced data sets that each  
213 had about 10 variable sites, which provides very little information for the estimation of evolutionary  
214 parameters. Additionally, due to the stochasticity of the simulation process, increased estimator  
215 variance between replicates is to be expected given the small number of variable sites. For the  
216 heterochronous data sets, GSS selected the heterochronous model with correct dates in at least 7  
217 out of 10 simulation replicates (fig. 2). Across the simulations with different clock models (40 in  
218 total), only in five heterochronous data sets did we find models with permuted sampling times to  
219 have the highest log marginal likelihoods. For NS, in 11 out of 40 simulations, either isochronous  
220 models or those with random sampling times were incorrectly selected when heterochronous data  
221 sets were analysed.

222

223 Log marginal likelihoods calculated using GSS tended to support models with sampling times  
224 (either permuted or those from the birth-death) for the isochronous data, whereas NS appeared to  
225 support all models with similar frequencies (supplementary fig. S2, Supplementary Material online).  
226 However, a critical feature of the results from the data sets with a low substitution rate is that the  
227 log marginal likelihoods for all models were more similar to one another than those for the data sets  
228 with high substitution rate (note that the log marginal likelihood scale in fig. 2 is smaller than that in  
229 fig. 1). As a case in point, for the isochronous data with  $\sigma=0.1$  there were log Bayes factors of about  
230 0.1 for the best model with birth-death sampling times relative to those with permuted sampling  
231 times. This result indicates that comparing models with permuted sampling times might be useful  
232 for determining whether the data are informative about a particular set of sampling times.

233

#### 234 *(iii) Simulations with High Substitution Rate and Narrow Sampling Window*

235 We conducted a set of simulations similar to those described in scenario (i) but where sequence  
236 sampling spanned only the last 10% of the age of the tree (0.5 units of time, compared with 5 units  
237 of time for the simulations with a wide sampling window). These conditions reflect those of  
238 organisms with deep evolutionary histories and for which samples are available for only a small  
239 portion of this time. Since in these trees the samples were collected over a narrower time window,  
240 we used a higher sampling probability to obtain about 100 samples, as in our other simulations (see  
241 examples of trees in supplementary fig. S3, Supplementary Material online). For these analyses we  
242 only considered heterochronous data because the isochronous case is the same as that in scenario  
243 (i).

244

245 Both GSS and NS showed excellent performance in detecting temporal signal in this scenario,  
246 almost always selecting models with correct sampling times. The exceptions to this pattern  
247 occurred for one data set with  $\sigma=0.5$  and for two data sets with  $\sigma=1.0$  for NS (fig. 3).

248

Differentiating between the strict clock and relaxed clock appeared somewhat more difficult,

249 particularly for NS, where the relaxed clock with correct sampling times yielded log marginal  
250 likelihoods very similar to those for the strict clock for data with low among-lineage rate variation  
251 ( $\sigma$  of 0.0 or 0.1). Although NS and GSS performed well in these simulations, the log Bayes factors  
252 for  $M_{\text{het}}$  relative to  $M_{\text{iso}}$  were much lower than those for data with a high substitution rate and a wide  
253 sampling window in (i). One obvious example is in the data with  $\sigma=0.0$ , where the mean log Bayes  
254 factors for  $M_{\text{het}}$  over  $M_{\text{iso}}$  using GSS was 203.15 with a wide sampling window (fig. 1), but only 35.77  
255 when sampling spanned a narrow time window (fig. 3).

256

#### 257 *(iv) Simulations with Low Substitution Rate and Narrow Sampling Window*

258 We considered data sets with a narrow sampling window, as in scenario (iii), and with a low  
259 substitution rate of  $10^{-5}$  subs/site/unit time, as in scenario (ii). We generated only heterochronous  
260 trees under these conditions, because the isochronous case would be the same as that in (ii).

261

262 Estimates of log marginal likelihoods with GSS and NS were very similar among models, with mean  
263 log Bayes factors among data sets of less than 1 for the two models with highest marginal  
264 likelihoods for GSS (fig. 4). In the data sets with  $\sigma=0.0$ , GSS and NS always preferred a  
265 heterochronous model. However, in a few cases (three for GSS and one for NS) the model with  
266 permuted sampling times was selected, indicating that temporal signal was not detected. As with  
267 the data sets with low substitution rate and constant sampling (ii), the relaxed clock was sometimes  
268 preferred over the strict clock, even when the data sets had no rate variation among lineages.

269

#### 270 *Comparison with Root-to-tip Regression*

271 Using a subset of the heterochronous data sets, we conducted root-to-tip regression using  
272 phylogenetic trees inferred using maximum likelihood in PhyML 3.1 (Guindon et al. 2010) with the  
273 same substitution model as in our BEAST analyses, and with the placement of the root chosen to  
274 maximize  $R^2$  in TempEst (Rambaut et al. 2016). We selected data sets generated with a high  
275 substitution rate and with both constant and narrow sampling windows. Because GSS and NS  
276 correctly detected temporal signal under these conditions, these regressions demonstrate the  
277 extent to which this informal regression assessment matches the BETS approach. We did not  
278 attempt to provide a thorough benchmarking of the two methods here.

279

280 All regressions had  $R^2$  values that matched our expectation from the degree of among-lineage rate  
281 variation, that is, higher values of  $\sigma$  corresponded to lower values of  $R^2$  (fig. 5). The data with a  
282 wide sampling window yielded regression slopes ranging from  $7.3 \times 10^{-3}$  to  $5.4 \times 10^{-3}$  subs/site/unit  
283 time, which is similar to the substitution rate values used to generate the data. Although the root-  
284 to-tip regression is sometimes used to assess temporal signal, it has no cut-off values to confirm  
285 temporal signal. This becomes critical when considering the data with a narrow sampling window,  
286 for which the  $R^2$  was between 0.13 and 0.02. For example, the regression for a data set with  $\sigma=1$   
287 and narrow sampling window had an  $R^2$  of 0.02, which is sometimes considered sufficiently low as  
288 to preclude molecular clock analyses (Rieux and Balloux 2016). However, BETS supported strong  
289 temporal structure under a relaxed clock in this data set, with log Bayes factors of 5.48 for this  
290 particular data set, which matches the simulation conditions. More importantly, even with such  
291 high rate variation, the substitution rate estimated using a relaxed clock and the correct sampling  
292 times included the true value used to generate the data ( $5 \times 10^{-3}$  subs/site/unit time), with a 95%  
293 highest posterior density (HPD) of between  $2.15 \times 10^{-3}$  and  $1.90 \times 10^{-2}$  subs/site/unit time, while the  
294 regression slope was  $2.22 \times 10^{-2}$  subs/site/unit time. A key implication of these comparisons is that  
295 BETS provides a formal assessment of temporal signal, unlike statistics computed from the  
296 regression. Moreover, the root-to-tip regression appears uninformative when the data have been  
297 sampled over a narrow time window and there is some rate variation among lineages.

298

## 299 *Analyses of Empirical Data Sets*

300 We analysed five empirical data sets with similar configurations of sampling times as in our  
301 simulation study (Table 1). Two data sets consisted of rapidly evolving pathogens: *A/H1N1 influenza*  
302 *virus* (Hedge et al. 2013) and *Bordetella pertussis* (Bart et al. 2014). We also analysed a data set with  
303 highly divergent sequences of coronaviruses (Wertheim et al. 2013), and two data sets with ancient  
304 DNA: *Hepatitis B virus* (Patterson Ross et al. 2018), and mitochondrial genomes of dog species  
305 (Thalmann et al. 2013). Due to the demonstrated higher accuracy of GSS over NS (Fourment et al.  
306 2019), we applied the BETS approach using the former method only.

307  
308 The *A/H1N1 influenza virus* data demonstrated clear temporal signal, with the strict clock and  
309 relaxed clock with the correct sampling times having the highest log marginal likelihoods, and a log  
310 Bayes factor of  $M_{\text{het}}$  with respect to  $M_{\text{iso}}$  of 150 (fig. 6). The strict clock had higher support than the  
311 relaxed clock for the correct sampling times (log Bayes factor 3.41). Broadly, this result is consistent  
312 with previous evidence of strong temporal signal and clocklike behaviour in this data set (Hedge et  
313 al. 2013). Using the strict clock with correct sampling times we estimated a substitution rate of  
314  $3.37 \times 10^{-3}$  subs/site/year (HPD:  $2.98 \times 10^{-3}$  to  $3.78 \times 10^{-3}$ ).

315  
316 We detected temporal signal in the *Bordetella pertussis* data set (fig. 6). The relaxed clock with the  
317 correct sampling times had the highest log marginal likelihood, with a log Bayes factor relative to  
318 the strict clock of 28.86. The log Bayes factor for  $M_{\text{het}}$  relative to  $M_{\text{iso}}$  was 47.40. These results echo  
319 previous assessments of these data using a date-randomization test (Duchene et al. 2016). We  
320 estimated a mean substitution rate using the best model of  $1.65 \times 10^{-7}$  subs/site/year (95% HPD:  
321  $1.36 \times 10^{-7}$  to  $2.00 \times 10^{-7}$ ).

322  
323 Our analyses did not detect temporal signal in the coronavirus data, for which the strict clock and  
324 relaxed clock with no sampling times had the highest log marginal likelihoods. The log Bayes factor  
325 of  $M_{\text{het}}$  relative to  $M_{\text{iso}}$  was -16.82, indicating strong support for the isochronous model. The relaxed  
326 clock was supported over the strict clock, with a log Bayes factor of 19.25 (fig. 7). Previous analyses  
327 of this data set suggested an ancient origin for this group of viruses, but here the lack of temporal  
328 signal precludes any interpretation of our estimates of substitution rates and timescales.

329  
330 The *Hepatitis B virus* data set included several human genotypes with complete genomes, where  
331 135 were modern sequences collected from 1963 to 2013 and two were ancient samples from  
332 human mummies from the 16<sup>th</sup> century. Previous studies have not found any temporal signal in  
333 these data using different approaches, despite the inclusion of ancient sequences. Our estimates of  
334 log marginal likelihoods were consistent with a lack of temporal signal, with a log Bayes factor of -  
335 101.51 for  $M_{\text{het}}$  relative to  $M_{\text{iso}}$ .

336  
337 The dog mitochondrial genome data contained samples from up to 36,000 years before the  
338 present. BETS detected temporal signal in these data, with a log Bayes factor of 38.77 for  $M_{\text{het}}$   
339 relative to  $M_{\text{iso}}$ ; this result is consistent with that of a date-randomization test in a previous study  
340 (Tong et al. 2018). The estimated substitution rate for these data using the best model had a mean  
341 of  $1.08 \times 10^{-7}$  subs/site/year (95% HPD:  $7.49 \times 10^{-8}$  to  $1.52 \times 10^{-7}$ ).

## 342 343 **Discussion**

344 We have proposed BETS, a method that explicitly assesses the statistical support for including  
345 sequence sampling times in a Bayesian framework. It is a test of the strength of the temporal signal  
346 in a data set, which is an important prerequisite for obtaining reliable inferences in phylodynamic  
347 analyses. BETS considers the model ensemble, such that the method can detect temporal signal  
348 using models that account for substitution rate variation among lineages. The results of our

349 analyses demonstrate that our method is effective in a range of conditions, including when the  
350 substitution rate is low or when the sampling window represents a small portion of the timespan of  
351 the tree.

352  
353 BETS does not require date permutations, which differentiates it from the widely used date-  
354 randomization test for temporal structure. Date-randomization tests address the question of  
355 whether a particular association between sequences and sampling times produces estimates  
356 different from those obtained from data sets with permuted sampling times (Duchene et al. 2015;  
357 Murray et al. 2015). However, such an approach is not a formal test of temporal signal in the data  
358 because the permutations do not necessarily constitute an appropriate null model. In contrast, our  
359 method does not require permutations and so has the benefit of being robust to using a small  
360 number of sampling times.

361  
362 Accurate calculations of marginal likelihoods are essential for BETS. In our simulation study, we  
363 found that GSS and NS correctly assessed the presence and absence of temporal signal in the data  
364 under most conditions. The correct clock model was also identified, although in a few instances NS  
365 preferred an overparameterized model. Conceivably, using different marginal-likelihood estimators  
366 might affect the actual model selected. Murray et al. (2015) also employed a Bayesian model-  
367 testing approach using the AICM to assess temporal signal. In their study, the AICM performed well  
368 in simulations, but failed to detect temporal signal in empirical data. We attribute this finding to the  
369 low accuracy of AICM relative to path-sampling methods (Baele et al. 2012, 2013), and suggest  
370 careful consideration of the marginal-likelihood estimator for tests of temporal signal.

371  
372 A key advantage of BETS is that the complete model is considered, unlike in simpler data-  
373 exploration methods such as root-to-tip regression. Specifically, root-to-tip regression is a visual  
374 tool for uncovering problems with data quality and to inspect clocklike behaviour, but the absence  
375 of appropriate statistics means that there is no clear way of determining whether the data contains  
376 temporal information. Consider the regressions in figure 5 for data with a high substitution rate and  
377 narrow sampling window. Even when among-lineage rate variation is low ( $\sigma=0.1$ ), the data points  
378 form a cloud, with a low  $R^2$  of 0.09. However, the apparent 'noise' around the regression line is  
379 probably the result of stochasticity in sequence evolution and of the narrow sampling window  
380 relative to the age of the root of the tree. In fact, for this particular data set the model with the  
381 highest log marginal likelihood is the strict clock with correct sampling times.

382  
383 In all of our analyses, we ensured that the priors for different models and configurations of sampling  
384 times were identical because, as with all Bayesian analyses, model comparison using marginal  
385 likelihoods can depend on the choice of prior (Oaks et al. 2019). For example, the tree prior can  
386 affect inferences of temporal signal, as it is part of the full model specification. Here we used an  
387 exponential-growth coalescent tree prior, which closely matches the demographic dynamics of the  
388 birth-death process under which the data were simulated. The effect of using an inappropriate tree  
389 prior on tests of temporal signal requires further investigation, but previous studies have suggested  
390 that there is only a small impact on estimates of rates and times if the sequence data are  
391 informative (Ritchie et al. 2017; Möller et al. 2018).

392  
393 An interesting finding is that statistical support for isochronous sampling times in truly isochronous  
394 data is lower than that for the correct sampling times in truly heterochronous data. This can  
395 potentially lead to an increased risk of incorrectly concluding the presence of temporal signal, but  
396 we only found this to be a problem in a small number of cases. In particular, in isochronous data  
397 simulated with a low substitution rate, and with very few variable sites, the best models were  
398 sometimes those that included sampling times, albeit with very low log Bayes factors (e.g.,



399 supplementary fig. S2, Supplementary Material online). This probably occurs because stochastic  
400 error associated with a small amount of evolution leads to low power for model selection.

401

402 Permuting sampling times led to poor model fit, as expected. This procedure requires substantial  
403 computing requirements, depending on the number of permutations that are performed, and we  
404 find that such date permutations are of limited value for model testing when the data are highly  
405 informative (e.g., figs. 1 and 3). However, in data sets with very low information content, such as  
406 those that were produced by simulation with a low substitution rate here, conducting a small  
407 number of date permutations might offer a conservative approach to determining whether model  
408 fit and parameter estimates are driven by a particular set of sampling times, as one would expect in  
409 the presence of temporal signal.

410

411 The nature of the BETS approach means that every parameter in the model has a prior probability,  
412 including the substitution rate. Because substitution rates and times are nonidentifiable, it is  
413 conceivable that an informative prior on the rate or on the age of an internal node might have a  
414 stronger effect than the sampling times on the posterior, for example if the samples span a very  
415 short window of time. Such analyses with informative substitution rate priors effectively include  
416 several simultaneous sources of calibrating information (i.e., sampling times, internal nodes, and an  
417 informative rate prior). Using sampling times in addition to other sources of calibration information  
418 might still be warranted if it improves the fit of the model, which can be tested using our proposed  
419 method.

420

421 Analyses with multiple calibrations can also allow uncertainty in sequence sampling times,  
422 especially in data sets that include ancient DNA, where sampling times can be treated as  
423 parameters in the model (Shapiro et al. 2011). BETS provides a coherent approach to assess  
424 temporal structure in these circumstances, unlike date-randomization tests that typically use point  
425 values for sampling times. In fact, BETS can be used as a means to validate whether a sample is  
426 modern or ancient.

427

428 In general, the uptake of Bayesian model testing in phylogenetics has great potential for improving  
429 our confidence in estimates of substitution rates and timescales. The test that we have proposed  
430 here, BETS, provides a coherent and intuitive framework to test for temporal information in the  
431 data.

432

## 433 **Materials and Methods**

434

### 434 **Simulations**

435 We simulated phylogenetic trees under a stochastic birth-death process using MASTER v6.1  
436 (Vaughan and Drummond 2013), by specifying birth rate  $\lambda=1.5$ , death rate  $\mu=0.5$ , and sampling  
437 rate  $\phi=0.5$ . This corresponds to an exponentially growing infectious outbreak with reproductive  
438 number  $R_0=1.5$  and a wide sampling window. We set the simulation time to 5 units of time, which  
439 corresponds to the time of origin of the process. For isochronous trees, we used similar settings, but  
440 instead of using the sampling rate, we sampled each tip with probability  $\rho=0.5$  when the process  
441 was stopped after 5 units of time (i.e.  $\mu=1$  and  $\phi=0$ ). Some of our analyses consisted of artificially  
442 specifying sampling times for isochronous trees, which we set to those that we would have  
443 obtained from a birth-death process with  $\mu=0.5$  and  $\phi=0.5$ .

444

445 In a second set of simulations of heterochronous trees, we generated trees with a narrow sampling  
446 window. We specified two intervals for  $\mu$  and  $\phi$ . The first interval spanned 4.5 units of time with  
447  $\mu=1.0$  and  $\phi=0$ , and the second interval 0.5 units of time with  $\mu=0.1$  and  $\phi=0.9$ . As a result, the  
448 process still had a constant become uninfected rate ( $\mu + \phi$ ), but samples were only collected in

449 the second interval. The high sampling rate in the second interval resulted in trees with similar  
450 numbers of tips to those with a wide sampling window, but where their ages only spanned 0.5 units  
451 of time.

452

453 We only considered the simulated trees that contained between 90 and 110 tips. The trees  
454 generated in MASTER are chronograms (with branch lengths in units of time), so we simulated  
455 substitution rates to generate phylograms (with branch lengths in units of subs/site). To do this we  
456 specified the uncorrelated lognormal relaxed clock with a mean rate of  $5 \times 10^{-3}$  or  $10^{-5}$  subs/site/unit  
457 time and a standard deviation  $\sigma$  of 0 (corresponding to a strict clock), 0.1, 0.5, or 1. We simulated  
458 sequence evolution along these phylograms under the HKY nucleotide substitution model  
459 (Hasegawa et al. 1985). We added among-site rate variation using a discretized gamma distribution  
460 (Yang 1994, 1996) using Phangorn v2.5 (Schliep 2011) to generate sequence alignments of 10,000  
461 nucleotides. We set the transition-to-transversion ratio of the HKY model to 10 and the shape of the  
462 gamma distribution to 1, which is similar to estimates of these parameters in influenza viruses  
463 (Duchene et al. 2014; Hedge and Wilson 2014). For each simulation scenario we generated 10  
464 sequence alignments.

465

#### 466 **Estimation of Marginal Likelihoods Using Nested Sampling**

467 We analysed the data in BEAST 2.5 using the matching substitution model, the exponential-growth  
468 coalescent tree prior, the strict clock or relaxed clock, and different configurations of sampling  
469 times. We chose the exponential-growth coalescent tree prior, instead of the birth-death tree prior,  
470 because it is conditioned on the samples instead of assuming a sampling process; this ensures that  
471 the marginal likelihoods for isochronous and heterochronous trees are comparable.

472

473 We specified proper priors on all parameters, which is essential for accurate estimation of marginal  
474 likelihoods (Baele et al., 2013). In our heterochronous analyses the prior on the substitution rate had  
475 a uniform distribution bounded between 0 and 1. We made this arbitrary choice to set a somewhat  
476 uninformative prior and because the default prior in BEAST 2.5 is a uniform distribution between 0  
477 and infinity, which is improper. Owing to the nonidentifiability of substitution rates and times,  
478 neither can be inferred in the absence of calibrating information, so in our isochronous analyses we  
479 fixed the value of the substitution rate to 1. The initial NS chain length was chosen so as to draw  
480 20,000 samples, with 20,000 steps, 32 particles, and a subchain length of 5,000 (note that NS is not  
481 equivalent to standard MCMC, nor is the definition of an iteration/step). The chain length and its  
482 accompanying sampling frequency were adjusted to obtain effective sample sizes for key  
483 parameters of at least 200 (computed in the NS output in BEAST 2.5). Examples of MASTER files  
484 and BEAST input files for NS are available online (supplementary data, Supplementary Material  
485 online).

486

#### 487 **Estimation of Marginal Likelihoods Using Generalized Stepping-Stone Sampling**

488 We used BEAST 1.10 with the same model specifications and priors as in BEAST2, except for the  
489 prior on the substitution rate, for which we used the approximate continuous-time Markov chain  
490 reference prior (Ferreira and Suchard 2008). Because our simulation analyses of GSS and NS differ  
491 in this prior, the marginal-likelihood estimates are not directly comparable, so for each simulation  
492 we report log Bayes factors of competing models instead of the log marginal likelihoods. The GSS  
493 implementation in BEAST 1.10 has two different working priors for the tree generative process: a  
494 matching tree prior and a product of exponentials. The latter approach is the most generally  
495 applicable and is the one that we used here (Baele et al. 2016).

496

497 We used an initial MCMC chain length of  $5 \times 10^7$  steps sampling every 5000 steps. After discarding  
498 10% of the samples obtained, the remaining samples were used to construct the working  
499 distributions for the GSS analysis. These comprised 100 path steps distributed according to

500 quantiles from a  $\beta$  distribution with  $\alpha=0.3$ , with each of the 101 resulting power posterior  
501 inferences running for  $5 \times 10^5$  iterations. We assessed sufficient sampling for the initial MCMC  
502 analysis by verifying that the effective sample sizes for key parameters were at least 200 in Coda  
503 v0.19 (Plummer et al. 2006). If this condition was not met, we doubled the length of the MCMC and  
504 reduced sampling frequency accordingly. Examples of MASTER files and BEAST input files for GSS  
505 are available online (supplementary data, Supplementary Material online).

506

### 507 **Analyses of Empirical Data Sets**

508 We downloaded sequence alignments from their original publications (Table 1): complete genomes  
509 of *A/H1N1 influenza virus* (Hedge et al. 2013), whole genome sequences of *B. pertussis* (Bart et al.  
510 2014; Duchene et al. 2016), RdRP sequences of coronaviruses (Wertheim et al. 2013), complete  
511 genomes of *Hepatitis B virus* (Patterson Ross et al. 2018), and dog mitochondrial genomes  
512 (Thalmann et al. 2013). The data and BEAST input files are available in the Supplementary Material  
513 online.

514

515 Briefly, we used similar settings as in our simulations to estimate marginal likelihoods using GSS.  
516 For sequence sampling times we considered the correct sampling times, no sampling times (i.e.,  
517 isochronous), and permuted sampling times. We also specified tree priors as follows: an  
518 exponential-growth coalescent for the *A/H1N2 influenza virus*, *Bordetella pertussis*, coronaviruses,  
519 and *Hepatitis B virus* data sets, and a constant-size coalescent for the dog mitochondrial genomes  
520 as used by Tong et al. (2018). We again chose the HKY+ $\Gamma$  substitution model, except in the analysis  
521 of *Hepatitis B virus* data, for which we used the GTR+ $\Gamma$  model (Tavaré 1986), and in the analysis of  
522 the dog data set for which we used the SRDo6 substitution model (Shapiro et al. 2006) for coding  
523 regions and the GTR+ $\Gamma$  for noncoding regions.

524

### 525 **Supplementary Material**

526 Supplementary data are available online.

527

### 528 **Funding**

529 SD was supported by an Australian Research Council Discovery Early Career Researcher Award  
530 (DE190100805) and an Australian National Health and Medical Research Council grant  
531 (APP1157586). PL acknowledges funding from the European Research Council under the European  
532 Union's Horizon 2020 research and innovation programme (grant agreement no. 725422-  
533 ReservoirDOCS) and the Research Foundation -- Flanders ('Fonds voor Wetenschappelijk  
534 Onderzoek -- Vlaanderen', Go66215N, GoD5117N and GoB9317N). SYWH was funded by the  
535 Australian Research Council (FT160100167). GB acknowledges support from the Interne Fondsen  
536 KU Leuven / Internal Funds KU Leuven under grant agreement C14/18/094. VD was supported by  
537 contract HHSN272201400006C from the National Institute of Allergy and Infectious Diseases,  
538 National Institutes of Health, U.S. Department of Health and Human Services, USA.

539

### 540 **Acknowledgements**

541 Pending.

542

### 543 **References**

544 Baele G., Lemey P. 2014. Bayesian model selection in phylogenetics and genealogy-based  
545 population genetics. In: Chen M.-H., Kuo L., Lewis P.O., editors. Bayesian Phylogenetics,  
546 Methods, Algorithms, and Applications. Boca Raton, Florida: CPC Press. p. 59–93.  
547 Baele G., Lemey P., Bedford T., Rambaut A., Suchard M.A., Alekseyenko A. V. 2012. Improving the  
548 accuracy of demographic and molecular clock model comparison while accommodating

- 549 phylogenetic uncertainty. *Mol. Biol. Evol.* 29:2157–2167.
- 550 Baele G., Lemey P., Suchard M.A. 2016. Genealogical working distributions for Bayesian model  
551 testing with phylogenetic uncertainty. *Syst. Biol.* 65:250–264.
- 552 Baele G., Li W.L.S., Drummond A.J., Suchard M.A., Lemey P. 2013. Accurate model selection of  
553 relaxed molecular clocks in bayesian phylogenetics. *Mol. Biol. Evol.* 30:239–243.
- 554 Bart M.J., Harris S.R., Advani A., Arakawa Y., Bottero D., Bouchez V., Cassiday P.K., Chiang C.-S.,  
555 Dalby T., Fry N.K. 2014. Global population structure and evolution of *Bordetella pertussis* and  
556 their relationship with vaccination. *MBio.* 5:e01074-14.
- 557 Biek R., Pybus O.G., Lloyd-Smith J.O., Didelot X. 2015. Measurably evolving pathogens in the  
558 genomic era. *Trends Ecol. Evol.* 30:306–313.
- 559 Bouckaert R., Vaughan T.G., Barido-Sottani J., Duchene S., Fourment M., Gavryushkina A., Heled J.,  
560 Jones G., Kuhnert D., de Maio N. 2019. BEAST 2.5: An Advanced Software Platform for  
561 Bayesian Evolutionary Analysis. *PLOS Comput. Biol.* 15:e1006650.
- 562 Bromham L., Duchêne S., Hua X., Ritchie A.M., Duchêne D.A., Ho S.Y.W. 2018. Bayesian molecular  
563 dating: Opening up the black box. *Biol. Rev.* 93:1165–1191.
- 564 Drummond A.J., Ho S.Y.W., Phillips M.J., Rambaut A. 2006. Relaxed phylogenetics and dating with  
565 confidence. *PLOS Biol.* 4:e88.
- 566 Drummond A.J., Pybus O.G., Rambaut A., Forsberg R., Rodrigo A.G. 2003. Measurably evolving  
567 populations. *Trends Ecol. Evol.* 18:481–488.
- 568 Duchene S., Duchene D.A., Geoghegan J.L., Dyson Z.A., Hawkey J., Holt K.E. 2018. Inferring  
569 demographic parameters in bacterial genomic data using Bayesian and hybrid phylogenetic  
570 methods. *BMC Evol. Biol.* 18:95.
- 571 Duchene S., Duchene D.A., Holmes E.C., Ho S.Y.W. 2015. The performance of the date-  
572 randomization test in phylogenetic analyses of time-structured virus data. *Mol. Biol. Evol.*  
573 32:1895–1906.
- 574 Duchene S., Holmes E.C., Ho S.Y.W. 2014. Analyses of evolutionary dynamics in viruses are  
575 hindered by a time-dependent bias in rate estimates. *Proc. R. Soc. London B.* 281:20140732.
- 576 Duchene S., Holt K.E., Weill F.-X., Le Hello S., Hawkey J., Edwards D.J., Fourment M., Holmes E.C.  
577 2016. Genome-scale rates of evolutionary change in bacteria. *Microb. Genomics.* 2:e000094.
- 578 Fan Y., Wu R., Chen M.-H., Kuo L., Lewis P.O. 2011. Choosing among partition models in Bayesian  
579 phylogenetics. *Mol. Biol. Evol.* 28:523–532.
- 580 Ferreira M.A.R., Suchard M.A. 2008. Bayesian analysis of elapsed times in continuous-time Markov  
581 chains. *Can. J. Stat.* 36:355–368.
- 582 Fourment M., Magee A.F., Whidden C., Bilge A., Matsen I. V, Frederick A., Minin V.N. 2019. 19  
583 dubious ways to compute the marginal likelihood of a phylogenetic tree topology. *Syst.*  
584 *Biol.*:syzo46.
- 585 Guindon S., Dufayard J.-F., Lefort V., Anisimova M., Hordijk W., Gascuel O. 2010. New algorithms  
586 and methods to estimate maximum likelihood phylogenies: assessing the performance of  
587 PhyML 3.0. *Syst. Biol.* 59:307–321.
- 588 Hasegawa M., Kishino H., Yano T. 1985. Dating of the human-ape splitting by a molecular clock of  
589 mitochondrial DNA. *J. Mol. Evol.* 22:160–174.
- 590 Hedge J., Lycett S.J., Rambaut A. 2013. Real-time characterization of the molecular epidemiology  
591 of an influenza pandemic. *Biol. Lett.* 9:20130331.
- 592 Hedge J., Wilson D.J. 2014. Bacterial phylogenetic reconstruction from whole genomes is robust to  
593 recombination but demographic inference is not. *MBio.* 5:e02158-14.
- 594 Hipsley C.A., Müller J. 2014. Beyond fossil calibrations: realities of molecular clock practices in  
595 evolutionary biology. *Front. Genet.* 5:138.
- 596 Ho S.Y.W., Duchene S. 2014. Molecular-clock methods for estimating evolutionary rates and time  
597 scales. *Mol. Ecol.* 23:5947–5975.
- 598 Kass R.E., Raftery A.E. 1995. Bayes factors. *J. Am. Stat. Assoc.* 90:773–795.
- 599 Korber B., Muldoon M., Theiler J., Gao F., Gupta R., Lapedes A., Hahn B.H., Wolinsky S.,

- 600           Bhattacharya T. 2000. Timing the ancestor of the HIV-1 pandemic strains. *Science*. 288:1789–  
601           1796.
- 602           Lartillot N., Philippe H. 2006. Computing Bayes factors using thermodynamic integration. *Syst.*  
603           *Biol.* 55:195–207.
- 604           Maturana P., Brewer B.J., Klaere S., Bouckaert R. 2019. Model selection and parameter inference in  
605           phylogenetics using Nested Sampling. *Syst. Biol.* 68:219–233.
- 606           Möller S., du Plessis L., Stadler T. 2018. Impact of the tree prior on estimating clock rates during  
607           epidemic outbreaks. *Proc. Natl. Acad. Sci. USA*. 115:4200–4205.
- 608           Murray G.G.R., Wang F., Harrison E.M., Paterson G.K., Mather A.E., Harris S.R., Holmes M.A.,  
609           Rambaut A., Welch J.J. 2015. The effect of genetic structure on molecular dating and tests for  
610           temporal signal. *Methods Ecol. Evol.* 7:80–89.
- 611           Nascimento F.F., dos Reis M., Yang Z. 2017. A biologist’s guide to Bayesian phylogenetic analysis.  
612           *Nat. Ecol. Evol.* 1:1446.
- 613           Newton M.A., Raftery A.E. 1994. Approximate Bayesian inference with the weighted likelihood  
614           bootstrap. *J. R. Stat. Soc. Ser. B*. 56:3–26.
- 615           Oaks J.R., Cobb K.A., Minin V.N., Leaché A.D. 2019. Marginal likelihoods in phylogenetics: a review  
616           of methods and applications. *Syst. Biol.* 68:681–697.
- 617           Patterson Ross Z., Klunk J., Fornaciari G., Giuffra V., Duchene S., Duggan A.T., Poinar D., Douglas  
618           M.W., Eden J.-S., Holmes E.C., Poinar H.N. 2018. The paradox of HBV evolution as revealed  
619           from a 16th century mummy. *PLoS Pathog.* 14:e1006887.
- 620           Plummer M., Best N., Cowles K., Vines K. 2006. CODA: Convergence diagnosis and output analysis  
621           for MCMC. *R News*. 6:7–11.
- 622           Raftery A., Newton M., Satagopan J., Krivitsky P. 2007. Estimating the integrated likelihood via  
623           posterior simulation using the harmonic mean identity. *Bayesian Statistics*. Oxford University  
624           Press. p. 1–45.
- 625           Rambaut A., Lam T.T., Carvalho L.M., Pybus O.G. 2016. Exploring the temporal structure of  
626           heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol.* 2:vev007.
- 627           Ramsden C., Holmes E.C.C., Charleston M.A.A. 2009. Hantavirus evolution in relation to its rodent  
628           and insectivore hosts: no evidence for codivergence. *Mol. Biol. Evol.* 26:143–153.
- 629           Rieux A., Balloux F. 2016. Inferences from tip-calibrated phylogenies: a review and a practical guide.  
630           *Mol. Ecol.* 25:1911–1924.
- 631           Ritchie A.M., Lo N., Ho S.Y.W. 2017. The impact of the tree prior on molecular dating of data sets  
632           containing a mixture of inter- and intraspecies sampling. *Syst. Biol.* 66:413–425.
- 633           Schliep K.P. 2011. phangorn: Phylogenetic analysis in R. *Bioinformatics*. 27:592–593.
- 634           Shapiro B., Ho S.Y.W., Drummond A.J., Suchard M.A., Pybus O.G., Rambaut A. 2011. A Bayesian  
635           phylogenetic method to estimate unknown sequence ages. *Mol. Biol. Evol.* 28:879–887.
- 636           Shapiro B., Rambaut A., Drummond A.J. 2006. Choosing appropriate substitution models for the  
637           phylogenetic analysis of protein-coding sequences. *Mol. Biol. Evol.* 23:7–9.
- 638           Skilling J. 2006. Nested sampling for general Bayesian computation. *Bayesian Anal.* 1:833–859.
- 639           Suchard M.A., Lemey P., Baele G., Ayres D.L., Drummond A.J., Rambaut A. 2018. Bayesian  
640           phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evol.* 4:vey016.
- 641           Tavaré S. 1986. Some probabilistic and statistical problems in the analysis of DNA sequences. *Lect.*  
642           *Math. Life Sci.* 17:57–86.
- 643           Thalmann O., Shapiro B., Cui P., Schuenemann V.J., Sawyer S.K., Greenfield D.L., Germonpré M.B.,  
644           Sablin M. V., López-Giráldez F., Domingo-Roura X. 2013. Complete mitochondrial genomes of  
645           ancient canids suggest a European origin of domestic dogs. *Science*. 342:871–874.
- 646           Tong K.J., Duchêne D.A., Duchêne S., Geoghegan J.L., Ho S.Y.W. 2018. A comparison of methods  
647           for estimating substitution rates from ancient DNA sequence data. *BMC Evol. Biol.* 18:70.
- 648           Trovão N.S., Baele G., Vrancken B., Bielejec F., Suchard M.A., Fargette D., Lemey P. 2015. Host  
649           ecology determines the dispersal patterns of a plant virus. *Virus Evol.* 1:vev016.
- 650           Vaughan T.G., Drummond A.J. 2013. A stochastic simulator of birth–death master equations with

- 651 application to phylodynamics. *Mol. Biol. Evol.* 30:1480–1493.
- 652 Wertheim J.O., Chu D.K.W., Peiris J.S.M., Pond S.L.K., Poon L.L.M. 2013. A case for the ancient  
653 origin of coronaviruses. *J. Virol.* 87:7039–7045.
- 654 Xie W., Lewis P.O., Fan Y., Kuo L., Chen M.H. 2011. Improving marginal likelihood estimation for  
655 Bayesian phylogenetic model selection. *Syst. Biol.* 60:150–160.
- 656 Yang Z. 1994. Estimating the pattern of nucleotide substitution. *J. Mol. Evol.* 39:105–111.
- 657 Yang Z. 1996. Among-site rate variation and its impact on phylogenetic analyses. *Trends Ecol. Evol.*  
658 11:367–372.
- 659 Zuckerkandl E., Pauling L. 1965. Evolutionary divergence and convergence in proteins. In: Bryson  
660 V., Vogel H., editors. *Evolving Genes and Proteins*. New York: Academic press. p. 97–166.
- 661

## 662 Figure Legends

663 **FIG. 1.** Log Bayes factors of heterochronous data simulated with a high substitution rate and wide  
664 sampling window. Each panel shows the results for data sets simulated with a different degree of  
665 among-lineage rate variation, governed by the standard deviation  $\sigma$  of a lognormal distribution. In  
666 each plot the x-axis depicts six analysis settings, with two clock models, strict clock (SC) and the  
667 uncorrelated relaxed clock with an underlying lognormal distribution (UCLN), and three settings for  
668 sampling times: generated under the birth-death process (BD), identical sampling times  
669 (Isochronous), and permuted (Permuted). The points have been jittered along the x-axis to facilitate  
670 visualization. The y-axis shows log Bayes factors relative to the best model. Red points correspond  
671 to estimates using generalized stepping-stone sampling and blue points correspond to estimates  
672 using nested sampling. We conducted 10 simulation replicates, with each replicate data set  
673 analysed under the six analysis settings and two marginal-likelihood estimators, such that  
674 stochastic error might cause differences in the preferred model. The number next to each cloud of  
675 points denotes the number of times (out of 10) that the corresponding model had the highest  
676 marginal likelihood with generalized stepping-stone sampling (red) and nested sampling (blue).  
677

678 **FIG. 2.** Log Bayes factors of heterochronous data simulated under a low substitution rate and a wide  
679 sampling window. Symbols and colours are the same as those in figure 1.

680  
681 **FIG. 3.** Log Bayes factors of heterochronous data simulated under a high substitution rate and  
682 narrow sampling window. Symbols and colours are the same as those in figure 1.

683  
684 **FIG. 4.** Log Bayes factors of heterochronous data simulated under a low substitution rate and  
685 narrow sampling window. Symbols and colours are the same as those in figure 1.

686  
687 **FIG. 5.** Root-to-tip regressions for a subset of data sets simulated with varying degrees of among-  
688 lineage rate variation (governed by the standard deviation  $\sigma$  of a lognormal distribution), using a  
689 high substitution rate and either a wide or narrow sampling window. The y-axis is the root-to-tip  
690 distance and the x-axis is the time from the youngest tip, where 0 is the present. Each point  
691 corresponds to a tip in the tree and the solid line is the best-fit linear regression using least-squares.  
692 The coefficient of determination,  $R^2$ , is shown in each case. For comparison, the log Bayes factors of  
693 the best heterochronous model relative the best isochronous model,  $BF(M_{\text{het}} - M_{\text{iso}})$ , are also shown.  
694

695 **FIG. 6.** Log marginal likelihoods estimated using generalized stepping-stone sampling for six  
696 analysis settings for sequence data from rapidly evolving pathogens, *A/H1N1 Human influenza virus*  
697 and *Bordetella pertussis*. The y-axis is the marginal likelihood and the x-axis shows the analysis  
698 settings, with two clock models, strict clock (SC) and the uncorrelated relaxed clock with an  
699 underlying lognormal distribution (UCLN), and three settings for sampling times: generated under  
700 the birth-death process (BD), identical sampling times (Isochronous), and permuted (Permuted).  
701 Solid points and dashed lines correspond to the log marginal-likelihood estimates. The asterisk  
702 denotes the model with the highest marginal likelihood.  
703

704 **FIG. 7.** Log marginal likelihoods estimated using generalized stepping-stone sampling for six  
705 analysis settings for data sets with ancient DNA or highly divergent sequences. The y-axis is the  
706 marginal likelihood and the x-axis shows the analysis settings, with two clock models, strict clock  
707 (SC) and the uncorrelated relaxed clock with an underlying lognormal distribution (UCLN), and  
708 three settings for sampling times: generated under the birth-death process (BD), identical sampling  
709 times (Isochronous), and randomized (Random). Solid points and dashed lines correspond to the log  
710 marginal-likelihood estimates. The asterisk denotes the model with the highest marginal likelihood.  
711

712

713  
714  
715  
716  
717  
718  
719  
720

## Tables

**Table 1.** Details of empirical data sets used in this study.

Data set	Number of sites (nucleotides)	Number of samples	Sampling time range	Reference
<i>A/H1N1 influenza virus</i>	13,154	329	10 months (March to December 2009)	Hedge et al. (2013)
<i>Bordetella pertussis</i>	$4.9 \times 10^6$	150	89 years (1920 to 2009)	Bart et al. (2014)
Coronaviruses	1,860	43	70 years (1941 to 2011)	Wertheim et al. (2013)
<i>Hepatitis B virus</i>	3,271	137	445 years (2103 to 1568)	Patterson Ross et al. (2018)
Dog mtDNA	14,596	50	36,000 years (to the present)	Thalmann et al. (2013)

721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733

## Supplementary Material

734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751

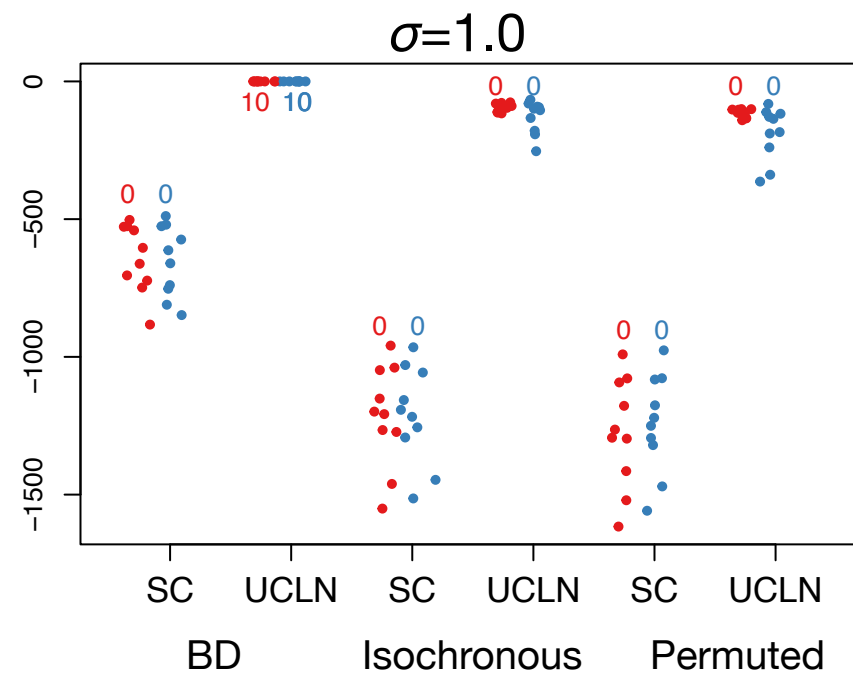
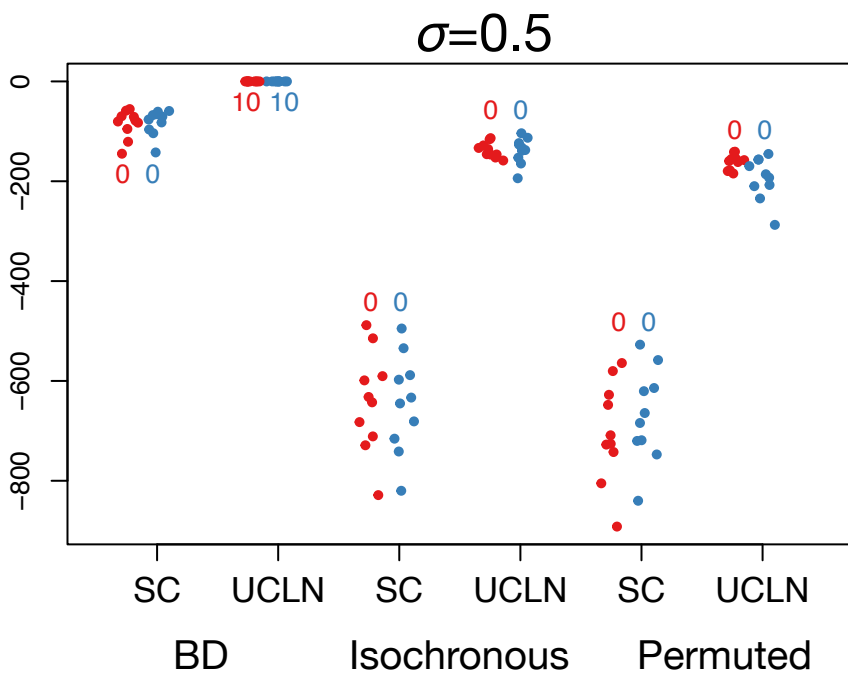
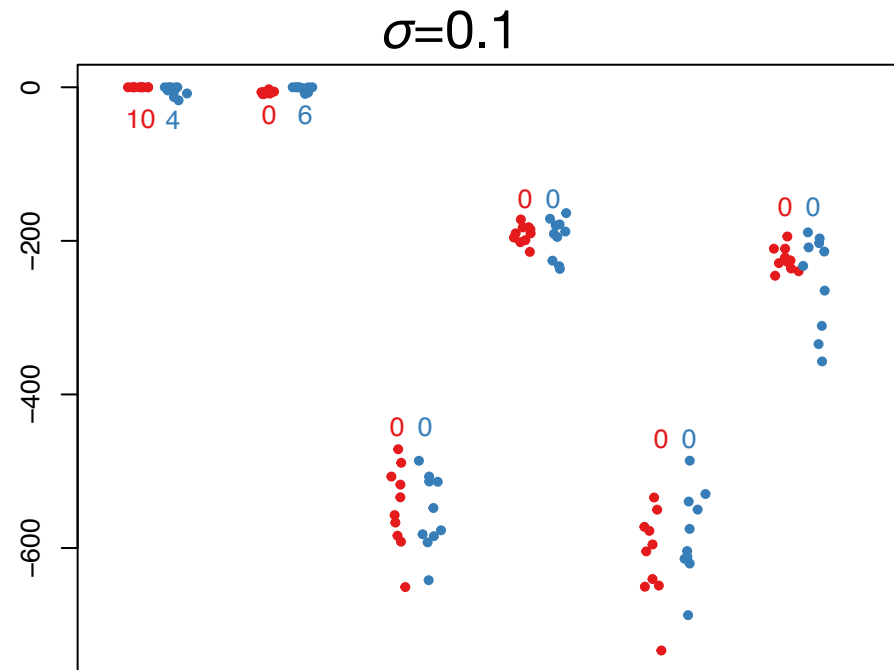
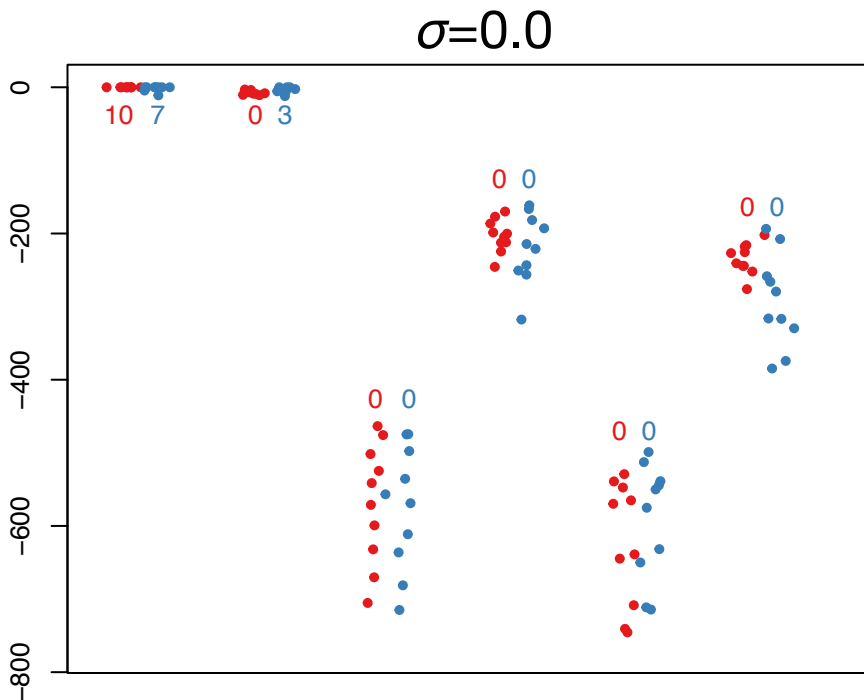
**FIG. S1.** Log Bayes factors of isochronous data simulated with a high substitution rate. Each panel shows the results for data sets simulated with a different degree of among-lineage rate variation, governed by the standard deviation  $\sigma$  of a lognormal distribution. The x-axis depicts six analysis settings, with two molecular clock models, strict clock (SC) and the uncorrelated relaxed clock with an underlying lognormal distribution (UCLN), and three settings for sampling times: generated under the birth-death process (BD), identical sampling times (Isochronous), and permuted (Permuted). The points have been jittered to facilitate visualization. The y-axis shows log Bayes factors relative to the best model. Red points correspond to estimates using generalized stepping-stone sampling and blue points correspond to estimates using nested sampling. We conducted 10 simulation replicates, with each replicate data set analysed under the six analysis settings and two marginal-likelihood estimators, such that stochastic error might cause differences in the preferred model. The number next to each cloud of points denotes the number of times (out of 10) that the corresponding model had the highest marginal likelihood with generalized stepping-stone sampling (red) and nested sampling (blue).

**FIG. S2.** Log Bayes factors of isochronous data simulated with a low substitution rate. Symbols and colours are the same as those in figure 1.



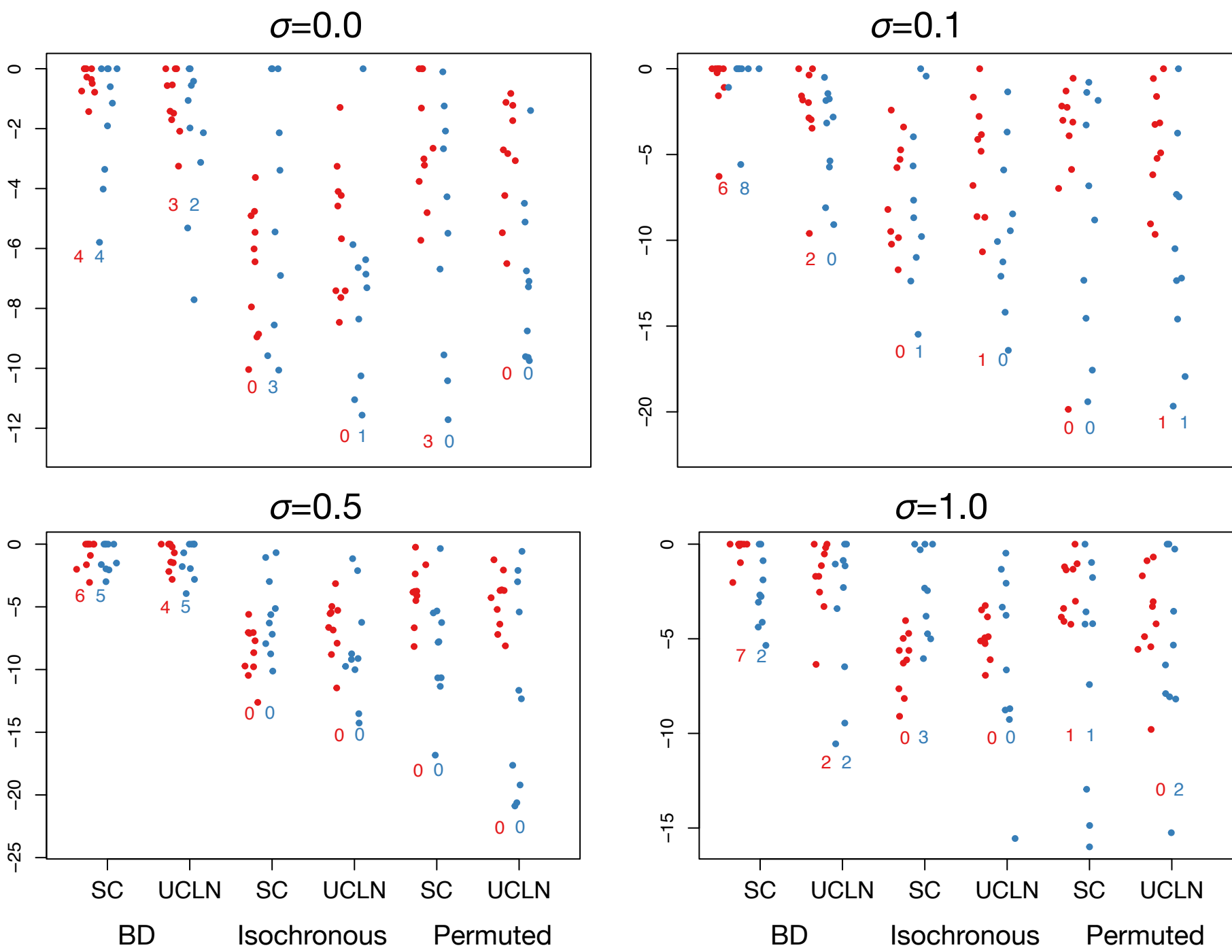
752 **FIG. S3.** Example of three phylogenetic trees used in simulations. Red dashed lines indicate the  
753 times of each of the tips and therefore represent the sampling process over time. All trees are  
754 simulated under a birth-death process with time of origin of 5, such that the sum of the tree height  
755 and the length of the stem branch leading to the root is always 5. In all trees, we set the birth rate  
756  $\lambda=1.5$ , and become uninfected rate  $\delta=1$ , where  $\delta=\mu+\phi$ , where  $\mu$  is the death rate and  $\phi$  is  
757 the sampling rate upon death. Thus, the population growth rate is constant and the same across all  
758 trees. The top tree assumes a constant sampling process and a wide sampling window ( $\phi=0.5$   
759 throughout the whole process), whereas in the second tree sampling starts after 4.5. Before this  
760 time the sampling rate,  $\phi_0$ , is zero. After 4.5 time units the sampling rate  $\phi_1$  is 0.9 (and thus  $\mu_1$   
761 = 0.1), resulting in a narrow sampling window. The bottom tree has samples drawn at a single point  
762 in time with a sampling probability at present,  $\rho$ , of 0.5 (and thus  $\phi=0$ ).

Log Bayes factors relative to best model



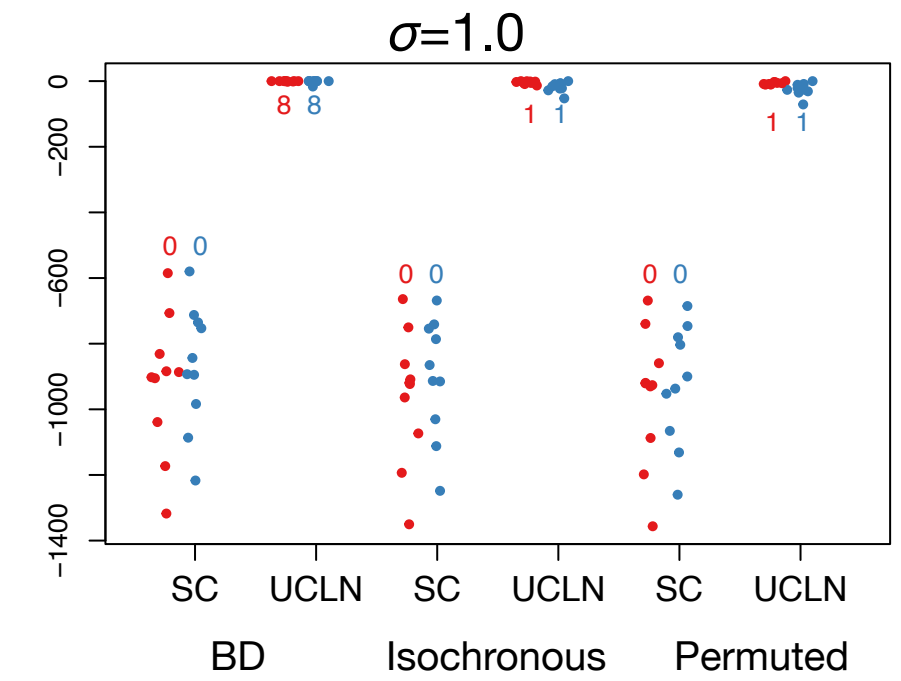
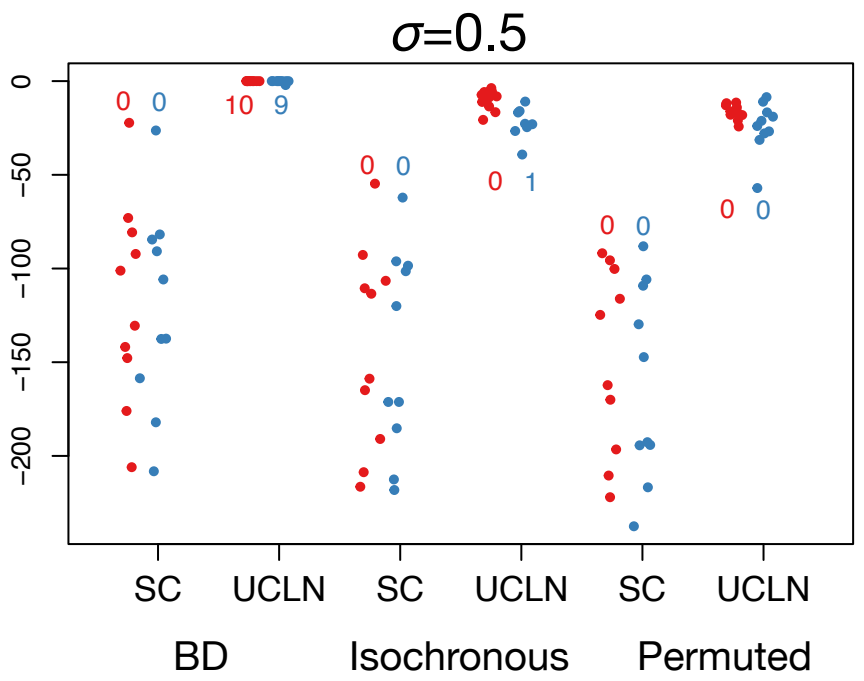
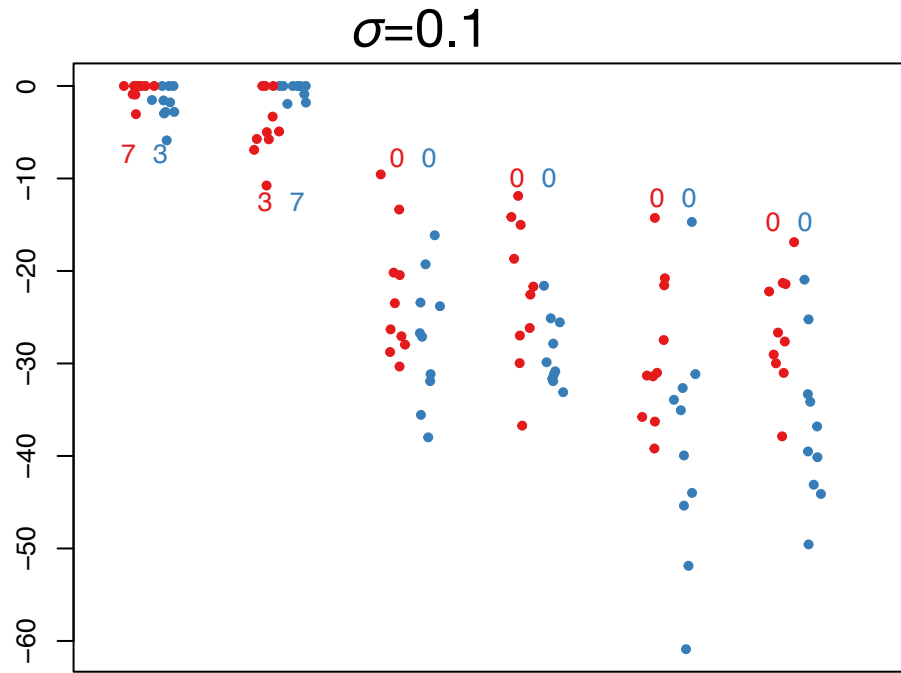
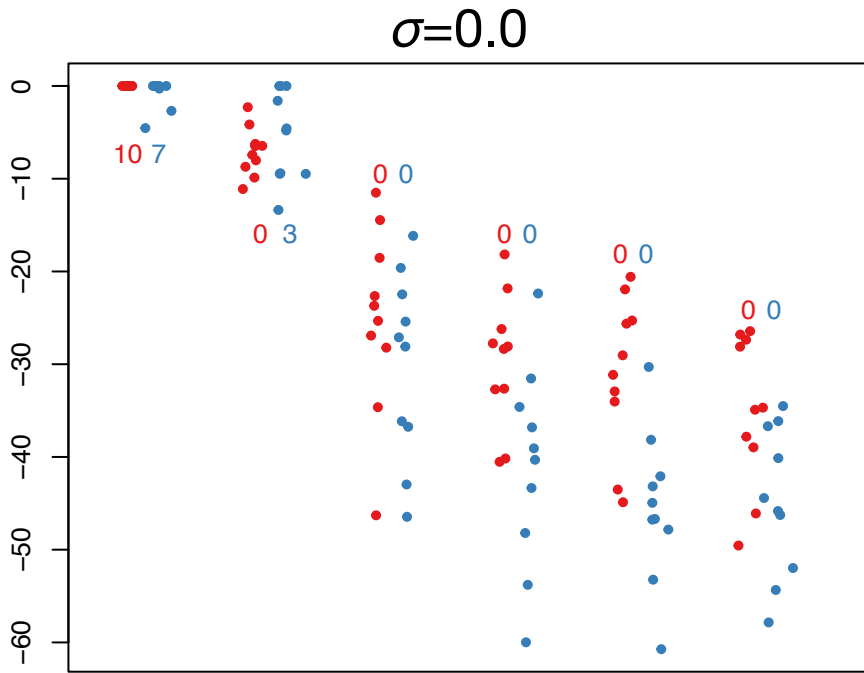
Analysis settings

Log Bayes factors relative to best model



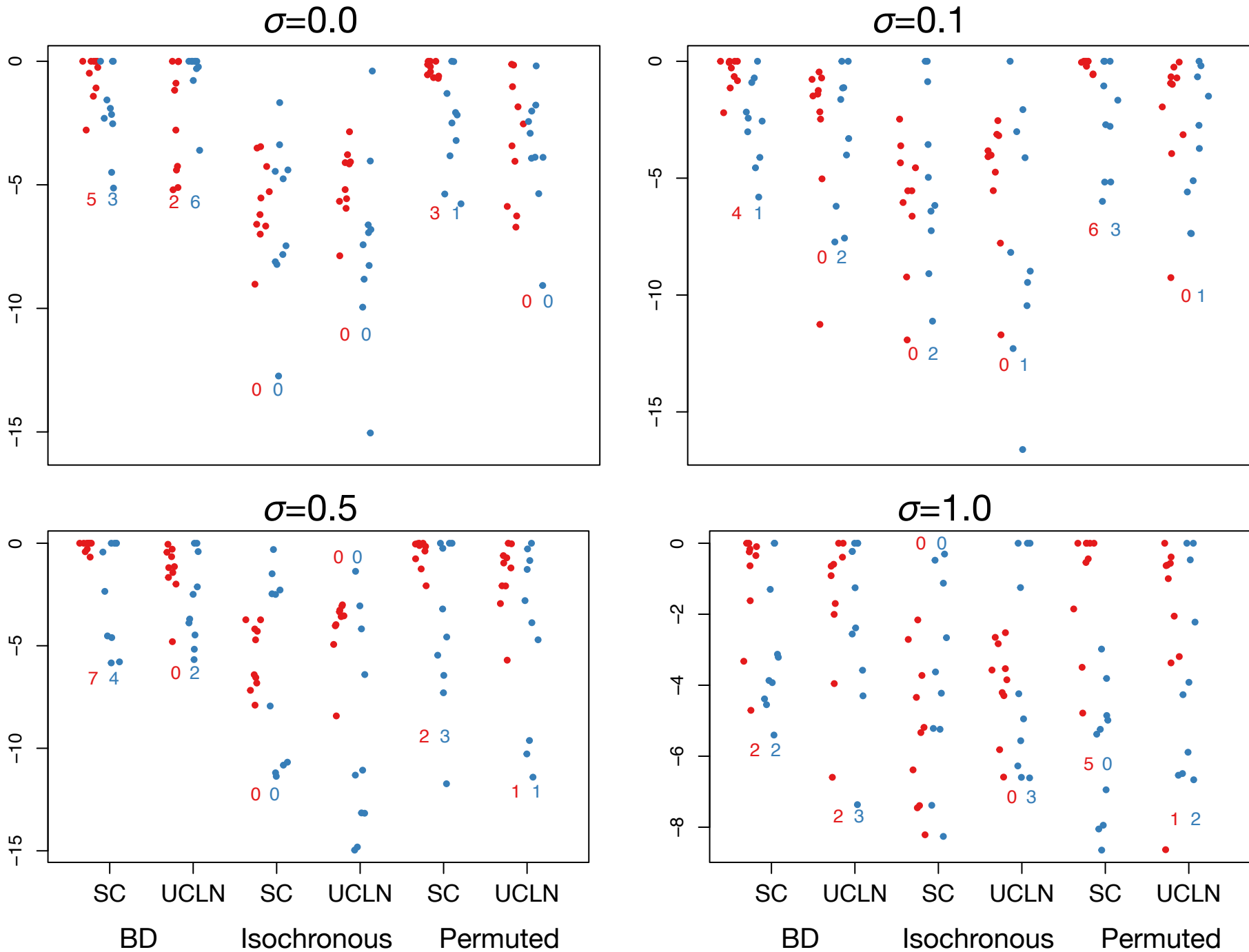
Analysis settings

Log Bayes factors relative to best model

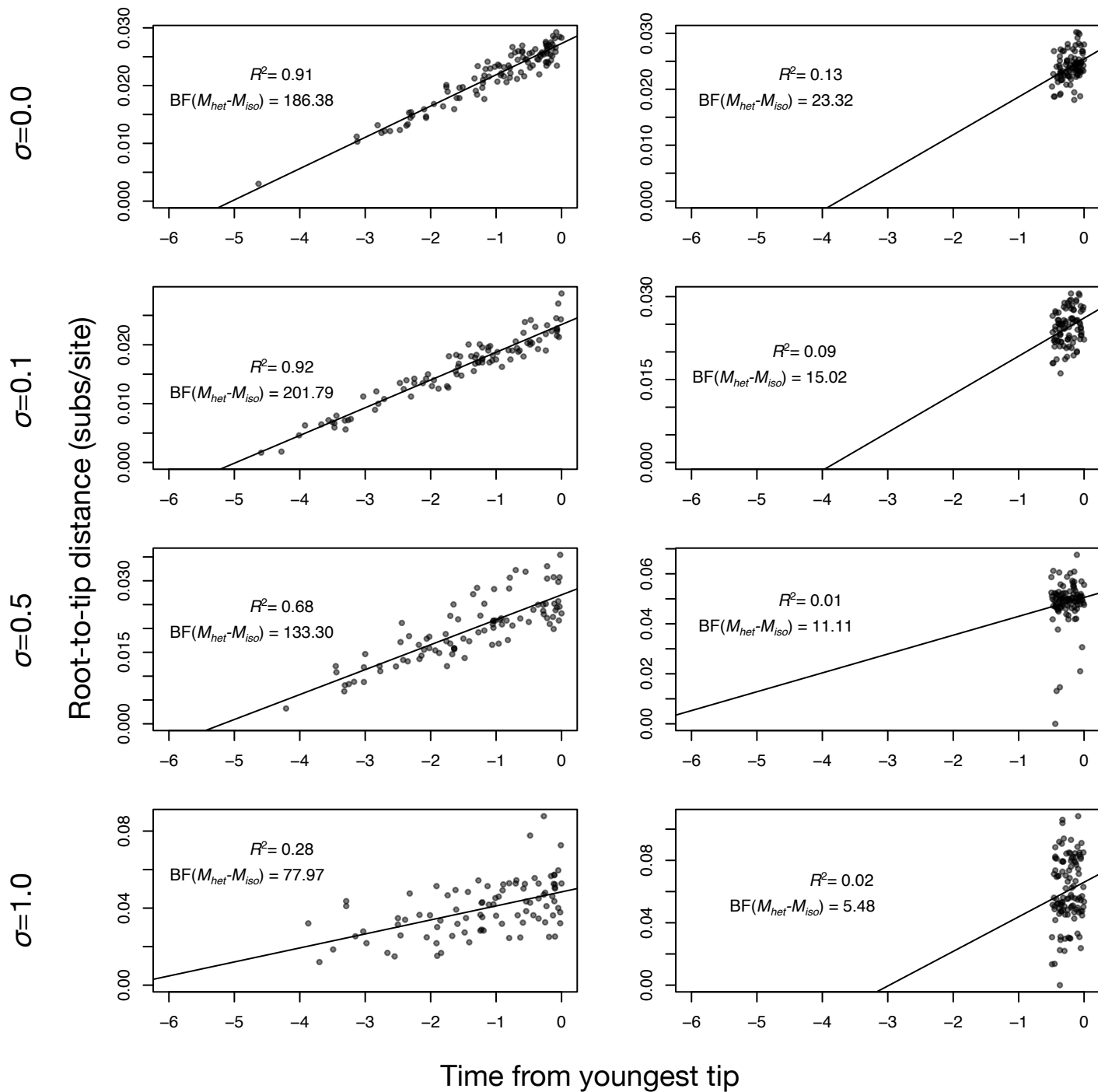


Analysis settings

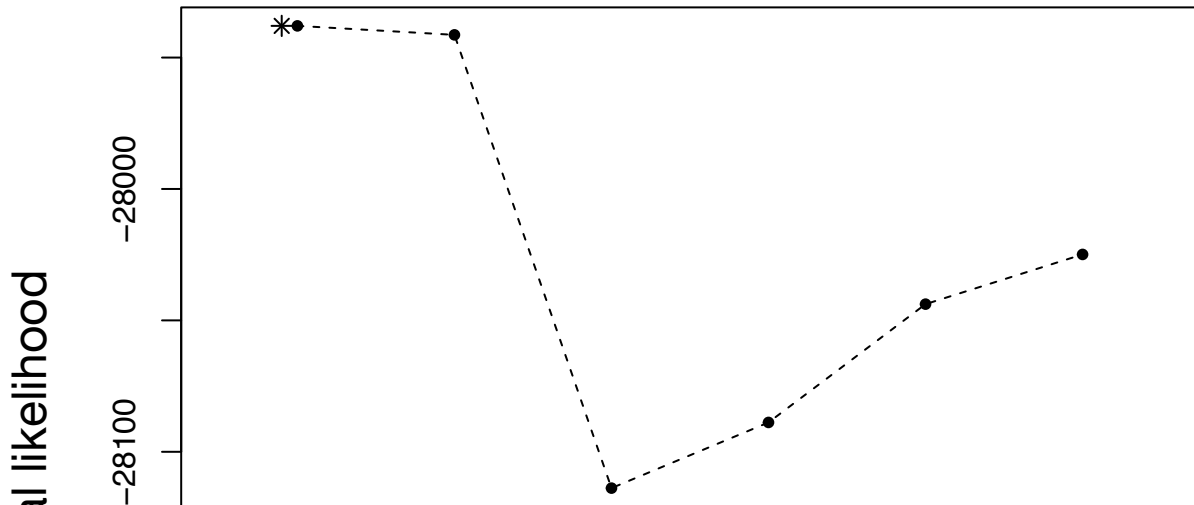
Log Bayes factors relative to best model



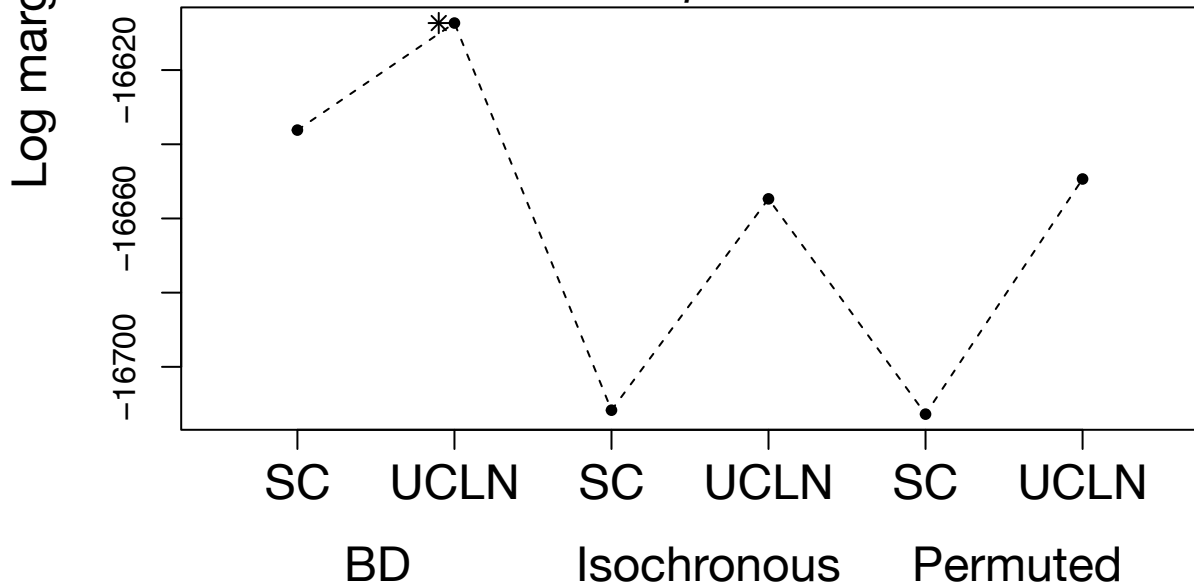
Analysis settings



*A/H1N1 Influenza virus*

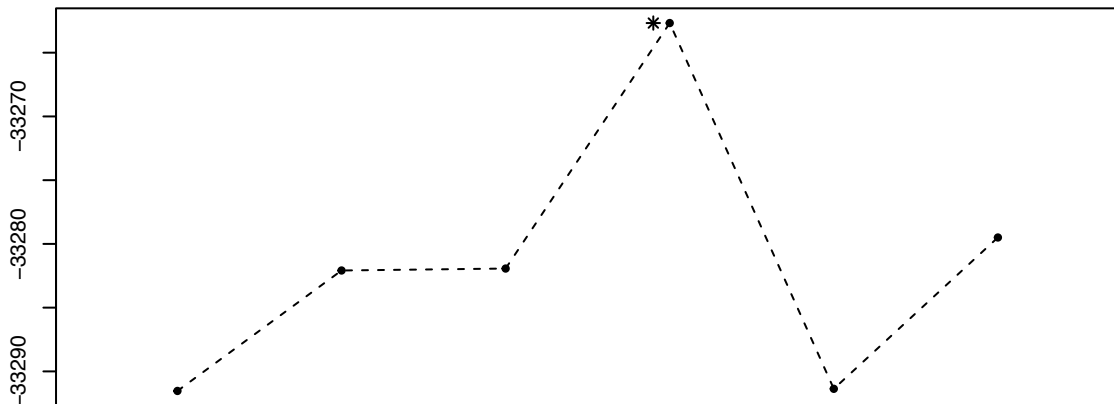


*Bordetella pertussis*

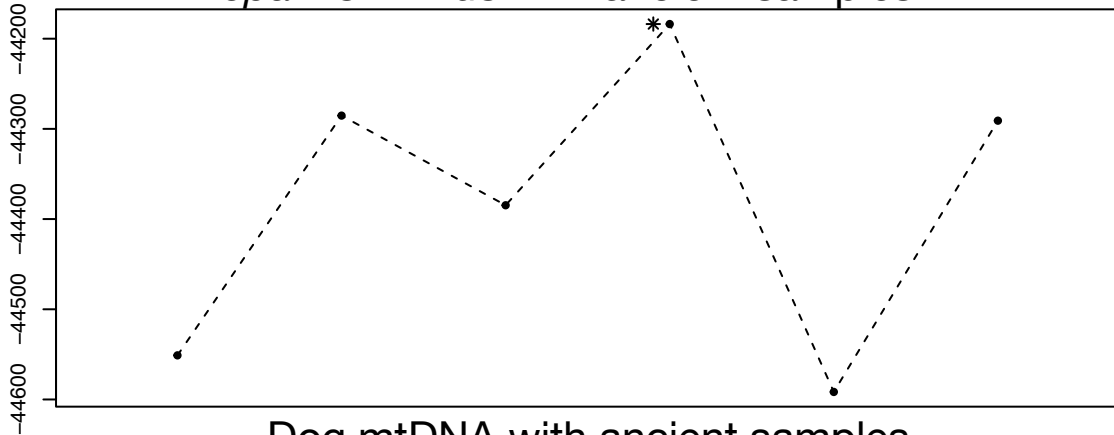


Log marginal likelihood

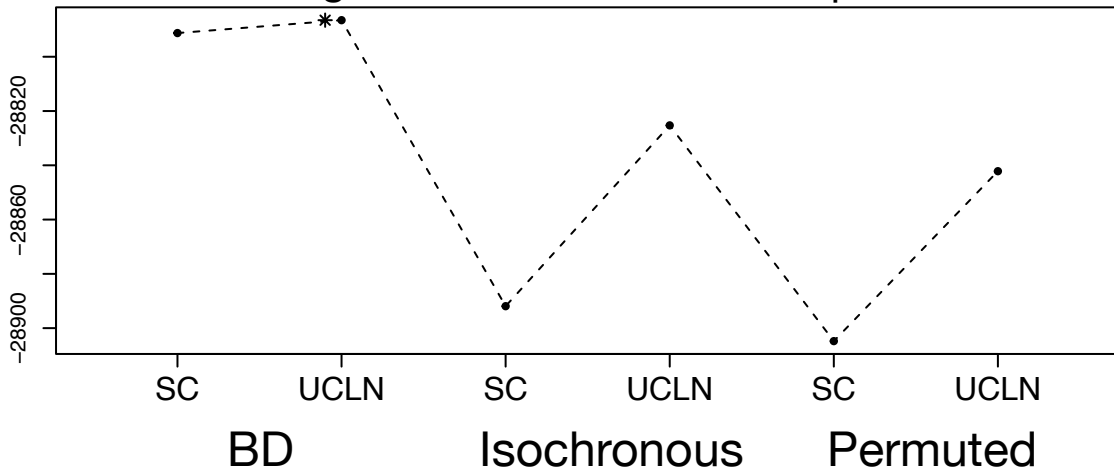
### Coronaviruses



### *Hepatitis B virus with ancient samples*

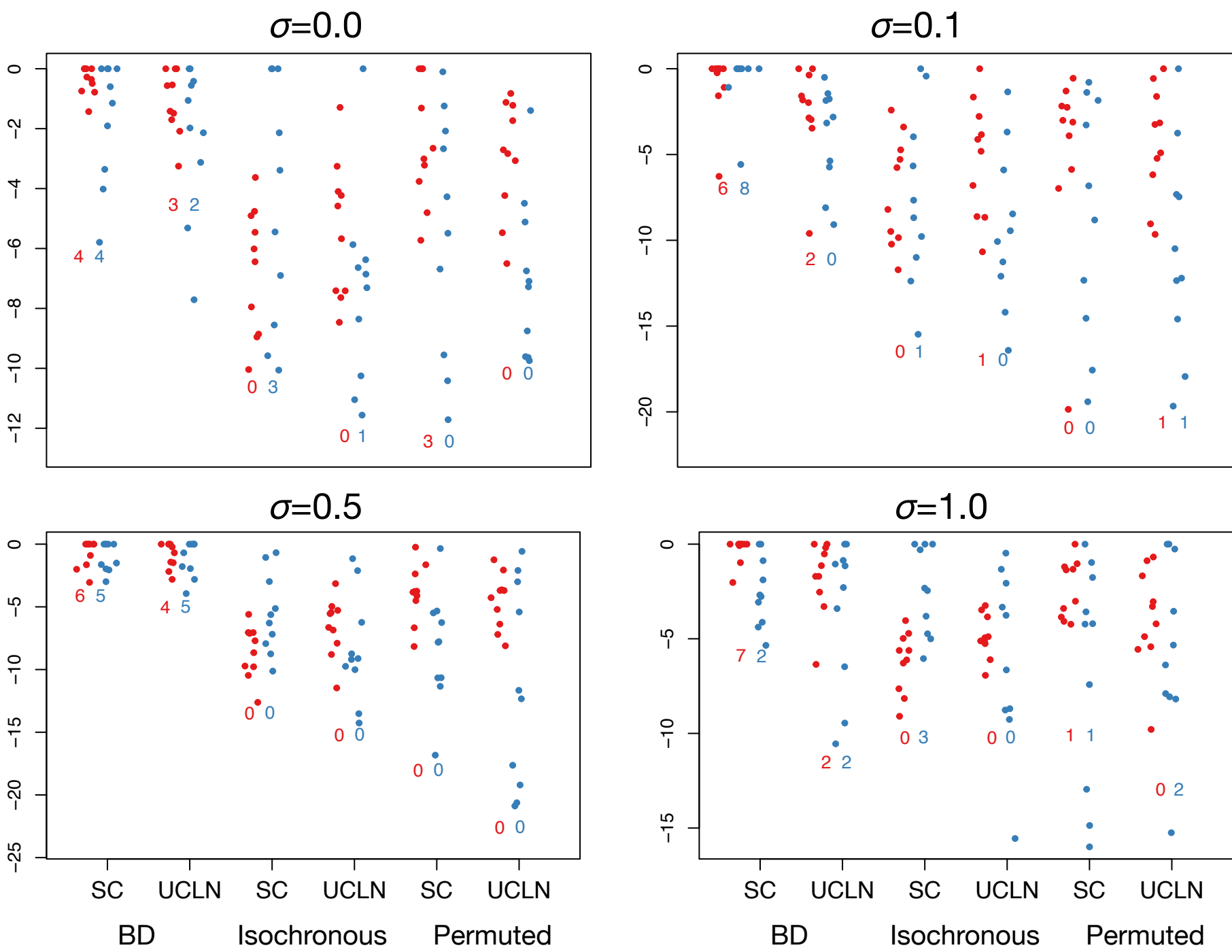


### Dog mtDNA with ancient samples



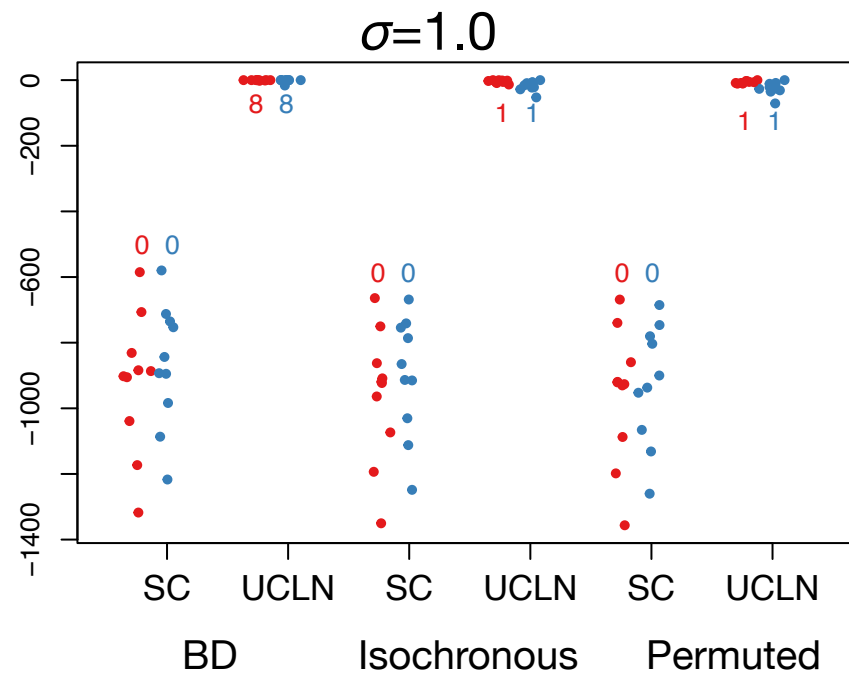
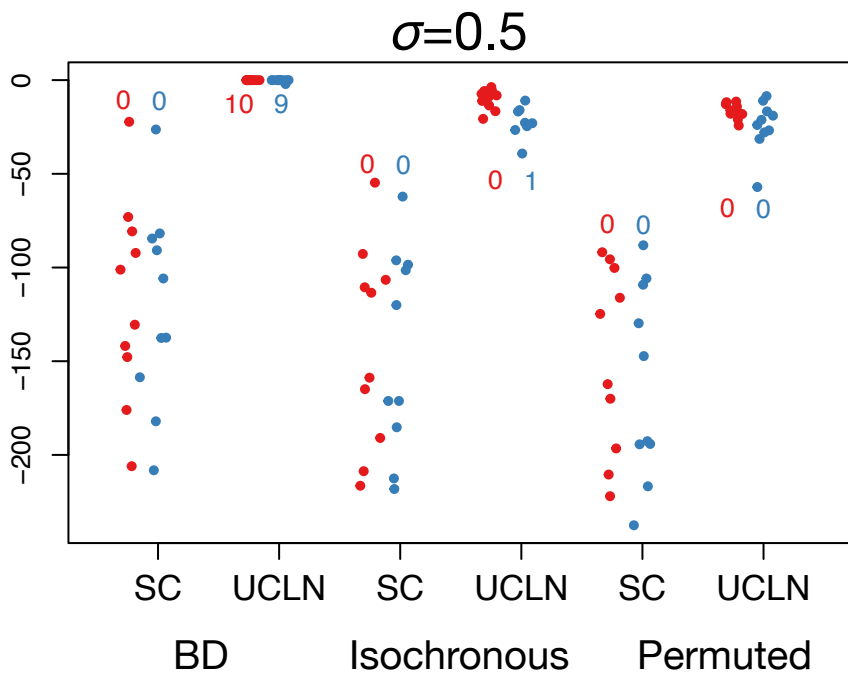
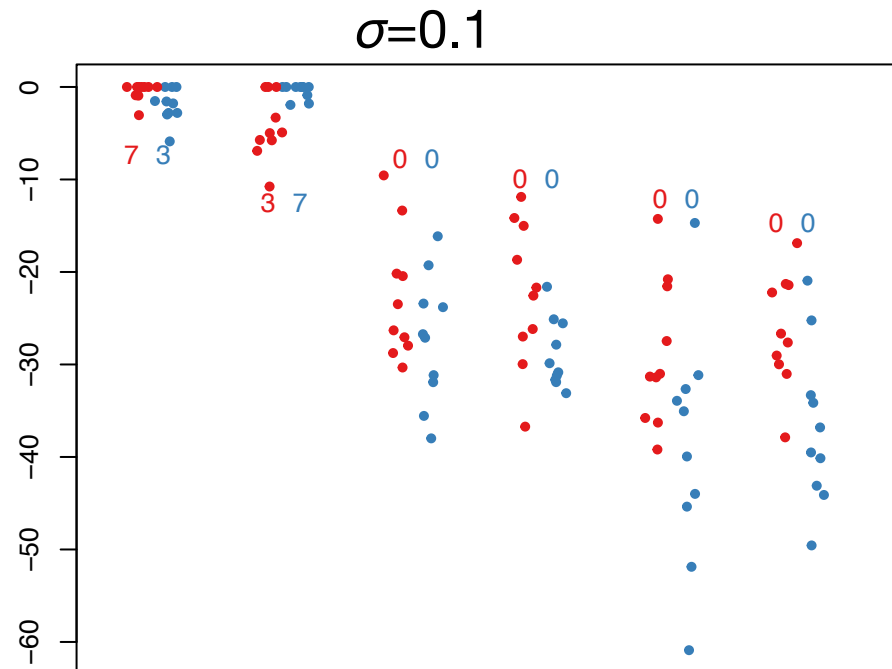
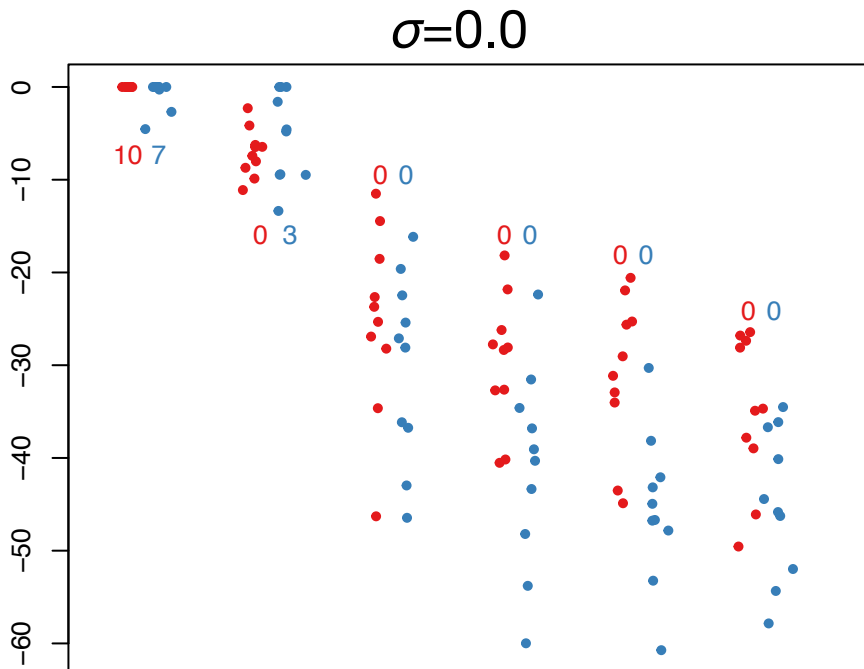


Log Bayes factors relative to best model



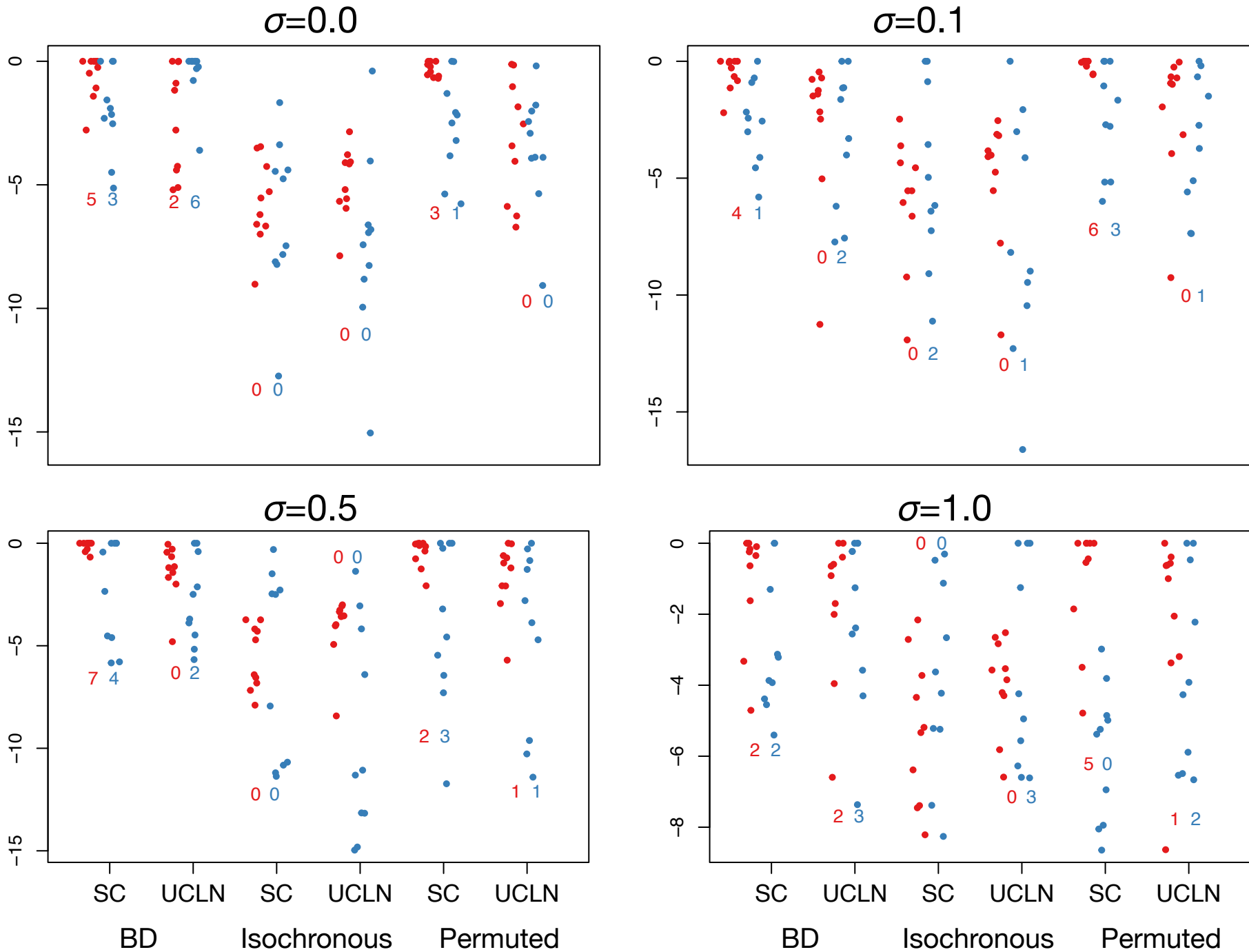
Analysis settings

Log Bayes factors relative to best model

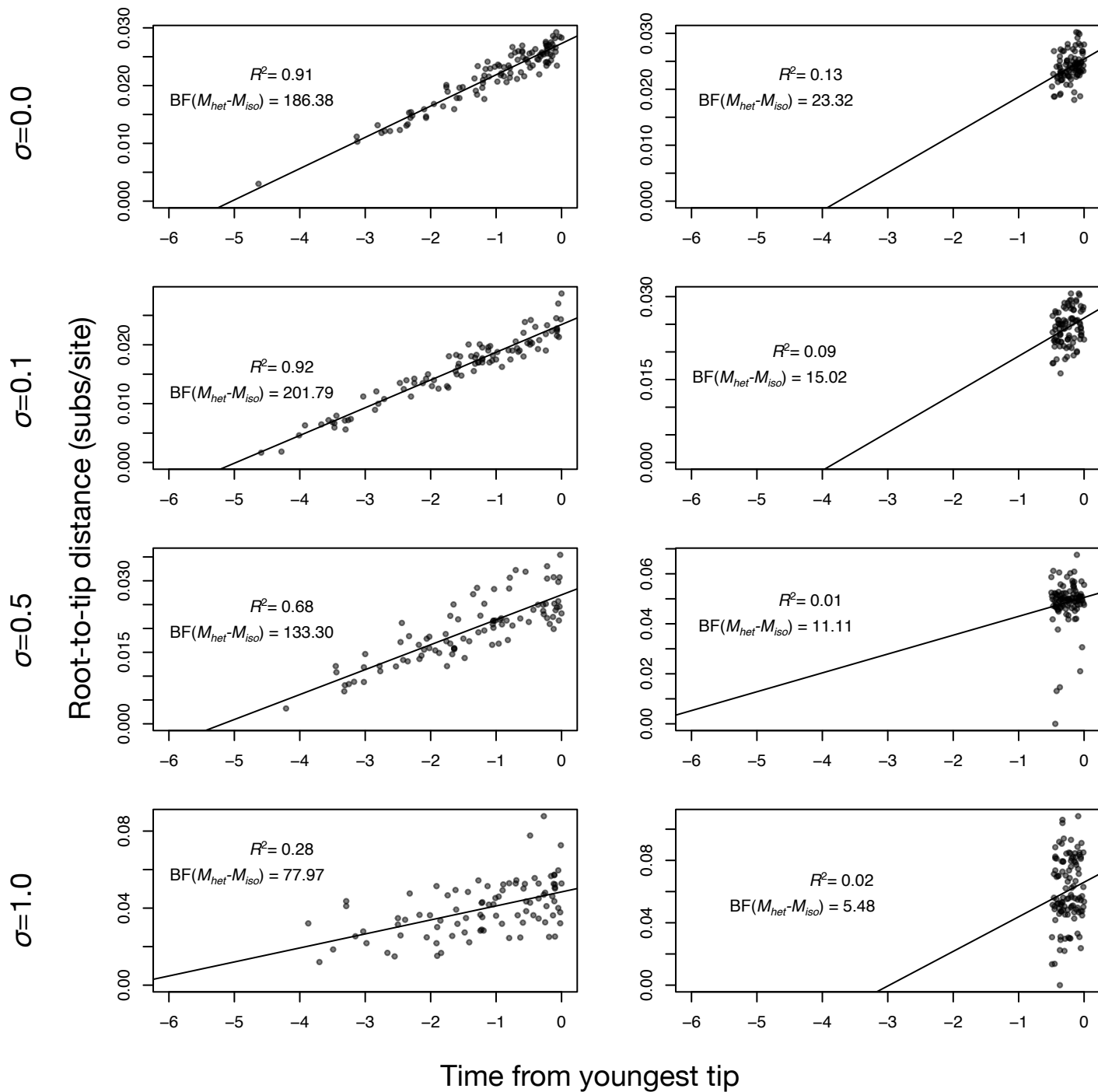


Analysis settings

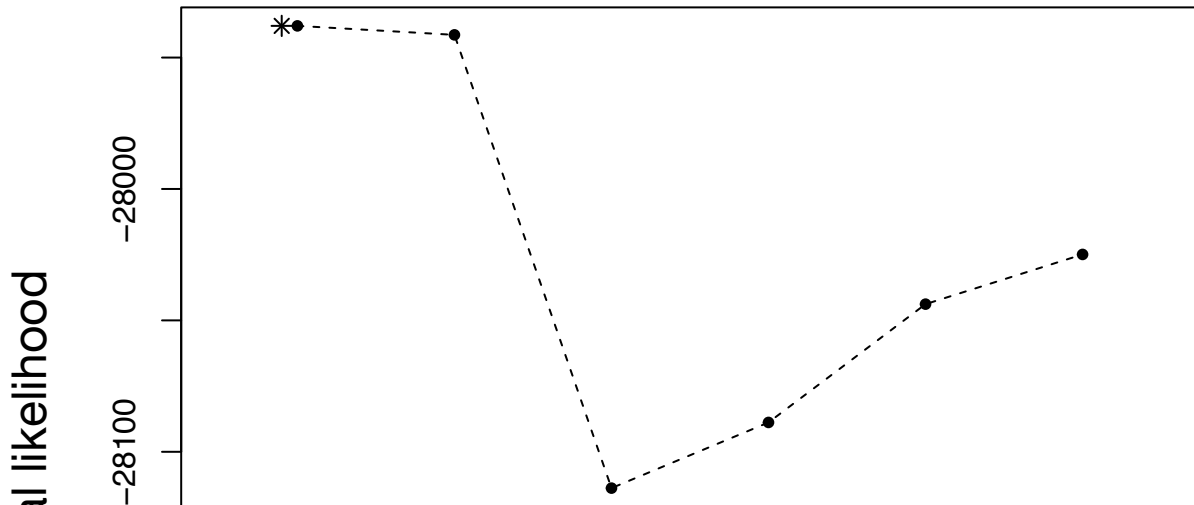
Log Bayes factors relative to best model



Analysis settings



*A/H1N1 Influenza virus*



*Bordetella pertussis*

