

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

## Inferring infectious disease phylodynamics with notification data

Sebastián Duchêne<sup>1\*</sup>, Francesca Di Giallonardo<sup>2</sup>, Edward C. Holmes<sup>3,5</sup>, Timothy G. Vaughan<sup>6,7</sup>.

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

<sup>2</sup>The Kirby Institute, University of New South Wales, Randwick, NSW, Australia.

<sup>3</sup>Marie Bashir Institute for Infectious Diseases and Biosecurity, The University of Sydney, Sydney, NSW, Australia.

<sup>4</sup>Charles Perkins Centre, School of Life and Environmental Sciences, The University of Sydney, Sydney, NSW, Australia.

<sup>5</sup>Sydney Medical School, The University of Sydney, Sydney, NSW, Australia.

<sup>6</sup>Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland.

<sup>7</sup>Swiss Institute of Bioinformatics (SIB), Switzerland.

\* Corresponding author

Email: [sebastian.duchene@unimelb.edu.au](mailto:sebastian.duchene@unimelb.edu.au)

### Abstract (150 words)

Genomic surveillance is increasingly common for infectious pathogens. Phylodynamic models can take advantage of pathogen genome sequence data to infer epidemiological dynamics, such as those based on the exponential growth coalescent and the birth-death process. Here we investigate the potential of including case notification data without associated genome sequences in such phylodynamic analyses. Using simulations, we demonstrate that birth-death phylodynamic models can capitalise on notification data to eliminate bias in estimates of the basic reproductive number,  $R_0$ , particularly when the sampling rate varies over time. In addition, an analysis of data collected from the 2009 pandemic H1N1 influenza virus demonstrates that using only samples from the prevalence peak results in biased estimates of the reproductive number over time, whereas using case notification data has a comparable accuracy to that achieved when using genome samples throughout the duration of the pandemic.

### Keywords

Phylodynamics, Notification data, Bayesian phylogenetics, Birth-death model, Coalescent model, Influenza virus.

### Main text (2000 words max.)

Outbreak investigations increasingly rely on genome sequencing of the causative pathogens. For example, it has been estimated that approximately 70% of Ebola cases that occurred in Sierra

43 Leone during the 2013-2016 West African Ebola virus outbreak have been sequenced (Stadler et al.  
44 2014). Phylodynamic methods can take advantage of these data to infer epidemiological dynamics  
45 (Grenfell et al. 2004). Recent sequencing technologies can generate such data very rapidly, such  
46 that phylodynamic inferences can be conducted in nearly real-time (Garday and Loman 2018;  
47 Hadfield et al. 2018; Grubaugh et al. 2019). The main appeal of phylodynamics is that the sequence  
48 data can inform on epidemiological dynamics for timescales prior to the earliest collected sample.  
49 Moreover, because phylodynamic inferences assume an underlying phylogenetic tree, the internal  
50 nodes and branches are informative about transmission.

51

52 Phylodynamic models describe a branching process. In Bayesian phylogenetic implementations the  
53 phylodynamic model is part of the prior and is sometimes referred to as the 'tree prior'. Internal  
54 nodes in the tree are assumed to be associated with transmission events while the tree tips  
55 represent sampling events, after which an individual is typically not infectious (du Plessis and  
56 Stadler 2015). The simplest models posit that the number of infected individuals increases  
57 exponentially over time. Although more sophisticated methods now exist (Kühnert et al. 2014;  
58 Poppinga et al. 2015; Kühnert et al. 2016; Rasmussen et al. 2017; Vaughan et al. 2017; Volz and  
59 Siveroni 2018), we focus our simulations on those that assume simple exponential growth that are  
60 appropriate for the early stages of an infectious disease outbreak.

61

62 Two commonly used phylodynamic models are the coalescent exponential and the birth-death,  
63 both of which assume that the infected population size,  $N$ , grows at a rate  $r$ ,  $N(t)=e^{rt}$ , where  $t$  is time  
64 after the origin. In the context of a branching process,  $r$ , is the difference between the transmission  
65 rate,  $\lambda$ , and the become uninfected rate,  $\delta$ , ( $r=\lambda-\delta$ ), and where  $1/\delta$  is the duration of infection. The  
66 basic reproductive number,  $R_0$ , is the average number of secondary infections in a fully susceptible  
67 population, estimated as  $R_0=\lambda/\delta$ . The exponential coalescent is a generalisation of the Wright-  
68 Fisher model where population size is a deterministic function of time (Griffiths and Tavaré 1994;  
69 Volz et al. 2009; Volz et al. 2013). The birth-death model typically assumes a stochastic process with  
70 sampling through time (Stadler 2010; Stadler et al. 2012; Stadler and Yang 2013), with  $\delta=\psi+\mu$ ,  
71 where  $\mu$  is the recovery rate and  $\psi$  is the sampling rate with recovery (the sampling proportion,  $p$ ,  
72 can be calculated as  $p=\psi/(\psi+\mu)$ ). This model can treat the time of origin of the outbreak as a  
73 parameter, which is not the case with the coalescent exponential. The individual parameters,  $\lambda$ ,  $\delta$ ,  
74 are non-identifiable because the tree likelihood in both models depends on two compound  
75 parameters,  $\lambda-\delta$  and  $\lambda\delta p$ , so prior information about any individual parameter is necessary to  
76 estimate the rest (Boskova et al. 2014).

77

78 Phylodynamic analyses typically require sequence data and sampling times (Rambaut 2000;  
79 Drummond et al. 2002; Drummond et al. 2003; Biek et al. 2015; Rieux and Balloux 2016). The  
80 number of samples and their times are also informative for the birth-death model because they are  
81 explicitly modelled (e.g. they inform  $\psi$ ) (Boskova et al. 2018). Although the amount of sequence  
82 data in outbreak investigations has increased, a key consideration is that sequencing efforts are  
83 often conducted only after a large number of cases are reported. For instance, the trees in Fig 1  
84 were simulated under an  $R_0$  of 2, a constant sampling effort, and over the course of 1 year. If  
85 sequencing was only conducted for samples collected after 0.75 years samples from the deep  
86 sections of the tree would be missed (*late sampling* in Fig 1). Such sampling bias can mislead

87 inferences of epidemiological dynamics because there are no data to inform inferences of the early  
88 stages of the outbreak.

89

90 Here we investigate bias in epidemiological parameters due to sampling heterogeneity, and we  
91 propose some effective approaches to reduce such bias. The first approach involves using a birth-  
92 death skyline model (Stadler et al. 2013), that requires an understanding of the sampling effort. For  
93 example, if there is knowledge that there was no attempt to collect samples early in the outbreak  
94 one can set two intervals for the  $\psi$  parameter. However, without knowledge of sampling effort this  
95 scenario is indistinguishable from one with a constant sampling effort but where initial prevalence  
96 was so low as to preclude obtaining any sequence data in the early stages of the outbreak. The  
97 second approach consists of including early case notification data in the analyses, where a  
98 notification is a clinically-confirmed case that was not sequenced (*notifications* scenario in Fig 1).  
99 Indeed, notifications are an inexpensive source of information traditionally used in epidemiology,  
100 such that they could be readily applied to leverage sequence data in outbreak investigations. In a  
101 Bayesian phylogenetic framework notification data can be incorporated by assigning a sampling  
102 time with no sequence data, and topological uncertainty is naturally incorporated into the analysis.  
103 An analogous approach can be used to coherently specify fossil data for molecular clock calibration  
104 (Heath et al. 2014; Heath and Moore 2014).

105

#### 106 *Simulation study*

107 We simulated phylogenetic trees under a birth-death process in MASTER v6.1 (Vaughan and  
108 Drummond 2013), with the following parameterisation;  $R_0=2$  or  $1.5$ ,  $\delta=91$ ,  $p=0.05$ , and an outbreak  
109 duration of one year ( $1/\delta = 0.011$  of one year for a duration of infection of about 4 days). The number  
110 of tips and their ages are naturally variable (from 100 to 150 tips). We assumed a strict molecular  
111 clock with an evolutionary rate of 0.01 substitutions per site per year (subs/site/year) and the HKY+ $\Gamma$   
112 substitution model to produce alignments of 13,000 nucleotides using NELSI (Ho et al. 2015) and  
113 Phangorn v2.4 (Schliep 2011). These settings are broadly similar to an influenza virus outbreak  
114 (Hedge et al. 2013). We then assumed three sampling scenarios: (i) *constant* sampling with all  
115 sequences from the simulation included (e.g. the sequence for every sample in the tree in Fig 1 is  
116 included), (ii) *late sampling* only with samples after time  $T_s$  (e.g. only sequences for samples after  
117 the dashed line in the tree in Fig 1), and (iii) *notifications* in which sequence data are available only  
118 after time  $T_s$  and the sampling time for samples before  $T_s$  are included with no sequence data (i.e.  
119 notifications). We set  $T_s$  at 0.75 and 0.9 years. For each parameter configuration we simulated 100  
120 sequence data sets which were subsampled according to the three scenarios above. We analysed  
121 the data in BEAST v2.5 (Bouckaert et al. 2014; Bouckaert et al. 2018), considering several  
122 phylodynamic models; a coalescent exponential and the birth-death. For the *late sampling* scenario  
123 we also considered the birth-death skyline with two intervals for the  $\psi$  parameter, with the interval  
124 time fixed at  $T_s$ . We matched the substitution and clock model to those used to generate the data  
125 and we used an informative prior on  $\delta$  using a  $\Gamma$  distribution with mean 91 and standard deviation of  
126 1.

127

128 Analyses of data sets with late sampling using a birth-death model produced inaccurate estimates  
129 of  $R_0$ . In only 11 of 100 simulations with  $R_0=2$  the 95% highest posterior density (HPD) for this  
130 parameter included the value used to generate the data (Table 1 and Fig 2). For simulations with  
131  $R_0=1.5$  this model was never able to recover the true  $R_0$  (Table S1 and Fig S1). The birth-death

132 skyline had much better performance, with 96 of 100 simulations estimating  $R_0$  accurately (i.e. the  
133 true value was within the HPD). The coalescent exponential had better performance than the birth-  
134 death, but it was still less accurate than the birth-death skyline, with 90 simulations producing  
135 accurate estimates. In general, for data sets with late sampling we observed that  $R_0$  tended to be  
136 overestimated with the birth-death and underestimated with the coalescent exponential (Fig 2).  
137 Interestingly, estimates of the evolutionary rate displayed a similar pattern to those of  $R_0$ , with the  
138 birth-death skyline and the birth-death being the most and least accurate, respectively.

139

140 As expected, analyses of the data with constant sampling were accurate in a majority of cases, with  
141 97 and 93 of 100 simulations being accurate for  $R_0$ , and 98 and 97 for the evolutionary rate, using the  
142 birth-death and the coalescent exponential, respectively (the correct model is the birth-death, such  
143 that it is expected to perform better than the coalescent). Estimates of  $R_0$  including notification  
144 data were similarly accurate as those with complete sampling under the birth-death model, where  
145 96 analyses correctly estimated this parameter, but this was not the case for the coalescent  
146 exponential, where only 84 analyses included the true value (Table 1). Evolutionary rate estimate  
147 with notification data were less accurate than those from the complete data, with 88 accurate  
148 analyses using the birth-death and 66 using the coalescent exponential. These results can be  
149 attributed to the fact that the birth-death treats sampling times as data, whereas the coalescent is  
150 conditioned on the number of samples and their ages (Stadler et al. 2015; Boskova et al. 2018).  
151 Notifications improve the accuracy of  $R_0$  in the birth-death and they pose an informative prior on  
152 the age of the tree height (Boskova et al. 2014), which can also improve the accuracy of the  
153 evolutionary rate relative to the coalescent exponential, but this estimate is unlikely to be as  
154 accurate as that with the complete sequence data because there is necessarily less molecular  
155 information.

156

157 The coalescent exponential appears to be more robust to the sampling process, with greater  
158 accuracy than the birth-death for the late sampling analyses. This model may be a good alternative  
159 when the sampling process is poorly understood, and with no reliable notification data. However,  
160 our simulations suggest that this comes at the expense of estimates that are less precise (with  
161 higher uncertainty) than those from the birth-death, as measured by the mean estimate divided by  
162 the HPD width (Table 1).

163

#### 164 Empirical case study: *A/H1N1 Influenza virus from North America*

165 To illustrate the accuracy of notification data relative to completely sequenced data sets we  
166 analysed 639 whole genome sequences sampled from the 2009 A/H1N1 influenza pandemic from  
167 North America that were downloaded from GenBank (Supplementary material). The sequences  
168 were collected from early April to October 2009. We chose this period of time because it  
169 corresponds to a densely sampled clade and captures the peak number of infections as reported in  
170 the FluView application (CDC 2019). We considered three data subsets for our analyses; (i)  
171 'complete sampling' with all of the genome sequences, (ii) 'notifications' with 104 sequences only  
172 from September and the remaining 535 samples treated as notifications, and (iii) 'late sampling'  
173 with only the 104 sequences from September 2009. Because we do not expect constant exponential  
174 growth for these data, we used a birth-death skyline model to estimate the reproductive number,  $R_e$   
175 (similar to  $R_0$ , but not assuming a fully susceptible population), on a monthly basis from January to  
176 October. We set the duration of infection at 4 days with a  $\Gamma$  prior on  $\delta$  with mean 91 and standard

177 deviation of 1, and we used the HKY+ $\Gamma$  substitution model and a strict molecular clock model. In all  
178 cases we set two intervals for  $\psi$ , to allow for a low sampling probability before the oldest sample.

179

180 The birth-death skyline plot revealed nearly identical trends for the complete data and that using  
181 notifications (Fig 3a). For example, the highest  $R_e$  was estimated in April, with a mean of 1.52 (HPD:  
182 1.40 – 1.66) for the complete sampling and a mean of 1.54 (HPD: 1.40 – 1.70) for the analysis using  
183 notifications. The estimate for the epidemic origin for the late sampling analyses was around late  
184 June, such that we cannot estimate  $R_e$  before this time. However, for July, August and September  
185 we found that late sampling resulted in substantially different estimates to those from the other  
186 analyses (Fig 3b-d), particularly in July where  $R_e$  with late sampling had a mean of 2.10 (HPD: 1.60 –  
187 2.54) and those with complete sampling and notifications were 0.79 (HPD: 0.72 – 0.85 and 0.72 –  
188 0.86, respectively). Importantly, estimates of  $R_e$  using complete sampling or notifications are overall  
189 similar in magnitude to those from large-scale epidemiological studies (Fraser et al. 2009;  
190 Biggerstaff et al. 2014), and previous estimates using genome sequence data (Hedge et al. 2013).

191

#### 192 *Notification data in empirical phylodynamic studies*

193 Our simulations and empirical data analyses reveal that notification data are a rich source of  
194 information for birth-death models that can dramatically improve the accuracy and precision in  
195 estimates of epidemiological parameters. A key consideration is that notifications should represent  
196 confirmed cases that would have been sequenced if sequencing effort had been constant.  
197 Combining notification and sequence data can be particularly useful in situations where it is  
198 unknown whether sequence sampling has been constant over time or where there exist several  
199 confirmed cases but a smaller number of sequences. For example, in recently emerging outbreaks  
200 combining both sources of data can provide timely insight about the recent evolution of the  
201 pathogen in question.

202

#### 203 **Acknowledgements**

204 SDG was supported by a Discovery Early Career Fellowship from the Australian Research Council  
205 (DE190100805) and a McKenzie fellowship from the University of Melbourne.

206

207

#### 208 **References**

- 209 Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. 2015. Measurably evolving pathogens in the genomic  
210 era. *Trends Ecol. Evol.* 30:306–313.
- 211 Biggerstaff M, Cauchemez S, Reed C, Gambhir M, Finelli L. 2014. Estimates of the reproduction  
212 number for seasonal, pandemic, and zoonotic influenza: a systematic review of the literature.  
213 *BMC Infect. Dis.* 14:480.
- 214 Boskova V, Bonhoeffer S, Stadler T. 2014. Inference of epidemiological dynamics based on  
215 simulated phylogenies using birth-death and coalescent models. *PLOS Comput Biol*  
216 10:e1003913.
- 217 Boskova V, Stadler T, Magnus C. 2018. The influence of phylodynamic model specifications on  
218 parameter estimates of the Zika virus epidemic. *Virus Evol.* 4:vex044.
- 219 Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, Suchard MA, Rambaut A, Drummond  
220 AJ. 2014. BEAST 2: a software platform for Bayesian evolutionary analysis. *PLOS Comput.*  
221 *Biol.* 10:e1003537.

- 222 Bouckaert R, Vaughan TG, Barido-Sottani J, Duchene S, Fourment M, Gavryushkina A, Heled J,  
223 Jones G, Kuhnert D, de Maio N. 2018. BEAST 2.5: An Advanced Software Platform for  
224 Bayesian Evolutionary Analysis. *bioRxiv*:474296.
- 225 CDC. 2019. Centres for Disease Control and Prevention. Available from:  
226 <https://gis.cdc.gov/grasp/fluview/>
- 227 Drummond AJ, Nicholls GK, Rodrigo AG, Solomon W. 2002. Estimating mutation parameters,  
228 population history and genealogy simultaneously from temporally spaced sequence data.  
229 *Genetics* 161:1307–1320.
- 230 Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG. 2003. Measurably evolving  
231 populations. *Trends Ecol. Evol.* 18:481–488.
- 232 Fraser C, Donnelly CA, Cauchemez S, Hanage WP, Van Kerkhove MD, Hollingsworth TD, Griffin J,  
233 Baggaley RF, Jenkins HE, Lyons EJ. 2009. Pandemic potential of a strain of influenza A (H1N1):  
234 early findings. *Science*. 324:1557–1561.
- 235 Gardy JL, Loman NJ. 2018. Towards a genomics-informed, real-time, global pathogen surveillance  
236 system. *Nat. Rev. Genet.* 19:9.
- 237 Grenfell BT, Pybus OG, Gog JR, Wood JLN, Daly JM, Mumford JA, Holmes EC. 2004. Unifying the  
238 epidemiological and evolutionary dynamics of pathogens. *Science* 303:327–332.
- 239 Griffiths RC, Tavaré S. 1994. Sampling theory for neutral alleles in a varying environment. *Philos.*  
240 *Trans. R. Soc. London. Ser. B Biol. Sci.* 344:403–410.
- 241 Grubaugh ND, Ladner JT, Lemey P, Pybus OG, Rambaut A, Holmes EC, Andersen KG. 2019.  
242 Tracking virus outbreaks in the twenty-first century. *Nat. Microbiol.* 4:10.
- 243 Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, Sagulenko P, Bedford T, Neher  
244 RA. 2018. Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* 34:4121–4123.
- 245 Heath TA, Huelsenbeck JP, Stadler T. 2014. The fossilized birth–death process for coherent  
246 calibration of divergence-time estimates. *Proc. Natl. Acad. Sci.* 111:E2957–E2966.
- 247 Heath TA, Moore BR. 2014. Bayesian inference of species divergence times. In: Chen M-H, Kuo L,  
248 Lewis PO, editors. *Bayesian Phylogenetics, Methods, Algorithms, and Applications*. Boca  
249 Raton, Florida: CRC Press. p. 277–318.
- 250 Hedge J, Lycett SJ, Rambaut A. 2013. Real-time characterization of the molecular epidemiology of  
251 an influenza pandemic. *Biol. Lett.* 9:20130331.
- 252 Ho SS, Duchêne S, Duchêne DA. 2015. Simulating and detecting autocorrelation of molecular  
253 evolutionary rates among lineages. *Mol. Ecol. Resour.* 15:688–696.
- 254 Kühnert D, Stadler T, Vaughan TG, Drummond AJ. 2014. Simultaneous reconstruction of  
255 evolutionary history and epidemiological dynamics from viral sequences with the birth–death  
256 SIR model. *J. R. Soc. Interface* 11:20131106.
- 257 Kühnert D, Stadler T, Vaughan TG, Drummond AJ. 2016. Phylodynamics with migration: A  
258 computational framework to quantify population structure from genomic data. *Mol. Biol. Evol.*  
259 33:2102–2116.
- 260 du Plessis L, Stadler T. 2015. Getting to the root of epidemic spread with phylodynamic analysis of  
261 genomic data. *Trends Microbiol.* 23:383–386.
- 262 Poppinga A, Vaughan T, Stadler T, Drummond AJ. 2015. Inferring epidemiological dynamics with  
263 Bayesian coalescent inference: the merits of deterministic and stochastic models. *Genetics*  
264 199:595–607.
- 265 Rambaut A. 2000. Estimating the rate of molecular evolution: incorporating non-contemporaneous  
266 sequences into maximum likelihood phylogenies. *Bioinformatics* 16:395–399.

- 267 Rasmussen DA, Kouyos R, Günthard HF, Stadler T. 2017. Phylodynamics on local sexual contact  
268 networks. PLoS Comput. Biol. 13:e1005448.
- 269 Rieux A, Balloux F. 2016. Inferences from tip-calibrated phylogenies: a review and a practical guide.  
270 Mol. Ecol. 25:1911–1924.
- 271 Schliep KP. 2011. phangorn: Phylogenetic analysis in R. Bioinformatics 27:592–593.
- 272 Stadler T. 2010. Sampling-through-time in birth-death trees. J. Theor. Biol. 167:696–404.
- 273 Stadler T, Kouyos R, von Wyl V, Yerly S, Böni J, Bürgisser P, Klimkait T, Joos B, Rieder P, Xie D.  
274 2012. Estimating the basic reproductive number from viral sequence data. Mol. Biol. Evol.  
275 29:347–357.
- 276 Stadler T, Kühnert D, Bonhoeffer S, Drummond AJ. 2013. Birth–death skyline plot reveals temporal  
277 changes of epidemic spread in HIV and hepatitis C virus (HCV). Proc. Natl. Acad. Sci. 110:228–  
278 233.
- 279 Stadler T, Kühnert D, Rasmussen DA, du Plessis L. 2014. Insights into the early epidemic spread of  
280 Ebola in Sierra Leone provided by viral sequence data. PLoS Curr. 6.
- 281 Stadler T, Vaughan TG, Gavryushkin A, Guindon S, Kühnert D, Leventhal GE, Drummond AJ. 2015.  
282 How well can the exponential-growth coalescent approximate constant-rate birth–death  
283 population dynamics? Proc. R. Soc. B Biol. Sci. 282:20150420.
- 284 Stadler T, Yang Z. 2013. Dating phylogenies with sequentially sampled tips. Syst. Biol. 62:674–688.
- 285 Vaughan TG, Drummond AJ. 2013. A stochastic simulator of birth–death master equations with  
286 application to phylodynamics. Mol. Biol. Evol. 30:1480–1493.
- 287 Vaughan TG, Leventhal GE, Rasmussen DA, Drummond AJ, Welch D, Stadler T. 2017. Directly  
288 Estimating Epidemic Curves From Genomic Data. bioRxiv:142570.
- 289 Volz E, Siveroni I. 2018. Bayesian phylodynamic inference with complex models. PLoS Comput.  
290 Biol. 14:e1006546.
- 291 Volz EM, Koelle K, Bedford T. 2013. Viral phylodynamics. PLoS Comput. Biol. 9:e1002947.
- 292 Volz EM, Pond SLK, Ward MJ, Brown AJL, Frost SDW. 2009. Phylodynamics of infectious disease  
293 epidemics. Genetics 183:1421–1430.

294

## 295 Figure legends

296

297 **Fig 1.** Example of a phylogenetic trees generated under a birth-death process with a basic  
298 reproductive number,  $R_0$ , of 2 a become uninfected rate,  $\delta$ , of 100 for three analysis scenarios. The  
299 solid line denotes the number of samples collected over time. In *constant* sampling samples are  
300 collected and sequenced at a rate  $\psi=5$  (i.e. sampling probability,  $p$ , of 0.05). In *late sampling*  
301 samples are collected and sequenced after time  $T_s$ , shown with the dashed line. In *notifications*  
302 samples are collected constantly over time, but only sequenced after time  $T_s$ , such that before  $T_s$ ,  
303 only notifications (sampling times with no sequence data) are included. Blue circles represent  
304 samples with sequence data, whereas those in orange correspond to notifications. In the  
305 *notifications* scenario a Bayesian analysis would integrate over their phylogenetic uncertainty. The  
306 solid line represents the number of samples collected over time. In *late sampling* there are no  
307 samples collected before  $T_s$ , such that assuming constant sampling can produce a bias in estimates  
308 of epidemiological dynamics.

309

310 **Fig 2.** Posterior densities for estimates of the basic reproductive number,  $R_0$ , and the evolutionary  
311 rate for 100 simulations with true  $R_0$  of 2 and an evolutionary rate of 0.01 subs/site/year. The bars

312 represent the 95% highest posterior density (HPD) and the points are the median. We analysed the  
313 data by sampling late in the outbreak only (i.e. after 0.75 of the tree height), with a constant  
314 sampling effort (with all samples sequenced), and by including notifications. The colours represent  
315 four different models; red for the coalescent exponential, blue for the birth-death skyline, and  
316 orange for the birth-death with constant sampling. For the data with sampling late in the outbreak  
317 only we use the birth-death skyline model with constant  $R_e$  and two intervals for the sampling rate,  
318  $\psi$ , before time 0.75. This model is not applicable to analyses with complete sampling or with  
319 notifications where sampling is constant. The dashed horizontal lines correspond to the true  
320 parameter value used to generate the data.

321

322 **Fig 3.** Estimates of the reproductive number,  $R_e$ , for empirical data of the 2009 A/H1N1 influenza  
323 pandemic in North America. **a.** Is a birth-death skyline plot where in which  $R_e$  is estimated per  
324 month from January to October 2009. Each line is a sampled trajectory from the posterior using  
325 each analysis, with red for that from the complete sampling (639 sequences), blue for notifications  
326 with 104 sequences from September (the month with largest number of infections) and 535  
327 notifications, and green is for only the 104 sequences from September. Note that the analysis using  
328 only sequence data from September has a more recent origin parameter, such that it is only  
329 possible to estimate  $R_e$  from June. The ticks along the x-axis represent the timing of sequences  
330 sampled in all analyses in black, and those treated as notifications (with no sequence data) in the  
331 'notifications' analysis. Panels **b.**, **c.** and **d.** show the posterior density for estimates of  $R_e$  in July,  
332 August, and September, respectively, for each analysis with colours matching those in panel **a.** (i.e.  
333 they are the densities for these months shown in **a.**).

334

335



336 **Tables**

337 **Table 1.** Results of the simulation study with  $R_0$  of 1.5 and evolutionary rate of 0.01 subs/site/year.  
338 The rows correspond to the seven treatments. The first two columns denote the number of  
339 simulations (out of 100) where the value used to generate the data was contained within the 95%  
340 highest posterior density (HPD). The last two columns are a measure of precision of the estimates  
341 calculated as the estimated mean estimate of  $R_0$  and the evolutionary rate divided by the 95% HPD  
342 width, such that large values imply low precision. Here we report the mean value over 100  
343 simulations.

344

345 **Supplementary material**

346 **Table S1.** Results of the simulation study with  $R_0=1.5$ , evolutionary rate of 0.01 subs/site/year, and  
347 late sampling starting at 0.9 years of a total time of 1 year. The rows correspond to the seven  
348 treatments. The first two columns denote the number of simulations (out of 100) where the value  
349 used to generate the data was contained within the 95% highest posterior density (HPD). The last  
350 two columns are a measure of precision of the estimates calculated as the estimated mean estimate  
351 of  $R_0$  and the evolutionary rate divided by the 95% HPD width, such that large values imply low  
352 precision. Here we report the mean value over 100 simulations.

353

354 **Fig S1.** Posterior densities for estimates of the basic reproductive number,  $R_0$ , and the evolutionary  
355 rate for 100 simulations with true  $R_0$  of 1.5 and an evolutionary rate of 0.01 subs/site/year. The bars  
356 represent the 95% highest posterior density (HPD) and the points are the median. We analysed the  
357 data by sampling late in the outbreak only (i.e. after 0.75 of the tree height), with a constant  
358 sampling effort (with all samples sequenced), and by including notifications. The colours represent  
359 four different models; red for the coalescent exponential, blue for the birth-death skyline, and  
360 orange for the birth-death with constant sampling. For the data with sampling late in the outbreak  
361 only we use the birth-death skyline model with constant  $R_0$  and two intervals for the sampling rate,  
362  $\psi$ , before time 0.75. This model is not applicable to analyses with complete sampling or with  
363 notifications where sampling is constant. The dashed horizontal lines correspond to the true  
364 parameter value used to generate the data.

365

366 **Supplementary data.** Zip file with input files to generate trees in MASTER and to analyse sequence  
367 data in BEAST according to the birth-death skyline, birth-death, and the coalescent exponential  
368 models. Accession numbers for empirical A/H1N1 Influenza virus data.

369

370

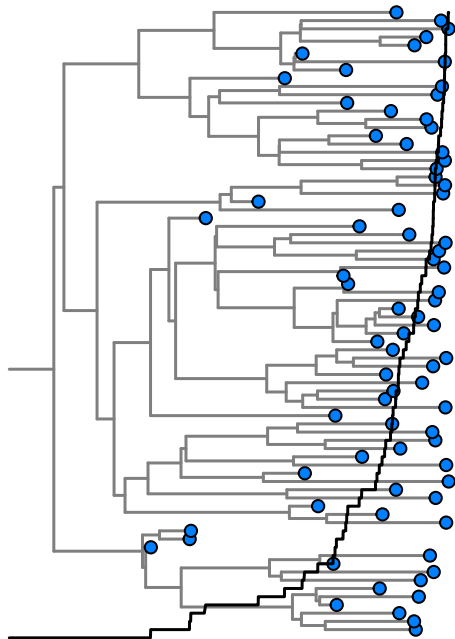
371

372

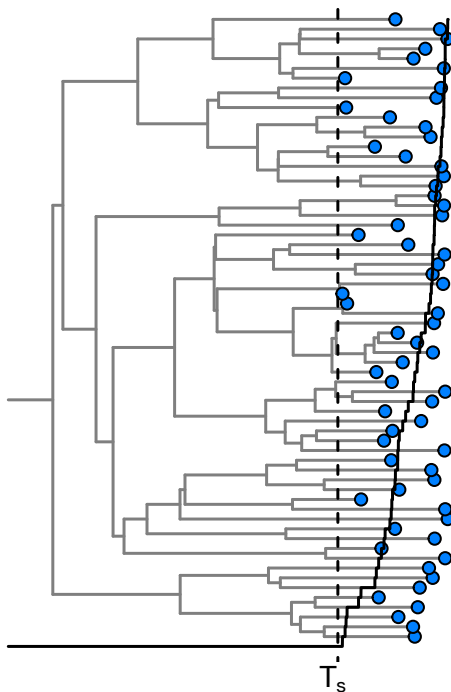
373

374

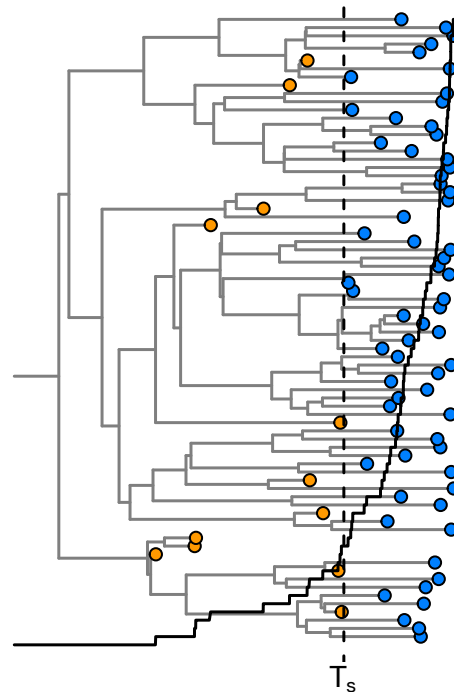
Constant sampling



Late sampling

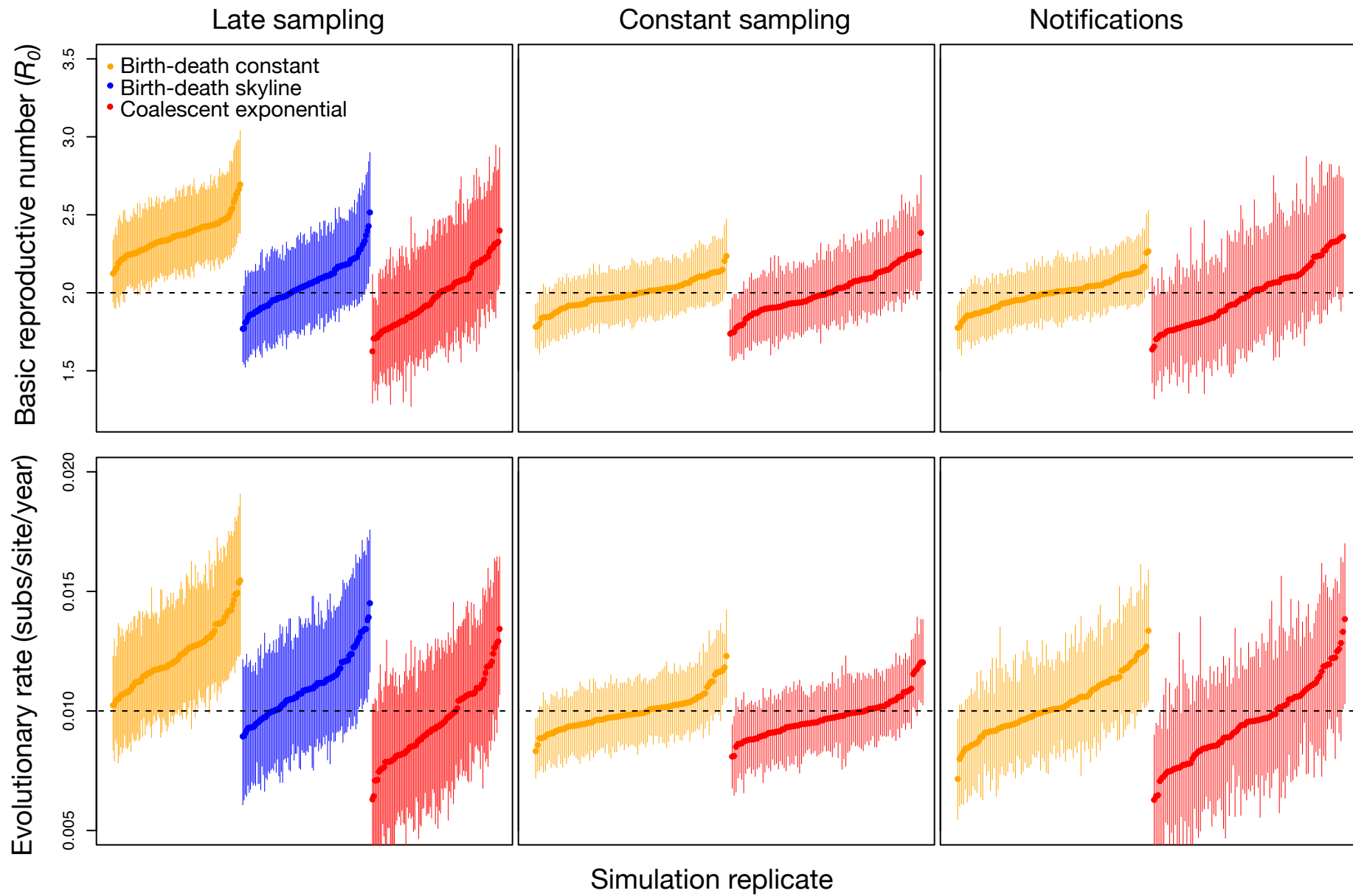


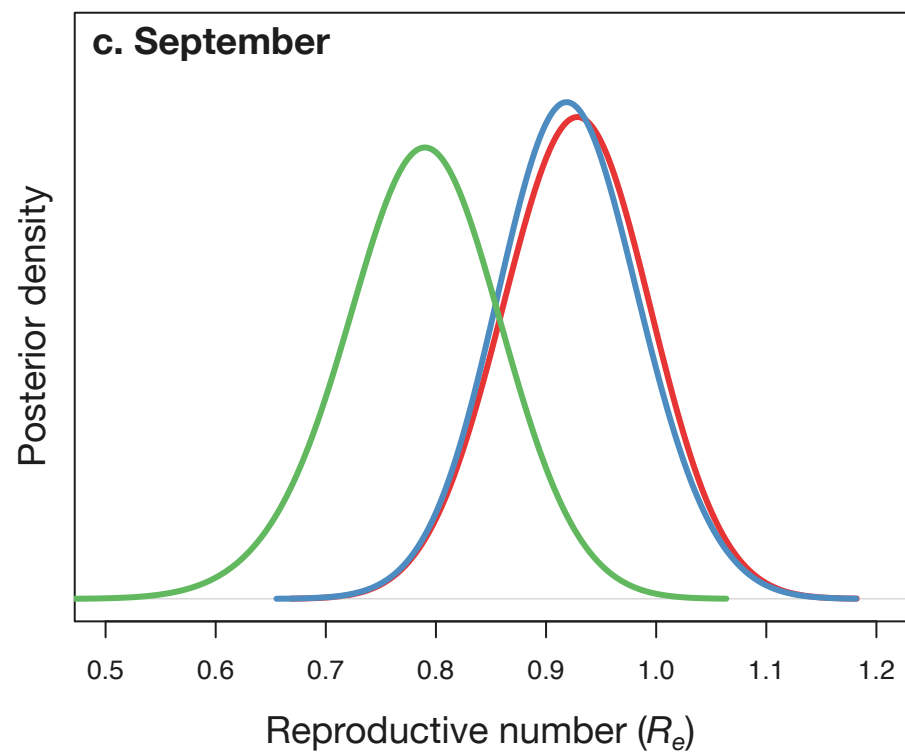
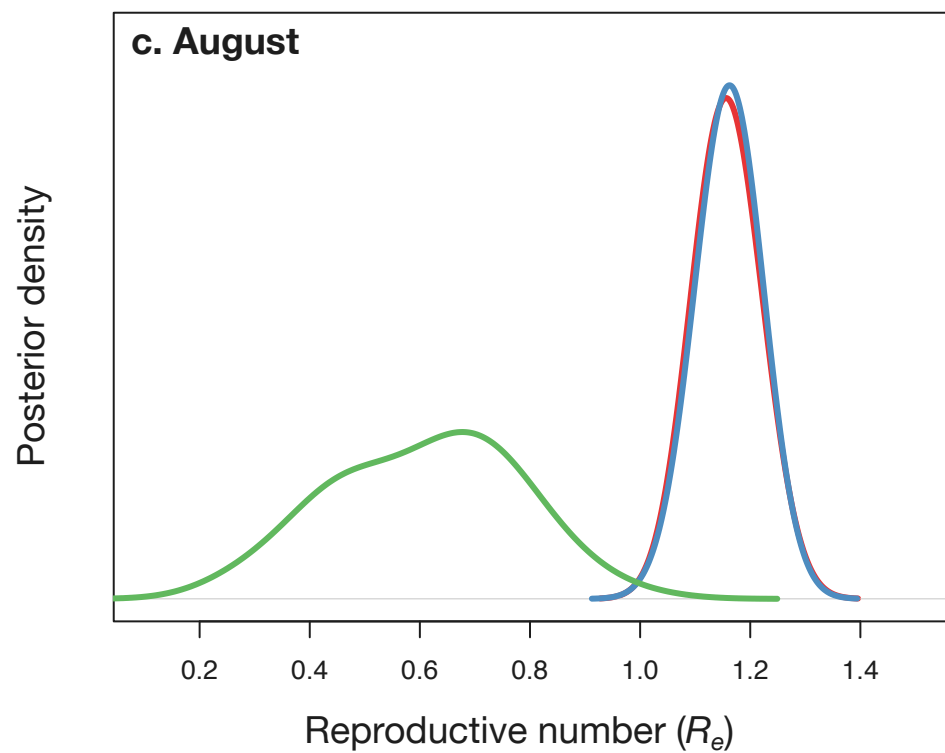
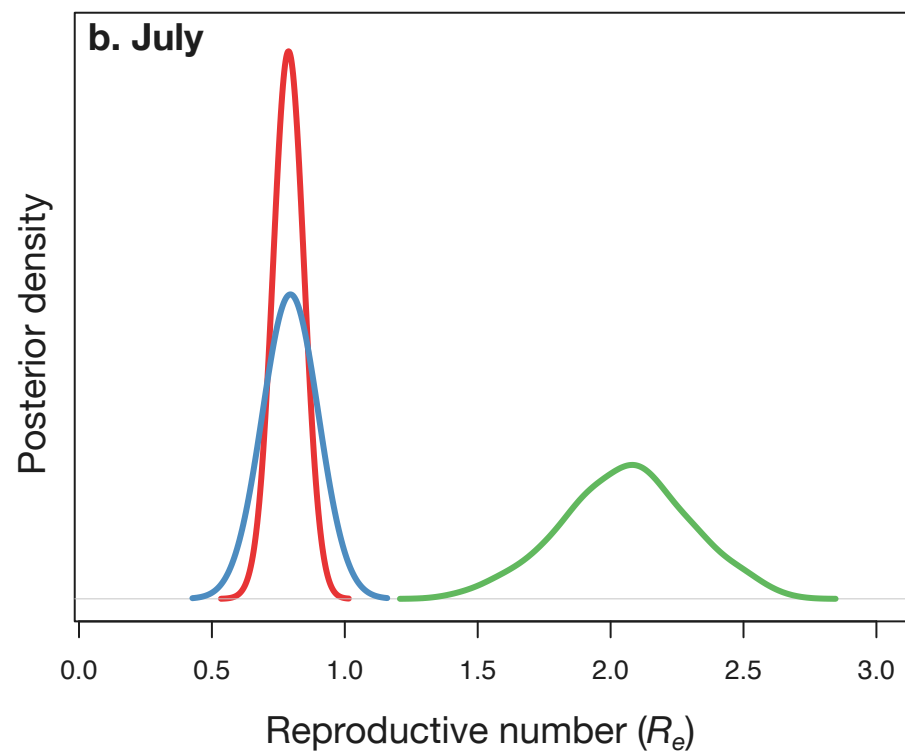
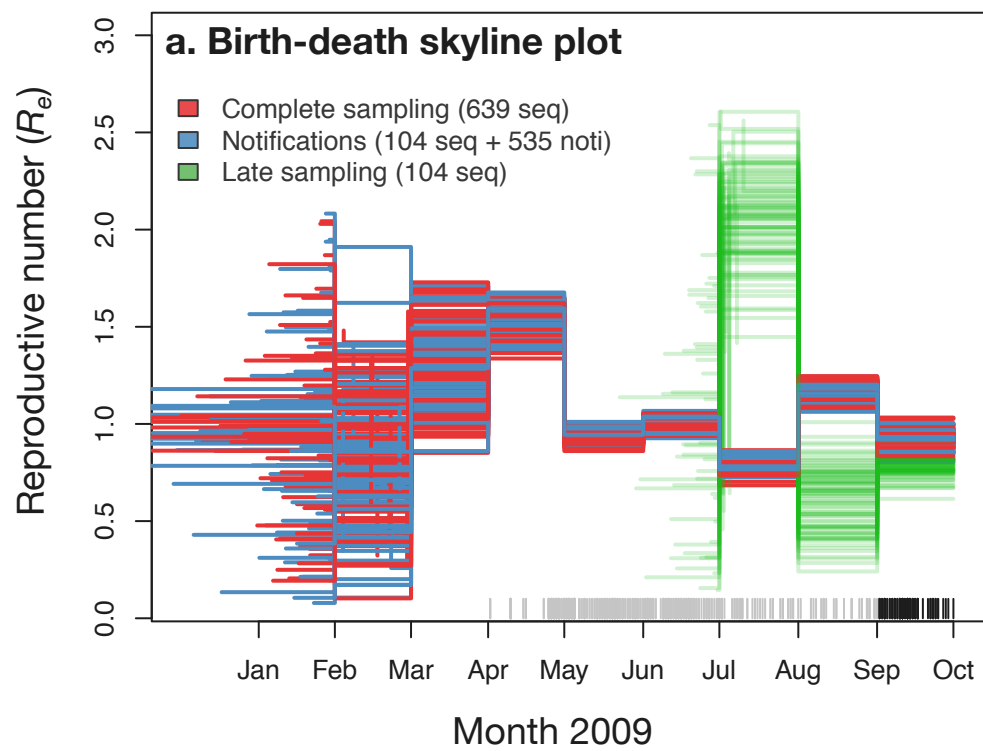
Notifications



● Sequenced samples

● Notifications (no sequence data available)





**Table 1.** Results of the simulation study with  $R_0=2$ , evolutionary rate of 0.01 subs/site/year, and late sampling starting at 0.75 years of a total time of 1 year. The rows correspond to the seven treatments. The first two columns denote the number of simulations (out of 100) where the value used to generate the data was contained within the 95% highest posterior density (HPD). The last two columns are a measure of precision of the estimates calculated as the estimated mean estimate of  $R_0$  and the evolutionary rate divided by the 95% HPD width, such that large values imply low precision. Here we report the mean value over 100 simulations.

	<b><math>R_0</math> within 95% HPD</b>	<b>Evol. rate within 95% HPD</b>	<b>Mean <math>R_0</math> / HPD width</b>	<b>Mean evol. rate / HPD width</b>
Late sampling BD const.	11	58	0.21	0.41
Late sampling BD skyline	96	96	0.28	0.51
Late sampling Coal. exp.	90	93	0.36	0.63
Constant sampling BD const.	97	98	0.18	0.27
Constant sampling Coal. exp.	93	97	0.23	0.29
Notifications BD const.	96	88	0.19	0.37
Notifications Coal. Exp	84	66	0.31	0.49