

Early phylodynamics analysis of the COVID-19 epidemics in France

Gonché Danesh¹, Baptiste Elie¹, Yannis Michalakis¹, Mircea T Sofonea¹,
Antonin Bal^{2,3}, Sylvie Behillil⁴,
Grégory Destras^{2,3}, David Boutolleau⁵, Sonia Burrel⁵, Anne-Geneviève Marcelin⁵,
Jean-Christophe Plantier⁶, Vincent Thibault⁷,
Etienne Simon-Lorier⁸, Sylvie van der Werf⁴, Bruno Lina^{2,3},
Laurence Josset^{2,3,*,+}, Vincent Enouf^{4,9,*,+}, Samuel Alizon^{1,*}
and the COVID SMIT PSL group

¹ Laboratoire MIVEGEC (UMR CNRS 5290, IRD, UM), Montpellier, France

² Laboratoire de Virologie des HCL, Institut des agents Infectieux, Centre National de Référence des virus des infections respiratoires (dont la grippe), Hôpital de la Croix-Rousse, Lyon, France

³ CIRI, Centre International de Recherche en Infectiologie, (Team VirPath), Univ Lyon, Inserm, U1111, Université Claude Bernard Lyon 1, CNRS, UMR5308, ENS de Lyon, F-69007, Lyon, France

⁴ National Reference Center for Respiratory Viruses, Molecular Genetics of RNA Viruses, CNRS-UMR 3569, The Institut Pasteur, Paris, France

⁵ Sorbonne Université, INSERM, Institut Pierre Louis d'Epidémiologie et de Santé Publique (iPLESP), AP-HP, Pitié Salpêtrière Hospital, Department of Virology, Paris, France

⁶ Laboratoire de Virologie, CHU Charles Nicolle, Rouen, France

⁷ Laboratoire de Virologie, CHU de Rennes, France

⁸ Evolutionary genomics of RNA viruses, Institut Pasteur, Paris, France

⁹ Mutualized Platform of Microbiology, Pasteur International Bioresources Network, The Institut Pasteur, Paris, France.

+ equal contribution

* authors for correspondence: laurence.josset@chu-lyon.fr, vincent.enouf@pasteur.fr, samuel.alizon@cnrs.fr

Abstract

France was one of the first countries to be reached by the COVID-19 pandemics. Here, we analyse 196 SARS-Cov-2 genomes collected between Jan 24 and Mar 24 2020, and perform a phylodynamics analysis. In particular, we analyse the doubling time, reproduction number (\mathcal{R}_t) and infection duration associated with the epidemic wave that was detected in incidence data starting from Feb 27. We show that a slowing down of the epidemic spread can be detected in Mar, which is consistent with the implementation of the national lock-down on Mar 17. The inferred distributions for the infection duration and \mathcal{R}_t are in line with those estimated from contact tracing data. Overall, this analysis shows the potential to use sequence genomic data to inform public health decisions in an epidemic crisis context.

1 Introduction

On Jan 8 2020, the Chinese Center for Disease Control announced that an outbreak of atypical pneumonia was caused by a novel coronavirus. The genetic sequence of what is now known as SARS-Cov-2 was released on Jan 10 (Liu et al., 2020). This was less than two weeks after the initial report of the outbreak by the Wuhan Health Commission, which took place on Dec 31 2019. Never has a novel pathogen been sequenced so rapidly.

The number of sequences in the databases grew rapidly thanks to an altruistic and international effort of virology departments all around the world gathered via the Global Initiative on Sharing All Influenza Data (GISAID, <https://www.gisaid.org/>). Early results allowed to better understand the origin of SARS-Cov-2 and identify a bat coronavirus (SARSr-CoV RaTG13) as its closest relative with more than 96% homology, as well as some potentially adaptive mutations (ICTV, 2020, Andersen et al., 2020, Xiao et al., 2020).

The available sequences were also analysed using the field of phylodynamics (Grenfell et al., 2004, Volz et al., 2013, Frost et al., 2015), which aims at inferring epidemiological processes from sequence data with known sampling dates. Most of these analyses were shared through the website virological.org. In particular, using 176 genomes from which he extracted 85 representative sequences (to avoid a potential cluster effect), Rambaut (2020) estimated the molecular clock to be approximately $8 \cdot 10^{-4}$ substitutions per position per year, with a 95% Highest Posterior Density (HPD) between $1.4 \cdot 10^{-4}$ and $1.3 \cdot 10^{-3}$ subst./pos./year, which yielded a date of origin of the outbreak mid-Nov 2019, with a 95% HPD spanning from Aug 27 to Dec 19. Further analysis with more recent

sequences found a median estimate of $1.1 \cdot 10^{-3}$ subst./pos./year with a similar HPD (Duchene et al., 2020). In their work, Scire et al. (2020) explored a variety of priors for the analysis and found similar orders of magnitude for the molecular clock estimate. They also applied a birth-death model to estimate several parameters among which the temporal reproduction number (\mathcal{R}_t) but a difficulty is that not all sequences originated from China and the sampling rate could also vary. Finally, Volz et al. (2020) performed one of the early analyses of the outbreak using coalescent models, allowing them to estimate the date of the origin of the epidemic in early Dec 2019 (with a 95% CI: between 6 Nov and 13 Dec 2019) and the doubling time of the epidemic to be 7.1 days (with a 95% CI: 3.0-20.5 days). These reports mention several caveats, which are due to the limited number of sequences, the limited amount of phylogenetic signal, the potentially unknown variations in sampling rates and the sampling across multiple countries.

Here, we focus on the COVID-19 epidemics in France by analysing 196 genomes sequenced from patients diagnosed in France that were available on Apr 4, 2020 thanks to the GISAID and to the collaborative effort of the French laboratories and National Reference Centers (CNR), who are listed in Supplementary Table S1.

The first COVID-19 cases were detected in France from Jan 24, 2020, mostly from travellers, but these remained isolated until Feb 27, when the national incidence curve of new COVID-19 cases started to increase steadily. Limited measures were announced on Feb 28, but schools were closed from Mar 16 and a nationwide lock-down was implemented from Mar 17. On Apr 19, the prime minister gave the first official estimate of the basic reproduction number (\mathcal{R}_0), which was 3.5, and of the temporal reproduction number after the lock-down, which was 0.5 (Salje et al., 2020).

Gambaro et al. (2020) provided a first picture of the general genomic structure of French COVID-19 epidemics using 97 genomes from samples collected in the north of France between Jan 24 and Mar 24, 2020. They identified several independent introductions of the virus in France but also found that the majority of the sequences belong to a major clade. This clade belongs to a larger clade labelled as G by GISAID, A2 by the nextstrain (<http://nextstrain.org/>) platform and B.1 following the dynamical taxonomy introduced by Rambaut et al. (2020). We refer to it as the clade related to the epidemic wave.

Here, we present early phylodynamics analyses on the French epidemic wave, focusing on the epidemic doubling time, the generation time and the temporal reproduction number $\mathcal{R}(t)$.

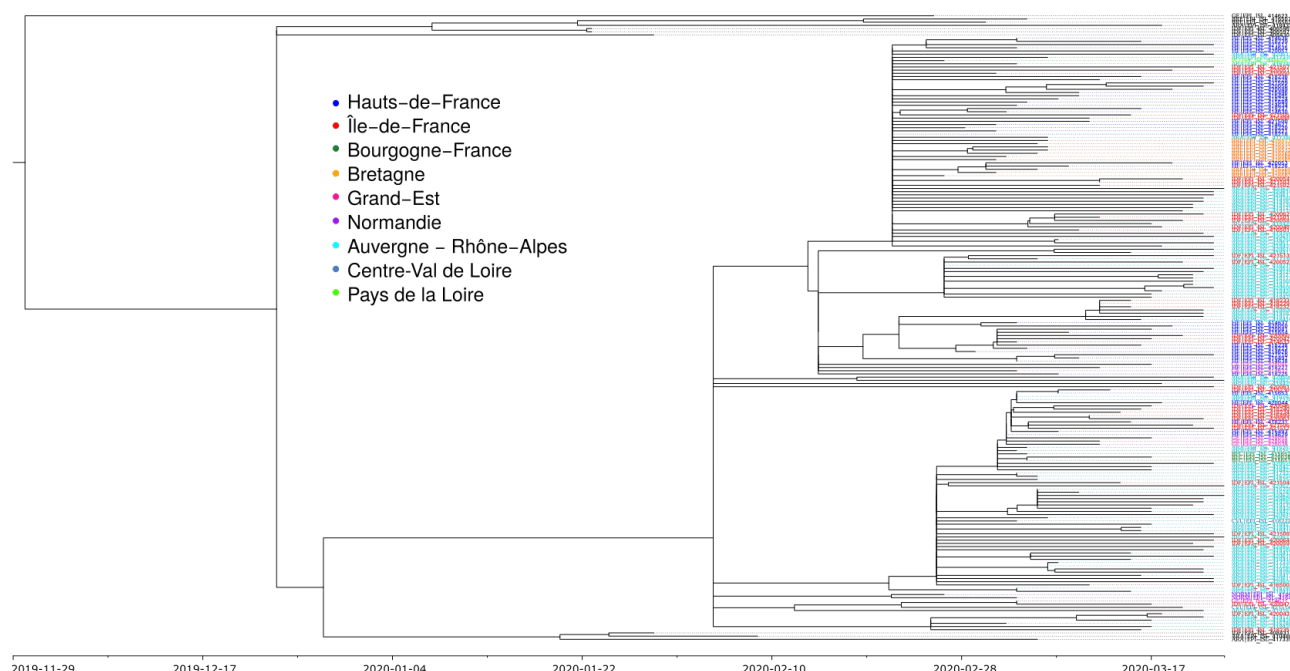


Figure 1 – **Phylogenetic structure of 196 SARS-Cov-2 genomes from France.** Color shows the French region of sampling. Sequences in black were removed from the analysis because outside the main clade corresponding to the epidemic wave.

2 Results

2.1 Phylogeny and regional structure

Figure 1 shows the regional structure of the French epidemics. Sequences corresponding to black leaves were ignored in the subsequent analyses because they do not belong to the main clade. Most of these originate from travelers isolated upon arrival in France, which explains their under-representation in the ongoing epidemic wave.

Focusing on the main clade, we see that all the leaves originate from a common branching event, which is approximately half-point of the phylogeny. The polytomy in this point can be due to a lack of phylogenetic signal. Another interpretation, could be independent introductions in France (up to 6 events). This could also be associated with a superspreading event, i.e. the infection of many people by the same person. Addressing this issue will require more sequences from the early stages of the epidemic wave since currently the earliest sequence in this major clade is from Feb 21, 2020.

Colors indicate the regional structure of the French epidemic. As expected, we see some regional clusters. We also see that sequences from the same region belong to different subclades of the major clade, which is consistent with multiple introductions or dispersal between regions. The general structure of the French epidemics will be the focus of a future study (see also the work by (Gambaro et al., 2020)).

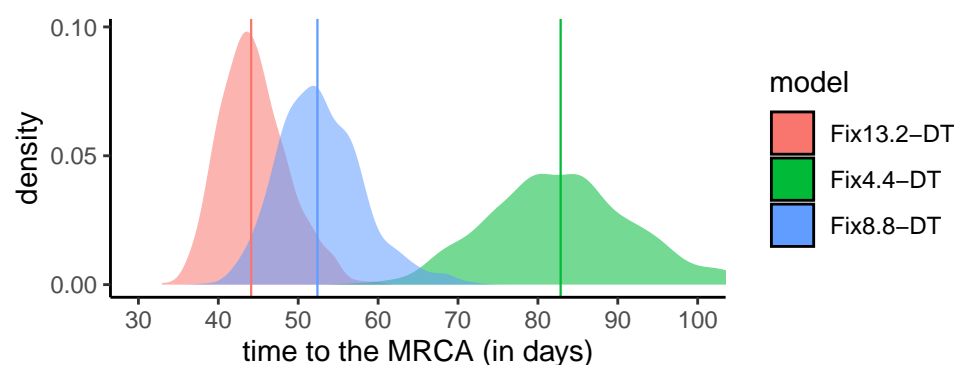


Figure 2 – Time to the origin of the French epidemic wave as a function of the molecular clock. This estimate was obtained assuming an exponential growth coalescent population model and a fix molecular clock (see Figure S4 for the BDSKY model). The slower the clock rate (in different colors), the further away in time the most recent common ancestor (MRCA). Vertical lines show the distribution medians. The most recent sample dates from Mar 24, 2020.

In the following of the work, we focus on the main clade associated with the epidemic wave.

2.2 Dating the epidemic wave

We first report the estimation of the time to the most recent common ancestor (TMRCA) of the 186 sequences that belong to the epidemic wave. Although this is the ancestor of the vast majority of the French sequences grouped in the B.1 clade (also referred to as G or A2 clade), the associated infection may have taken place outside France because the epidemic wave may be due to multiple introduction events (although from infections caused by similar viruses given the clustering).

Estimates of SARS-Cov-2 molecular clock should be treated with care given the limited amount of phylogenetic signal (Rambaut, 2020, Duchene et al., 2020). This is particularly true in our case since we are analysing a small subset of the data. In Appendix, we present the analysis of the temporal signal in the data using the TempEst software (Rambaut et al., 2016) and show that it strongly relies on early estimates that do not belong to the epidemic wave clade (Figure S2).

As shown in Figure 2, the molecular clock value directly affected the time to the most recent common ancestor for the coalescent model. This was also true for the BDSKY model, where the prior shape for the recovery rate had little impact (Figure S3). For both models, sampling of the 122 sequences amongst the 186 has a much smaller impact (Figure S5).

Table 1 shows the dates for models with different evolution rates and different population models (exponential coalescent or BDSKY). Note that smaller dataset may not include the most recent samples.

For most of our datasets and models, the origin for the clade corresponding to the sequences from

the French epidemic wave is dated between mid-Jan and early Feb. This large interval is due to the scarcity of "old" sequences (the first one collected in this clade dates from Feb 21) and on the fact that this clade averages the epidemics in several regions of France, which could have been seeded by independent introductions from outside France. The date provided by the slowest molecular clock (Fix4.4-DT) seems at odds with the data as we will see below.

To evaluate the effect of a potential sampling bias, we also estimate the time to the MRCA for 10 different sets of 122 sequences (Figure S5). We found similar median values for 9 of these 10 random datasets. Notice that the value of their parent dataset (France186), was slightly larger. For the BDSKY model, the effect was even less pronounced (Figure S4).

Overall, these dates (except for the slowest molecular clock) are consistent with those obtained by Rambaut (2020) regarding the beginning of the epidemic in China, which is dated November 17, 2019 with a confidence interval between Aug 27 and Dec 19, 2020. This interval is highly dependent on the number of available sequences as there are documented (but unsequenced) cases of COVID-19 in China early Dec 2019 (Li et al., 2020).

2.3 Doubling time

Using a coalescent model with exponential growth and serial sampling (Drummond et al., 2002), we can estimate the doubling time, which corresponds to the number of days for the epidemic wave to double in size. This parameter is key to calculate the basic reproduction number \mathcal{R}_0 (Wallinga and Lipsitch, 2007).

In Figure 3, we show this doubling time for four of our datasets that cover the whole (France122a), the first three quarters (France81), the first half (France61-1) and the second half (France61-2) of the time period. Since the first dataset includes more recent sequences than the second, which itself

Table 1 – Date of the most recent common ancestor of the clade corresponding to the French epidemic wave. Unless specified otherwise, the year is 2020. The "model" indicates the value of the molecular clock and the population dynamics model used (DT or BDSKY).

model	size	most recent sample	median value	95% HPD
Fix8.8-DT	122a	24 Mar	31 Jan	[19 Jan - 9 Feb]
Fix8.8-BDSKY	122a	24 Mar	31 Jan	[20 Jan - 11 Feb]
Fix13.2-DT	122a	24 Mar	8 Feb	[30 Jan - 15 Feb]
Fix4.4-DT	122a	24 Mar	1 Jan	[11 Dec 2019 - 17 Jan]
Fix8.8-DT	81	17 Mar	2 Feb	[17 Jan - 11 Feb]
Fix8.8-DT	61-1	12 Mar	03 Feb	[21 Jan - 12 Feb]
Fix8.8-DT	61-2	24 Mar	08 Feb	[25 Jan - 17 Feb]

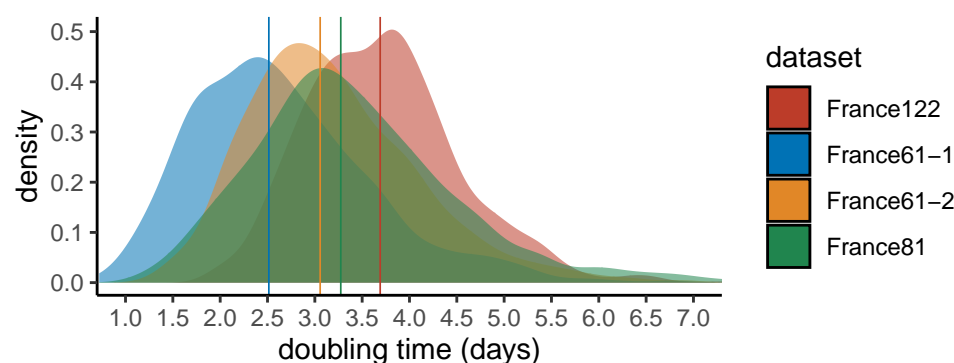


Figure 3 – Epidemic doubling time. We assume an exponential growth coalescent model with a fixed molecular clock. The four datasets differ in the sequences analysed (see the Methods). Vertical lines show the distribution medians.

includes more recent sequences than the third, our hypothesis is that we can detect variations in doubling time over the course of the epidemic.

Adding more recent sequence data indeed leads to an increase in epidemic doubling time. Initially, with the first 61 sequences (which run from Feb 21 to Mar 12), the epidemic spreads rapidly, with a median doubling time of 2.5 days. With the addition of sequences sampled between Mar 12 and 17, the doubling time increases to 3.3 days. Finally, by adding sequences sampled between Mar 17 and 24, the doubling time rises to 3.7 days.

Again, we explore the effect of sampling by estimating the doubling time on 10 different sets of 122 sequences and find a limited effect on the median value (Figure S7). Notice that the value of the parent dataset (France186), was slightly larger.

We also studied the effect of the molecular clock on the doubling time (Figure S6). As expected, the slower the molecular clock, the higher the doubling time. However, for our realistic molecular clocks, the effect is limited: the median is 3.4 days assuming a high value for the molecular clock and 3.7 days for our default (medium) value. The low value of the molecular clock led to a high median doubling time of 5.6 days. This is at odds with the incidence data in France, which indicates an exponential growth rate of 0.23 days^{-1} which corresponds to a doubling time of 3 days, suggesting that our default molecular clock is more realistic.

In comparison, phylodynamic inferences made from data from China with 86 genomes (Rambaut, 2020) found a median doubling time of about 7 days with a confidence interval between 4.7 and 16.3 days). One reason for the slower growth rate of the epidemic compared to ours is that we have focused on one rapidly expanding clade of the epidemic and neglected the smaller clades. Another possibility could be related to the timing of the sampling (early or late in the infection).

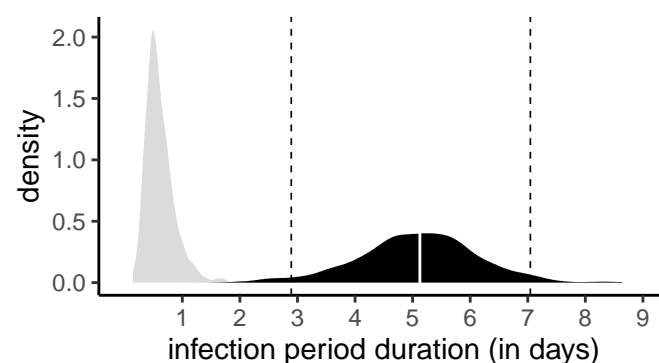


Figure 4 – **Distribution of infection duration.** The prior distribution is shown in gray, and the posterior distribution in black. The white line shows the distribution median and the dashed line the 95% highest posterior density (HPD).

2.4 Duration of infectious period

The birth-death skyline (BDSKY) model (Stadler et al., 2013) allows us to estimate the duration of contagiousness and the reproduction number of the epidemic (i.e. the number of secondary infections caused by an infected host). The exponential growth coalescent model described above cannot distinguish between these two quantities. However, the BDSKY model requires more parameter values to be estimated.

From this rate, we can obtain the duration of contagiousness, knowing that it is also necessary to account for the sampling rate because patients whose infections are sequenced can be assumed not to transmit the infection after this detection. The sampling rate after Feb 21 (it is set to 0 before that date) is estimated at 0.093 days^{-1} with a (wide) 95% confidence interval between 0.006 and 0.627 days^{-1} . If we analyse this in days, the median value of the distribution yields 10.8 days and is consistent with the fact that in the French epidemics most of the screening for SARS-Cov-2 is done on severe cases upon hospital admission.

The distribution of infectious periods is obtained by taking the inverse of the sum of the sampling rate and the contagiousness end rate. The median of this distribution is 5.12 days and 95% of its values are between 2.89 and 7.05 days (Figure 4).

In Supplementary Figure S8, we show that the estimate for the infectious period duration is sensitive to the shape of the prior assumed for the recovery rate. Indeed, if we use a, less informative, uniform prior, then the median sampling rate estimate is larger and the median infectious period estimated is shorter.

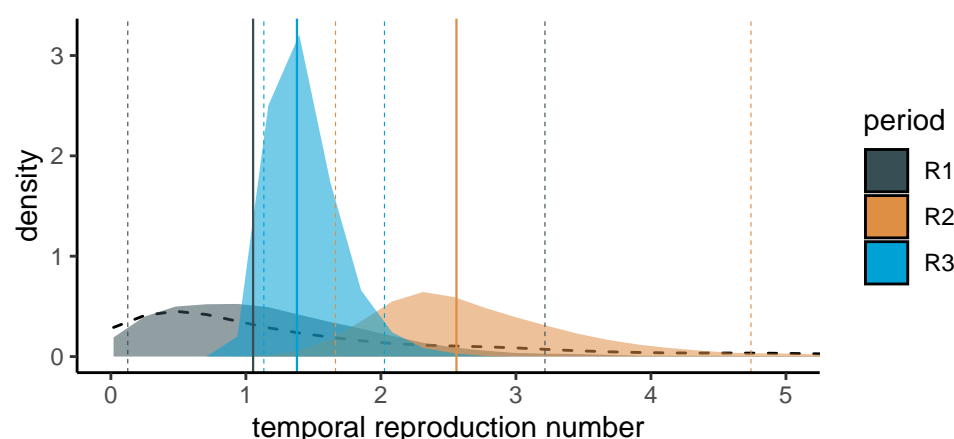


Figure 5 – **Temporal reproduction numbers inferred using the BDSKY model.** These results are obtained for the France186 dataset. The black dashed curve show the prior distribution, and the posterior distributions are in color. Vertical plain lines show distribution medians, while vertical dashed lines indicate the 95% highest posterior density (HPD).

2.5 Reproduction number

With the BDSKY model, we can estimate the temporal reproductive number, noted $\mathcal{R}(t)$, since the onset of the epidemic wave. Here, given the limited temporal signal, we only divided the time into 3 intervals to estimate three reproduction numbers: \mathcal{R}_1 before Feb 19, \mathcal{R}_2 between Feb 19 and Mar 7, and \mathcal{R}_3 between Mar 7 and Mar 24.

These results are very consistent with those obtained for the doubling time, even if the time periods are different. For the period before Feb 19, the estimate is the least accurate with values of \mathcal{R}_1 with a median of 1.05 but at 95% Highest Posterior Density (HPD) between 0.13 and 3.22. The lack of information can be seen in Figure 5 as the posterior distribution (gray area) is very similar to the prior (dashed curve). This is consistent with the fact that the oldest sequence dates from Feb 21, while the tree root is estimated at the beginning of Feb. Over the second time period (in orange), the distribution shape is similar to that of the prior but the median is very different and rapid growth is detected with a median value of \mathcal{R}_2 of 2.56 (95% HPD between 1.66 and 4.74). Finally, the most recent period after Mar 7 is the most accurate and detects a slowing down of the epidemic with a \mathcal{R}_3 of 1.38 (95% HPD between 1.13 and 2.03).

In Appendix, we show that these estimates for \mathcal{R}_t are robust to the prior used for the recovery rate (Figure S9). They are also robust to the sampling of 122 of the 186 sequences (Figure S10).

3 Discussion

Analysing SARS-Cov-2 genome sequences with a known date of sampling allows to infer phylogenies of infections and to estimate the value of epidemiological parameters of interest (Volz et al., 2013, Frost et al., 2015). We performed this analysis based on the 196 sequences sampled in France and available on Apr 4, 2020. We focus in particular on the largest clade regrouping 186 of the most recent sequences and likely corresponding to the epidemic wave that peaked in France early Apr 2020.

Before summarizing the results, we prefer to point out several limitations of our analysis. First, the French clade we analysed is in fact an international clade: although most French sequences appear to be grouping into two main subclades within this clade, it is possible that the variations in epidemic growth that we detect are more due to European than French control policies. Second, some French regions (e.g. Auvergne-Rhône-Alpes) are more represented than others (e.g. Occitanie is absent), which could bias the analysis at the national level. Finally, the molecular clock had to be set in this analysis because we do not have sufficient sampling during the month of Feb in France.

Despite these limitations, our results obtained early Apr confirm a slowing down of the epidemic in France, where the epidemic peak in terms of ICU admissions was reached on Apr 1. Indeed, by adding sequences sampled between Mar 12 and 24 to the phylogeny, the doubling time of the epidemic estimated by an exponential growth coalescent model increased by 48%. This slowdown is more clearly detected using a birth death model via the temporal reproduction number $\mathcal{R}(t)$: the median value decreased by 41% after Mar 12. This is consistent with the implementation of strict control measures in France as of Mar 17. These variations and even these orders of magnitude are consistent with our estimates based on the time series of incidence of new hospitalizations and deaths (unpublished reports, <http://covid-ete.ouvaton.org>).

Finally, the BDSKY model also provides us with an estimate of the duration of the infectious period, which is essential in the calculation of \mathcal{R}_0 (Wallinga and Lipsitch, 2007). The result we obtain, with a 95% Highest Posterior Distribution between 3 and 7 days and a median of 5.2, is highly relevant biologically and comparable to results obtained using contact tracing data. For instance, Ferretti et al. (2020) estimated a serial interval with a median of 5 days and a standard deviation of 1.9 days. To date, there is no estimate of the serial interval in France.

By increasing the number of SARS-CoV-2 genomic sequences from the French epidemic (and the number of people working on the subject), in particular sequences collected at the beginning of the epidemic, it would be possible to better estimate the date at which the epidemic wave took

210 off in France, improve the estimate for the infection generation time and the reproduction number, better understand the spread between the different French regions, and estimate the number of virus introductions into the country.

Finally, it is important to set these results into their context. As acknowledge in the introduction, the French state only acknowledged the magnitude of the COVID-19 epidemics on the last days of 215 Feb 2020 and these genomes were mostly collected between Feb 21 and Mar 24. Most of this analysis was published on Apr 6. At this time, the epidemic peak was barely noticeable in the incidence data. Furthermore, the serial interval, which is used to estimate the generation time of the infection and classically measured from contact tracing data, is still unknown in France. These results illustrate the contribution phylodynamics can make to public health during a crisis.

220 **Materials and Methods**

3.1 Data and quality check

On Apr 4, 196 sequences were available from samples originating from France via the Global Initiative on Sharing All Influenza Data (GISAID, <https://www.gisaid.org/>) thanks to the work of the two Centre National de Référence and local virology laboratories. These sequences only provide a 225 partial view of the epidemics as they originate from 8 the 18 French regions (Figure S1). Sequences were aligned and cleaned using the Augur pipeline developed by nextstrain (Hadfield et al., 2018). One sequence was removed due to low quality. The list of the sequences used is shown in Supplementary Table S1.

We screened the dataset with RDP4 (Martin et al., 2015), which did not detect any recombination 230 events.

3.2 Phylogenetic inference

We first performed a maximum likelihood inference of the phylogeny using SMS (Lefort et al., 2017) and PhyML (Guindon and Gascuel, 2003). The mutation model inferred by SMS was GTR.

We then used Beast 1.8.3 (Drummond et al., 2012) to perform inference using a bayesian ap- 235 proach. More specifically, we assumed an exponential coalescent for the population model (Drummond et al., 2002). We used the default settings for the model, which correspond to a gamma distribution for the growth rate prior $\Gamma(0.001, 1000)$ and an inverse prior for the population size $1/x$ (see

Supplementary Methods S2).

We also used Beast 2.3 (Bouckaert et al., 2014) to estimate key parameters using the birth-death skyline (BDSKY) model (Stadler et al., 2013, Kühnert et al., 2014). One of these parameters is the temporal reproduction number (\mathcal{R}_t) and we here assume three periods in the epidemics (which means we estimate 3 values \mathcal{R}_1 , \mathcal{R}_2 and \mathcal{R}_3). Another parameter is the rate at which the infectiousness ends in absence of treatment. The final key parameter is the sampling rate, the inverse of which corresponds to the average number of days until an infected person is sampled. The ratio between the sampling rate and the sum of the sampling and end of infection rates indicates the fraction of infections that are actually sampled. By sampled, we mean that the patient is identified and the virus population causing the infection is sequenced. Note that we assume sampled hosts are not infectious anymore. We considered multiple priors for the rate of end of the infectious period by setting a lognormal prior LogNorm(90,0.5) and a uniform prior Unif(5,350). We assumed a beta prior $\beta(1.0, 1.0)$ for the sampling rate (see Supplementary Methods S2). As previous models (Stadler et al., 2013), we set the sampling rate to 0 before the first infected host is sampled (here on Feb 21, 2020).

For both analyses in Beast, we assumed a GTR mutation model, following the results of SMS. We also assumed a uniform prior $U(0, 1)$ for the nucleotide frequencies and a lognormal prior for parameter κ , LogNorm(1, 1.25).

Regarding the molecular clock, earlier studies have reported a limited amount of phylogenetic signal in the first sequences from the COVID-19 pandemics. Given that we here focus on a subset of these sequences, we chose to fix the value of the strict molecular clock to $8.8 \cdot 10^{-4}$ substitutions/position/year, following the analysis by Rambaut (2020). In Appendix, we study the influence of this value on the results by setting it to a lower ($4.4 \cdot 10^{-4}$ subst./pos./year) or a higher ($13.2 \cdot 10^{-4}$ subst./pos./year) value. Finally, we also estimate this parameter assuming a strict molecular clock. The most recent estimates suggest that the intermediate and high value are the most realistic ones (Duchene et al., 2020).

3.3 Data subsets

We analysed subsets of the whole data set. Our largest subset, excluded 10 sequences that did not belong to the French epidemic wave clade and therefore contained 186 sequences. Figure S1 shows the sampling date and French region of origin for each sequence.

Some sampling dates are over-represented in the dataset, which could bias the estimation of di-

vergence times (Seo et al., 2002, Stadler et al., 2012). To correct for this, we sampled 6 sequences for each of the days where more than 6 sequences were available. This was done 10 times to generate 10 datasets with 122 sequences (France122a to France122h).

To investigate temporal effects using the coalescent model, we created three other subsets of the France122a dataset: "France61-1" contains the 61 sequences sampled first (i.e. from Feb 21 to Mar 12), "France61-2" contains the 61 sequences sampled more recently (i.e. from Mar 12 to Mar 24), and "France81" contains the 81 sequences sampled first (i.e. from Feb 21 to Mar 17).

With the exponential coalescent model (denoted DT for "doubling time"), we analysed all subsets of data (Fra61-1, Fra61-2, Fra81, and all the 10 Fra122 datasets), whereas for the BDSKY model we show the main dataset (Fra186) and analyse the 10 subsets with 122 leaves in Appendix.

Acknowledgments

We thank the two Centre Nationaux de Référence (Nord and Sud) in France for generating the sequence data and sharing it via GISAID.

Gonché Danesh is supported by a doctoral grant from the Fondation pour la Recherche Médicale (FRM grant number ECO20170637560). We thank the ETE modelling group for discussion. We also thank the patients, nurses, doctors and all the French laboratories who made this work possible by generating and sharing the virus genome sequences.

We acknowledge the IRD itrop high-performance computer (South Green Platform) at IRD Montpellier for providing resources that have contributed to the results presented in this work (more details on bioinfo.ird.fr).

Preliminary versions of this report were posted on Apr 9 (in French on <http://covid-ete.ouvaton.org>) and on Apr 21 (in English on <http://virological.org/>).

References

- Andersen, K. G., A. Rambaut, W. I. Lipkin, E. C. Holmes and R. F. Garry. 2020. The proximal origin of SARS-CoV-2. *Nat Med* pp. 1–3. Publisher: Nature Publishing Group.
- Bouckaert, R., J. Heled, D. Kühnert, T. Vaughan, C.-H. Wu, D. Xie, M. A. Suchard, A. Rambaut and A. J. Drummond. 2014. BEAST 2: A Software Platform for Bayesian Evolutionary Analysis. *PLoS Computational Biology* 10(4):e1003537.
- Drummond, A. J., G. K. Nicholls, A. G. Rodrigo and W. Solomon. 2002. Estimating mutation parameters, population history and genealogy simultaneously from temporally spaced sequence data. *Genetics* 161(3):1307–20.
- Drummond, A. J., M. A. Suchard, D. Xie and A. Rambaut. 2012. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol Biol Evol* 29(8):1969–73.
- Duchene, S., L. Featherstone, M. Haritopoulou-Sinanidou, A. Rambaut, P. Lemey and G. Baele. 2020. Temporal signal and the phylodynamic threshold of SARS-CoV-2. *bioRxiv* p. 2020.05.04.077735. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- Ferretti, L., C. Wymant, M. Kendall, L. Zhao, A. Nurtay, L. Abeler-Dörner, M. Parker, D. Bonsall and C. Fraser. 2020. Quantifying SARS-CoV-2 transmission suggests epidemic control with digital contact tracing. *Science*. Publisher: American Association for the Advancement of Science Section: Research Article.
- Frost, S. D., O. G. Pybus, J. R. Gog, C. Viboud, S. Bonhoeffer and T. Bedford. 2015. Eight challenges in phylodynamic inference. *Epidemics* 10:88–92.
- Gambaro, F., A. Baidaliuk, S. Behillil, F. Donati, M. Albert, A. Alexandru, M. Vanpeene, M. Bizard, A. Brisebarre, M. Barbet, F. Derrar, S. v. d. Werf, V. Enouf and E. Simon-Loriere. 2020. Introductions and early spread of SARS-CoV-2 in France. *bioRxiv* p. 2020.04.24.059576. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- Grenfell, B. T., O. G. Pybus, J. R. Gog, J. L. Wood, J. M. Daly, J. A. Mumford and E. C. Holmes. 2004. Unifying the epidemiological and evolutionary dynamics of pathogens. *Science* 303(5656):327–32.
- Guindon, S. and O. Gascuel. 2003. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst. Biol.* 52(5):696–704.

Hadfield, J., C. Megill, S. M. Bell, J. Huddleston, B. Potter, C. Callender, P. Sagulenko, T. Bedford and R. A. Neher. 2018. Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* 34(23):4121–4123.

ICTV 2020. The species Severe acute respiratory syndrome-related coronavirus : classifying 2019-nCoV and naming it SARS-CoV-2. *Nat Microbiol* pp. 1–9.

Kühnert, D., T. Stadler, T. G. Vaughan and A. J. Drummond. 2014. Simultaneous reconstruction of evolutionary history and epidemiological dynamics from viral sequences with the birth-death SIR model. *J R Soc Interface* 11(94):20131106.

Lefort, V., J.-E. Longueville and O. Gascuel. 2017. SMS: Smart Model Selection in PhyML. *Molecular Biology and Evolution* 34(9):2422–2424.

Li, Q., X. Guan, P. Wu, X. Wang, L. Zhou, Y. Tong, R. Ren, K. S. M. Leung, E. H. Y. Lau, J. Y. Wong, X. Xing, N. Xiang, Y. Wu, C. Li, Q. Chen, D. Li, T. Liu, J. Zhao, M. Liu, W. Tu, C. Chen, L. Jin, R. Yang, Q. Wang, S. Zhou, R. Wang, H. Liu, Y. Luo, Y. Liu, G. Shao, H. Li, Z. Tao, Y. Yang, Z. Deng, B. Liu, Z. Ma, Y. Zhang, G. Shi, T. T. Y. Lam, J. T. Wu, G. F. Gao, B. J. Cowling, B. Yang, G. M. Leung and Z. Feng. 2020. Early Transmission Dynamics in Wuhan, China, of Novel Coronavirus–Infected Pneumonia. *New England Journal of Medicine*.

Liu, Y., R. M. Eggo and A. J. Kucharski. 2020. Secondary attack rate and superspreading events for SARS-CoV-2. *The Lancet* 395(10227):e47.

Martin, D. P., B. Murrell, M. Golden, A. Khoosal and B. Muhire. 2015. RDP4: Detection and analysis of recombination patterns in virus genomes. *Virus Evolution* 1(1):vev003.

Rambaut, A. (2020). Phylodynamic Analysis | 176 genomes | 6 Mar 2020. Library Catalog: virological.org.

Rambaut, A., E. C. Holmes, V. Hill, A. O’Toole, J. T. McCrone, C. Ruis, L. d. Plessis and O. G. Pybus. 2020. A dynamic nomenclature proposal for SARS-CoV-2 to assist genomic epidemiology. *bioRxiv* p. 2020.04.17.046086. Publisher: Cold Spring Harbor Laboratory Section: New Results.

Rambaut, A., T. T. Lam, L. Max Carvalho and O. G. Pybus. 2016. Exploring the temporal structure of heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evolution*. Publisher: Oxford Academic.

Salje, H., C. T. Kiem, N. Lefrancq, N. Courtejoie, P. Bosetti, J. Paireau, A. Andronico, N. Hoze, J. Richet, C.-L. Dubost, Y. L. Strat, J. Lessler, D. Bruhl, A. Fontanet, L. Opatowski, P.-Y. Boëlle and S. Cauchemez. (2020). Estimating the burden of SARS-CoV-2 in France.

Scire, J., T. G. Vaughan and T. Stadler. 2020. Phylodynamic analyses based on 93 genomes.

350 Seo, T.-K., J. L. Thorne, M. Hasegawa and H. Kishino. 2002. A viral sampling design for testing the molecular clock and for estimating evolutionary rates and divergence times. *Bioinformatics* 18(1):115–23.

Stadler, T., D. Kühnert, S. Bonhoeffer and A. J. Drummond. 2013. Birth-death skyline plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc Natl Acad Sci*
355 USA 110(1):228–33.

Stadler, T., R. Kouyos, V. von Wyl, S. Yerly, J. Böni, P. Bürgisser, T. Klimkait, B. Joos, P. Rieder, D. Xie, H. F. Günthard, A. J. Drummond, S. Bonhoeffer and the Swiss HIV Cohort Study. 2012. Estimating the Basic Reproductive Number from Viral Sequence Data. *Mol Biol Evol* 29(1):347–357.

360 Volz, E., M. Baguelin, S. Bhatia, A. Boonyasiri, A. Cori, Z. Cucunubá, G. Cuomo, C. A. Donnelly, I. Dorigatti, R. FitzJohn, H. Fu, K. Gaythorpe, A. Ghani, W. Hinsley, N. Imai, D. Laydon, G. Nedjati-Gilani, L. Okell, S. Riley, S. van, H. Wang, Y. Wang, X. Xi and N. M. Ferguson. 2020. Report 5: Phylogenetic analysis of SARS-CoV-2. p. 7.

Volz, E. M., K. Koelle and T. Bedford. 2013. Viral phylodynamics. *PLoS Comput Biol*
365 9(3):e1002947.

Wallinga, J. and M. Lipsitch. 2007. How generation intervals shape the relationship between growth rates and reproductive numbers. *Proc. R. Soc. Lond. B* 274:599–604.

Xiao, K., J. Zhai, Y. Feng, N. Zhou, X. Zhang, J.-J. Zou, N. Li, Y. Guo, X. Li, X. Shen, Z. Zhang, F. Shu, W. Huang, Y. Li, Z. Zhang, R.-A. Chen, Y.-J. Wu, S.-M. Peng, M. Huang, W.-J. Xie, Q.-H. Cai, F.-H. Hou, W. Chen, L. Xiao and Y. Shen. 2020. Isolation of SARS-CoV-2-related coronavirus
370 from Malayan pangolins. *Nature* pp. 1–7. Publisher: Nature Publishing Group.

Appendix

S1 Sequences used

See attached .csv file.

375 S2 BEAST priors

MCMC chains were run for $5 \cdot 10^8$ iterations. The first 10% runs were discarded as a burn in and convergence was assessed using Effective Sample Size (ESS). All parameters had ESS greater than 200.

Original XML files cannot be shared due to the GISAID agreement.

Table S1 – Prior summary for the exponential coalescent model

Parameter	Value
Molecular clock	fixed
Evolution model	GTR
kappa	LogNormal(1,1.25)
frequencies	Uniform(0,1]
popsiz	1/x
growth rate	Gamma(0.001,1000)

Table S2 – Prior summary for the BDSKY model

Parameter	Value
Molecular clock	fixed
Evolution model	GTR
kappa	LogNormal(1,1.25)
frequencies	Uniform[0,1]
Rate of end of infection	Uniform(1.2, ∞) or LogNormal(0,1.2)
Sampling rate	Beta(1,1)
Reproduction number	LogNormal(0,1.2) with maximum at 10

380 S3 Supplementary figures

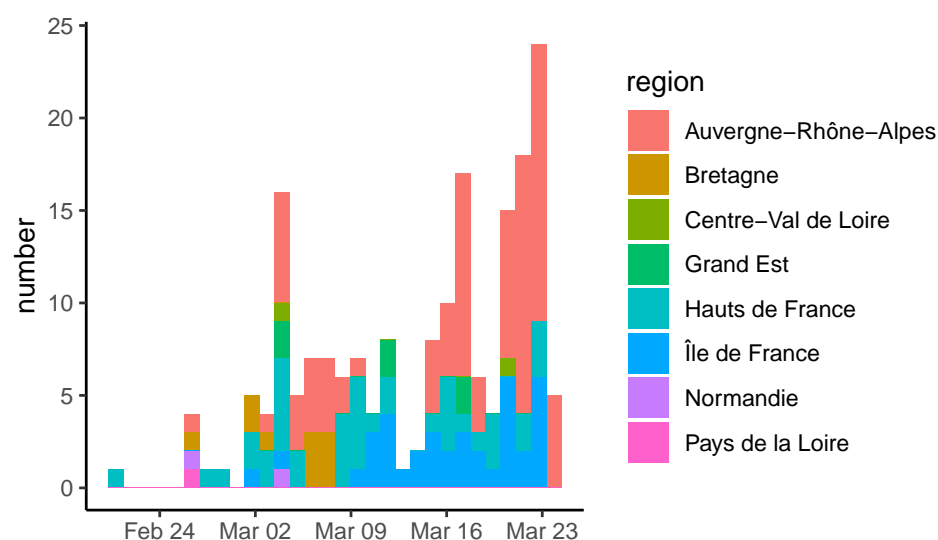


Figure S1 – **Sampling date and region.** List of samples collected, analysed and shared via GISAID by the two French National Reference Centers (CNR) as of Apr 4, 2020.

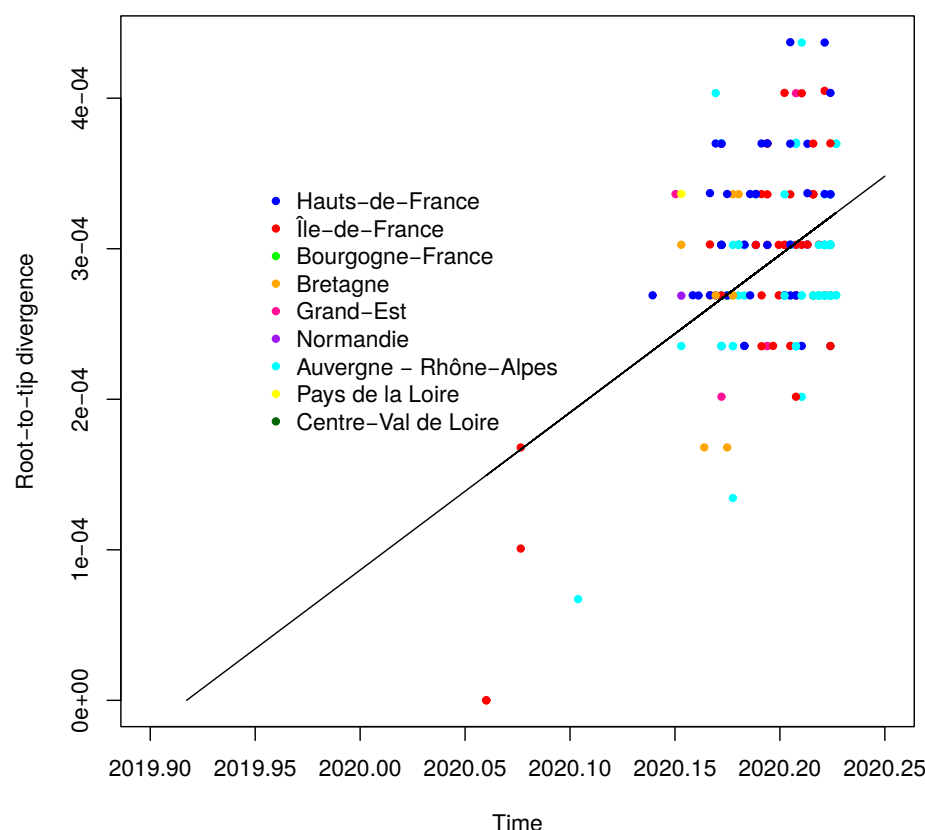


Figure S2 – **Root-to-tip correlation.** We analyse a phylogeny based on all 196 French sequences (i.e. not only that from the epidemic wave). The four earliest cases in Jan and early Feb were all isolated and belong to another clade than the rest of the sequences. The figure was obtained using TempEst (Rambaut et al., 2016).

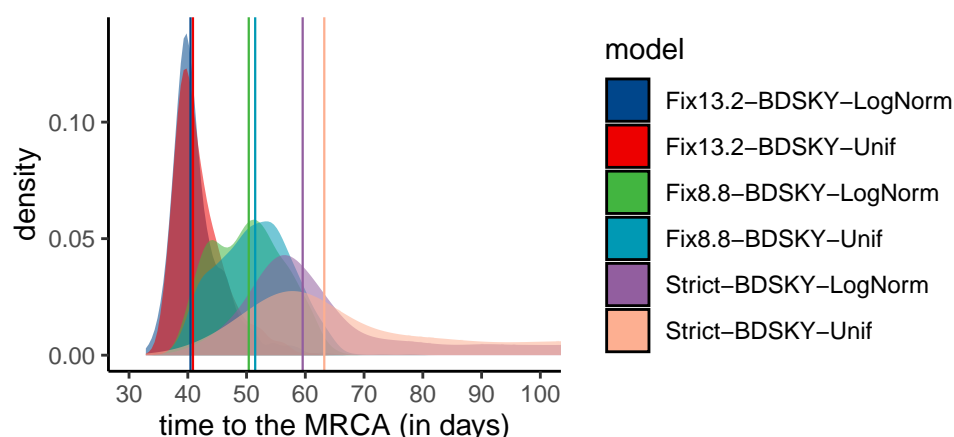


Figure S3 – **Time to the MRCA as a function of the molecular clock and of the recovery rate prior.** Here we assume BDSKY model.

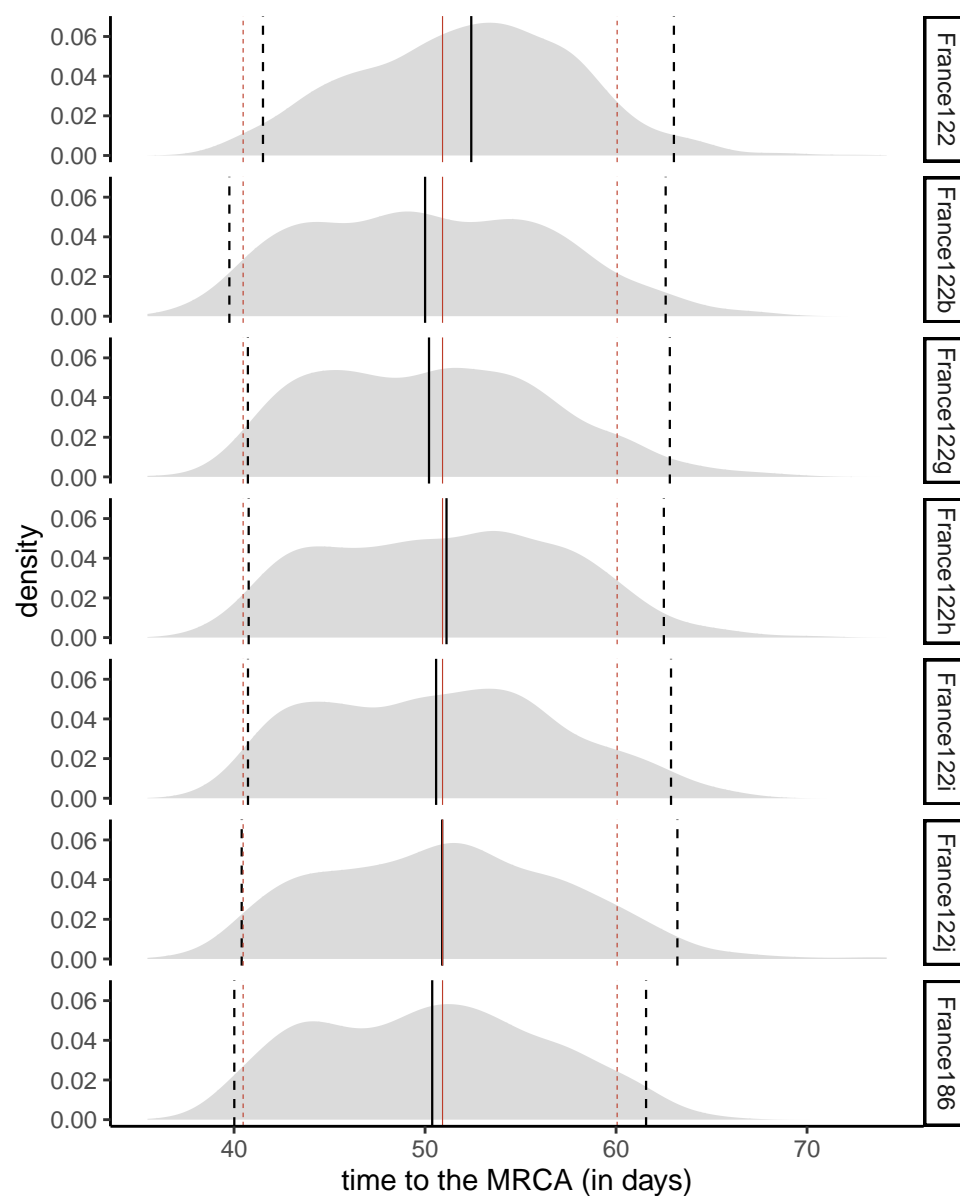


Figure S4 – **Time to the MRCA for the France186 dataset and subsets with 122 sequences.** Here we assume BDSKY model.

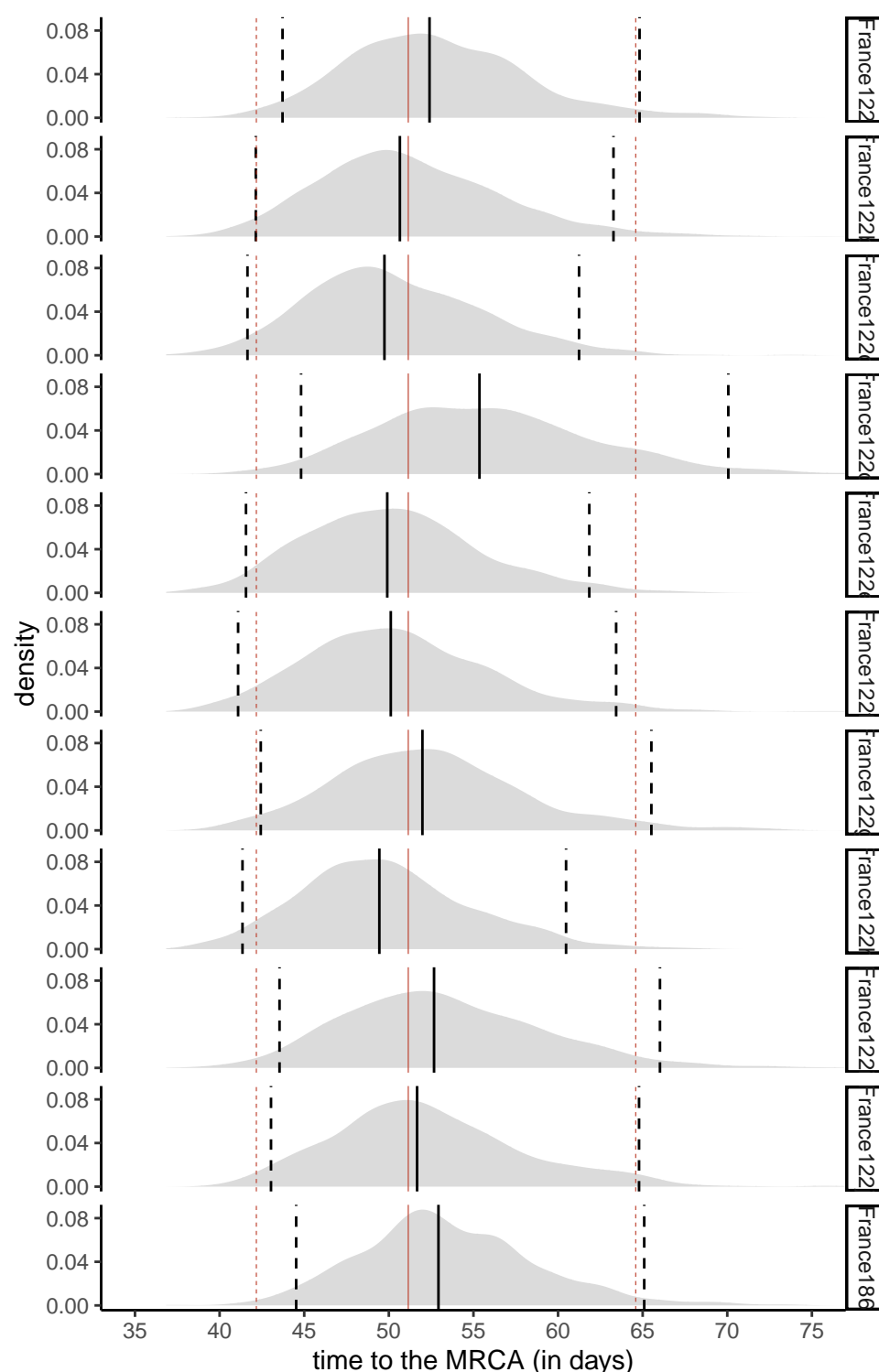


Figure S5 – **Time to the MRCA the 10 trees with 122 leaves sampled.** The red lines show the quantiles (0.025, 0.5 and 0.975) for the average of the 10 datasets. The black line shows the quantile for each dataset. The last panel shows the largest phylogeny with 186 leaves (which slightly overestimates the time to the MRCA).

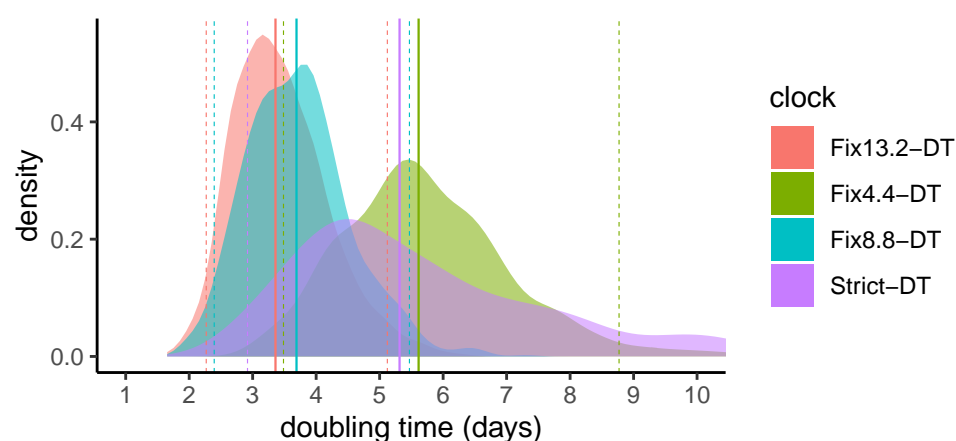


Figure S6 – **Effect of the molecular clock on the doubling time.** Note that in the "Strict" model we infer the value using a strict molecular clock but convergence is limited. Here, an exponential growth coalescent model is assumed.

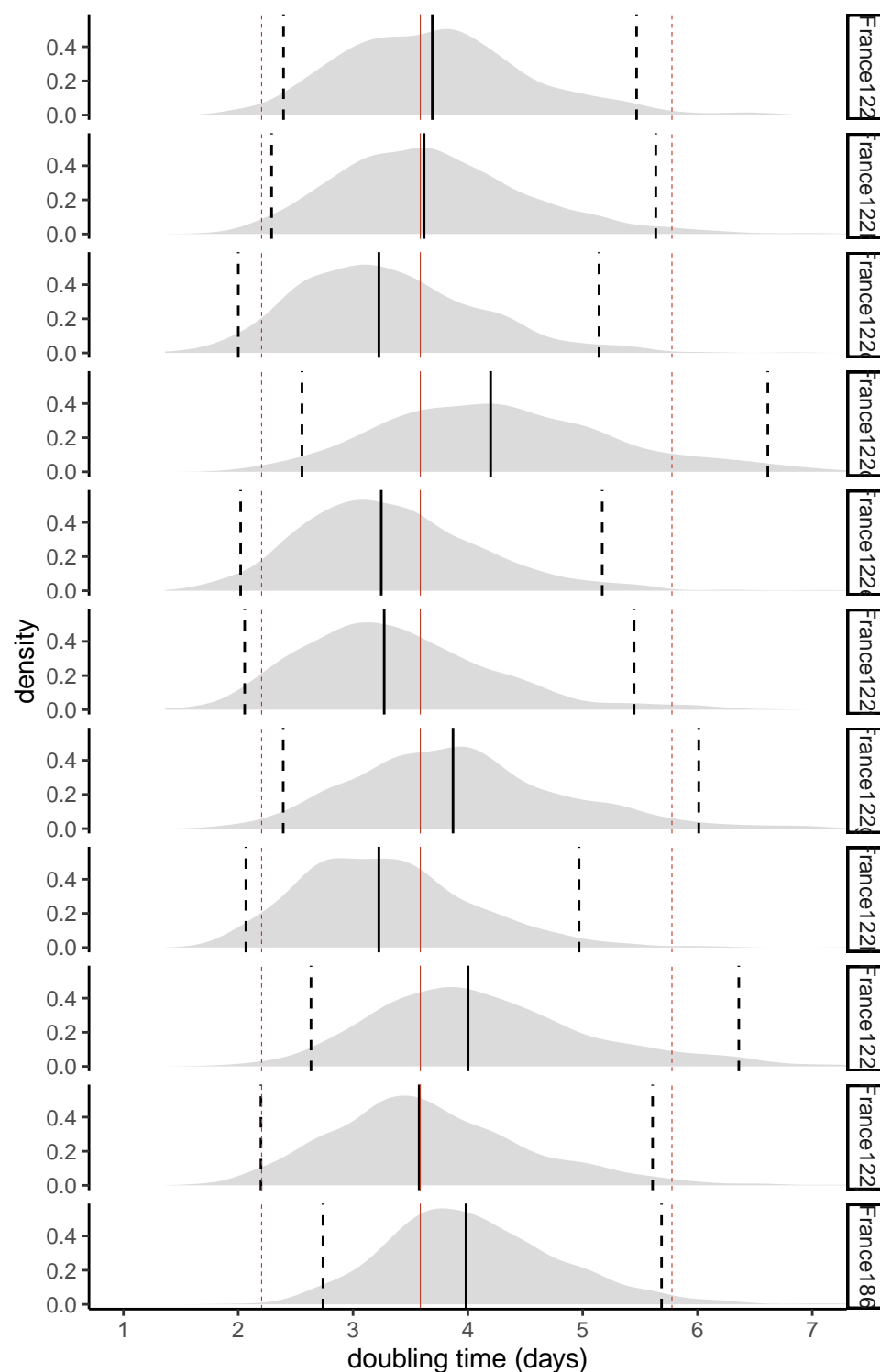


Figure S7 – **Doubling time for the 10 trees with 122 leaves sampled.** The red lines show the quantiles (0.025, 0.5 and 0.975) for the average of the 10 datasets. The black line shows the quantile for each dataset. The last panel shows the largest phylogeny with 186 leaves (which slightly overestimates the doubling time).

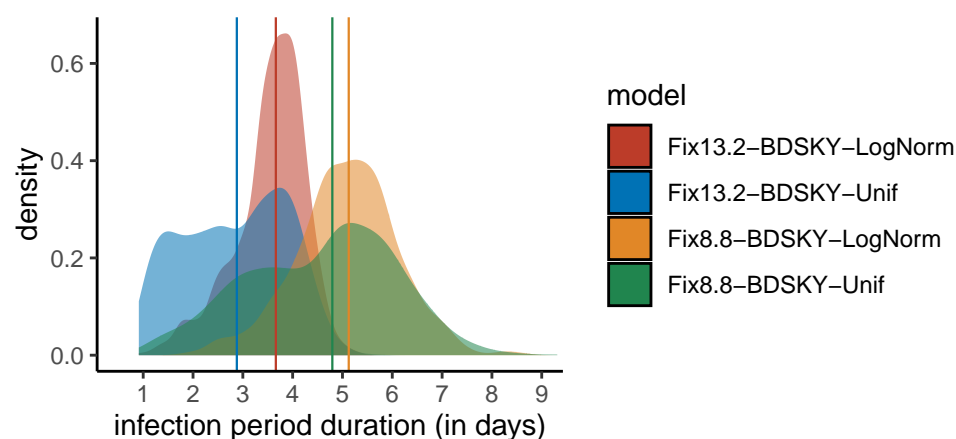


Figure S8 – **Effect of the prior shape and of the molecular clock on the infection duration estimate.** The molecular clock value has a stronger effect than the prior shape.

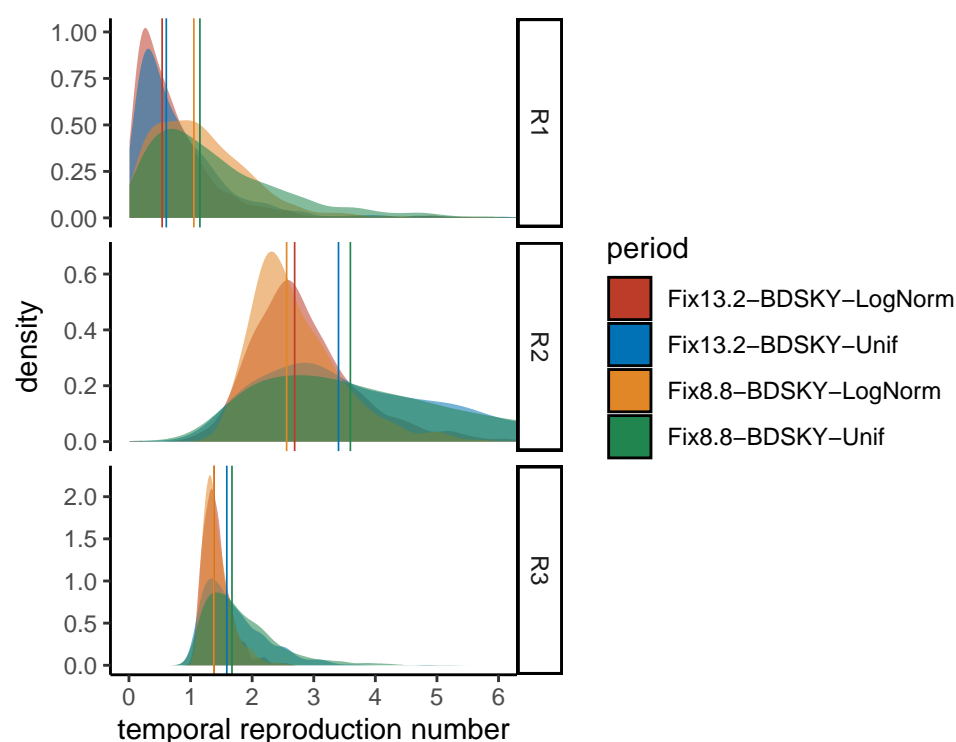


Figure S9 – **Effect of the prior shape and of the molecular clock on the temporal reproduction number estimate.** The molecular clock value has a stronger effect than the prior shape except for \mathcal{R}_3 , where the effect is limited.

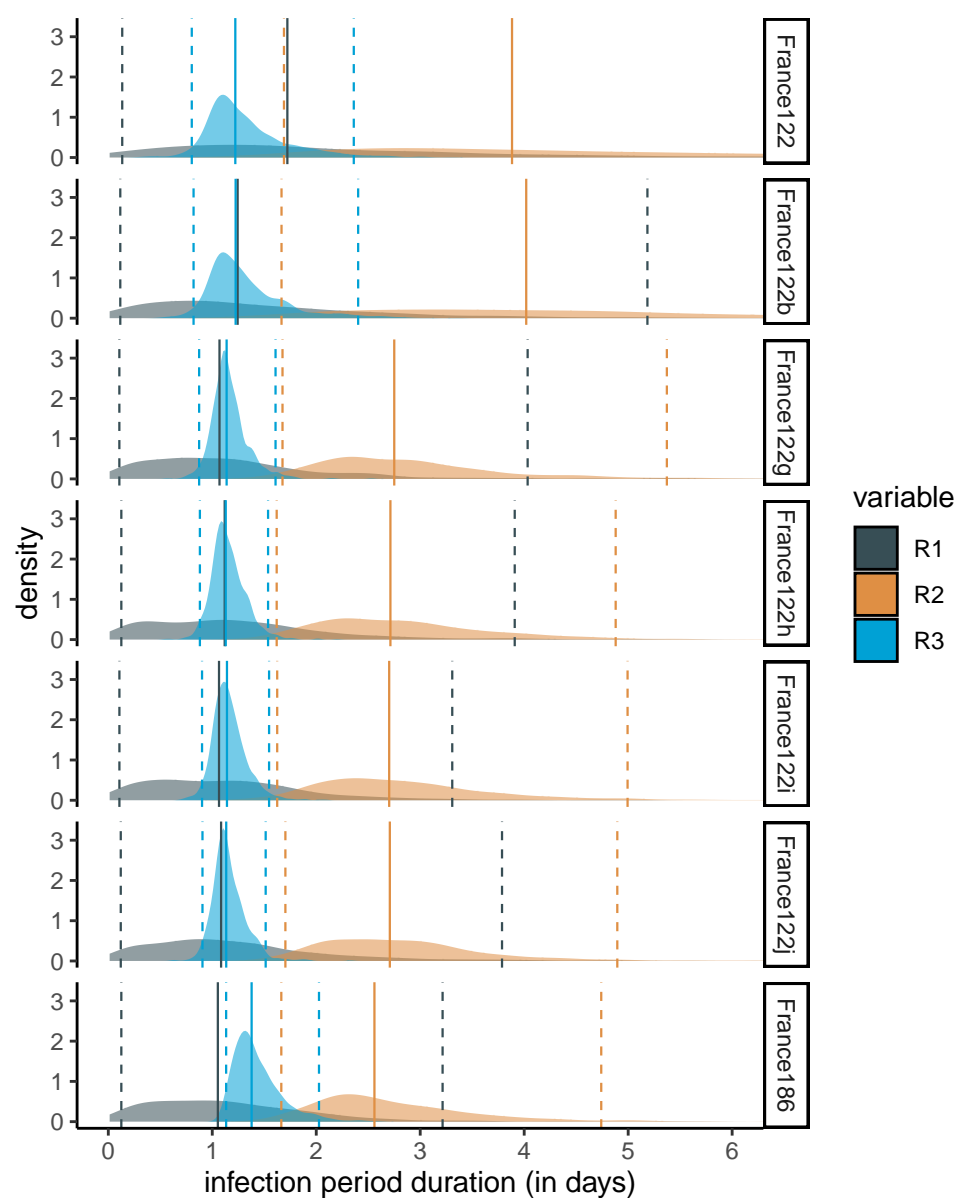


Figure S10 – **Reproduction numbers for the France186 dataset and subsets with 122 sequences.** Here we assume BDSKY model.

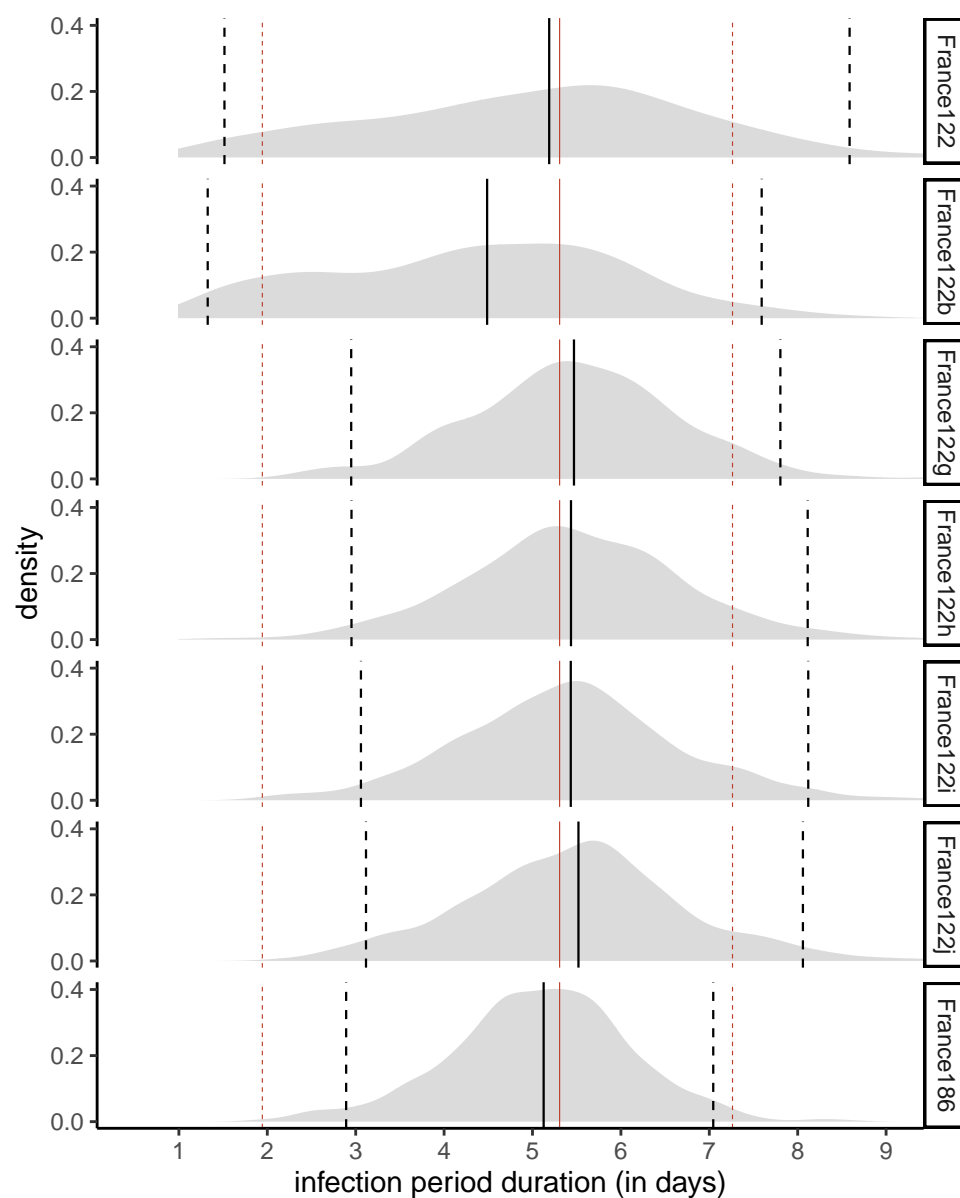


Figure S11 – **Infection duration for the France186 dataset and subsets with 122 sequences.** Here we assume BDSKY model.