

Using early detection data to estimate the date of emergence of an epidemic outbreak

Jijón, S.¹, Czuppon, P.², Blanquart, F.³ and Débarre, F.¹

¹Institute of ecology and environmental sciences of Paris (iEES-Paris, UMR 7618),
Sorbonne Université, CNRS, UPEC, IRD, INRAE, Paris 75005, France

²Institute for Evolution and Biodiversity, University of Münster, Münster 48149, Germany

³Center for Interdisciplinary Research in Biology, CNRS, Collège de France,
PSL Research University, Paris 75005, France

Abstract

While the first infection by an emerging disease is often unknown, information on early cases can be used to date it, which is of great interest to trace the disease's origin and understand early infection dynamics. In the context of the COVID-19 pandemic, previous studies have estimated the date of emergence (e.g., first human SARS-CoV-2 infection in Wuhan, emergence of the Alpha variant in the UK) using mainly genomic data. Another dating attempt only relied on case data, estimating a date of emergence using a non-Markovian stochastic model and considering the first case detection. Here, we extend this stochastic approach to use available data of the whole early case dynamics. Our model provides estimates of the delay from the first infection to the N^{th} reported case. We first validate our model using data concerning the spread of the Alpha SARS-CoV-2 variant in the UK. Our results suggest that the first Alpha infection occurred on (median) August 20 (95% interquartile range across retained simulations, IqR: July 20–September 4), 2020. Next, we apply our model to data on the early reported cases of COVID-19. We used data on the date of symptom onset up to mid-January, 2020. We date the first SARS-CoV-2 infection in Wuhan at (median) November 26 (95%IqR: October 31–December 7), 2019. Our results fall within ranges previously estimated by studies relying on genomic data. Our population dynamics-based modelling framework is generic and flexible, and thus can be applied to estimate the starting time of outbreaks, in contexts other than COVID-19, as long as some key parameters (such as transmission and detection rates) are known.

Keywords. Epidemic emergence; Early dynamics; Stochastic simulations; COVID-19; SARS-CoV-2.

1 Introduction

Dating the first infection of an emerging infectious disease is a step towards tracing the disease's origin and understanding early epidemic dynamics. Beyond the early transmission of a new pathogen, estimating the date of first infections is also of interest while studying the initiations of local clusters in naïve populations. For example, this happens when the pathogen is first introduced to a new location, but also when the pathogen evolves to distinct genotypes such as emerging variants of concern (VOCs).

Various attempts have been made to date the first human infections by SARS-CoV-2 that led to the COVID-19 pandemic (noting that earlier spillovers, leading to dead-ends, may have occurred). Using a stochastic model for the epidemic spread coupled with genomic data allowing to trace transmission at the individual level, Pekar et al. [1] estimated that the first human infection took place between late October and early December 2019. This estimate resulted from revising previous findings yielding an emergence date between mid-October and mid-November 2019 [2], notably after updating the dates of the first case reported [3]. Another modeling study dated the first COVID-19 case between early October and mid-November 2019 by adapting a technique used in conservation science [4], but was also based on outdated data. Other studies have focused on the introductions of SARS-CoV-2 to different countries. For instance, studies using molecular clock analyses relying on genomic data to determine the time of most recent common ancestor (tMRCA) of lineages introduced in a focal country have been conducted in the context of France [5], the United States [6] and the United Kingdom (UK) [7]; while another study used a stochastic non-Markovian approach relying on mortality data to estimate the date of SARS-CoV-2 introduction to France [8]. One study focusing on the emergence of the 'EU1' SARS-CoV-2 variant (B.1.1778) circulating among European countries during the summer of 2020 used genomic data to date most introductions to June, 2020 [9]. Dating attempts have also been done for the 'Alpha' variant (B.1.1.7), whose date of emergence was estimated at early August 2020 using a stochastic, non-Markovian approach relying on the date of the first observed case [10], and whose tMRCA was estimated at late August 2020 [11].

The tMRCA however does not necessarily approximate the emergence date [12], which can have taken place earlier. Infection times occurring earlier than the tMRCA can be estimated thanks to mathematical models. Moreover, modeling studies have helped unveiling other unobserved indicators during the early stages of epidemics, such as the epidemic size at the time of first detection [10, 13]. In particular, because infection numbers are low, stochastic approaches are key to studying early dynamics. Hence, methodological developments of stochastic models to study the early stages of infectious diseases remain of great interest in the field of mathematical epidemiology.

The main objective of our study is the estimation of the date of the first infection leading to a sustained epidemic (hereafter named the date of epidemic/outbreak *emergence*), using available data on the first N detected cases. To this end, we build a stochastic model and designed a simulation framework extending previous work [10]. This approach provides the time elapsed between epidemic emergence and the N first observed cases, as well as the proportion of the epidemic that remains undetected.

Here, we present in detail the construction of our model and its applications to two epidemiological contexts. First, we use data on the spread of the Alpha variant in the UK, and validate our extension of the model presented in [10]. Next, we parameterize our model to reproduce the dynamics of the early outbreak in Wuhan to estimate a range of probable dates for the emergence of the COVID-19 pandemic.

2 Results

2.1 Modeling the early dynamics of an epidemic outbreak

We develop a stochastic epidemic model that estimates the time elapsed between the first infection and N reported cases, using available data and estimates on key epidemiological parameters. We model infectious disease transmission with a general branching process starting from a single infectious individual. This implies that times from infection to transmission events are not exponentially distributed as assumed by ordinary differential equation models. We then model the detection of infected individuals, which constitutes the modeled time series of cases. Both infection and detection processes follow distributions with known fixed parameters; cf. Table 2. Importantly, we assume fixed parameters as estimated by previous studies (namely, the detection probability), and we later examine the robustness of our findings to these exact values. We consider the time-series of infections and detections up to the day of occurrence of the N^{th} case. We calibrate our model to reproduce the observed epidemic, using available data on disease cases. For more details on the model, we refer to the Methods section.

The main outcome of our model is the time series of cases, from which we deduce the delay between emergence (first infection) and N^{th} case. By first infection, we refer to successful epidemic outbreaks only; that is, we do not account for the first infections that may have led to epidemics that went extinct. In addition, by keeping track of the whole epidemic, we retrieve the time series of infections and the number of secondary cases produced by each infectious individual. We can thereby compute, for instance, the epidemic size at the time where the N^{th} case is reported, and deduce the proportion of detected infections (i.e., cases). This proportion is impacted by detection delays and stochasticity and thus, it is not straightforwardly obtained from the probability of detection considered in the simulations.

We run as many numerical simulations as needed to obtain 5 000 successful epidemics, i.e., epidemics that were sustained after a predetermined period of time and that verified the calibration conditions imposed by the data from a certain epidemiological context (details in the Methods section). We apply our model to two epidemiological contexts: the emergence of the ‘Alpha’ variant in the UK and the emergence of SARS-CoV-2 in Wuhan. Both applications and the corresponding results are described in more detail below.

2.2 Estimating the date of the first infection with the Alpha variant in the UK

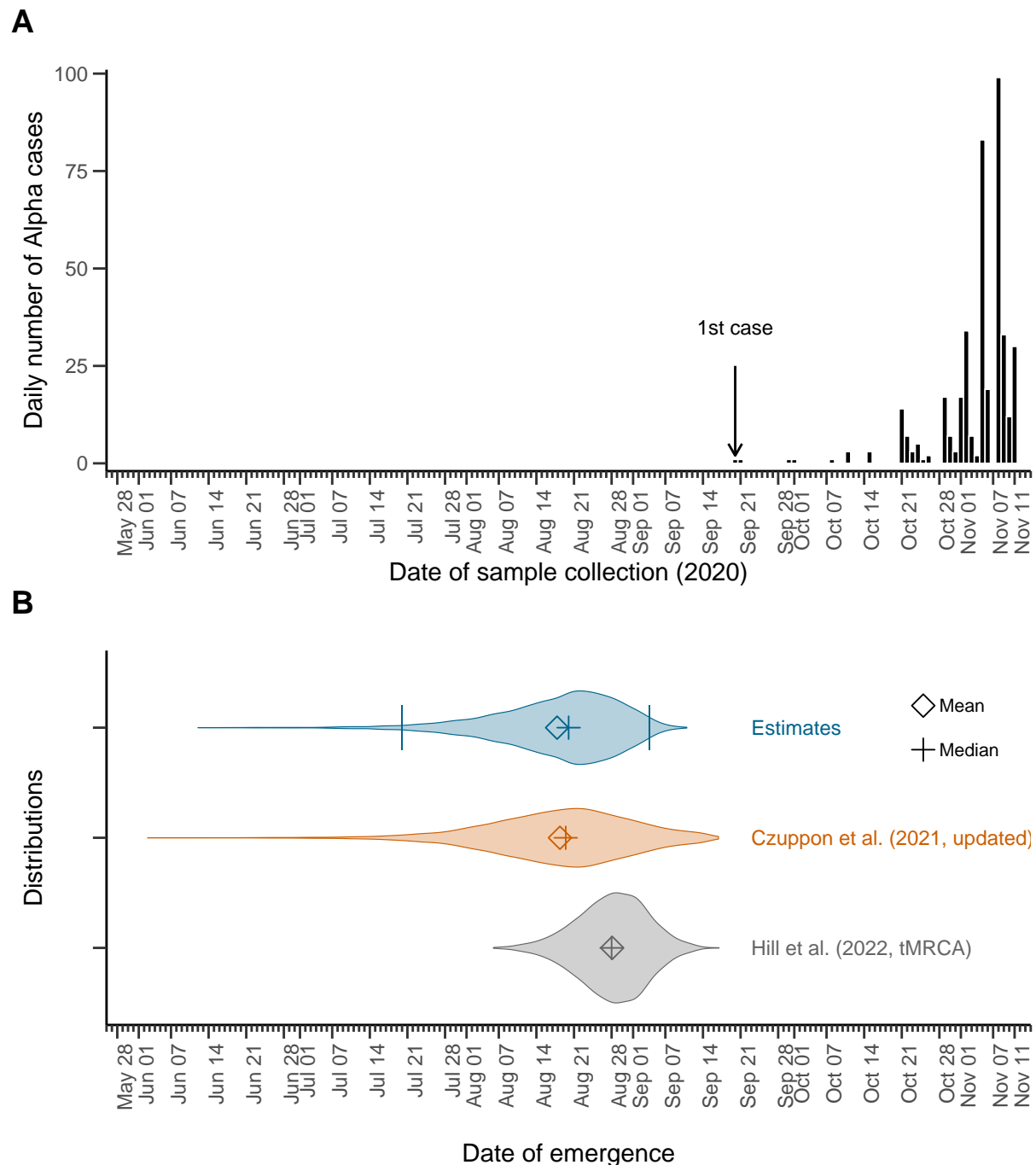
As a first validation of our model, we estimate the emergence of the Alpha variant in the UK, which was the main result of the numerical applications presented in [10]. We applied our model to a dataset of $N = 406$ samples carrying the Alpha variant collected and sequenced between September 20 and November 11, 2020 [14]. Table 2 summarizes the parameter values used in our simulations. The 5 000 simulated epidemics that we analyse below result from model calibration (i.e., epidemics arising from a single infectious individual and verifying the calibration constraints; details in the Methods section), and represent 25% of all simulations run with the input parameters. The cumulative cases of the accepted epidemics are depicted in Supplementary Figure S3.

Hereafter, we summarize our results using median values and 95% interquartile ranges (95%IQR; values between the 2.5th and the 97.5th percentiles) from the distributions of the different epidemiological indicators obtained from the 5 000 simulated epidemics, similar to a posterior distribution obtained in an Approximate Bayesian Computation framework (see Methods). We estimated the number of days between the 1st infection and the N^{th} case at 83 (95%IQR: 68–114), dating the emer-

gence of the Alpha variant in the UK at August 20, 2020 (95%IqR: July 20–September 4), and not earlier than June 12, 2020. Alpha transmissions were ongoing about 26 (95%IqR: 9–56) days before the date at which the first known case was sampled and sequenced. Table S1 in the Appendix provides other calibration metrics like the delay between the 1st and N^{th} sequenced samples in the simulations. Figure 1 depicts our estimates of the date of emergence along with the epidemic curve (i.e., the daily number of sequenced samples; by sampling date), for context, as well as previous estimates, for comparison. In particular, we ran an updated version of the model presented in [10] (the distribution of the number of secondary cases is negative-binomial instead of Poisson previously, and its mean, R , is now equal to 1.9 instead of 1.5 in [10]). We also compare our results to tMRCA estimates [11] (personal communication of the distributions). Our median estimates for the emergence date fall within a very close, slightly narrower range than that found by running an updated version of [10], while falling ~ 1 week earlier than the estimated tMRCA [11]. These comparisons are summarized in Supplementary Table S2.

We further estimated the epidemic size at the date of infection of the N^{th} case at about 90 700 (95%IqR: 80 400–102 000). The simulated detected cases thus represent a proportion of about 0.48% (95%IqR: 0.43%–0.53%) of the total number of infections. Note that a case is an infected individual who underwent a PCR test and whose sample was sequenced and contained the Alpha variant, which accounts for the low probability of detection. Our results are summarized in Table 1. For the results of running our model for $N = 1$ and the comparison of our results with those of [10], we refer to Supplementary Figure S5 and Supplementary Table S3.

Figure 1: Estimates of the date of emergence of the Alpha variant in the United Kingdom (UK). (A) The epidemic curve corresponding to the data used to calibrate the model, for context. A total of $N = 406$ samples carrying the Alpha variant were collected between September 20 and November 11, 2020 [14]. (B) Violin plots for the distributions of the date of emergence. We estimated the emergence of the Alpha variant in the UK at August 20 (95%IqR: July 20–September 4), 2020 (top, blue; upper and lower bound of the 95%IqR depicted by bars). For comparison, we also show the distributions of the estimates from an updated version of the model developed in [10]—which relies on a single observation on September 20—where we set $R = 1.9$ and a negative-binomial distribution for the number of the secondary cases (middle, orange). The distribution for the estimated time of most recent ancestor (tMRCA) [11] is also shown (bottom, gray).



2.3 Estimating the date of emergence of SARS-CoV-2 in Wuhan

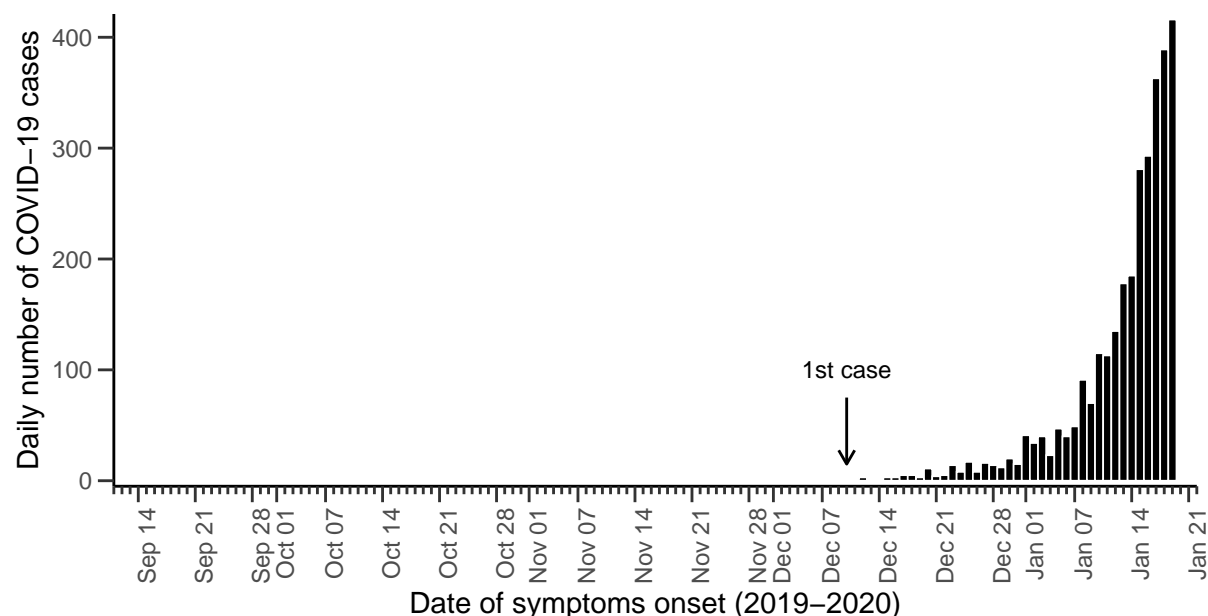
Next, we applied our model to the dataset of the early cases of COVID-19 considered in [1] (personal communication). A total of 3 072 confirmed COVID-19 cases were reported to have had symptoms onset between December 10, 2019 and January 19, 2020, the day before the first public declaration of human-to-human transmission, shortly before the first lockdown interventions [15]. We parameterized our model using estimates from the literature; cf. Table 2. The 5 000 selected simulations represent 13% of all runs (cf. Supplementary Table S1, Appendix) and are depicted in Supplementary Figure S4.

Our simulations yield an estimated median number of days between the 1st SARS-CoV-2 infection to the N^{th} symptomatic COVID-19 case recorded of 54 (95%IQR: 43–80) days, dating the emergence (i.e., the first sustained human infection) of SARS-CoV-2 to November 26 (95%IQR: October 31–December 7), 2019, and not earlier than September 28, 2019. This also implies that the epidemic remained completely undetected (i.e., no detected infections) for about 9 (95%IQR 3–20) days. These findings are depicted in Figure 2, along with the observed epidemic curve (i.e., COVID-19 cases dataset) as well as recently published estimates of the date of emergence of the COVID-19 pandemic [1] (personal communication of the distributions), for comparison. Our estimates fall remarkably close to those previously found in [1] (cf. Supplementary Table S2, Appendix).

We further estimate the median number of infections on the day of infection of the N^{th} case at about 63 600 (95%IQR: 57 300–69 800), which results in a median proportion of detected infections of 5.31% (95%IQR: 5.07%–5.56%). Our results are summarized in Table 1.

Figure 2: Estimates of the emergence of SARS-CoV-2 in Wuhan. (A) Observed epidemic curve, for context. A total of $N = 3072$ COVID-19 cases with symptom onset between December 10, 2019 and January 19, 2020, the day before the first public statement on human-to-human transmission. (B) Violin plots for the distributions of the date of emergence. We estimated the median date of the first SARS-CoV-2 infection (i.e., emergence) at November 26 (95%IqR: October 31–December 7), 2019 (top, blue; upper and lower bound of the 95%IqR depicted by bars). For comparison, the distribution for the estimates from [1] are also shown (bottom, purple).

A



B

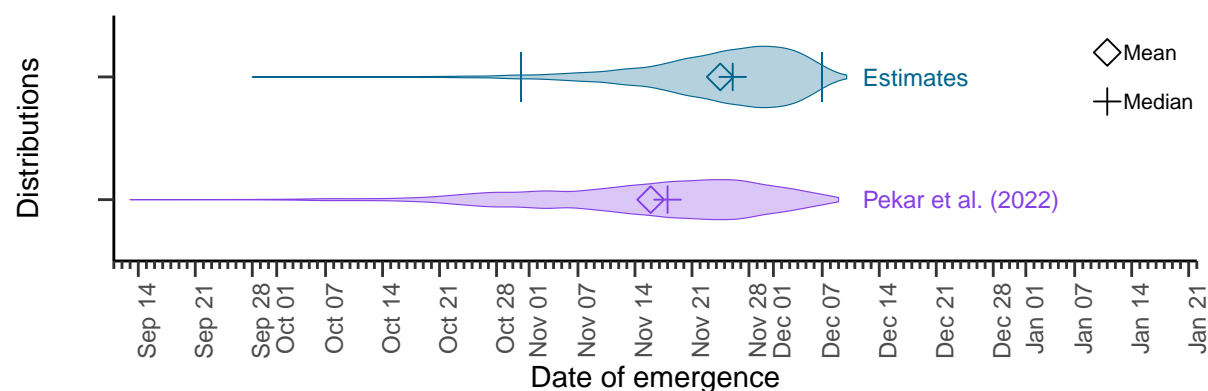


Table 1: Estimates of the date of emergence and other epidemiological indicators resulting from the calibrated model. Here we summarize the estimates obtained from the numerical application of our model to two epidemiological contexts: the Alpha variant infections in UK, and the first COVID-19 cases reported in Wuhan. The estimated time elapsed between the first infection to the N^{th} observed case yields the estimated emergence date. In addition, we estimate the epidemic size at the date of detection of the N^{th} case. The proportion of detected infections and mean secondary cases are retrieved for comparison with the input epidemic parameters. Median and 95% interquartile ranges are shown, unless stated otherwise.

Epidemiological indicator	Alpha (UK, 2020)	COVID-19 (Wuhan, 2019)
Number of days from 1 st infection to N^{th} case	83 (68–114)	54 (43–80)
Date of first infection	Aug 20 (Jul 20–Sep 4), 2020	Nov 26 (Oct 31–Dec 7), 2019
Date of earliest infection	Jun 12, 2020	Sep 28, 2019
Epidemic size at day of infection of the N^{th} case	90 700 (80 400–102 000)	63 600 (57 300–69 800)
Proportion of detected infections	0.48% (0.43%–0.53%)	5.31% (5.07%–5.56%)
Mean number of secondary infections	1.90 (1.88–1.92)	2.51 (2.43–2.59)

2.4 Sensitivity analyses

Here, we evaluate the impact of uncertainty around the main model parameters on the results, by running our simulations while varying the input values of the expected number of secondary cases (R), the overdispersion parameter (κ) and the probability of detecting an infection (p_{detect}). We find that increasing both p_{detect} and R implies a reduction in the number of days between the first infection and the N^{th} case, meaning that the epidemic emerges later. In addition, we find that increasing R (i.e., higher number of per-day transmission events) yields distributions with more pronounced skewness, reflecting the occurrence of larger super-spreading events; cf. the Methods section. Median estimates and 95%IqR are summarized in the Supplementary Table S2 in the Appendix.

We further evaluate the impact of the tolerances chosen for accepting the simulated epidemics on our results (cf. Methods for details). We run our simulations while varying the tolerance for the difference between the simulated and the observed daily number of infections, δ_{tol}^Y , and the tolerance on the difference between the simulated and the observed period of time between the 1st and the N^{th} case, δ_{tol}^τ , for both applications. Our estimates are robust with respect to these variations. However, our model seems to be slightly more sensitive to the choice of the minimum length of the time interval of cases.

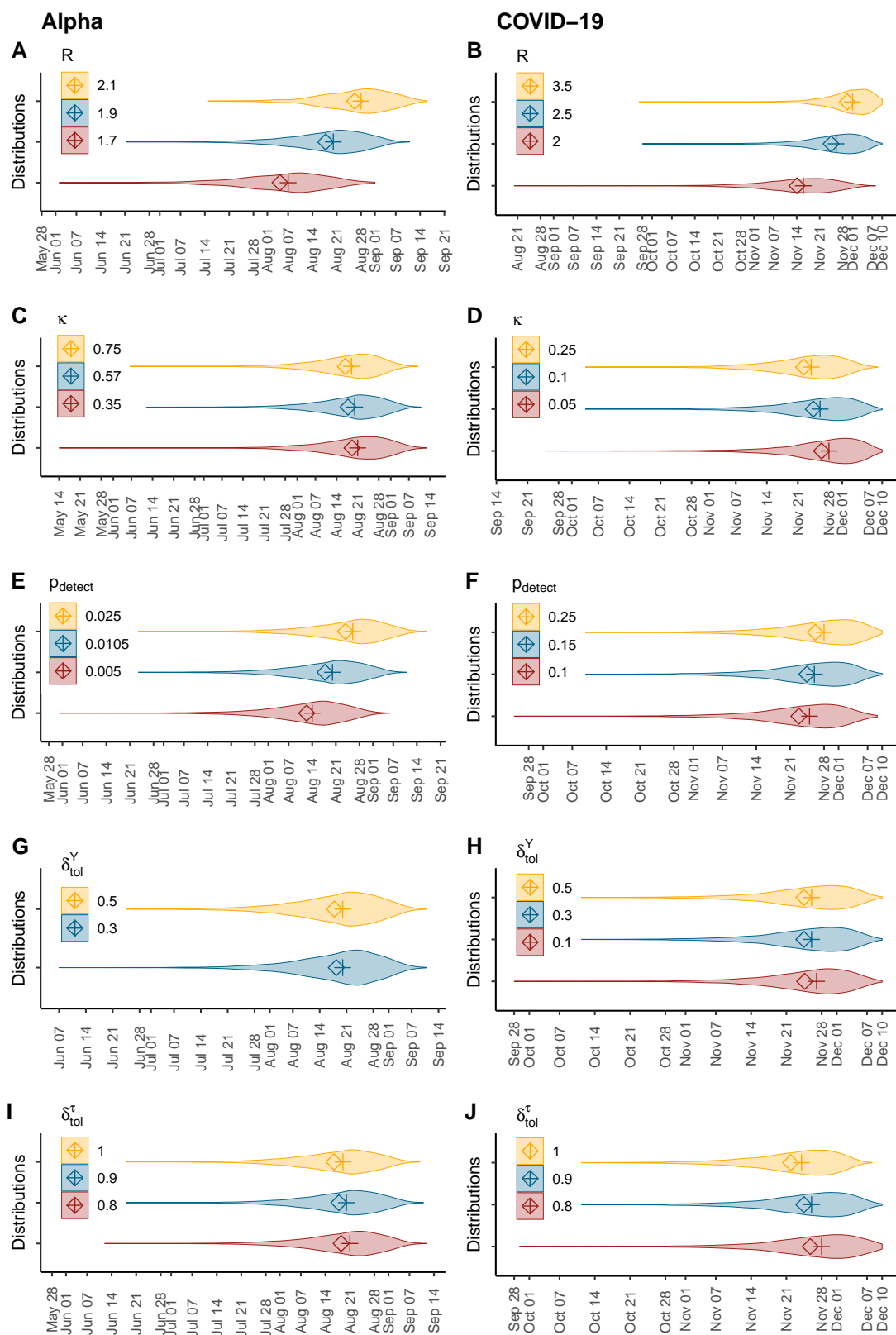


Figure 3: Sensitivity analyses. Distributions for the date of emergence of the Alpha variant in the UK (left) and the COVID-19 epidemic in Wuhan (right) obtained from running our simulations setting different input values for key model parameters: the reproduction number, R (panels A and B), the overdispersion parameter, κ , (panels C and D), the probability of detection, p_{detect} (panels E and F), the tolerance for the difference between the simulated and the observed daily number of infections, δ_{tol}^Y (panels G and H) and the tolerance for the difference between the lengths of the simulated and the observed cases time periods, δ_{tol}^τ (panels I and J). The absence of a violin plot for $\delta_{\text{tol}}^Y = 0.1$ in panel G results from the absence of selected simulations in > 3 million runs. Median values are depicted by crosses and mean values by diamonds. Baseline parameterization of the model is depicted in blue (middle).

Uncertainty around very early case declaration

We ran additional sensitivity analyses for the application to the emergence of the COVID-19 cases in Wuhan, to evaluate the impact of the uncertainty around very early case declaration, by applying our model to different datasets; cf. Section 2.3 of the Appendix. We first used a shorter dataset, with data on COVID-19 cases with symptoms onset only up to December 31, 2019, the day of the first public declaration of a cluster of pneumonia of unknown etiology [15] ($N = 169$). We also used an outdated, later corrected dataset published by the WHO [16] ($N = 202$) where the first case reported symptoms onset was on December 2, 2019. We found that using the shorter dataset dates the epidemic emergence 2–3 days earlier than our main results, and selects simulations with a slightly higher mean number of secondary infections (c.f. Table 1 and Supplementary Table S4), which may reflect a change in the population behavior following the first public announcement. However, using the outdated dataset [16] dates the epidemic emergence about a week prior to our main results, while the mean number of secondary infections remained about the same, which reflects the impact of the date of first detection on the results.

3 Discussion

We estimate the date of emergence of an epidemic outbreak, defined as the first infection leading to a sustained transmission chain, relying on estimates of key epidemiological parameters as well as available data on the first N observed cases. With our population-dynamics approach, we recover estimates very close to those of previous studies [1, 11], which were obtained using information from whole genome sequences.

Our model was conceived as an extension of the numerical application presented in [10]. This methodology relies on a general branching process to model disease transmission and detection. Our model is informed by available data on the first N observed cases (unlike [10], who used the date of first detection only). We further account for super-spreading events using a negative-binomial distribution for the generation of secondary infections. This assumption may reduce the time elapsed between the first infection and the first detection or increase its variance, in comparison to ignoring super-spreading events by considering other distributions (e.g. Poisson) [10] or by ignoring individual heterogeneity in infectiousness [8], as well as by using deterministic approaches [8].

We first study the emergence of the Alpha variant as model validation. Our results suggest that the Alpha variant emerged in the UK around August 20, 2020 and not earlier than June 12, 2020. Our results fall indeed within the same ranges as those of the approach presented in [10] when updated to match our parameterization; and fall shortly earlier than previous tMRCA estimates of the Alpha variant [11]. This result makes sense: tMRCA does not necessarily yield the date of first infection, which may have occurred before the most recent common ancestor. Next, we apply our model to data on the early COVID-19 cases in Wuhan, estimating the date of the first SARS-CoV-2 infection around November 26 (95%IQR: Oct 31–Dec 7), 2019, and not earlier than September 18, 2019. These ranges also fall remarkably close to—slightly later than—recently published estimates [1]. To the best of our knowledge, the novelty of these findings rests on using exclusively a population-dynamics approach, unlike previous studies addressing the subject. The median estimate in Pekar et al. [1] (which is a first infection, and not a tMRCA) falls ~ 8 days before ours, which can be explained by the fact that our approach ignores genomic information. Namely, with our approach, having two cases with the exact same infecting virus, or two cases with viruses two mutations away, are treated the same way, as our approach only uses case numbers. Using genomic information would however yield an earlier date of emergence if the infecting viruses are genetically more distant. Our findings are thus also compatible with a previous study that found no evidence of widespread transmission in Wuhan before December 2019, using serological data [17]. Hence, in accordance to previous discussions [1, 2], our results suggest that widespread SARS-CoV-2 circulation (and even more so, international spread) earlier than the end of 2019 is unlikely. Assuming an origin of the pandemic in China [18], claims of large early (i.e., before January 2020) circulation outside of China [e.g. 19, and references therein] would be therefore extraordinary, and require extraordinary evidence, excluding potential false positive by setting appropriate controls.

Quantifying the time that emerging epidemics remain undetected before detecting the first cases is particularly important in the context of emergent pathogens such as SARS-CoV-2, where very early cases may remain unidentified, especially if a high proportion of the infections are asymptomatic [20] (N.B.: most COVID-19 cases with symptoms onset up to December 31, 2019 were declared retrospectively [21]). There is also evidence that SARS-CoV-2 may have been introduced in other countries for some time before the first reported cases [5, 6, 7, 22]. Accordingly, our approach shows that early SARS-CoV-2 infections of the transmission chain that was first detected at Wuhan's Huanan market [18] remained undetected for more than a week and up to 3 weeks before the first case of symptomatic COVID-19 was observed.

The impact of uncertainty on our results is assessed by varying the main transmission and de-

tection parameters, as well as the rejection criteria for the simulated epidemics. These sensitivity analyses unveil the robustness of our approach regarding the choice of the rejection criteria. On the other hand, we find that the main model parameters (expected number of secondary infections, R , and probability of detection, p_{detect}) have a greater impact on the model outcomes: as expected, both higher transmissibility and higher detection shorten the time between emergence and N detected infections. This variation in the results is particularly true for the application on the Alpha variant data, probably due to the notably smaller case dataset we use.

Our study has several limitations. First, our results depend heavily on input data, while access to good quality data on the early stages of an epidemic outbreak may be challenging. Datasets may be scarce, they may face reporting delays, early cases may be detected retrospectively and detection protocols may change. Early Wuhan COVID-19 cases with—severe—symptoms onset before December 30 (the day of issue of the emergency notice from the Wuhan Municipal Health Commission) [23, 15] were diagnosed clinically before tests for SARS-CoV-2 infection were available [3]. In particular, our methods rely on the date of first infection (cf. the Methods section, condition (C3)). See Supplementary Figure S6 and Supplementary Table S4, where emergence estimates change. Second, our model requires early estimates of the distributions for key epidemiological indicators such as the mean number of secondary infections, the secondary infection generation time, the probability of case detection and the incubation period, which depend themselves on the quality of early observed data and may not be available for new emerging infectious diseases. In particular, it can be challenging to estimate the probability of case detection (or ascertainment rate) during early stages of an emergent infectious diseases, and it may vary between countries [24]. Third, we model epidemic spread starting from a single infectious individual, thus neglecting scenarios of multiple introductions. This impedes, for instance, the application of our model to contexts such as SARS-CoV-2 importation to France [5]. This may also be a limitation in the context of epidemics emerging from multiple spillover events, such as has been concluded by [1] relying on data on the early SARS-CoV-2 lineages. Fourth, the forward simulations of our model do not allow to consider time-dependent parameters. Hence, we are constraint to use data on relatively short periods of time to ensure that epidemiological parameters remain nearly constant over the study period. This may not reflect early epidemic dynamics, where public outbreak alerts may provoke, on the one hand, an increase in testing efforts and thus, rapid changes in the probability of detection and, on the other hand, changes in individual behaviors that may impact the probability of disease transmission and thus, the expected number of secondary infections.

The numerous, fast and free availability of genomic data for the COVID-19 epidemic is unprecedented. Here, we built our model in a parsimonious, generic and flexible manner intended to be applied to contexts other than COVID-19, provided that key epidemiological parameters are known and transmission chains arise from a single infectious individual. Further developments of our model need to include genomic data on top of case data, as it is likely that sequencing will remain as intensive for other infectious diseases as it has been for COVID-19.

4 Methods

4.1 Model

We extend the methods presented in the numerical applications of [10], where the time of emergence was estimated from data on the first reported case only, using a stochastic population-dynamics approach. We model the early stages of the epidemic and use dates of the first N observed cases to estimate the date of emergence. We use the term ‘*case*’ to refer to infections that are ascertained and reported: time series of cases may thus correspond to one of the following types of time series: infection detection, sample sequences, symptoms onset declarations, etc. We use the term ‘*probability of detection*’ for the probability of such ascertainment to occur.

Our model is defined by a non-Markovian branching process to model the transmission of an infectious disease, starting from a single infectious individual in a fully susceptible population [8, 10]. Since we study early epidemic dynamics, i.e., for a relatively short period of time, all distribution parameters are assumed to be constant during the modeled time period. At any time t , infected individuals may transmit the disease. We account for super-spreading events by assuming that the number of secondary cases follows a negative binomial distribution,

$$\text{NegBinom}\left(\text{number of failures} = \kappa, \text{probability of transmission} = \frac{\kappa}{\kappa + R}\right), \quad (1)$$

where R denotes the expected number of secondary infections (i.e., the effective reproduction number).

Then, we model the ascertainment of infections by drawing the per-day number of observed cases from a Binomial distribution

$$\text{Binom}\left(\text{number of trials} = I(d_k), \text{probability of success} = p_{\text{detect}}\right), \quad (2)$$

where $I(d_k)$ is the number of incident infections at day d_k , with $k = 1, 2, \dots$, and p_{detect} is the probability of infection detection. The generation time of each new infection, $\{t_i\}_{i=1,2,\dots}$, and the time from infection to detection of a case, $\{\tau_j\}_{j=1,2,\dots,N}$ are drawn from a Gamma distribution

$$\text{Gamma}(\text{shape} = \omega_x, \text{scale} = \theta_x), \quad (3)$$

with $x \in \{t, \tau\}$, respectively.

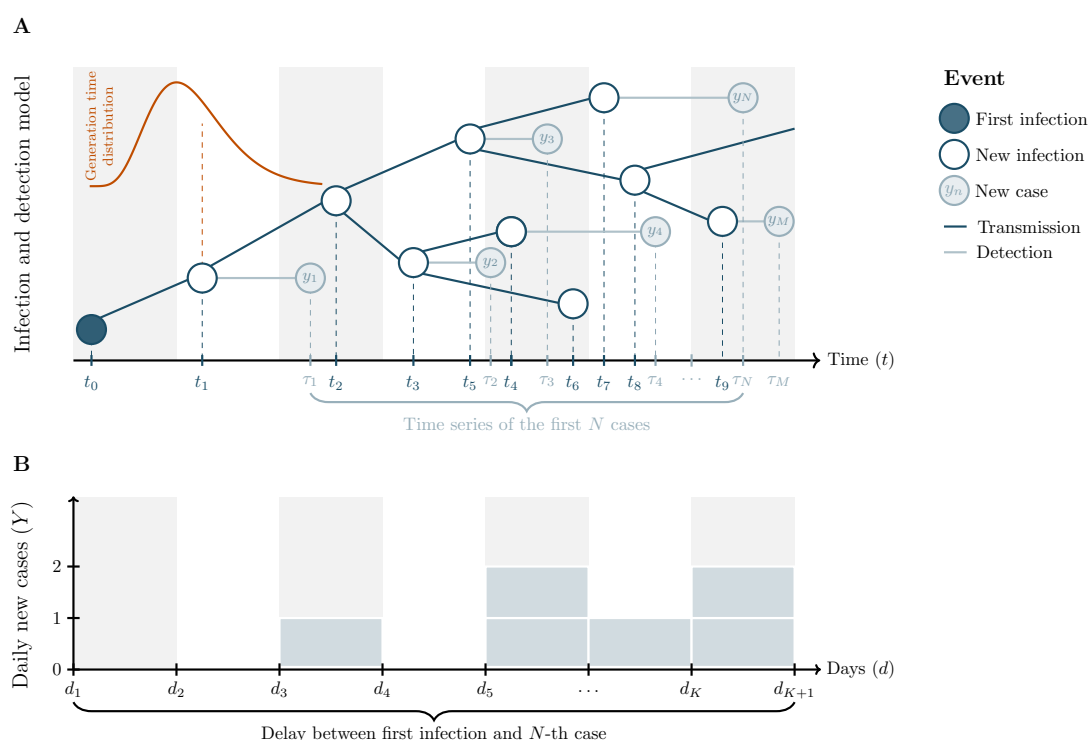
4.2 Estimation of the date of first infection

In our simulations, we discretize time by $\Delta t = 0.1$ days. Note that times to infection and detection are still on a continuous scale. We then aggregate epidemiological indicators such as new infections and cases by day (denoted by d_k , where $k = 1, 2, \dots$), since this is the time scale at which most data are presented. We run stochastic simulations forward in time, from $t = t_0$, the time of occurrence of the first (undetected) infection (i.e., $I(t_0) = 1$), until the end of the day of detection of the N^{th} case, d_K , which is determined by $d_K \leq \tau_N < d_{K+1}$. Note that while the N^{th} case is a stopping criterion of our model, our analyses still deal with all $M \geq N$ cases occurring up to the end of day d_K ; that is,

$$Y(d_k) \equiv \sum_{\substack{i \in \{1, \dots, M\} \\ \text{s.t. } d_k \leq \tau_i < d_{k+1}}} y_i, \quad k = 1, 2, \dots, K, \quad (4)$$

where y_i denotes the i^{th} detected infection (i.e., the i^{th} case) and $Y(d_k)$ denotes the number of infections detected on day k . A depiction of our infectious disease transmission and detection model is shown in Figure 4.

Figure 4: Model diagram. (A) We model infectious disease transmission (dark blue elements) starting from a single infectious individual (full dark blue dot), using a general branching process. The generation time of secondary infections, $\{t_{i+1} - t_i\}_{i=0,1,\dots}$, follows a Gamma distribution (shown in orange, above the first transmission event). In addition, we model the detection of infected individuals (light blue elements), which yields the time series of observed cases, $\{\tau_j\}_{j=1,\dots,N}$. The number of cases are aggregated daily. Days are denoted by $\{d_k\}_{k=1,\dots,K}$ and depicted by alternating gray and white bands. Our algorithm stops the day at which the N^{th} case is observed, d_K , but our analyses deal with the set of all cases detected on day d_K , $\{y_j\}_{j=1}^M$, where $M \geq N$. (B) Resulting epidemic curve (cases per day). Our model is calibrated so that the simulated epidemic curve, $\{Y_k\}_{k=1}^K$, reproduces the observed number of cases per day. The main outcome of our model is the number of days elapsed between the first infection and the N^{th} observed case, d_K . NB. the time scale in the figure is not representative of our simulations: infection and detection delays in the simulations usually span multiple days.



The goal is to estimate the date of the first infection. To this end, we follow a strategy that is similar to an approximate Bayesian computation (ABC) algorithm. In ABC, one would define a (typically uniform) prior distribution for the date of the first infection (i.e., date of emergence), randomly draw emergence dates from this distribution, simulate stochastic epidemics and compute a distance between the data (the observed time series of cases) and each simulation. One simple algorithm (*rejection* algorithm) consists in accepting the small fraction of trajectories that are closest to the data. The distribution of emergence dates from the accepted simulations then approximates the posterior distribution of the date of first infection.

Here, we do not follow this computationally intensive strategy, yet our inference method is very close to an ABC algorithm and yields an approximate posterior distribution of the date of emergence. Let us assume for simplicity that we want to infer with ABC a single parameter, the date of emergence, defining to this end a distance metric based only on the date d_K when the simulations reach the N^{th} case. We only accept simulations with identical date (distance 0) to that observed in the data. If the prior distribution for the dates of emergence is uniform, the set of accepted simulations is identical to a set of simulations with unlabelled date of emergence, where simulation time is shifted to absolute date such that the date d_K is indeed identical to that observed in the data. This follows from our assumption that all parameters are constant in time, implying that the epidemio-

logical dynamics do not depend on absolute date but solely on time elapsed since the first infection. Thus, all simulations can actually be retained and shifted to absolute dates in such a way that they correspond to a sample from the posterior distribution in this simple ABC algorithm. The actual rejection criterion is actually slightly more complex than described, as we not only match exactly the date of the N^{th} case, but also additionally require a set of calibration conditions ensuring that the stochastic epidemiological dynamics are similar to those observed in the data. This procedure is repeated until we simulated 5 000 accepted epidemics. To compute the dates of emergence from these simulations, we take the final date of the observed case data and subtract the duration of each simulated epidemic, which produces the posterior distribution of the date of emergence.

We now describe our calibration conditions in detail. Simulated epidemics ('sim') are calibrated to reproduce the observed ('obs') epidemic, via four conditions. First we constrain our model to accept the simulations where an epidemic occurs, i.e., if the number of infections is high enough that N cases are ensured (here, five times higher than the expected number of cases):

$$\sum_i I(t_i) > 5 \times N / p_{\text{detect}}. \quad (\text{C1})$$

Second, we require that the first simulated infection predates the first observed case:

$$t_1^{\text{sim}} \leq \tau_1^{\text{obs}}. \quad (\text{C2})$$

Third, the length of the time period (in days) between the first and the N^{th} case must be similar to the observed time period in the case data set, under a certain tolerance $\delta_{\text{tol}}^{\tau}$:

$$\tau_N^{\text{sim}} - \tau_1^{\text{sim}} \geq \delta_{\text{tol}}^{\tau} (\tau_N^{\text{obs}} - \tau_1^{\text{obs}}), \quad (\text{C3})$$

where $\{\tau_j^{\text{obs}}\}_{j=1,\dots,N}$ is the time series of observed cases (i.e., the data set of reported cases). The fourth and last constraint concerns the epidemic curve; the daily number of cases of the simulated epidemic is required to resemble (under a certain tolerance δ_{tol}^Y) to the observed behavior:

$$\max_k \left| \sum_{j=1}^k Y^{\text{obs}}(d_j) - \sum_{j=1}^k Y^{\text{sim}}(d_j) \right| \leq \delta_{\text{tol}}^Y N, \quad k = K, K-1, K-2, \dots, \quad (\text{C4})$$

where $\sum_{j=1}^k Y(d_j)$ denotes the cumulative number of cases at day d_k . The daily case count in the accepted simulations thus depends on N , i.e. on the epidemiological context.

For more details on the numerical application of the model, please refer to the pseudo-algorithm in the Appendix. The simulations were run in Julia [25] version 1.8, and the results figures were generated in R [26] version 4.1.2, using the ggplot2 package [27], version 3.3.6. All data and codes needed for reproducibility of our results and the corresponding figures are available at a public Github repository: <https://github.com/sjijon/estimate-emergence-from-data>.

4.3 Applications

4.3.1 Alpha variant in the UK

The first application concerns the early stages of the spread of the Alpha SARS-CoV-2 variant in the UK, and thus serves as a validation of our extended version of the model presented in [10].

The first reported sequence of the SARS-CoV-2 Alpha variant of concern, was collected on September 20, 2020 [28]. Here, we define a case as a sequenced sample carrying the Alpha variant, and we

define the probability of a case detection as the probability of sampling and sequencing such variant. The parameterization of our model is as in [10], except for the expected number of secondary infections, R , which was updated to match the hypotheses made in [11], to ensure comparability of results; cf. Table 2. The data on early Alpha cases were retrieved from the Global Initiative on Sharing Avian Influenza Data (GISAID) [29], available at doi.org/10.55876/gis8.230104xg (Appendix). We use the data on the sequences submitted to GISAID up to November 30, 2020, and used only the sample collected up to November 11, 2020. This choice was done to overcome reporting delays; cf. Supplementary Figure S1.

4.3.2 COVID-19 in Wuhan

The second application concerns the early COVID-19 cases reported in Wuhan, China. Here, we define a case as a confirmed COVID-19 infection, in many cases determined retrospectively [21]. We use the dataset considered in [1] (personal communication), comprising 3 072 COVID-19 cases with symptoms onset between December 10, 2019 and January 19, 2020, the day before the first public statement on human-to-human transmission [15]; cf. Supplementary Figure S2 in the Appendix. The input parameters are summarized in Table 2.

Table 2: Input parameters. Dates for the first and N^{th} observed cases correspond to the dates in the data set. A case in the context of the Alpha variant in the UK is defined as a sequenced sample, whereas a case in the context of COVID-19 in Wuhan is defined as a confirmed, symptomatic case. The total numbers of observed cases correspond to the size of the data set used to inform the model. The key epidemiological parameters are obtained from available literature.

Parameter	Symbol	Alpha	SARS-CoV-2
Total number of observed cases	N	406	3072
Date of first reported case	d_1^{obs}	Sep 20, 2020	Dec 10, 2019
Date of N^{th} observed case	$d_{K'}^{\text{obs}}$	Nov 11, 2020	Jan 19, 2020
Expected number of secondary cases	R	1.90 [11]	2.50 [30]
	κ	0.57 [10]	0.10 [31]
Secondary infection generation time	κ_t	0.83 [10]	0.83 [10]
	θ_t	6.60 [10]	6.60 [10]
Probability of case detection	p_{detect}	0.01 [10]	0.15 [30]
Time to detection	κ_τ	0.58 [10]	1.04 [32]
	θ_τ	12.0 [10]	6.25 [32]
Tolerance on the length of the cases time period*	δ_{tol}^τ	0.9	0.9
Tolerance on the daily number of cases*	δ_{tol}^Y	0.3	0.3

*Used to select the simulations that resemble the observed data; cf. the Methods section.

References

- [1] Jonathan E. Pekar, Andrew Magee, Edyth Parker, Niema Moshiri, Katherine Izhikevich, Jennifer L. Havens, Karthik Gangavarapu, Lorena Mariana Malpica Serrano, Alexander Crits-Christoph, Nathaniel L. Matteson, Mark Zeller, Joshua I. Levy, Jade C. Wang, Scott Hughes, Jungmin Lee, Heedo Park, Man-Seong Park, Katherine Zi Yan Ching, Raymond Tzer Pin Lin, Mohd Noor Mat Isa, Yusuf Muhammad Noor, Tetyana I. Vasylyeva, Robert F. Garry, Edward C. Holmes, Andrew Rambaut, Marc A. Suchard, Kristian G. Andersen, Michael Worobey, and Joel O. Wertheim. The molecular epidemiology of multiple zoonotic origins of SARS-CoV-2. *Science*, page eabp8337, July 2022. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.abp8337.
- [2] Jonathan Pekar, Michael Worobey, Niema Moshiri, Konrad Scheffler, and Joel O. Wertheim. Timing the SARS-CoV-2 index case in Hubei province. *Science*, 372(6540):412–417, April 2021. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.abf8003.
- [3] Michael Worobey. Dissecting the early COVID-19 cases in Wuhan. *Science*, 374(6572):1202–1204, December 2021. doi: 10.1126/science.abm4454.
- [4] David L. Roberts, Jeremy S. Rossman, and Ivan Jarić. Dating first cases of COVID-19. *PLOS Pathogens*, 17(6):e1009620, June 2021. ISSN 1553-7374. doi: 10.1371/journal.ppat.1009620.
- [5] Fabiana Gámbaro, Sylvie Behillil, Artem Baidaliuk, Flora Donati, Mélanie Albert, Andreea Alexandru, Maud Vanpeene, Meline Bizard, Angela Brisebarre, Marion Barbet, Fawzi Derrar, Sylvie van der Werf, Vincent Enouf, and Etienne Simon-Loriere. Introductions and early spread of SARS-CoV-2 in France, 24 January to 23 March 2020. *Eurosurveillance*, 25(26), July 2020. ISSN 1560-7917. doi: 10.2807/1560-7917.ES.2020.25.26.2001200.
- [6] Michael Worobey, Jonathan Pekar, Brendan B. Larsen, Martha I. Nelson, Verity Hill, Jeffrey B. Joy, Andrew Rambaut, Marc A. Suchard, Joel O. Wertheim, and Philippe Lemey. The emergence of SARS-CoV-2 in Europe and North America. *Science*, 370(6516):564–570, October 2020. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.abc8169.
- [7] Louis du Plessis, John T. McCrone, Alexander E. Zarebski, Verity Hill, Christopher Ruis, Bernardo Gutierrez, Jayna Raghwan, Jordan Ashworth, Rachel Colquhoun, Thomas R. Connor, Nuno R. Faria, Ben Jackson, Nicholas J. Loman, Áine O’Toole, Samuel M. Nicholls, Kris V. Parag, Emily Scher, Tetyana I. Vasylyeva, Erik M. Volz, Alexander Watts, Isaac I. Bogoch, Kamran Khan, COVID-19 Genomics UK (COG-UK) Consortium, David M. Aanensen, Moritz U. G. Kraemer, Andrew Rambaut, and Oliver G. Pybus. Establishment and lineage dynamics of the SARS-CoV-2 epidemic in the UK. *Science*, 371(6530):708–712, February 2021. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.abf2946.
- [8] Thomas Beneteau, Baptiste Elie, Mircea T. Sofonea, and Samuel Alizon. Estimating dates of origin and end of COVID-19 epidemics. *Peer Community Journal*, 1:e70, December 2021. ISSN 2804-3871. doi: 10.24072/pcjournal.63.
- [9] Philippe Lemey, Nick Ruktanonchai, Samuel L. Hong, Vittoria Colizza, Chiara Poletto, Frederik Van den Broeck, Mandev S. Gill, Xiang Ji, Anthony Levasseur, Bas B. Oude Munnink, Marion Koopmans, Adam Sadilek, Shengjie Lai, Andrew J. Tatem, Guy Baele, Marc A. Suchard, and Simon Dellicour. Untangling introductions and persistence in COVID-19 resurgence in Europe. *Nature*, 595(7869):713–717, July 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03754-2.
- [10] Peter Czippon, Emmanuel Schertzer, François Blanquart, and Florence Débarre. The stochastic dynamics of early epidemics: probability of establishment, initial growth rate, and infection cluster size at first detection. *Journal of The Royal Society Interface*, 18(184):20210575, November 2021. ISSN 1742-5662. doi: 10.1098/rsif.2021.0575.
- [11] Verity Hill, Louis Du Plessis, Thomas P Peacock, Dinesh Aggarwal, Rachel Colquhoun, Alesandro M Carabelli, Nicholas Ellaby, Eileen Gallagher, Natalie Groves, Ben Jackson, J T McCrone, Áine O’Toole, Anna Price, Theo Sanderson, Emily Scher, Joel Southgate, Erik Volz, Wendy S Barclay, Jeffrey C Barrett, Meera Chand, Thomas Connor, Ian Goodfellow, Ravindra K Gupta, Ewan M Harrison, Nicholas Loman, Richard Myers, David L Robertson, Oliver G Pybus, and Andrew Rambaut. The Origins and Molecular Evolution of SARS-CoV-2 Lineage B.1.1.7 in the UK. *Virus Evolution*, page veac080, August 2022. ISSN 2057-1577. doi: 10.1093/ve/veac080.

- [12] James B. Pettengill. The Time to Most Recent Common Ancestor Does Not (Usually) Approximate the Date of Divergence. *PLOS ONE*, 10(8):e0128407, August 2015. ISSN 1932-6203. doi: 10.1371/journal.pone.0128407.
- [13] Thibaut Jombart, Kevin van Zandvoort, Timothy W. Russell, Christopher I. Jarvis, Amy Gimma, Sam Abbott, Sam Clifford, Sebastian Funk, Hamish Gibbs, Yang Liu, Carl A. B. Pearson, Nikos I. Bosse, Centre for the Mathematical Modelling of Infectious Diseases COVID-19 Working Group, Rosalind M. Eggo, Adam J. Kucharski, and W. John Edmunds. Inferring the number of COVID-19 cases from recently reported deaths. *Wellcome Open Research*, 5:78, April 2020. ISSN 2398-502X. doi: 10.12688/wellcomeopenres.15786.1.
- [14] Andrew Rambaut, Edward C. Holmes, Áine O'Toole, Verity Hill, John T. McCrone, Christopher Ruis, Louis du Plessis, and Oliver G. Pybus. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nature Microbiology*, 5(11):1403–1407, November 2020. ISSN 2058-5276. doi: 10.1038/s41564-020-0770-5. Number: 11 Publisher: Nature Publishing Group.
- [15] WHO. Listings of WHO's response to COVID-19, June 2020. URL <https://www.who.int/news/item/29-06-2020-covidtimeline>. [Accessed: Nov 25, 2022].
- [16] The Novel Coronavirus Pneumonia Emergency Response Epidemiology Team. The epidemiological characteristics of an outbreak of 2019 novel coronavirus diseases (covid-19) – china. *China CDC Weekly*, 2:113, 2020. ISSN 2096-7071. doi: 10.46234/ccdcw2020.032.
- [17] Le Chang, Lei Zhao, Yan Xiao, Tingting Xu, Lan Chen, Yan Cai, Xiaojing Dong, Conghui Wang, Xia Xiao, Lili Ren, and Lunan Wang. Serosurvey for SARS-CoV-2 among blood donors in Wuhan, China from September to December 2019. *Protein & Cell*, page pwac013, May 2022. ISSN 1674-800X, 1674-8018. doi: 10.1093/procel/pwac013.
- [18] Michael Worobey, Joshua I. Levy, Lorena Malpica Serrano, Alexander Crits-Christoph, Jonathan E. Pekar, Stephen A. Goldstein, Angela L. Rasmussen, Moritz U. G. Kraemer, Chris Newman, Marion P. G. Koopmans, Marc A. Suchard, Joel O. Wertheim, Philippe Lemey, David L. Robertson, Robert F. Garry, Edward C. Holmes, Andrew Rambaut, and Kristian G. Andersen. The Huanan Seafood Wholesale Market in Wuhan was the early epicenter of the COVID-19 pandemic. *Science*, page abp8715, July 2022. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.abp8715.
- [19] Antonella Amendola, Marta Canuti, Silvia Bianchi, Sudhir Kumar, Clara Fappani, Maria Gori, Daniela Colzani, Sergei L. Kosakovsky Pond, Sayaka Miura, Melissa Baggieri, Antonella Marchi, Elisa Borghi, Gianvincenzo Zuccotti, Mario C. Raviglione, Fabio Magurano, and Elisabetta Tanzi. Molecular evidence for SARS-CoV-2 in samples collected from patients with morbilliform eruptions since late 2019 in Lombardy, northern Italy. *Environmental Research*, 215:113979, December 2022.
- [20] Mercedes Yanes-Lane, Nicholas Winters, Federica Fregonese, Mayara Bastos, Sara Perlman-Arrow, Jonathon R. Campbell, and Dick Menzies. Proportion of asymptomatic infection among COVID-19 positive persons and their transmission potential: A systematic review and meta-analysis. *PLOS ONE*, 15(11):e0241536, November 2020. ISSN 1932-6203. doi: 10.1371/journal.pone.0241536.
- [21] WHO. WHO-convened global study of origins of SARS-CoV-2: China Part. Technical report, WHO, 2021. URL <https://www.who.int/publications/i/item/who-convened-global-study-of-origins-of-sars-cov-2-china-part>.
- [22] Francesca Di Giallonardo, Sebastian Duchene, Ilaria Puglia, Valentina Curini, Francesca Profeta, Cesare Cammà, Maurilia Marcacci, Paolo Calistri, Edward Holmes, and Alessio Lorusso. Genomic Epidemiology of the First Wave of SARS-CoV-2 in Italy. *Viruses*, 12(12):1438, December 2020. ISSN 1999-4915. doi: 10.3390/v12121438.
- [23] Congressional Research Service. COVID-19 and China: A Chronology of Events (December 2019-January 2020), 2020. URL <https://crsreports.congress.gov/product/pdf/r/r46354>. [Accessed on: December 15, 2022].

- [24] Timothy W. Russell, Nick Golding, Joel Hellewell, Sam Abbott, Lawrence Wright, Carl A. B. Pearson, Kevin van Zandvoort, Christopher I. Jarvis, Hamish Gibbs, Yang Liu, Rosalind M. Eggo, W. John Edmunds, Adam J. Kucharski, CMMID COVID-19 working group, Arminder K. Deol, C. Julian Villabona-Arenas, Thibaut Jombart, Kathleen O'Reilly, James D. Munday, Sophie R. Meakin, Rachel Lowe, Amy Gimma, Akira Endo, Emily S. Nightingale, Graham Medley, Anna M. Foss, Gwenan M. Knight, Kiesha Prem, Stéphane Hué, Charlie Diamond, James W. Rudge, Katherine E. Atkins, Megan Auzenberg, Stefan Flasche, Rein M. G. J. Houben, Billy J. Quilty, Petra Klepac, Matthew Quaife, Sebastian Funk, Quentin J. Leclerc, Jon C. Emery, Mark Jit, David Simons, Nikos I. Bosse, Simon R. Procter, Fiona Yueqian Sun, Samuel Clifford, Katharine Sherratt, Alicia Rosello, Nicholas G. Davies, Oliver Brady, Damien C. Tully, and Georgia R. Gore-Langton. Reconstructing the early global dynamics of under-ascertained COVID-19 cases and infections. *BMC Medicine*, 18(1):332, December 2020.
- [25] Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B. Shah. Julia: A Fresh Approach to Numerical Computing. *SIAM Review*, 59(1):65–98, January 2017. ISSN 0036-1445, 1095-7200. doi: 10.1137/141000671.
- [26] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2022. URL <https://www.R-project.org/>.
- [27] Hadley Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN 978-3-319-24277-4. URL <https://ggplot2.tidyverse.org>.
- [28] Andrew Rambaut, Nick Loman, Oliver Pybus, Wendy Barclay, Jeff Barrett, Alesandro Carabelli, Tom Connor, Tom Peacock, David L Robertson, Erik Volz, and on behalf of COVID-19 Genomics Consortium UK (CoG-UK). Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in the UK defined by a novel set of spike mutations - SARS-CoV-2 coronavirus / nCoV-2019 Genomic Epidemiology, December 2020. URL <https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563>. [Accessed: Nov 25, 2022].
- [29] Shruti Khare, Céline Gurry, Lucas Freitas, Mark B Schultz, Gunter Bach, Amadou Diallo, Nancy Akite, Jose Ho, Raphael TC Lee, Winston Yeo, Gisaid Core Curation Team, Sebastian Maurer-Stroh, GISAID Global Data Science Initiative (GISAID), Munich, Germany, Bioinformatics Institute, Agency for Science Technology and Research, Singapore, Oswaldo Cruz Foundation (FIOCRUZ), Rio de Janeiro, Brazil, Institut Pasteur de Dakar, Dakar, Senegal, National Institutes of Biotechnology Malaysia, Selangor, Malaysia, Smorodintsev Research Institute of Influenza, St. Petersburg, Russia, Genome Institute of Singapore, Agency for Science Technology and Research, Singapore, China National GeneBank, Shenzhen, China, A*STAR Infectious Disease Labs (ID Labs), Singapore, National Public Health Laboratory, National Centre for Infectious Diseases, Ministry of Health, Singapore, and Department of Biological Sciences, National University of Singapore, Singapore. GISAID's Role in Pandemic Response. *China CDC Weekly*, 3(49):1049–1051, 2021. ISSN 2096-7071. doi: 10.46234/ccdcw2021.255.
- [30] Xingjie Hao, Shanshan Cheng, Degang Wu, Tangchun Wu, Xihong Lin, and Chaolong Wang. Reconstruction of the full transmission dynamics of COVID-19 in Wuhan. *Nature*, 584(7821):420–424, August 2020. ISSN 1476-4687. doi: 10.1038/s41586-020-2554-8.
- [31] Akira Endo, Centre for the Mathematical Modelling of Infectious Diseases COVID-19 Working Group, Sam Abbott, Adam J. Kucharski, and Sebastian Funk. Estimating the overdispersion in COVID-19 transmission using outbreak sizes outside China. *Wellcome Open Research*, 5:67, July 2020. ISSN 2398-502X. doi: 10.12688/wellcomeopenres.15842.3.
- [32] Jantien A Backer, Don Klinkenberg, and Jacco Wallinga. Incubation period of 2019 novel coronavirus (2019-nCoV) infections among travellers from Wuhan, China, 20–28 January 2020. *Eurosurveillance*, 25(5), February 2020. ISSN 1560-7917. doi: 10.2807/1560-7917.ES.2020.25.5.2000062.

Acknowledgements

SJ's postdoctoral fellowship was funded by a grant from the MODCOV19 platform of the National Institute of Mathematical Sciences and their Interactions (Insmi, CNRS) to FD. FD was funded by ANR-19-CE45-0009 (TheoGeneDrive). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript. We also thank Verity Hill [11] and Jonathan Pekar [1] for sharing their data and results for comparison with ours. We gratefully acknowledge all data contributors, i.e., the Authors and their Originating laboratories responsible for obtaining the specimens, and their Submitting laboratories for generating the genetic sequence and metadata and sharing via the GISAID Initiative, on which this research is based.