

RESEARCH ARTICLE

A Bayesian inference method to estimate transmission trees with multiple introductions; applied to SARS-CoV-2 in Dutch mink farms

Bastiaan R. Van der Roest^{1*}, Martin C. J. Bootsma^{1,2}, Egil A. J. Fischer³, Don Klinkenberg⁴, Mirjam E. E. Kretzschmar^{1,4}

1 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University, Utrecht, Netherlands, **2** Department of Mathematics, Faculty of Science, Utrecht University, Utrecht, Netherlands, **3** Department of Population Health Sciences, Faculty of Veterinary Medicine, Utrecht University, Utrecht, Netherlands, **4** Centre for Infectious Disease Control, National Institute for Public Health and the Environment (RIVM), Bilthoven, Netherlands

* b.r.vanderroest-2@umcutrecht.nl



OPEN ACCESS

Citation: Van der Roest BR, Bootsma MCJ, Fischer EAJ, Klinkenberg D, Kretzschmar MEE (2023) A Bayesian inference method to estimate transmission trees with multiple introductions; applied to SARS-CoV-2 in Dutch mink farms. *PLoS Comput Biol* 19(11): e1010928. <https://doi.org/10.1371/journal.pcbi.1010928>

Editor: Claudio José Struchiner, Fundação Getúlio Vargas: Fundacao Getulio Vargas, BRAZIL

Received: February 6, 2023

Accepted: November 12, 2023

Published: November 27, 2023

Copyright: © 2023 Van der Roest et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: There are no primary data in the paper. All data can be found at the GISAID database (<https://gisaid.org>) and as supplementary information included by the paper from Lu et al (2021) (<https://www.nature.com/articles/s41467-021-27096-9#Sec20>). The *phybreak* package is available at the github page <https://github.com/bastiaanvdroest/phybreak>. All code written in support of this publication, including the version of *phybreak* used, is publicly

Abstract

Knowledge of who infected whom during an outbreak of an infectious disease is important to determine risk factors for transmission and to design effective control measures. Both whole-genome sequencing of pathogens and epidemiological data provide useful information about the transmission events and underlying processes. Existing models to infer transmission trees usually assume that the pathogen is introduced only once from outside into the population of interest. However, this is not always true. For instance, SARS-CoV-2 is suggested to be introduced multiple times in mink farms in the Netherlands from the SARS-CoV-2 pandemic among humans. Here, we developed a Bayesian inference method combining whole-genome sequencing data and epidemiological data, allowing for multiple introductions of the pathogen in the population. Our method does not a priori split the outbreak into multiple phylogenetic clusters, nor does it break the dependency between the processes of mutation, within-host dynamics, transmission, and observation. We implemented our method as an additional feature in the R-package *phybreak*. On simulated data, our method correctly identifies the number of introductions, with an accuracy depending on the proportion of all observed cases that are introductions. Moreover, when a single introduction was simulated, our method produced similar estimates of parameters and transmission trees as the existing package. When applied to data from a SARS-CoV-2 outbreak in Dutch mink farms, the method provides strong evidence for independent introductions of the pathogen at 13 farms, infecting a total of 63 farms. Using the new feature of the *phybreak* package, transmission routes of a more complex class of infectious disease outbreaks can be inferred which will aid infection control in future outbreaks.

available at https://github.com/bastiaanvdroest/phybreak_multiple_introductions.

Funding: This work was funded as part of the research program of the Netherlands Center for One Health (www.ncoh.nl). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Author summary

Information about transmission routes is essential to design effective control measures for infectious disease outbreaks. A key question is whether outbreaks are caused by single imported index cases, or by multiple introductions. We introduce the concept of a ‘history host’ in a Bayesian model for simultaneous inference of phylogenetic and transmission trees. This artificial host—representing a source population—merges transmission trees by acting as the infector of all index cases, creating one large transmission tree. We applied the model to data from a SARS-CoV2 outbreak among mink farms in the Netherlands. We conclude that introductions from humans were an important factor in the outbreak, which implies that culling of infected farms alone was insufficient to prevent newly infected farms.

Introduction

Knowledge of who infected whom during an infectious disease outbreak is an important source of information. Characteristics of the outbreak, such as the generation time distribution, are derived from data on these transmission events [1]. Moreover, risk factors for transmission, such as distance between individuals or time lag since infection, can be more accurately quantified, if the infection chain is known. Several methods exist that use data on the time of symptom onset, contacts, or other proximity information, to reconstruct the most likely transmission links between cases [2–4]. Currently, genetic data is increasingly incorporated into epidemiological inference as an additional source of information to infer individual transmission events, transmission clusters, and even complete transmission trees [5–9]. The use of both genetic data (i.e., differences in nucleotides between different samples of the pathogen) and epidemiological data (e.g., time of sampling, contacts, and geographic distance) increases the evidence on who infected whom. Moreover, high-risk contacts and superspreaders can be identified when a model is based on both types of data [10, 11]. Therefore, several statistical methods have been developed which take both transmission and evolutionary dynamics of the pathogen into account [12–15].

Most methods assume a single introduction to the population of interest. However, there are many outbreaks where this assumption does not hold, e.g., *Staphylococcus aureus* or *Pseudomonas aeruginosa* are often introduced multiple times on a hospital ward when infected patients are admitted [16], highly pathogenic avian influenza (HPAI) outbreaks among farms are initiated multiple times by wild birds [17], and Foot and Mouth Disease (FMD) can be introduced multiple times from outside a district [18]. Control measures focusing on transmission between hosts may be less effective if there are also external introductions.

Currently, several methods to infer transmission trees from both genetic and epidemiological data are available. A method designed by Worby et al. [9] allows for multiple introductions, but it only has phenomenological distributions of genetic distances. There is no underlying mechanistic mutation model for the genetic difference within and between transmission trees. The *outbreaker2* package in R [19] also allows for multiple introductions, but there is only a phenomenological distribution of the genetic distances between trees. Moreover, *outbreaker2* assumes mutation at transmission, thereby ignoring within-host evolution of the virus. A method that uses a phylogenetic tree and within-host evolution is *Transphylo* [20], although transmission links are placed on a fixed phylogenetic tree. Both *outbreaker2* and *Transphylo* can deal with unsampled cases within the population, which can be used to link transmission clusters, although this is different than inferring introductions from an exogenous population.

To model multiple introductions from an exogenous population, Mollentze et al. [21] extended the transmission model of Morelli et al. [22], which simultaneously infers a transmission and phylogenetic tree. Here, the within-host evolution was modeled by the use of a binary tree, making the use of multiple samples per host problematic. Moreover and most importantly, there is no publicly available software to use the method.

To make optimal use of genetic and epidemiological data while allowing for multiple introductions of a pathogen, we propose a method to simultaneously infer introductions and transmissions consistent with an explicit phylogeny describing the genetic history of all samples. This extended version of the method developed by Klinkenberg et al. [23] aims to infer the transmission dynamics of an outbreak, i.e., who infected whom, from both genetic data of the pathogen and epidemiological data, such as the time of sampling and culling. Inference of the transmission tree and the phylogenetic tree is done simultaneously, concerning four processes: genetic diversity (within and between transmission trees), within-host diversity, transmission, and case observation, i.e., sampling time of hosts. Samples from posterior distributions of the model parameters are taken using a Markov-Chain Monte Carlo (MCMC) method. These samples provide information on how likely certain infection times and infectors of hosts are.

To address the possibility of multiple introductions, we relax the assumption of a single index case. We add an artificial host to the set of sampled hosts, which serves as an infector for all index cases (Fig 1). For this artificial host, we introduce the term ‘history host’, referring to the representation of the history of the lineages within the index cases. Using the history host, multiple outbreaks of a pathogen in the same population are merged into a single phylogenetic tree.

After evaluation of the performance on simulated outbreaks with single and multiple introductions, we illustrate the application of our method with an analysis of an outbreak of the SARS-CoV-2 virus in the Dutch mink farm industry. From April to November 2020, 63 mink farms tested positive for SARS-CoV-2. To investigate whether the virus was introduced on several farms, we estimated the number of introductions and compared the resulting

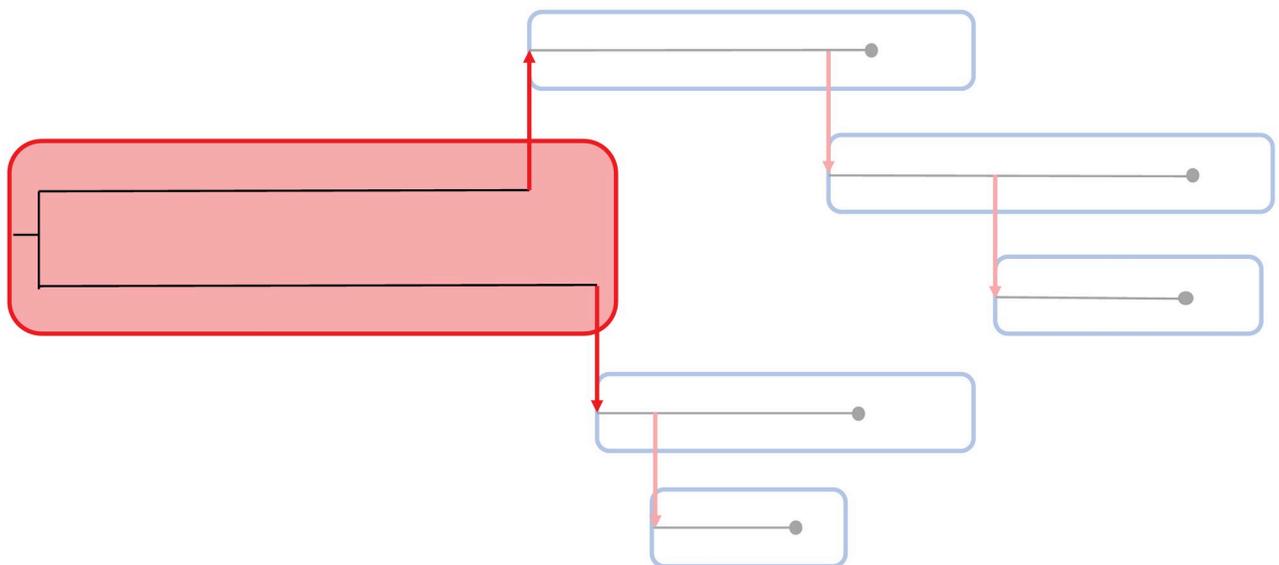


Fig 1. Overview of an outbreak with five sampled hosts and two introductions. The index cases of the sampled hosts (blue squares) are connected via the history host (red square). Coalescence of lineages happens at a different rate in the history host than in the sampled hosts. The black lines give the phylogenetic tree of the outbreak and the red arrows indicate transmissions between hosts.

<https://doi.org/10.1371/journal.pcbi.1010928.g001>

farm-to-farm transmission tree and phylogenetic tree to the phylogenetic tree obtained in [24, 25]. To describe the generation time distribution of infected farms, we used a within-farm model of time since infection, that takes measures to reduce the spread and culling of all animals into account. Furthermore, we implemented the possibility to include multiple sequences per host.

Results

Modelling with the history host

To infer the transmission tree of an infectious disease outbreak, we developed a Bayesian method in which four processes define the likelihood of a tree. Mutation events are modeled with a mutation rate μ . For the within-host dynamics, we make a distinction between the history host and the sampled hosts. The history host represents either a different population of the same host species, a different host species (e.g., zoonotic infection), or an environmental source. Therefore, it contains the evolution of the pathogen in the source population, with coalescence happening on a different time scale than within the sampled hosts (see Fig 1). Coalescence, i.e. lineages merging backward in time, is thus described by two rates: rate $1/w(\tau, r) = 1/r\tau$, with τ the time since infection, for the coalescence events in the sampled hosts, and rate $1/w(\tau, r_{\text{history}}) = 1/r_{\text{history}}\tau$ for the coalescence events in the history host. Timing of transmission is described by a generation time distribution, the time between infection and transmission. The generation time distribution follows in the default model a gamma distribution with mean m_G and shape a_G and for the analysis of the mink farm data we used the generation time described in the methods. Sampling time intervals, i.e., the time between infection and sampling, as a representation of case observations, are also described by a gamma distribution with mean m_S and shape a_S .

Improving the efficiency of the MCMC

The posterior is sampled by MCMC, with proposals that simultaneously change the phylogenetic and transmission trees. In case there are many introductions, starting the MCMC in overdispersed starting points lead to entrapment of the chain in local optima (S1 Fig). The history host contains many tips, i.e. the introductions, and therefore many branches and possible configurations, of which some are rarely proposed by the update steps of the MCMC currently implemented in phybreak. We solved this problem by (1) initializing the MCMC run by making each host an introduction and using the neighbor joining tree (NJ tree) for the phylogenetic tree in the history host, and (2) implementing the paralleled Metropolis Coupled Monte Carlo Markov Chain (p(MC³)) algorithm to give more freedom to the chain [26]. The p(MC³) algorithm allows multiple peaks in the landscape of trees to be more readily explored, and eases convergence without the cost of increased runtime. We tested for convergence by comparing the likelihood reached by each algorithm, to the likelihood reached by an MCMC run starting with the simulated (true) phylogenetic and transmission trees. It turned out that the NJ initialization and the p(MC³) algorithm always led to optimal convergence, whereas starting from a random tree and using MCMC sometimes ended up in a local optimum, especially when the number of introductions is high (Table A in S1 Results). We concluded that the configuration of the history host is a bottleneck for performance with initialized randomly, and decided to run all analyses with the NJ tree initialization, even though this breaks some assumptions of MCMC diagnostics.

Varying number of introductions and the coalescent rate

Before assessing in detail the method's performance to identify the correct introductions and infectors, we compared its performance in relation to different priors. Outbreaks with 20 hosts were simulated with 5 introductions and a set of default parameters (see [Materials and methods](#)). The outbreak-size of 20 hosts represents small real-life outbreaks, although we simulate larger outbreaks later on. The outbreaks were analyzed with uninformative priors on all parameters, informative priors on the mutation rate and mean generation and sampling intervals, and with all parameters set to their true values. The results were compared with respect to identifying the correct infectors, infection times, and parameter values. Only small differences were found between the results of each set of priors for the outbreaks with 5 introductions (Table B in [S1 Results](#)). For example, the mean numbers of correctly identified infectors were 15, 15, and 15.7, with increasing prior information.

Next, we simulated outbreaks with varying numbers of introductions and varying coalescent rates of the history host. While fixing the number of sampled hosts at 20, we simulated outbreaks with either 1, 2, 5, 10, 15, or 20 introductions. For each number of introductions, we used coalescent rates of 0.004, 0.02, and 0.1 coalescence events per day in the history host, against the background of a mean generation time interval of 1 day for transmission events. Thereby we changed the genetic variability of the index cases. Low coalescent rates in the history host result in long branch lengths and thus large differences in sequences of the index cases. Vice versa, high coalescent rates result in short branch lengths and small differences. Each combination of a number of introductions and coalescent rate was used for 25 simulated outbreaks, resulting in 450 outbreaks. We analyzed the simulated data with informative priors (Table B in [S1 Results](#)), as in outbreak research most of the time there is some prior information about the generation time and mutation rate.

Analyzing simulated outbreaks with 1 introduction resulted in a mean number (of 25 posterior medians) of 1 introduction, see [Fig 2A](#). This result did not change with the coalescent rate, because there is no coalescence in the history host. With 2 or 5 introductions, the estimated medians were still close to the simulated number. However, with 10 or more introductions the estimated medians were lower than the simulated number of introductions, and a high coalescent rate increased this gap. When all hosts are simulated as an introduction, no more than 40% of all introductions were truly identified by the inference method. This indicates that simulated clusters were merged due to the low genetic variability.

Approximately 70% of all hosts have correctly identified infectors when there was 1 introduction, and more than 95% of the hosts had their true infectors present in the 95% support set ([Fig 2B](#)). This is the set of infectors for a host with cumulative support of at least 95%, with infectors added by decreasing support. For more introductions and low coalescent rates, more infectors were correctly identified, whereas for higher coalescent rates the number of correctly identified infectors decreased.

Several types of incorrectly identified infectors can be distinguished. We define a transmission cluster as the set of hosts derived from one index case. We separate the errors into two classes: involving a single transmission cluster in both the simulated and estimated tree (single, S) or involving multiple transmission clusters in the simulated and/or estimated tree (multiple, M). The simulated or identified infector is then in a different transmission cluster than the case in the simulated or estimated tree. Both classes of error can be subdivided into three subclasses: both simulated and identified infectors are other cases in the data set (case to case, $C \rightarrow C$), the simulated infector is the history host and the identified infector is a case (history to case, $H \rightarrow C$), and the simulated infector is a case and the identified infector is the history host (case to history, $C \rightarrow H$) (see [S1 Fig](#)). In our analysis, we find that for small numbers of

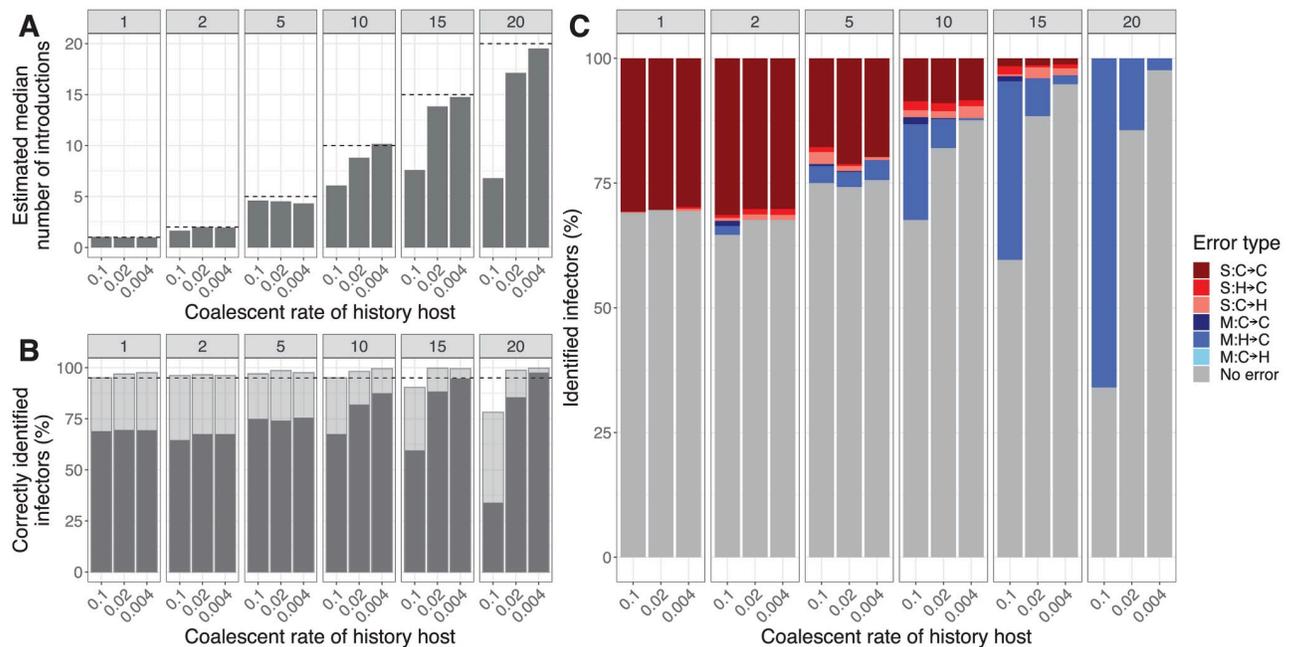


Fig 2. Analysis of simulated outbreaks with a varying number of introductions and coalescent rate (r_{history}) in the history host. The facets give the results for either 1, 2, 5, 10, 15, or 20 simulated introductions. (A) **The mean estimated median number of introductions.** The black line indicates the simulated number of introductions. (B) **Percentage of correctly identified infectors.** The grey bar indicates cases for which the true infector has the highest posterior weight. The transparent bar indicates cases for which the true infector is contained in the smallest set of candidate infectors with at least 95% of the posterior weight. (C) **Classification of the falsely identified infectors based on the highest support.** The grey bars indicate the correctly identified infectors. S: single transmission cluster involved, M: multiple transmission clusters involved. For the infector of a host: C2C: case becomes case, H2C: history becomes case, C2H: case becomes history.

<https://doi.org/10.1371/journal.pcbi.1010928.g002>

introductions, i.e. 1, 2, and 5, almost all errors are within a single transmission cluster and do not involve an index case (single none). For 10 introductions, this is around half of the errors, while the other half are merges of transmission clusters (multiple simulated). Larger numbers of introductions, i.e. 15 and 20, mostly lead to merged transmission clusters. With the number of introductions approaching the number of sampled hosts, there are only very few transmission events, such that it is hard to estimate the mutation rate or the coalescent rate in the history host correctly. An overestimation of the mutation rate, or an underestimation of the coalescent rate, makes it more likely that index cases are placed in the same cluster, causing merges. Fewer index cases imply more transmission events to estimate the correct parameter values. However, even if all parameters were fixed at their true value, an incorrect infector sometimes has the highest posterior probability (S2 Fig).

So, for low numbers of introductions, in these simulations up to 5, the model can reliably infer the number of introductions when informative priors are given for the model parameters. The number of introductions tends to be underestimated if there are many, due to the merging of clusters.

SARS-CoV-2 in mink farms: Analysis of simulated data

In 2020, an outbreak of SARS-CoV-2 occurred among mink farms in the Netherlands. Symptomatic infections in minks first occurred two months after the virus was introduced into the Dutch human population, which suggests that the outbreak was a spillover from humans to mink. To investigate whether there were multiple introductions of the virus into the mink

Table 1. Summary statistics of simulated SARS-CoV-2 outbreaks in mink farms.

	Number of simulated introductions					
	1	2	5	10	20	30
Estimated number of introductions	1.2	2.1	4.5	7	12.7	14.3
Correct infectors with highest support	75%	75%	71%	74%	74%	66%
Correct infectors in 95% CI	96%	97%	96%	97%	97%	92%

<https://doi.org/10.1371/journal.pcbi.1010928.t001>

farm population, we applied our extended method to sequence data collected from minks together with their time of sampling. Culling times of the farms were also known. To assess the accuracy of our method on outbreaks with sizes similar to the SARS-CoV-2 outbreak, we simulated and analyzed outbreaks with comparable settings of the number of hosts, mutation and coalescence rates, and prior distributions for the sampling time interval and the generation time interval (Material and Methods). Sequences and sampling times were simulated for these outbreaks. Again, we tested different numbers of introductions, for which 10 outbreaks each were simulated and analyzed. The results are shown in Table 1. Compared to the percentages of correctly identified infectors for outbreaks with 20 hosts, the model performs equally well for the larger outbreak size of 63 hosts. Around 70–75% of all infectors are correctly identified with the highest support, and the true infector of a host is present in the 95% CI set for at least 95% of all hosts. Only for a high number of introductions (e.g., 20, or 30 introductions), the performance decreases, due to merged clusters, with 5–10% (S3 Fig). The genome of the SARS-CoV-2 virus is found to have conserved regions and regions with higher mutation rates [27]. To see what the implications of these different mutation rates are for the results, analyzed with a single mutation rate over the complete genome, we simulated outbreaks for which only 50 base pairs are under mutation and analyzed these with our model (S5 Fig). If many mutations happen in only a small part of the genome, and possibly multiple mutations on the same positions, the model will estimate slightly more introductions. Furthermore, it will infer slightly more false infectors, although the false infectors are in the same transmission chain as the true infector. Therefore, it is important to see which farms are in the set of infectors of a farm, as the farm with the highest support will not always be the true infector. The performance of the model, however, is not notably different.

SARS-CoV-2 in mink farms: Analysis of the Dutch outbreak

During the first and second wave of SARS-CoV-2 infections in the Netherlands (starting in March 2020 and September 2020, respectively), 63 out of a total of 126 mink farms in the Netherlands were sampled positive for the virus. From the end of April 2020 to November 2020, genetic and epidemiological data were collected on these farms, including viral sequences, sampling times, and culling times. A phylogenetic analysis of the viral sequences showed 5 distinct genetic clusters of farms, based on their separation by sequences from human samples [24]. Classification by PANGO lineages [28] showed that each cluster contained one PANGO lineage, with 2 clusters containing the same lineage (S1 Data). One farm, NB-EMC-8, contained samples from 2 different clusters and is therefore divided into NB-EMC-8a and NB-EMC-8b in our analysis. While phylogenetic analysis could distinguish five clusters based on human intermediate samples, suggesting five introductions, it could not rule out multiple introductions within each cluster. For an estimate of the number of introductions without the need for intermediate samples from the source population, we analyzed this outbreak with our extended version of phylbreak. We set the following priors on the model parameters: $\mu_\mu = 3 \cdot$

Table 2. Summary statistics of SARS-CoV-2 outbreak in mink farms from real data.

Parameter inference	median (95% quantile range) of posterior
μ	$5.5 \cdot 10^{-6}$ ($4.7 \cdot 10^{-6}$; $6.4 \cdot 10^{-6}$)
m_S	11.9 (10.2; 14.1)
r_{history}	30.5 (17.2; 53.6)
L	1.0 (0.6; 1.5)
Tree inference	
Number of introductions	13 (11; 14)
Time of first coalescent event in history	-51.7 (-87.4; -27.9)

<https://doi.org/10.1371/journal.pcbi.1010928.t002>

10^{-6} substitutions per nucleotide per day, $\sigma_\mu = 1 \cdot 10^{-6}$ [29] and the mean $r_{\text{history}} = 20$, i.e., coalescence between any pair of lineages in a host happens with a rate $1/20$, with shape equal to 3 (see [Materials and methods](#)). The mean of the prior distribution of the introduction rate is $5/180$, as five genetic clusters were reported within 180 days, with a shape equal to 3. Finally, we set the prior mean sampling time μ_S at 10 days, with standard deviation $\sigma_S = 2$, as infection is expected to happen 1–2 weeks before sampling [30].

The method estimated the time of the first coalescent event in the history host on March 4th, 2020 ([Table 2](#)). The reduction factor of infectiousness after sampling L was estimated at 1, meaning that the method did not find an influence of sampling on infectiousness. We find 13 introductions in the maximum parent credibility tree (see [Fig 3](#)), of which 11 have minimal support above 0.5. The median number of introductions in all cycles was 13, with the first and third quartile being 11 and 14 introductions respectively ([S8 Fig](#)). Six introductions initiated a

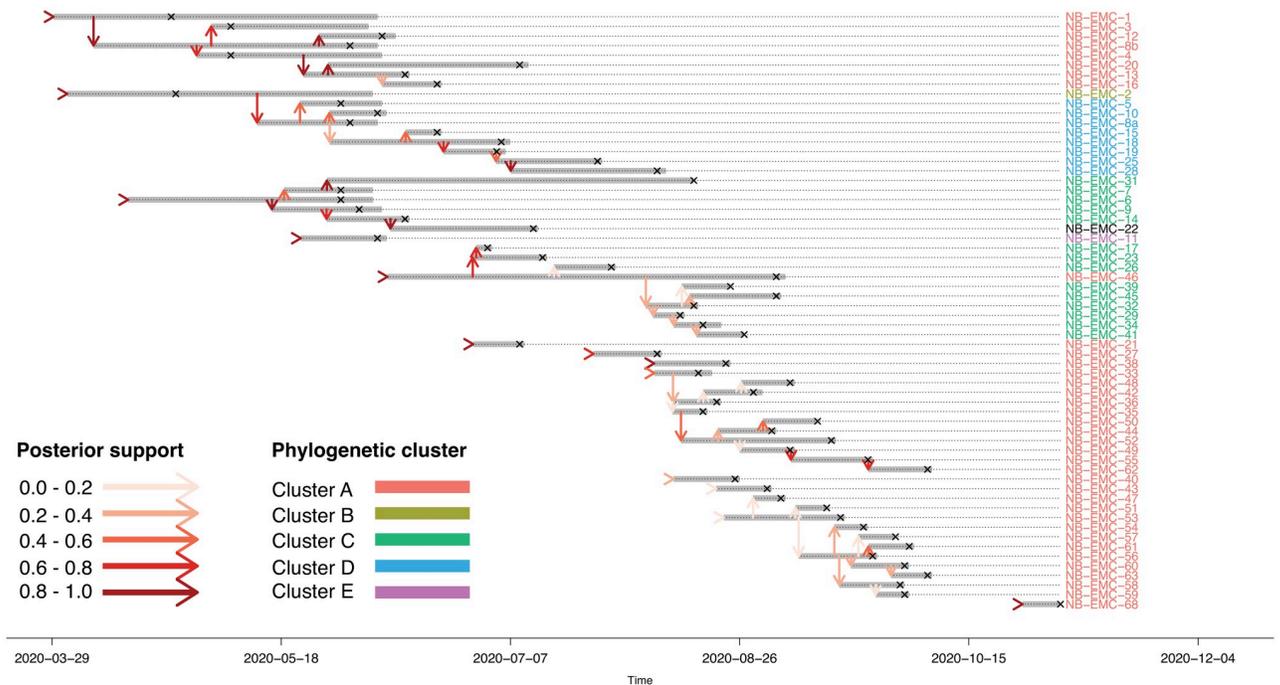


Fig 3. Maximum parent credibility transmission tree of a SARS-CoV-2 outbreak in mink farms. In total 13 introductions are found in the outbreak. Vertical arrows represent transmission links and all arrows are colored according to the support in the posterior distribution. The grey bars show the infectiousness of the hosts and hosts are sampled at the crosses. Host labels are colored according to phylogenetic clusters found by Lu et al. [24].

<https://doi.org/10.1371/journal.pcbi.1010928.g003>

transmission chain, whereas the other 7 were single cases. By coloring the host labels, we see that the method divided the hosts into subtrees similar to the phylogenetic clusters found by Lu et al. [24]. Two genetic clusters, i.e. cluster B and cluster D, were merged into a single transmission cluster, and with a genetic distance of only 4 nucleotides they belong to the same PANGO lineage. Genetic cluster C is split into two transmission clusters, with NB-EMC-46 as the index case of one of them. NB-EMC-46 was placed in genetic cluster A, but its samples were found to belong to multiple PANGO lineages, including the lineage of genetic cluster C. This indicates that farm NB-EMC-46 is infected multiple times. The large genetic cluster A is separated into multiple transmission clusters, meaning that not all genetically clustered farms are linked by one transmission chain. We find that the single cases which are part of this phylogenetic cluster have common ancestors with cases in the human population (S6 Fig). Time of infection and genetic distance made it less likely that the single farms were part of the transmission cluster of farms. In the later stage of the outbreak, there are two larger transmission chains, for which the exact index case is less certain (S9 Fig). There is support for the scenario that these transmission clusters are merged into one. The multivariate potential scale reduction factor (PSRF) of 3 p(MC³) chains was 1.48 after a burn-in of 10,000 cycles and sampling of 10,000 cycles. Only the PSRF of the loglikelihood was above 1.1, namely 1.75, although trace-plots show no further converging of the loglikelihood (S10 Fig). In conclusion, by using a phylogenetic model combining the phylogenetic history of the samples with the transmission history between the farms, we were able to distinguish farm-to-farm transmission routes within a group of farms with a common introduction from the human population.

Our extensions are implemented in the package *phybreak* [23] for the R software [31] and can be found at <https://github.com/bastiaanvdroest/phybreak>. The package version used, together with the code for the analyses, is found at https://github.com/bastiaanvdroest/phybreak_multiple_introductions.

Discussion

The method presented enables for the first time to simultaneously estimate the phylogenetic tree and the transmission tree of an outbreak in the case where there may have been multiple introductions. The inference is done without breaking the dependencies between mutations, within-host dynamics, transmission, and observation. By modeling the history of lineages infecting index cases through a phylogenetic tree in a history host, we can distinguish between single and multiple introductions. As an extension to the model of Klinkenberg et al. [23], we now have an easily accessible method for transmission tree inference, with the possibility to assess multiple introductions.

From analyses of simulated outbreaks, we conclude that the model can infer the true number of introductions if there are few introductions compared to the total outbreak size. For an increasing number of introductions, the model increasingly underestimated the number of introductions, but the posterior distribution did include the actual number of introductions. The simulated index cases which were incorrectly identified as non-index cases did have support as an index in the posterior trees. This means that interpretation of the transmission trees should take into account the support as index for cases.

The ability to infer multiple introductions in the analysis of an outbreak is not only useful for finding transmission clusters but also gives valuable information on how to respond to an outbreak. In the case of multiple introductions, measures aimed at reducing transmission events need to be complemented by preventing introduction from outside the target population. Therefore it is of great importance to distinguish between single and multiple introductions of a pathogen in a population. With simulated data sets, we showed that our method is a

useful tool to make this distinction: outbreaks with a single introduction are almost always inferred to have a single index case, and outbreaks with multiple introductions are seldom inferred to have a single introduction.

Although the model can distinguish between single or multiple introductions, the accuracy strongly depends on genetic variability. High genetic variability makes it easier to distinguish clusters of hosts and thus gives more weight to the true number of introductions in the posterior. Low genetic variability, however, will cause sub-trees to be merged and therefore will lead to an underestimation of the number of introductions. The accuracy of the results cannot be determined in advance due to the variability in the external source population. This variation is determined by the mutation rate and the effective population size of the history host. When available, strong priors on the mutation rate and coalescent rate in the history host will increase the accuracy, although even with the true values of the model parameters sub-trees will not always be separated. In that case, there is too little information in the genetic and epidemiological data to find all introductions.

Transmission clusters of an infectious disease outbreak in a population are often derived using phylogenetic analyses. However, with closely related index cases, defining clusters may become arbitrary. If obtainable sequences sampled outside of the study population may help to discriminate the clusters by acting as ‘missing links’ between clusters, but discrimination is not so likely if clusters are closely connected. As with the SARS-CoV-2 outbreak in minks, low genetic variability may cause transmission clusters to be merged in the phylogenetic tree, thereby underestimating the number of introductions. We have shown that our method can be used as an alternative approach, which only depends on the genetic data from the study population. Moreover, with the addition of epidemiological data, e.g. sampling times and culling times, it can differentiate genetically similar transmission clusters.

Application of the model to a SARS-CoV-2 outbreak in the Dutch mink farms led to confirmation of previously found phylogenetic clusters, although the phylogenetic clusters are broken down into multiple transmission clusters. These transmission clusters are composed of individual infections along with a larger transmission tree. We split farm NB-EMC-8 based on the genetic clustering of the samples taken on this farm. Without this split, a transmission cluster would have been formed containing multiple PANGO lineages and always having NB-EMC-8 separating the two genetic clusters within that transmission cluster. Farm NB-EMC-46 is also likely to be infected multiple times, as in our results it is the index case of a transmission cluster containing samples from a different genetic cluster than NB-EMC-46. Currently, our method does not allow for multiple infections of a host with different strains, and therefore these clusters could not be separated by the estimation procedure. Extending the method to allow multiple infections of the same host is a challenge for future development. Clustering of the farms could be performed by setting a cutoff for the SNP distance between samples of the farms (Table D in [S1 Results](#)). A 3 SNP cutoff results in 6 genetic clusters, but it excludes relations between farms found in the phylogeny of Lu et al. [24]. A 2 SNP cutoff will separate the clusters found with a 3 SNP cutoff, but it will produce too many introductions compared to our results. The usage of samples of humans around and on the farms in building the phylogeny makes it possible to derive the genetic clusters. Here we show that we come to similar conclusions, but do not need samples of the source population to distinguish transmission clusters. Often such data is not available, for example, with introductions from other countries, the general population is a case of non-notifiable diseases or from wildlife.

The possibility to distinguish multiple introductions of a pathogen into a host population opens up a new avenue for the analysis of outbreaks. Currently, more methods can link transmission clusters, as described in the introduction, but a software package to infer introductions

from an exogenous population was not yet publicly available. However, the method assumes a large population of which a small part gets infected and where contact is equally likely for all pairs of hosts. An outbreak on, for instance, a hospital ward does not meet this assumption with its small population size, in and outflow of patients, and spatial distance between patients. To address these assumptions, the population size has to be accounted for, and contact data, i.e., possible (in)direct contacts between hosts, as well as the geographical location of hosts give a probability of the contact between hosts. Transmission routes can be excluded based on these data sources, such that the certainty of the results increases.

In conclusion, we developed a new method for transmission tree inference which makes it possible to estimate the number of introductions of a pathogen during an outbreak. The analysis of the SARS-Cov-2 outbreak in Dutch mink farms shows multiple introductions of the virus, indicating that even with fully controlling farm-to-farm transmission, newly infected farms would arise by new introductions from the human population. Our method opens the way to evaluate outbreaks in such a way that information about new introductions can be derived; knowledge that is useful for policy-making.

Materials and methods

Tree inference model

The transmission and phylogenetic tree inference model describes the likelihood of observing an infectious disease outbreak based on the epidemiological and genetic links between hosts and samples. The outbreak dynamics are described by four processes: incidence of new cases by introduction from outside or transmission by existing cases, the observation of the pathogen through sampling, the dynamics of the pathogen within infected hosts and the history host, and genetic mutations in the pathogen. The inference is done by a Bayesian analysis, using Markov-Chain Monte Carlo (MCMC) to obtain samples from the posterior distributions of the phylogenetic and transmission tree, along with all outbreak parameters, formed by prior distributions and four likelihood functions for the four processes. We will briefly summarize the likelihood functions, the posterior distributions, and the update steps in the MCMC run.

The incidence of cases after the first introduction is modeled by two independent processes: additional introduction from outside the study population and transmission between hosts. Additional introductions occur with exponentially distributed waiting times with a rate λ_{intro} , after the first introduction until the last sample time. We denote by T the time between the first introduction and the last sample time, and k the number of introductions. Transmission occurs with a dynamic rate, depending on the times since infection of infected hosts, described by the generation time distribution. This is a Gamma distribution with shape a_G and mean m_G . By the use of vector \mathbf{I} of all infection times, including introductions, and the numeric vector \mathbf{M} indicating the infectors of all hosts and 0 for introductions, the probability density function of the generation time of a host i , with $M_i \neq 0$, is $d_{\Gamma(a_G, m_G)}(I_i - I_{M_i})$. The likelihood for the transmission tree is therefore:

$$\Pr(\mathbf{I}, \mathbf{M} | a_G, m_G) = \lambda_{\text{intro}}^{k-1} \cdot e^{(-\lambda_{\text{intro}} * T)} \cdot \prod_{i | M_i > 0} d_{\Gamma(a_G, m_G)}(I_i - I_{M_i}) \quad (1)$$

For sampling, we assume that all hosts are detected and sampled at random times after they were infected, according to a Gamma distribution with shape a_S and mean m_S . The likelihood

uses the vector S of sampling times of all hosts and is therefore:

$$\Pr(\mathbf{S}|\mathbf{I}, a_s, m_s) = \prod_i d_{\Gamma(a_s, m_s)}(S_i - I_i) \tag{2}$$

The phylogenetic tree P describes the evolutionary history of all sampled sequences and is built from the phylogenetic mini-trees for each host, connected through the transmission links. The introductions are connected by a phylogenetic tree in a separate ‘history host’. Each mini-tree has tips formed by samples and lineages from secondary cases, and a single root which is a tip in the mini-tree of the infector. Coalescent processes form mini-trees. In (observed) hosts, a rate $1/w(\tau, r)$ describes coalescence between any pair of lineages within the host going backward in time; in the history host, the rate is constant over time: r_{history} . In our analysis, we use $w(\tau, r) = r\tau$, the linearly increasing within-host effective population size of the pathogen at forward time τ since infection of the host. A consequence of this function is a complete bottleneck: only one lineage is transmitted between hosts. In the phylogenetic tree P of the outbreak with the set of nodes V , there are three sets of nodes: sampling nodes V_S , i.e. the tips of the tree where sampling took place, coalescent nodes V_C and transmission nodes V_T , where a lineage goes from the infector to its infectee. For node x , τ_x gives the time of the node since infection of the host. The number of lineages in host i at time τ is then denoted by $L_i(\tau)$:

$$L_i(\tau) = 1 + \sum_{\substack{x \in P_i \\ x \in V_C}} (u(\tau - \tau_x)) - \sum_{\substack{x \in P_i \\ x \in (V_T \cup V_S)}} (u(\tau - \tau_x)) \tag{3}$$

where $u(\tau)$ is the heaviside step function, i.e. $u(\tau) = 0$ if $\tau < 0$, and $u(\tau) = 1$ if $\tau \geq 0$. The likelihood of each host’s tree is then

$$\Pr(P_i|S_i, \mathbf{I}, \mathbf{M}, r) = \exp\left(-\int_0^\infty \binom{L_i(\tau)}{2} \frac{1}{w(\tau, r)} d\tau\right) \cdot \prod_{\substack{x \in P_i \\ x \in V_C}} \frac{1}{w(\tau_x, r)} \tag{4}$$

with $\binom{0}{2} \equiv \binom{1}{2} \equiv 0$. Here, the first term is the probability of having no coalescent event during the intervals in which there are two or more lineages, and the second term is the product of coalescent rates at the coalescent nodes. The prior distribution of the slope r is Gamma distributed with shape a_r and rate b_r . Those were set to $a_r = b_r = 3$ for all analyses. For the history host, we assume that the coalescent rate is constant over time, so $w(\tau, r_{\text{hist}}) = r_{\text{hist}}$. The total likelihood of the within-host dynamics is the product of all hosts’ likelihoods:

$$\Pr(P|\mathbf{S}, \mathbf{I}, \mathbf{M}, r) = \Pr(P_0|\mathbf{I}, \mathbf{M}, r_{\text{history}}) \cdot \prod_{i|j>0} \Pr(P_i|S_i, \mathbf{I}, \mathbf{M}, r) \tag{5}$$

Mutations are described by a Jukes-Cantor model, stating that any of the four nucleotides have equal probability to mutate to, with a fixed mutation rate μ for all sites in the set of sequences \mathbf{G} . For all coalescent and transmission nodes x , which occur at time t_x with parent

node v_x , the mutation likelihood is:

$$\Pr(\mathbf{G}|P, \mu) = \prod_{loci \in \{A,C,G,T\}^{3n-1}} \sum_x \prod \left(\frac{1}{4} - \frac{1}{4} \exp(-\mu(t_x - t_{v_x})) \right)^{\mathcal{I}_{mut}(1-N)} \cdot \left(\frac{1}{4} + \frac{3}{4} \exp(-\mu(t_x - t_{v_x})) \right)^{(1-\mathcal{I}_{mut})(1-N)} \quad (6)$$

Here, \mathcal{I}_{mut} indicates if a mutation occurred on the branch between x and v_x , and N indicates if a branch ends with a tip with an unknown nucleotide ('n' in the sequence). We use here a strict molecular clock model, i.e. one mutation rate for all branches of the phylogenetic tree, because, on the time frame of most outbreaks, there will not be any effect of different mutation rates. In the history, changes in mutation rates are met by the coalescent rate of the history host. The likelihood is calculated using Felsenstein's pruning algorithm [32].

The transmission tree and its parameters are inferred by a Bayesian analysis, using Markov-Chain Monte Carlo (MCMC). From the MCMC we obtain samples from the posterior distributions of the model parameters, the infectors, and the infection times of all hosts. The posterior distribution, with θ the set of model parameters, is given by

$$\Pr(\mathbf{I}, \mathbf{M}, P, \theta | \mathbf{S}, \mathbf{G}) \propto \Pr(\mathbf{G}|P, \theta) \cdot \Pr(P|\mathbf{S}, \mathbf{I}, \mathbf{M}, \theta) \cdot \Pr(\mathbf{S}|\mathbf{I}, \theta) \cdot \Pr(\mathbf{I}, \mathbf{M}|\theta) \cdot \Pr(\theta) \quad (7)$$

MCMC sampling

An MCMC run is run to get the posterior distribution of the model parameters, together with the transmission and phylogenetic tree of the outbreak. The MCMC runs were initialized by first choosing the means of priors for the parameters (except for μ), then constructing the transmission and phylogenetic trees, and finally computing a value for μ . The trees were constructed by first sampling infection times from the observed sampling times and sample time distribution. All cases were assumed to be index cases (other options are possible within the package), and the topology of the phylogenetic tree was made with the neighbor-joining algorithm using the first sequence of each host. The times of the coalescent nodes were simulated with the coalescent model. This guaranteed an optimized tree topology in the history host, not needing to be reached by sampling in the MCMC run. The parameter μ was for the initial state set to be the tree parsimony (the number of mutations on the tree) divided by the sum of all branch lengths and the genome size. The default prior distributions for the model parameters are found in Table 3. The priors for m_G and m_S are translated into a prior for the rate parameter in the Gamma distribution. More detail about the prior and posterior distributions is included in S1 Methods. Per iteration cycle, each host is picked once in random order as the focal host. A new infection time I'_i is proposed for focal host i and consecutive steps are made according to this new infection time. At the start of a proposal, there are two main ways of updating: within a sub-tree, by following all hosts with a common index case along their transmission links, or between sub-trees. Here we will describe the proposal step for updating between sub-trees, as this is the step where the number of introductions can be altered. The update steps within a sub-tree are as in the original *phybreak* package and can be found in S1 Methods.

Three situations describe the possibility to update the transmission tree between sub-trees (see Fig 4):

1. The focal host i is the history host. In this case, new coalescent times are proposed. Optionally, a new phylogenetic mini-tree can be proposed.

Table 3. Prior distributions of the model parameters.

Parameter	Description	Type of distribution	Distribution parameters
$\log_{10}(\mu)$	Mutation rate	$N(\mu, \sigma)$	$\mu_{\log_{10}}(\mu) = -4; \sigma_{\log_{10}}(\mu) = 0.5$
m_G	Mean generation time interval	$D(\mu_{m_G}, \sigma_{m_G})$	$\mu_{m_G} = 1; \sigma_{m_G} = \infty$
m_S	Mean sampling time interval	$D(\mu_{m_S}, \sigma_{m_S})$	$\mu_{m_S} = 1; \sigma_{m_S} = \infty$
r	Within-host coalescent rate	$\Gamma(a_r, b_r)$	$a_r = 3; b_r = 3/1$
r_{history}	History host coalescent rate	$\Gamma(a_{r_{\text{history}}}, b_{r_{\text{history}}})$	$a_{r_{\text{history}}} = 1; b_{r_{\text{history}}} = 1/100$
r_{intro}	Introduction rate	$N(\mu_{r_{\text{intro}}}, \sigma_{r_{\text{intro}}})$	$\mu_{r_{\text{intro}}} = 1; \sigma_{r_{\text{intro}}} = \infty$

<https://doi.org/10.1371/journal.pcbi.1010928.t003>

2. The focal host i is an index case. An infection time I'_i is proposed. If this I'_i is before the first transmission from host i , a new infector M'_i is proposed out of the hosts which are infectious at time I'_i . Two situations are now possible:
 - a If $M'_i = 0$, then host i remains an index case, with infection time I'_i .
 - b If $M'_i \neq 0$, then host i is no longer an index case, and there is one fewer introduction. Host i and its descendants will be merged as a branch to another sub-tree.
3. The focal host i is not an index case. An infection time I'_i is proposed. If this I'_i is before the first transmission from host i , a new infector M'_i is proposed out of the hosts which are infectious at time I'_i . Two situations are now possible:

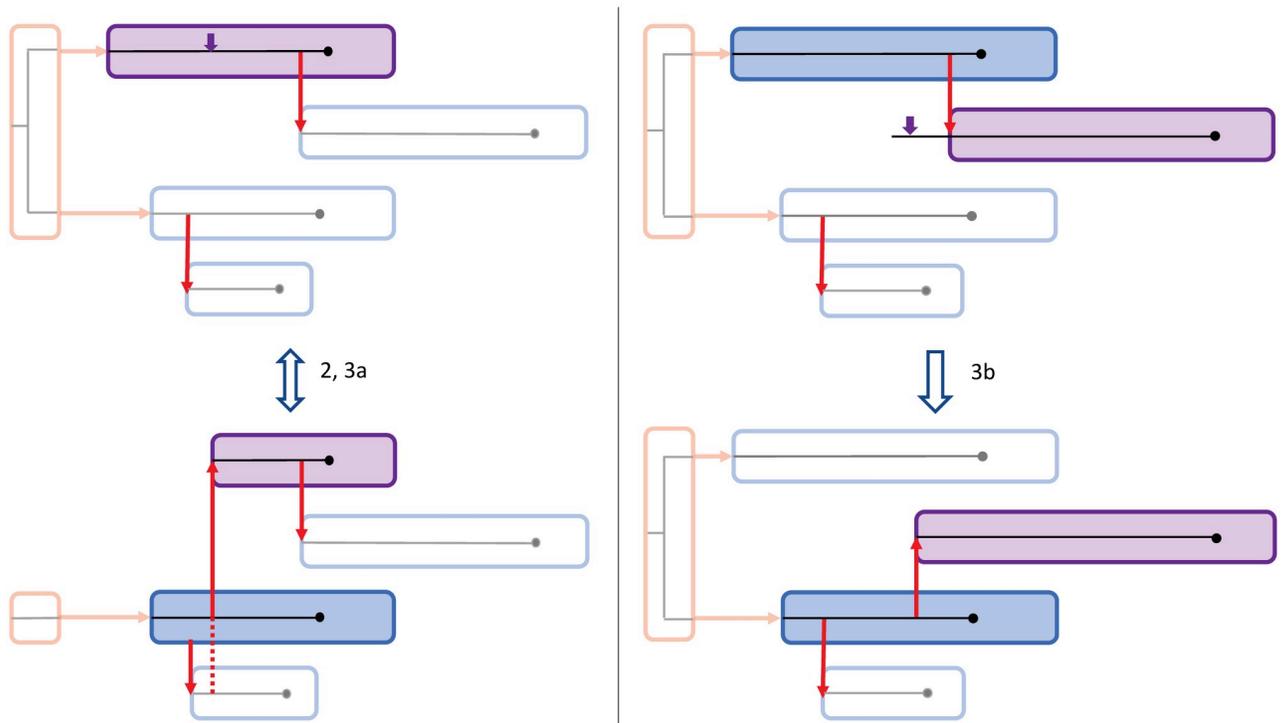


Fig 4. Proposal steps for updates between sub-trees. In purple is the focal host, with the purple arrow indicating the proposed infection time I'_i . The red arrows indicate the transmission events and the history host is colored red, with the introductions as transmission from the history host. 2: Losing an introduction by proposing a new infector $M'_i \neq 0$ for an index case. 3a: The reverse of 2, by proposing a new infector $M'_i = 0$ for a non-index case. 3b: Switching sub-trees by proposing a new infector $M'_i \neq 0$ on a different sub-tree. Situation 3b is also possible within the same sub-tree.

<https://doi.org/10.1371/journal.pcbi.1010928.g004>

- a If $M'_i = 0$, then host i will become an index case, and there is one extra introduction. The new sub-tree consists of host i and all of its descendants.
- b If $M'_i \neq 0$, then host i either switch to another branch in its sub-tree or switch to another sub-tree. There is no change in the number of introductions.

Each proposal step is followed by proposing new phylogenetic mini-trees for all hosts involved. The proposal distributions and acceptance probabilities of all steps are described in [S1 Methods](#). The MCMC run is run according to the (MC)³ algorithm described by Altekar et al. [26] to more readily explore multiple peaks in the posterior tree landscape, decreasing the probability to get entrapped in a local optimum. The chains consisted of 35,000 cycles of which the first 10,000 were used as burn-in.

Construction and analysis of simulated outbreaks

To verify the implementation of multiple introductions in the model, we simulated outbreaks including one or more index cases, and analyzed them by MCMC runs. The simulation of an outbreak starts with the simulation of a transmission tree:

1. Set an observation size, i.e. the number of hosts, the number of introductions k , and the duration of the outbreak T .
2. Calculate the optimal population size in which to simulate the outbreak from the set parameter R_0 and the observation size.
3. Sample $k - 1$ introduction times from the exponential waiting time distribution with rate λ_{intro} . The introduction time of the first index case will be 0, and other introductions are at cumulative waiting times from the first index.
4. For the index cases, sample the number of secondary cases from a Poisson distribution with parameter R_0 .
5. The generation time between two hosts is Gamma distributed with shape a_G and mean m_G . After infection, the sampling of a host takes place after a Gamma distributed time with shape a_S and mean m_S .
6. Repeat steps 3 and 4 for the complete population size, where the infection time for a host is not after T . Remove non-index cases without any links.
7. Repeat 3–6 till the desired observation size was given.
8. Add the history host and connect the index cases to this host.

After the simulation of the transmission tree, the phylogenetic tree is constructed by simulating phylogenetic mini-trees for each host. Coalescent times are sampled according to the given coalescent rate $1/w(\tau, r)$. Edges between sample, coalescent, and transmission nodes are made backward in time. In the history host, coalescence events occur with a constant rate $1/r_{\text{history}}$.

For the sequences, we sample the number of mutations from a Poisson distribution with parameter equal to $\lambda = \mu \cdot \text{sequencelength} \cdot \text{totallengthofalldges}$, where μ is the mutation rate. The mutations are distributed over the edges, with weights the lengths of the edges. For each mutation, a uniform random locus is changed to a uniform random nucleotide.

We simulated outbreaks with a basic set of parameter values, the same as in Klinkenberg et al. [23], ($m_G = 1$, $a_G = 10$, $m_S = 1$, $a_S = 10$, $R_0 = 1.5$, $r = 1$, a sequence length of 10^4

nucleotides and a mutation rate of $\mu = 10^{-4}$), with new parameters at $\lambda_{\text{intro}} = 1$. The number of introductions varied between the simulations to assess the performance of the model. MCMC chains were run following the $(MC)^3$ algorithm, with 3 parallel chains with heats 1, 0.5, and 0.333. The chains are 35,000 cycles long, of which the first 10,000 cycles are used as burn-in. Posterior distributions for infectors, infection times, and model parameters are collected from the remaining 25,000 cycles.

Analysis of SARS-CoV-2 in mink farms: Simulated data

As an application of the method, we analyzed the SARS-CoV-2 outbreak in the Dutch mink industry in 2020 [24]. We described the outbreak by taking the farms as hosts. The prior distributions of the model parameters are set as follows: we set the mean sampling time interval $m_S = 10$ days (with a shape $a_S = 3$), as the time between infection and detection was estimated to be 1–2 weeks [30]. We set the mean introduction rate to 5/180 per day (with a shape of 3), as five different clusters were found during the outbreak, which lasted for approximately 180 days, by Lu et al. [24]. The coalescent rate parameter r_{history} was set to 20, i.e., 1/20 coalescent events per pair of lineages. With an expected number of 5 introductions, this rate represented the introduction of the virus in the Netherlands two months before the first positive mink sample. The other prior distributions were set to default.

As the hosts are farms here, we introduced an infectiousness function describing the growth and circulation of the virus within the mink population of a farm. This function replaced the gamma distribution for the likelihood that one farm infected another. We assumed that infectiousness follows a logistic curve, with a reduced level after detection at time T_s , and exponential decay after culling at time T_c :

$$I = \begin{cases} \frac{1}{1 + ae^{-gt}} & t < T_s \\ \frac{L}{1 + ae^{-gt}} & T_s < t < T_c \\ \frac{L}{1 + a \cdot e^{-gT_c}} \cdot e^{-C(t-T_c)} & t > T_c \end{cases} \quad (8)$$

Here, $a = 1 \cdot 10^{-4}$ is the initial part of the mink population at a farm being infected, g is the growth rate, and t is the time after infection of the farm. Parameter L is estimated to see if there was some reduction of infectiousness after detection, and C is a fixed value. The logic behind the exponential decay after culling is to have a delay in clearing of infectiousness, for instance due to environmental contamination. Because the values for T_s and T_c differ per farm, the infectiousness curves differ between the farms. Therefore we normalize the curves, such that the average infectiousness function integrates to 1, while accounting for higher total infectivity of longer infected farms. Another addition used for the mink farms was to include multiple samples per farm. Phylogenetic mini-trees are then built with multiple lineages within a farm, increasing the amount of genetic data. For the sampling time distribution, only the first sample of each host is used.

To test the new model, with a similar history host, and sampling time distribution, we simulated outbreaks with the same parameters as before but with the new infectiousness curve. Culling times were set 15 days after infection, such that the hosts have a fixed infectiousness curve. As for the outbreak size, we used 63 hosts with 1 sample per host. Prior distributions were set with the same parameter values as the analysis of the real data. We set C to 5, such that in 5 days after culling the infectiousness of a farm was effectively 0. We varied the number of introductions, from 1, 2, 5, 10, 20, up to 30 introductions. As the genome of the SARS-CoV-2

virus has different mutation rates for certain regions, we also tested the performance of the model with a single mutation rate on simulations where only a region of 50 bp is mutated. For this, we used the transmission and phylogenetic trees of the simulated outbreaks of the minks with 10 introductions. Sequences of 50 bp were simulated with the same mutation rate as previous, adjusted for the genome length in order to get the same amount of mutations. When mutations were acquired, sequences were extended to their original length and the outbreaks were analyzed with the same parameter settings as before. Results of the simulated SARS-CoV-2 outbreak were obtained by running three parallel chains, with 25,000 cycles each, according to the $(MC)^3$ algorithm.

Analysis of SARS-CoV-2 in mink farms: Real data

We collected the full viral genomes in minks at 63 farms from GISAID (gisaid.org) and aligned them with MUSCLE [33]. The alignment contains 326 sequences of 29,775 nucleotides long. All positions with N in all 326 sequences are removed because we do not know if there is a mutation at such a position. This left us with 326 sequences of 16,289 nucleotides long. Each farm is sampled at least once, and we have an average number of 5 samples per farm, each farm sampled on a single day. Besides the date of sampling, we also have the date of culling, which is between 1 day and 45 days after sampling, with an average of 4 days. The first 5 farms found to be infected had more than 30 days between sampling and culling, but for the rest of the farms, this was no more than 10 days. Results were obtained by running three independent $p(MC^3)$ chains with 3 heats, with 25,000 cycles after a burn-in period of 10,000 cycles. The maximum parent credibility tree is used for visualization, computation of the number of introductions, and comparison to the phylogenetic tree [24].

Supporting information

S1 Results. Tables with additional results.

(PDF)

S1 Methods. Extensive treatment of model and MCMC updating steps.

(PDF)

S1 Data. PANGO lineage classifications of the sampled sequences.

(XLSX)

S1 Fig. Comparison of MCMC and $p(MC^3)$ with and without the neighbour-joining tree initialization step. A: For low numbers of introductions (5 of the 20 hosts), there is no difference between methods in the posterior log-likelihood distribution. B: Higher numbers of introductions (15 of the 20 hosts), performance of MCMC with a random tree as initialization of the history host is inferior to either $p(MC^3)$, neighbour-joining tree initialization of the history host or the combination of both. Moreover, the simulated outbreak has a log-likelihood (the vertical black line) that is higher than the log-likelihood distribution of MCMC with a random tree as initialization. The latter gives the highest likelihood distribution and is chosen as default option in all analyses. ‘random’ is random tree initialization, ‘nj’ is neighbour-joining tree initialization, ‘2’ is MCMC and ‘3’ is $p(MC^3)$. The black lines are the log-likelihood values of the simulated outbreaks.

(TIF)

S2 Fig. Type of errors in the estimated transmission tree. The left figure represents the transmission tree of a simulated outbreak with 5 cases; there are 2 introductions (clusters) and 3 transmission events. The right figure represents possible estimates of the transmission tree of

the simulated outbreak. The vertical ordering of cases in the left and the right figures is identical. The upper right figure shows errors in which an incorrect infector is identified, but the incorrect infector belongs to the same cluster as the true infector (type A errors), the lower right figure represents incorrect identifications of the infector in which the incorrect infector belongs to a different cluster as the true infector (type B errors). In Type 1 errors neither the true infector nor the incorrect identified infector is an index case. For type 2 errors, the host is an index case in the simulated outbreak but not in the estimated outbreak. For type 3 errors, the host is not an index case in the simulated outbreak but is an index case in the estimated outbreak.

(TIF)

S3 Fig. Analysis of simulated outbreaks with fixed model parameters in the MCMC runs, while varying the number of introductions and coalescent rate in the history host.

The model parameters are fixed at the simulation values. (A) The mean estimated median number of introductions. The black line indicates the simulated number of introductions. (B) Percentage of correctly identified infectors. The grey bar indicates cases for which the true infector has the highest posterior weight. The transparent bar indicates cases for which the true infector is contained in the smallest set of candidate infectors with at least 95% of the posterior weight. (C) Classification of the incorrectly identified infectors in the maximum credibility tree. The grey bars indicate the correctly identified infectors. S: single transmission cluster involved, M: multiple transmission clusters involved. C→C: simulated and inferred infectors are cases, H→C: simulated infector was history host, inferred infector is case, C→H: simulated infector was case, inferred infector is history host.

(TIF)

S4 Fig. Analysis of simulated outbreaks with similar parameter values as the SARS-CoV-2 outbreak in mink farms.

(A) The mean estimated median number of introductions. The black line indicates the simulated number of introductions. (B) Percentage of correctly identified infectors. The grey bar indicates cases for which the true infector has the highest posterior weight. The transparent bar indicates cases for which the true infector is contained in the smallest set of candidate infectors with at least 95% of the posterior weight. (C) Classification of the falsely identified infectors based on highest support. (C) Classification of the falsely identified infectors based on highest support. The grey bars indicate the correctly identified infectors. S: single transmission cluster involved, M: multiple transmission clusters involved. For the infector of a host: C→C: case becomes case, H→C: history becomes case, C→H: case becomes history.

(TIF)

S5 Fig. Analysis of simulated outbreaks with similar parameter values as the SARS-CoV-2 outbreak in mink farms, with only 50 base pairs of the genome under mutation.

(A) The mean estimated median number of introductions. The black line indicates the simulated number of introductions. (B) Percentage of correctly identified infectors. The grey bar indicates cases for which the true infector has the highest posterior weight. The transparent bar indicates cases for which the true infector is contained in the smallest set of candidate infectors with at least 95% of the posterior weight. (C) Classification of the falsely identified infectors based on highest support. (C) Classification of the falsely identified infectors based on highest support. The grey bars indicate the correctly identified infectors. S: single transmission cluster involved, M: multiple transmission clusters involved. For the infector of a host: C→C: case becomes case, H→C: history becomes case, C→H: case becomes history.

(TIF)

S6 Fig. Maximum parent credibility transmission tree with with-host phylogenetic trees for SARS-CoV-2 outbreak in mink farms. The farms are colored according to the clusters found by Lu et al. (2021): cluster A: red; cluster B; yellow, cluster C: green; cluster D: blue, cluster E: purple, cluster unknown: black. Cluster A is divided into 5 smaller clusters, with cluster A1 introduced in NB-EMC-1 and cluster A2 introduced in NB-EMC-46.

(TIF)

S7 Fig. Maximum parent credibility phylogenetic tree for SARS-CoV-2 outbreak in mink farms. The history host is shown as the most-left red line, and the hosts are given in alternating colors. The black boxes represent the clusters in the transmission tree, with the lowest box the assumed bigger cluster with index case NB-EMC-46.

(TIF)

S8 Fig. Histogram of number of introductions for the mink farms.

(TIF)

S9 Fig. Posterior support of infectors of all hosts. There is a high certainty of the index cases (infected with the history host as infector) in the beginning of the outbreak. Transmission clusters with index cases NB-EMC-33 and NB-EMC-53 show more variation of the infectors, even outside their transmission cluster. Posterior support is shown from 0 (white) to 1 (blue). Hosts are ordered by transmission cluster and infection time. The grey bars show the transmission clusters.

(TIF)

S10 Fig. Traceplots of loglikelihood of 3 MCMC chains. Traceplots are shown of the loglikelihood of 3 MCMC chains analyzing the SARS-CoV-2 outbreak in Dutch mink farms.

(TIF)

Acknowledgments

This work was performed as part of the research program of the Netherlands Centre for One Health (www.ncoh.nl). We thank Bas Oude Munnik, Francisca Velkers and the 'One Health mink outbreak investigation consortium' for providing a prepublication of the mink data.

Author Contributions

Conceptualization: Bastiaan R. Van der Roest, Martin C. J. Bootsma, Egil A. J. Fischer, Don Klinkenberg, Mirjam E. E. Kretzschmar.

Formal analysis: Bastiaan R. Van der Roest.

Investigation: Bastiaan R. Van der Roest, Don Klinkenberg.

Methodology: Bastiaan R. Van der Roest, Martin C. J. Bootsma, Egil A. J. Fischer, Don Klinkenberg, Mirjam E. E. Kretzschmar.

Software: Bastiaan R. Van der Roest, Don Klinkenberg.

Supervision: Martin C. J. Bootsma, Egil A. J. Fischer, Don Klinkenberg, Mirjam E. E. Kretzschmar.

Visualization: Bastiaan R. Van der Roest.

Writing – original draft: Bastiaan R. Van der Roest.

Writing – review & editing: Bastiaan R. Van der Roest, Martin C. J. Bootsma, Egil A. J. Fischer, Don Klinkenberg, Mirjam E. E. Kretzschmar.

References

1. Zhao S, Tang B, Musa SS, Ma S, Zhang J, Zeng M, et al. Estimating the generation interval and inferring the latent period of COVID-19 from the contact tracing data. *Epidemics*. 2021; 36:100482. <https://doi.org/10.1016/j.epidem.2021.100482> PMID: 34175549
2. Haydon DT, Chase-Topping M, Shaw DJ, Matthews L, Friar JK, Wilesmith J, et al. The construction and analysis of epidemic trees with reference to the 2001 UK foot-and-mouth outbreak. *Proceedings Biological sciences*. 2003; 270(1511):121–7. <https://doi.org/10.1098/rspb.2002.2191> PMID: 12590749
3. Cauchemez S, Boelle PY, Donnelly CA, Ferguson NM, Thomas G, Leung GM, et al. Real-time estimates in early detection of SARS. *Emerging infectious diseases*. 2006; 12(1):110–3. <https://doi.org/10.3201/eid1201.050593> PMID: 16494726
4. Cauchemez S, Ferguson NM. Methods to infer transmission risk factors in complex outbreak data. *Journal of the Royal Society, Interface*. 2012; 9(68):456–69. <https://doi.org/10.1098/rsif.2011.0379> PMID: 21831890
5. Fraser C, Donnelly CA, Cauchemez S, Hanage WP, Van Kerkhove MD, Hollingsworth TD, et al. Pandemic Potential of a Strain of Influenza A (H1N1): Early Findings. *Science*. 2009; 324(5934):1557–1561. <https://doi.org/10.1126/science.1176062> PMID: 19433588
6. Harris SR, Feil EJ, Holden MTG, Quail MA, Nickerson EK, Chantratita N, et al. Evolution of MRSA During Hospital Transmission and Intercontinental Spread. *Science*. 2010; 327(5964):469–474. <https://doi.org/10.1126/science.1182395> PMID: 20093474
7. Mutreja A, Kim DW, Thomson NR, Connor TR, Lee JH, Kariuki S, et al. Evidence for several waves of global transmission in the seventh cholera pandemic. *Nature*. 2011; 477(7365):462–465. <https://doi.org/10.1038/nature10392> PMID: 21866102
8. Ruan Y, Wei CL, Ling AE, Vega VB, Thoreau H, Se Thoe SY, et al. Comparative full-length genome sequence analysis of 14 SARS coronavirus isolates and common mutations associated with putative origins of infection. *The Lancet*. 2003; 361(9371):1779–1785. [https://doi.org/10.1016/S0140-6736\(03\)13414-9](https://doi.org/10.1016/S0140-6736(03)13414-9) PMID: 12781537
9. Worby CJ, O'Neill PD, Kyraios T, Robotham JV, De Angelis D, Cartwright EJP, et al. Reconstructing transmission trees for communicable diseases using densely sampled genetic data. *Annals of Applied Statistics*. 2016; 10(1):395–417. <https://doi.org/10.1214/15-aos898> PMID: 27042253
10. Kenah E. Semiparametric Relative-risk Regression for Infectious Disease Transmission Data. *Journal of the American Statistical Association*. 2015; 110(509):313–325. <https://doi.org/10.1080/01621459.2014.896807> PMID: 26146425
11. Kenah E, Britton T, Halloran ME, Longini IM. Molecular Infectious Disease Epidemiology: Survival Analysis and Algorithms Linking Phylogenies to Transmission Trees. *PLoS computational biology*. 2016; 12(4):e1004869. <https://doi.org/10.1371/journal.pcbi.1004869> PMID: 27070316
12. Didelot X, Gardy J, Colijn C. Bayesian inference of infectious disease transmission from whole-genome sequence data. *Molecular Biology and Evolution*. 2014; <https://doi.org/10.1093/molbev/msu121> PMID: 24714079
13. Numminen E, Chewapreecha C, Sirén J, Turner C, Turner P, Bentley SD, et al. Two-phase importance sampling for inference about transmission trees. *Proceedings of the Royal Society B: Biological Sciences*. 2014; 281(1794). <https://doi.org/10.1098/rspb.2014.1324> PMID: 25253455
14. Ypma RJF, van Ballegooijen WM, Wallinga J. Relating phylogenetic trees to transmission trees of infectious disease outbreaks. *Genetics*. 2013; 195(3):1055–1062. <https://doi.org/10.1534/genetics.113.154856> PMID: 24037268
15. Hall M, Woolhouse M, Rambaut A. Epidemic Reconstruction in a Phylogenetics Framework: Transmission Trees as Partitions of the Node Set. *PLoS Computational Biology*. 2015; 11(12):1–36. <https://doi.org/10.1371/journal.pcbi.1004613>
16. Pham TM, Kretzschmar M, Bertrand X, Bootsma M. Tracking *Pseudomonas aeruginosa* transmissions due to environmental contamination after discharge in ICUs using mathematical models. *PLOS Computational Biology*. 2019; 15(8):e1006697. <https://doi.org/10.1371/journal.pcbi.1006697> PMID: 31461450
17. Si Y, de Boer WF, Gong P. Different environmental drivers of highly pathogenic avian influenza H5N1 outbreaks in poultry and wild birds. *PloS one*. 2013; 8(1):e53362. <https://doi.org/10.1371/journal.pone.0053362> PMID: 23308201
18. Kerfua SD, Shirima G, Kusiluka L, Ayebazibwe C, Mwebe R, Cleaveland S, et al. Spatial and temporal distribution of foot-and-mouth disease in four districts situated along the Uganda-Tanzania border:

- Implications for cross-border efforts in disease control. *The Onderstepoort journal of veterinary research*. 2018; 85(1):e1–e8. <https://doi.org/10.4102/ojvr.v85i1.1528> PMID: 30198279
19. Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N. Bayesian Reconstruction of Disease Outbreaks by Combining Epidemiologic and Genomic Data. *PLoS Computational Biology*. 2014; 10(1). <https://doi.org/10.1371/journal.pcbi.1003457> PMID: 24465202
 20. Didelot X, Fraser C, Gardy J, Colijn C, Malik H. Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks. *Molecular Biology and Evolution*. 2017; 34(4):997–1007. <https://doi.org/10.1093/molbev/msw275> PMID: 28100788
 21. Mollentze N, Nel LH, Townsend S, le Roux K, Hampson K, Haydon DT, et al. A bayesian approach for inferring the dynamics of partially observed endemic infectious diseases from space-time-genetic data. *Proceedings of the Royal Society B: Biological Sciences*. 2014; 281 (1782). <https://doi.org/10.1098/rspb.2013.3251> PMID: 24619442
 22. Morelli MJ, Thébaud G, Chadœuf J, King DP, Haydon DT, Soubeyrand S. A Bayesian Inference Framework to Reconstruct Transmission Trees Using Epidemiological and Genetic Data. *PLoS Computational Biology*. 2012; 8(11):e1002768. <https://doi.org/10.1371/journal.pcbi.1002768> PMID: 23166481
 23. Klinkenberg D, Backer JA, Didelot X, Colijn C, Wallinga J. Simultaneous inference of phylogenetic and transmission trees in infectious disease outbreaks. *PLoS Computational Biology*. 2017; 13(5). <https://doi.org/10.1371/journal.pcbi.1005495> PMID: 28545083
 24. Lu L, Sikkema RS, Velkers FC, Nieuwenhuijse DF, Fischer EAJ, Meijer PA, et al. Adaptation, spread and transmission of SARS-CoV-2 in farmed minks and associated humans in the Netherlands. *Nature Communications*. 2021; 12(1). <https://doi.org/10.1038/s41467-021-27096-9> PMID: 34815406
 25. Munnink BBO, Sikkema RS, Nieuwenhuijse DF, Molenaar RJ, Munger E, Molenkamp R, et al. Transmission of SARS-CoV-2 on mink farms between humans and mink and back to humans. *Science*. 2021; 371(6525):172–177. <https://doi.org/10.1126/science.abe5901>
 26. Altekar G, Dwarkadas S, Huelsenbeck JP, Ronquist F. Parallel Metropolis coupled Markov chain Monte Carlo for Bayesian phylogenetic inference. *Bioinformatics*. 2004; 20(3):407–415. <https://doi.org/10.1093/bioinformatics/btg427> PMID: 14960467
 27. Abbasian MH, Mahmanzar M, Rahimian K, Mahdavi B, Tokhanbigli S, Moradi B, et al. Global landscape of SARS-CoV-2 mutations and conserved regions. *Journal of Translational Medicine*. 2023; 21(1):152. <https://doi.org/10.1186/s12967-023-03996-w> PMID: 36841805
 28. O'Toole A, Scher E, Underwood A, Jackson B, Hill V, McCrone JT, et al. Assignment of Epidemiological Lineages in an Emerging Pandemic Using the Pangolin Tool. *Virus Evolution*. 2021; <https://doi.org/10.1093/ve/veab064> PMID: 34527285
 29. Amicone M, Borges V, Alves MJ, Isidro J, Zé-Zé L, Duarte S, et al. Mutation rate of SARS-CoV-2 and emergence of mutators during experimental evolution. *Evolution, Medicine, and Public Health*. 2022; 10(1):142–155. <https://doi.org/10.1093/emph/eoac010> PMID: 35419205
 30. Hammer AS, Quaade ML, Rasmussen TB, Fonager J, Rasmussen M, Mundbjerg K, et al. SARS-CoV-2 Transmission between Mink (*Neovison vison*) and Humans, Denmark. *Emerging Infectious Diseases*. 2021; 27(2):547–551. <https://doi.org/10.3201/eid2702.203794> PMID: 33207152
 31. R Core Team. R: A Language and Environment for Statistical Computing; 2022. Available from: <https://www.r-project.org/>.
 32. Felsenstein J. Evolutionary trees from DNA sequences: A maximum likelihood approach. *Journal of Molecular Evolution*. 1981; 17(6):368–376. <https://doi.org/10.1007/BF01734359> PMID: 7288891
 33. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research*. 2004; 32(5):1792–1797. <https://doi.org/10.1093/nar/gkh340> PMID: 15034147