

RESEARCH ARTICLE

Simultaneous inference of phylogenetic and transmission trees in infectious disease outbreaks

Don Klinkenberg^{1*}, Jantien A. Backer¹, Xavier Didelot², Caroline Colijn³, Jacco Wallinga^{1,4}

1 Department of Epidemiology and Surveillance, National Institute for Public Health and the Environment, Bilthoven, The Netherlands, **2** Department of Infectious Disease Epidemiology, Imperial College London, London, United Kingdom, **3** Department of Mathematics, Imperial College London, London, United Kingdom, **4** Department of Medical Statistics and Bio-Informatics, Leiden University Medical Center, Leiden, The Netherlands

* don.klinkenberg@rivm.nl



OPEN ACCESS

Citation: Klinkenberg D, Backer JA, Didelot X, Colijn C, Wallinga J (2017) Simultaneous inference of phylogenetic and transmission trees in infectious disease outbreaks. *PLoS Comput Biol* 13 (5): e1005495. <https://doi.org/10.1371/journal.pcbi.1005495>

Editor: Mark M. Tanaka, University of New South Wales, AUSTRALIA

Received: August 16, 2016

Accepted: April 3, 2017

Published: May 18, 2017

Copyright: © 2017 Klinkenberg 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: Data are in the supplementary file "S1 Data."

Funding: XD received funding from the UK National Institute for Health Research (grant PRU-2012-10080) and the UK Medical Research Council (grant MR/N010760/1). Caroline Colijn received funding from the UK Engineering and Physical Sciences Research Council (grant EP/K026003/1). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Abstract

Whole-genome sequencing of pathogens from host samples becomes more and more routine during infectious disease outbreaks. These data provide information on possible transmission events which can be used for further epidemiologic analyses, such as identification of risk factors for infectivity and transmission. However, the relationship between transmission events and sequence data is obscured by uncertainty arising from four largely unobserved processes: transmission, case observation, within-host pathogen dynamics and mutation. To properly resolve transmission events, these processes need to be taken into account. Recent years have seen much progress in theory and method development, but existing applications make simplifying assumptions that often break up the dependency between the four processes, or are tailored to specific datasets with matching model assumptions and code. To obtain a method with wider applicability, we have developed a novel approach to reconstruct transmission trees with sequence data. Our approach combines elementary models for transmission, case observation, within-host pathogen dynamics, and mutation, under the assumption that the outbreak is over and all cases have been observed. We use Bayesian inference with MCMC for which we have designed novel proposal steps to efficiently traverse the posterior distribution, taking account of all unobserved processes at once. This allows for efficient sampling of transmission trees from the posterior distribution, and robust estimation of consensus transmission trees. We implemented the proposed method in a new R package *phybreak*. The method performs well in tests of both new and published simulated data. We apply the model to five datasets on densely sampled infectious disease outbreaks, covering a wide range of epidemiological settings. Using only sampling times and sequences as data, our analyses confirmed the original results or improved on them: the more realistic infection times place more confidence in the inferred transmission trees.

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

Author summary

It is becoming easier and cheaper to obtain (whole genome) sequences of pathogen samples during outbreaks of infectious diseases. If all hosts during an outbreak are sampled, and these samples are sequenced, the small differences between the sequences (single nucleotide polymorphisms, SNPs) give information on the transmission tree, i.e. who infected whom, and when. However, correctly inferring this tree is not straightforward, because SNPs arise from unobserved processes including infection events, as well as pathogen growth and mutation within the hosts. Several methods have been developed in recent years, but often for specific applications or with limiting assumptions, so that they are not easily applied to new settings and datasets. We have developed a new model and method to infer transmission trees without putting prior limiting constraints on the order of unobserved events. The method is easily accessible in an R package implementation. We show that the method performs well on new and previously published simulated data. We illustrate applicability to a wide range of infectious diseases and settings by analysing five published datasets on densely sampled infectious disease outbreaks, confirming or improving the original results.

This is a *PLOS Computational Biology Methods* paper.”

Introduction

As sequencing technology becomes easier and cheaper, detailed outbreak investigation increasingly involves the use of pathogen sequences from host samples [1]. These sequences can be used for studies ranging from virulence or resistance related to particular genes [1, 2], to the interaction of epidemiological, immunological and evolutionary processes on the scale of populations [3, 4]. If most or all hosts in an outbreak are sampled, it is also possible to use differences in nucleotides, i.e. single-nucleotide polymorphisms (SNPs), to resolve transmission clusters, individual transmission events, or complete transmission trees. With that information it becomes possible to identify high risk contacts and superspreaders, as well as characteristics of hosts or contacts that are associated with infectiousness and transmission [5, 6]. Much progress has been made in recent years in theory and model development, but existing methods either include assumptions that do not take the full uncertainty in the evolutionary process into account [7, 8], are designed for specific datasets, with fit-for-purpose code for data analysis [9–11], or make limiting assumptions about the relation between sampling times and infectivity [12]. An easily accessible method without these restrictions and with the flexibility to cover a wide range of infections is currently lacking, and would bring analysis of outbreak sequence data within reach of a much broader community.

The interest in easily applicable methods for sequence data analysis in outbreak settings is demonstrated by the community’s widespread use of the Outbreaker package in R [8, 13–15]. However, the model in Outbreaker assumes that mutations occur at the time of transmission, which does not take the pathogen’s in-host population dynamics into account, nor the fact that mutations occur within hosts. The publications by Didelot et al [7] and Ypma et al [11] revealed that within-host evolution is crucial to relate sequence data to transmission trees, as is illustrated in Fig 1A: there are four unobserved processes, i.e. the time between subsequent infections, the time between infection and sampling, the pathogen dynamics within hosts, and

- a. infection to transmission
- b. infection to sampling
- c. coalescence
- d. mutation

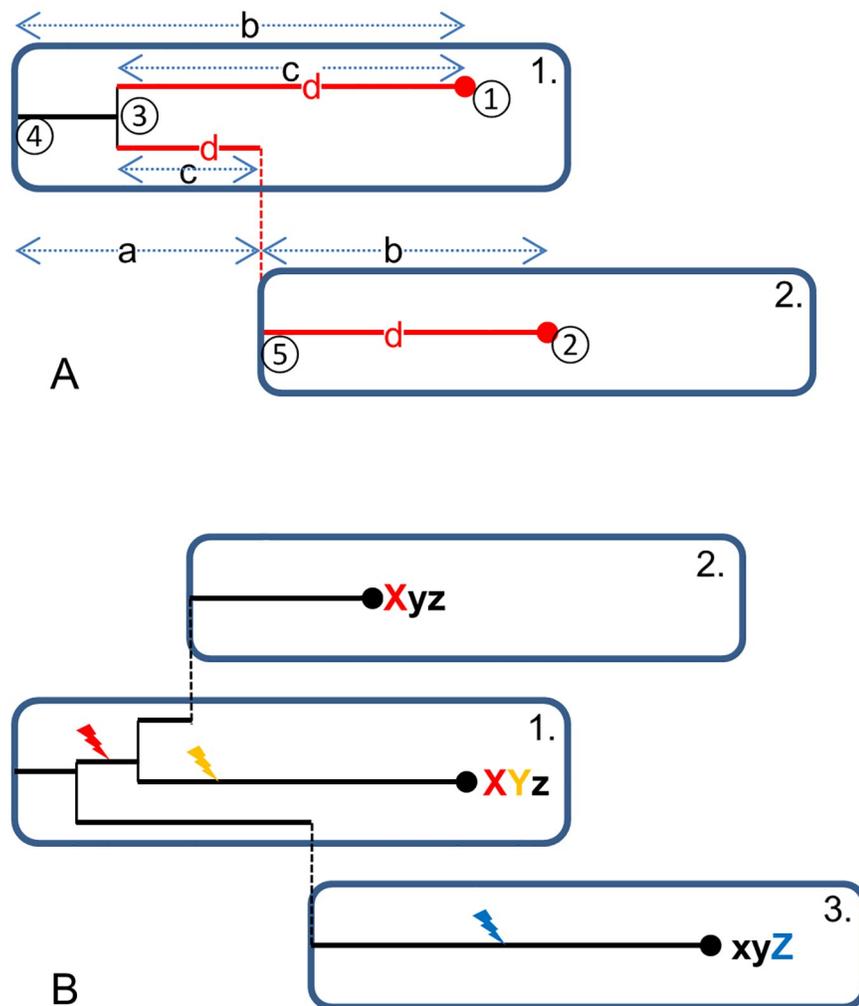


Fig 1. Sketch of stochastic processes involved in data generation process. (A) The four processes indicated by host 1 infecting host 2, together leading to a difference between the sampled sequences of hosts 1 and 2: (a) transmission, (b) sampling, (c) coalescence, and (d) mutation. As described in Methods, the node ID numbers are indicated in circles as if this was a complete outbreak. (B) Examples of differences in sequences for host 1 infecting both hosts 2 and 3. Host 1 is infected by the original sequence xyz, and the lightning indicates when mutations take place.

<https://doi.org/10.1371/journal.pcbi.1005495.g001>

mutation. The difference in sequences between host 2 and infector 1 result from all of these processes. As a result, a host's sample can have different SNPs from his infector's (Fig 1B: hosts 1 and 2); a host can even be sampled earlier than his infector with fewer SNPs (Fig 1B: hosts 1 and 3).

Several recently published methods do allow mutations to occur within the host, but make other assumptions not fully reflecting the above-described process, such as using a phenomenological model for pairwise genetic distances [16], presence of a single dominant strain in

which mutations can accumulate [9, 17], or absence of a clearly defined infection time [18]. To take the complete process into account, Didelot et al [7] and Numminen et al [10] took a two-step approach: first, phylogenetic trees were built, and second, these trees were used to infer transmission trees. Didelot et al [7] used the software BEAST [19, 20] to make a timed phylogenetic tree, and used a Bayesian MCMC method to colour the branches such that changes in colour represent transmission events. Numminen et al [10] took the most parsimonious tree topology, and accounted for unobserved hosts by a sampling model (which is an additional complication). This two-step approach is likely to work better if the phylogenetic tree is properly resolved (unique sequences with many SNPs), but less so if there is uncertainty in the phylogenetic tree. However, also in that case construction of the phylogenetic tree is done without taking into account that only lineages in the same host can coalesce, and that these go through transmission bottlenecks during the whole outbreak. That is likely to result in incorrect branch lengths and consequently incorrect infection times.

Hall et al [12] implemented a method in BEAST for simultaneous inference of transmission and phylogenetic trees. BEAST allows for much flexibility when it comes to phylogeny and population dynamics reconstruction (for which it was originally developed [19, 20]), e.g. by allowing variation in mutation rates between sites in the genome, between lineages, and in time. However, datasets of fully observed outbreaks often do not contain sufficient information for reliable inference: they typically cover only a few months up to at most several years (as in Didelot et al [7], with tuberculosis) and do not contain many SNPs (usually of the same order of magnitude as the number of samples). A more important limitation is that the transmission model implemented in BEAST is rather specific: it allows for transmission only during an infectious period informed by positive and negative samples, during which infectiousness is assumed to be constant. This may put prior constraints on the topology and order of events in the transmission and phylogenetic trees, which is undesirable if the primary aim is to reconstruct the transmission tree with little or no prior information about when hosts were infectious.

Previously, Ypma et al [11] had also developed a method for simultaneous inference of transmission and phylogenetic trees, albeit with rather specific assumptions on the within-host pathogen dynamics and the time and order of transmission events, and with no available implementation. However, their view on the phylogenetic and transmission trees was quite different. Instead of a phylogenetic tree with transmission events, they regarded it as a hierarchical tree. The top level is the transmission tree, with hosts having infected other hosts according to an epidemiological transmission model. The lower level consists of phylogenetic “mini-trees” within each host. A mini-tree describes the within-host micro-evolution. It is rooted at the infection time and has tips at transmission and sampling events; in its simplest form it is only a single branch from infection to sampling. The complete phylogenetic tree then consists of all these mini-trees, connected through the transmission tree. That description allowed them to develop new MCMC updating steps, some changing the transmission tree, some the phylogenetic mini-trees.

We built further on that principle to reconstruct the transmission trees of outbreaks, in a new model and estimation method. The method requires data on pathogen sequences and sampling times. The model includes all four underlying stochastic processes (Fig 1A), each described in a flexible and generic way, such that we avoid putting unnecessary prior constraints on the order of unobserved events (Fig 1B). This allows for application of the method to a wide range of infectious diseases, including new emerging infections where we have little quantitative information on the infection cycle. The method is implemented in R, in a package called *phybreak*. We illustrate the performance of the method by applying it to both new and previously published simulated datasets. We demonstrate the range of applicability by applying

the model to five datasets on densely sampled infectious disease outbreaks, covering a wide range of epidemiological settings.

Results

Outline of the method

The method infers infection times and infectors of all cases in an outbreak. The data consist of sampling times and sequences of all cases, where some of the sequences may be non-informative if no sequence is available. Using simple models for transmission, sampling, within-host dynamics and mutation, samples are taken from the posterior distributions of model parameters and transmission and phylogenetic trees, by a Markov-Chain Monte Carlo (MCMC) method. The main novelty of our method lies in the proposal steps for the phylogenetic and transmission trees that are used to generate the MCMC chain. It makes use of the hierarchical tree perspective, in which the phylogenetic tree is described as a collection of phylogenetic mini-trees (one for each host), connected through the transmission tree (see [Methods](#) for details).

The posterior samples are summarized by medians and credible intervals for parameters and infection times, and by consensus transmission trees. Consensus transmission trees are based on the posterior support for infectors of each host, defined as the proportion of posterior trees in which a particular infector infects a host. The Edmonds' consensus tree takes for each host the infector with highest support, and uses Edmonds' algorithm to resolve cycle and multiple index cases [21], whereas the Maximum Parent Credibility (MPC) tree is the one tree among the posterior trees with maximum product of supports [12].

The models and parameters used for inference are as follows:

- transmission: assuming that all cases are sampled and the outbreak is over, the mean number of secondary infections must be 1. The transmission model therefore consists only of a Gamma distribution for the generation interval, i.e. the time interval between a primary and a secondary case. This transmission model is equivalent to a Kermack-McKendrick type renewal model with a generalized infectiousness function, with basic reproduction ratio $R_0 = 1$ in an infinite population [22]. The model contains two parameters: the shape a_G , which we fixed at 3 during our analyses, and the mean m_G , which is estimated and has a prior distribution with mean μ_G and standard deviation σ_G . In an uninformative analysis, $\mu_G = 1$, and $\sigma_G = \infty$.
- sampling: the sampling model consists of a Gamma distribution for the sampling interval, which is the interval between infection and sampling of a case. The model contains two parameters: the shape a_S , which is fixed during the analysis, and the mean m_S , which is estimated and has a prior distribution with mean μ_S and standard deviation σ_S . In an uninformative analysis, $\mu_S = 1$, and $\sigma_S = \infty$, but a_S is chosen to reflect prior information (the coefficient of variation is $a_S^{-1/2}$); in a naïve analysis we additionally set $a_S = 3$.
- within-host dynamics: the within-host model determines the genetic relation between the tips of the within-host phylogenetic mini-tree (at sampling and transmission) through a coalescence model, assuming that samples are clonal lineages. The within-host model describes a linearly increasing pathogen population size $w(\tau) = r\tau$, at time τ since infection of a host. This within-host model results in a bottleneck at transmission of 1 lineage. The slope r has a Gamma distributed prior distribution with shape a_r and rate b_r . In an uninformative analysis, we used $a_r = b_r = 3$.
- The mutation model is a site-homogeneous Jukes-Cantor model, with per-site mutation rate μ . The prior distribution for $\log(\mu)$ is uniform.

Analysis of newly simulated datasets

We generated 25 new simulated datasets of 50 cases with the above model, which we modified by taking a population of 86 individuals and a basic reproduction number $R_0 = 1.5$ (instead of an infinite population with $R_0 = 1$). Parameters were $a_G = a_S = 10$, $m_G = m_S = r = 1$, $\mu = 10^{-4}$ and sequence length 10^4 , resulting in 1 genome-wide mutation per mean generation interval of one year.

Table 1 shows some summary measures on performance of the method (see S1 Results for additional measures and results for more simulations). A 5,000 cycle burn-in followed by sampling a single chain of 25,000 MCMC cycles took about 30 minutes on a 2.6 GHz CPU (Linux). Four sets of results are shown, all with an uninformative prior for μ : one with all parameters other than μ fixed at their correct value, and three with uninformative priors for m_G and r , and different levels of prior knowledge on m_S : informative with correct mean,

Table 1. Performance of the method: analysis of 25 newly simulated datasets of 50 cases, with shape parameters $a_S = a_G = 10$.

	Reference ^a	Level of prior information on m_S		
		Informative Correct ^b	Uninformative ^c	Informative Wrong ^d
MCMC sampling				
Continuous parameter samples (95% interval of ESS)				
M	7402; 9591	1047; 2343	218; 1664	1037; 1924
m_G		1330; 3399	873; 2873	1003; 2615
m_S		314; 779	38; 124	401; 888
r		251; 390	195; 383	197; 335
t_{inf}	1526; 8484	895; 4978	257; 969	531; 4013
phylogenetic tree topology ^e	1307; 3108	1313; 3035	1311; 3111	1315; 2787
Infectors (% Fisher's exact tests accepted)				
between chains	98.6%	97.8%	98.3%	98.6%
autocorrelation	95.6%	97.1%	95.1%	95.6%
Parameter inference (95% percentile interval of posterior medians)				
$\log_{10}(\mu)$	-4.11; -3.90	-4.10; -3.90	-4.16; -3.88	-4.27; -4.08
m_G		0.74; 1.04	0.72; 1.08	0.68; 1.01
m_S		0.95; 1.07	0.64; 1.54	1.90; 1.98
r		0.49; 1.11	0.60; 1.18	0.33; 0.74
Tree inference				
Infection times (coverage: % of 95% CIs containing the true value)				
	95.8%	96.3%	95.4%	46.2%
Infection time bias (median)	0.01 yr	0.01 yr	0.09 yr	-0.84 yr
95% interval of medians	-0.45; 0.62	-0.46; 0.61	-0.78; 0.82	-1.52; -0.15
Infectors (number correct/number identified)				
Edmonds'	34.9/50	34.4/50	34.0/50	34.9/50
MPC	33.3/50	33.1/50	33.0/50	30.6/50
≥50% support	27.8/33.2	28.0/33.8	28.2/34.0	21.6/24.1
≥80% support	15.2/15.8	15.4/15.9	15.4/16.1	7.8/7.9

Results are based on two MCMC chains of 25,000 samples each; ESS, effective sample size; CI, credible interval; MPC, maximum parent credibility.

^a $m_G, m_S, r = 1$

^b $\mu_S = 1, \sigma_S = 0.1$

^c $\mu_S = 1, \sigma_S = \infty$

^d $\mu_S = 2, \sigma_S = 0.1$

^e the path-distance approximate ESS [23]

<https://doi.org/10.1371/journal.pcbi.1005495.t001>

uninformative, and informative with incorrect mean. The top of the table shows effective sample sizes (ESSs) for all parameters and for the infection times to evaluate mixing of continuous parameter samples. The path-distance approximate topological ESS [23] was calculated to assess phylogenetic tree mixing. To evaluate mixing across and within chains of infectors per host, we tested for differences between the chains and for dependency within the chains by Fisher's exact tests: the proportion of accepted tests ($P > 0.05$) is a measure of mixing. The MCMC mixing is generally good for tree inference and model parameters, as most ESSs are above 200 and an expected 95% of Fisher's tests is accepted; the only exceptions being m_S with an uninformative prior.

The bottom part of Table 1 shows the results on tree inference. Infection times (using all MCMC samples) are well recovered if the mean sampling interval does not have a strong incorrect prior: coverage of 95% credible intervals is good, and medians may only be slightly positively biased (later than true infection time) if uninformative priors are used. For this simulation scenario, consensus transmission trees contained almost 70% (35 out of 50) correct infectors, as determined by counting infectors and resolving multiple index cases and cycles in the tree (Edmonds' method [21]) and slightly fewer when choosing the maximum parent credibility (MPC) tree [12] among the 50,000 posterior trees. Infectors with high support are more likely correct: 84% (28 out of 33) are correct if the support is above 50%, and 96% (15.2 out of 15.8) are correct if the support is above 80%. These numbers are similar in smaller outbreaks (S1 Results). If sampling and generation interval distributions are wider, the sampling times contain less information on the order of infection, which reduces the accuracy of transmission tree reconstruction (S1 Results). Using prior information on the mean sampling interval did not improve on this, but if prior information is incorrect, fewer hosts have a strongly supported infector, which makes inference more uncertain. In conclusion, the method is fast and efficient if applied to simulated data fitting the model. In that case, no informative priors are needed for transmission tree inference, though correct estimation of the infection time is aided by some information.

For comparison, we analysed the same datasets with the *Outbreaker* package in R [8], which uses the assumption of mutation at transmission, and with the *TransPhylo* package [7, 24], which requires input of a phylogenetic tree that we created in BEAST v2 [19] with a constant population coalescent model and Jukes-Cantor substitution model. Both *Outbreaker* and *TransPhylo* require input of a generation and sampling interval distribution, for which we supplied the distributions used to simulate the data. Thus, the results are best compared to the results of the reference scenario of our model (Table 1, left-most column). The numbers of correctly identified infectors (Edmonds' consensus tree [21]) were smaller with both alternative methods: in the 25 outbreaks of Table 1 (50 cases, $a_G = a_S = 10$), *Outbreaker* identifies on average 27.5 out of 50 infectors correctly, *TransPhylo* 32.2, and *phybreak* 34.9. Also in smaller outbreaks or with different generation and sampling interval distributions, *phybreak* identified 8–22% more infectors correctly (S1 Results).

We also analysed the simulated results with 20% of the cases removed from the dataset, to assess performance if outbreaks are not completely observed. Table 2 shows the results with reference (parameters fixed and correct) and uninformative analyses, in comparison with the reference scenario and all data observed. With some of the cases removed, some of the remaining cases did not have their infector in the dataset anymore; these cases are referred to as orphans in Table 2. Infection time estimation was less accurate, with only 85% of credible interval containing the correct value, and more infection times estimated too early in the outbreak. Surprisingly, this was not only the case with orphans, for which this may have been expected with their infector not present in the data. It turns out that infectors are correctly identified about 20% less accurately, for all threshold levels of support. However, when

Table 2. Tree inference with incomplete data, with 25 newly simulated datasets of 50 cases, of which 40 observed, simulated with shape parameters $a_S = a_G = 10$.

	Complete data	Incomplete data	
	Reference ^a	Reference ^a	Uninformative ^b
Tree inference			
Infection times (coverage: % of 95% CIs containing the true value)			
	95.8%	84.9%	85.4%
Infection time bias (median)			
> all cases	0.01 yr (-0.45; 0.62)	-0.04 yr (-1.47; 0.56)	-0.04 yr (-1.62; 0.70)
> infectors in data ^c		-0.03 yr (-1.46; 0.56)	-0.02 yr (-1.61; 0.68)
> orphans ^d		-0.10 yr (-1.50; 0.56)	-0.09 yr (-1.62; 0.68)
Infectors (number correct/number identified)			
Edmonds'			
> all cases	34.9/50 (70%)	23.5/40 (59%)	23.0/40 (58%)
> infectors in data ^c		23.5/32.3 (73%)	23.0/32.3 (71%)
> orphans (ancestors) ^d		3.4/7.7 (45%)	3.6/7.7 (46%)
>50% support			
> all cases	27.8/33.2 (84%)	19.2/28.1 (68%)	18.6/27.7 (67%)
> infectors in data ^c		19.2/22.6 (85%)	18.6/22.3 (84%)
> orphans (ancestors) ^d		3.2/5.5 (58%)	3.2/5.4 (59%)
>80% support			
> all cases	15.2/15.8 (97%)	11.2/14.2 (79%)	10.6/13.1 (80%)
> infectors in data ^c		11.2/11.8 (95%)	10.6/11.0 (96%)
> orphans (ancestors) ^d		1.9/2.4 (81%)	1.8/2.1 (88%)

Results are based on two MCMC chains of 25,000 samples each; CI, credible interval.

^a $m_G, m_S, r = 1$

^b $\mu_S = 1, \sigma_S = \infty$

^c only cases whose infector was in the dataset

^d only orphans (infector not in the dataset), counting identified ancestors of the true infector.

<https://doi.org/10.1371/journal.pcbi.1005495.t002>

correcting for presence of the infector in the data, infectors are identified with the same accuracy as in the complete dataset. We also checked how frequently the identified infector of orphans was in fact an earlier ancestor in the transmission tree, i.e. the infector's infector in most cases. It turned out that ancestors were often identified as infector, but not as accurately as the true infector identification in complete datasets (Table 2).

Analysis of previously published simulated data

We applied the method to previously published outbreak simulations [12]. Briefly, a spatial susceptible-exposed-infectious-recovered (SEIR) model was simulated in a population of 50 farms, with a latent period (exposed) of two days and a random infectious period with mean 10 days and standard deviation 1 day, at the end of which the farm was sampled. Two mutation rates were used with an HKY substitution model, either *Slow Clock* or *Fast Clock*, equivalent to 1 or 50 genome-wide mutations per generation interval of one week, respectively.

Table 3 shows performance of the method with naïve and informative prior information on the sampling interval distribution (see S1 Results for uninformative). Effective sample sizes of parameters and phylogenetic trees are a bit smaller than with analysis of the new simulations, but in most cases still good for infection times, whereas sampling of infectors was excellent. The low variance of the sampling interval distribution caused some problems in efficient

Table 3. Performance on 25 published simulated datasets in populations of size 50 [12].

Prior information	Slow Clock simulations		Fast Clock simulations	
	Naïve ^a	Informative ^b	Naïve ^a	Informative ^b
MCMC sampling				
Continuous parameter samples (95% interval of ESS)				
μ	168; 578	401; 951	62; 570	90; 916
m_G	259; 936	1948; 4531	798; 1652	2999; 12748
m_S	42; 145	43; 87	199; 1737	170; 461
r	189; 357	220; 369	43; 143	184; 613
t_{inf}	166; 1532	222; 559	216; 2264	295; 2796
phylogenetic tree topology ^c	580; 1764	1501; 4331	104; 972	83; 796
Infectors (% Fisher's exact tests accepted)				
between chains	95.8%	97.4%	92.5%	97.6%
autocorrelation	94.6%	96.4%	87.7%	96.0%
Parameter inference (95% interval of posterior medians)				
$\log_{10}(\mu)$	-4.85; -4.67	-4.95; -4.78	-3.26; -3.15	-3.20; -3.16
m_G	2.1; 5.0	3.7; 5.7	4.2; 6.1	4.7; 6.1
m_S	6.5; 10.2	11.2; 12.6	9.7; 13.9	11.4; 12.6
r	0.75; 1.3	0.36; 0.77	0.82; 2.1	0.30; 1.2
Tree inference				
Infection times (coverage: % of 95% CIs containing the true value)				
	76.1%	97.8%	94.6%	94.2%
Infection time bias (median)	3.97 days	0.10 days	0.25 days	0.01 days
95% interval of medians	(-3.20; 8.80)	(-1.98; 2.28)	(-6.49; 5.86)	(-2.00; 1.97)
Infectors (number correct/number identified)				
Edmonds'	29.0/49.3	30.8/49.3	32.8/49.3	45.5/49.3
MPC	25.4/49.3	29.9/49.3	30.8/49.3	45.5/49.3
≥50% support	12.9/14.4	25.2/30.9	23.4/28.9	45.4/48.7
≥80% support	3.0/3.1	18.5/19.7	4.4/5.0	40.9/42.1

Results are based on two MCMC chains of 25,000 samples each. The mean outbreak size was 49.3 cases; ESS, effective sample size; CI, credible interval; MPC, maximum parent credibility.

^a $a_S = 3, \mu_S = 1, \sigma_S = \infty$

^b $a_S = 144, \mu_S = 12, \sigma_S = 1$

^c path-distance approximate ESS [23]

<https://doi.org/10.1371/journal.pcbi.1005495.t003>

sampling of m_S because of its high correlation with the associated infection times, but it also caused problems in the burn-in phase if inference starts with parameter values far from their actual values (not shown). This was especially the case in the uninformative *Slow Clock* analyses, resulting in unreliable estimates of the mean sampling interval and infection times (S1 Results). With the *Fast Clock* analyses there were no such problems, as long as the full set of proposal paths in the MCMC chain was used (see Methods for details). Posterior median mutation rates are slightly higher than used for simulation, which could be due to different rates for transition and transversion in the simulation model [12].

Consensus trees with uninformative and informative settings were almost as good as in the original publication [12], in which spatial data were used and in which it was known that there was a latent period and that hosts could not transmit after sampling. In the *Slow Clock* simulations about 62% of infectors were correct, and in the *Fast Clock* simulations about 92%. Infection times were also well recovered in most cases, but not in the uninformative *Slow Clock*

analysis (S1 Results). In the naïve analyses, the *Slow Clock* consensus trees were only slightly less good (but not the infection times), whereas the *Fast Clock* consensus trees became much worse, with only 65% of infectors correct. In conclusion, the method performs well if applied to data simulated with a very different model. Good prior information on the sampling interval can significantly improve both MCMC mixing and transmission tree inference, especially if the genetic data contain many SNPs.

Analysis of published datasets

We finally applied the method to five published datasets on outbreaks of *Mycobacterium tuberculosis* (Mtb, [7]), Methicillin-resistant *Staphylococcus aureus* (MRSA, [25]), Foot-and-mouth disease (FMD2001 and FMD2007, [9, 11, 26, 27]), and H7N7 avian influenza (H7N7, [12, 28–30]).

The results for the four smaller datasets are shown in Table 4, which shows that mixing of the MCMC chains was generally good. Fig 2 shows the Edmond’s consensus trees (full details in S1 Results), with each host’s estimated infection time and most likely infector, colour coded to indicate posterior support. Fig 3 shows one sampled tree for each dataset (from the posterior set of 50,000), matching the MPC consensus tree topology.

The Mtb data were analysed with naïve prior information, which resulted in a median sampling interval of 419 days (similar to estimated incubation times [31]), a median generation interval of 107 days, and a mutation rate equivalent to 0.3–1.3 mutations per genome per year, as estimated before [32, 33]. The Edmonds’ consensus transmission tree (Fig 2A) shows low support for most infectors, which is a reflection of the low number of SNPs, but also of the long sampling interval relative to the generation interval, which makes the sampling time less

Table 4. Summary statistics for four published datasets.

	Mtb	MRSA	FMD2001	FMD2007
Prior information	Naïve^a	Informative^b	Naïve^a	Naïve^a
MCMC sampling				
Continuous parameter samples (ESS)				
μ	273	3558	393	1007
m_G	145	551	1614	1536
m_S	38	560	362	292
r	591	428	866	1144
t_{inf} (range of ESSs)	143; 1134	76; 2255	289; 3401	686; 3405
phylogenetic tree ^c	1126	600	1002	3071
Infectors (% Fisher’s exact tests accepted)				
between chains	31/33	34/36	14/15	10/11
autocorrelation	31/33	36/36	14/15	11/11
Parameter inference (95% interval of posterior medians)				
$\log_{10}(\mu)$	-9.4 (-9.7; -9.1)	-8.1 (-8.3; -8.0)	-4.4 (-4.5; -4.3)	-4.6 (-4.8; -4.4)
m_G	107 (49; 177)	22 (16; 33)	14 (10; 21)	7.5 (4.8; 13)
m_S	419 (185; 619)	31 (21; 44)	14 (8; 22)	9.9 (5.5; 18)
r	0.88 (0.23; 2.3)	0.85 (0.21; 2.2)	1.1 (0.40; 2.5)	0.91 (0.27; 2.4)

Results are based on two MCMC chains of 25,000 samples each; ESS, effective sample size;

^a $a_S = 3, \mu_S = 1, \sigma_S = \infty$

^b $a_S = 1, \mu_S = 15, \sigma_S = 5$

^c path-distance approximate ESS [23]

<https://doi.org/10.1371/journal.pcbi.1005495.t004>

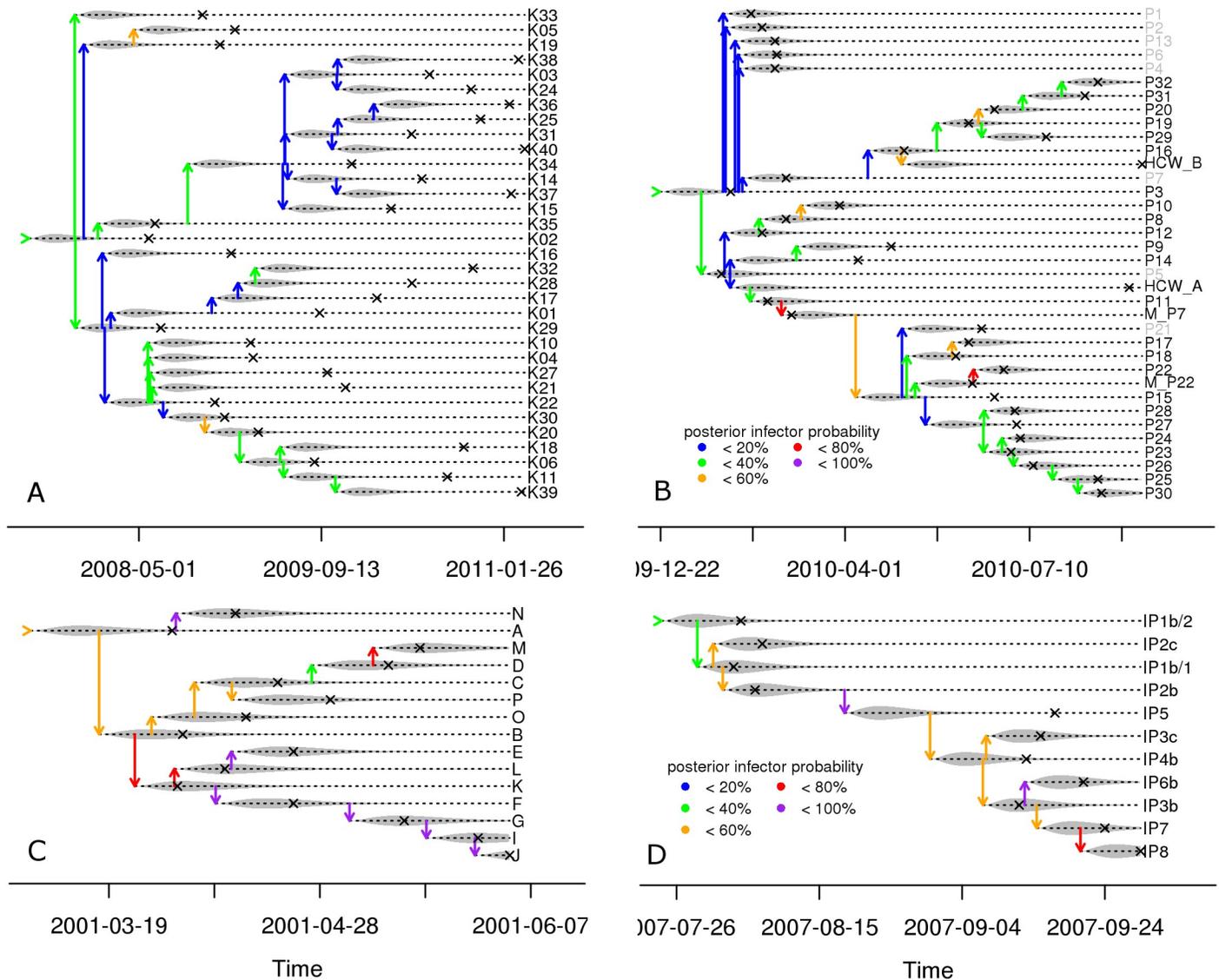


Fig 2. Consensus Edmonds' transmission trees for four of the five analysed datasets. Crosses indicate sampling days, coloured links indicate most likely infectors, with colours indicating the posterior support for that infector. (A) Mtb data [7]; (B) MRSA data [25]; (C) FMD2001 data [26]; (D) FMD2007 data [27].

<https://doi.org/10.1371/journal.pcbi.1005495.g002>

informative of the order of infection. However, the same index case K02 and three clusters as identified in Didelot et al [7] are distinguished: one starting with K22, one with K35, and the remaining cases starting with K16 or infected by the index case. The main difference compared to the original analysis lies in the shape of the phylogenetic tree and the estimated infection times (Fig 3A). Whereas the topology is very similar, the timing of the branching events is different: in the original tree, internal branches decrease in length when going from root to tips. That shape is consistent with a coalescent tree based on a single panmictic population but also reflects the fact that three mutations separate the two clades after the root node, whereas the posterior median genome-wide mutation rate is estimated at 0.48 per year (mutation rate \times sequence length). By taking into account the fact that coalescent events take place within individual hosts, our analysis shows branch lengths that are more balanced in length across the

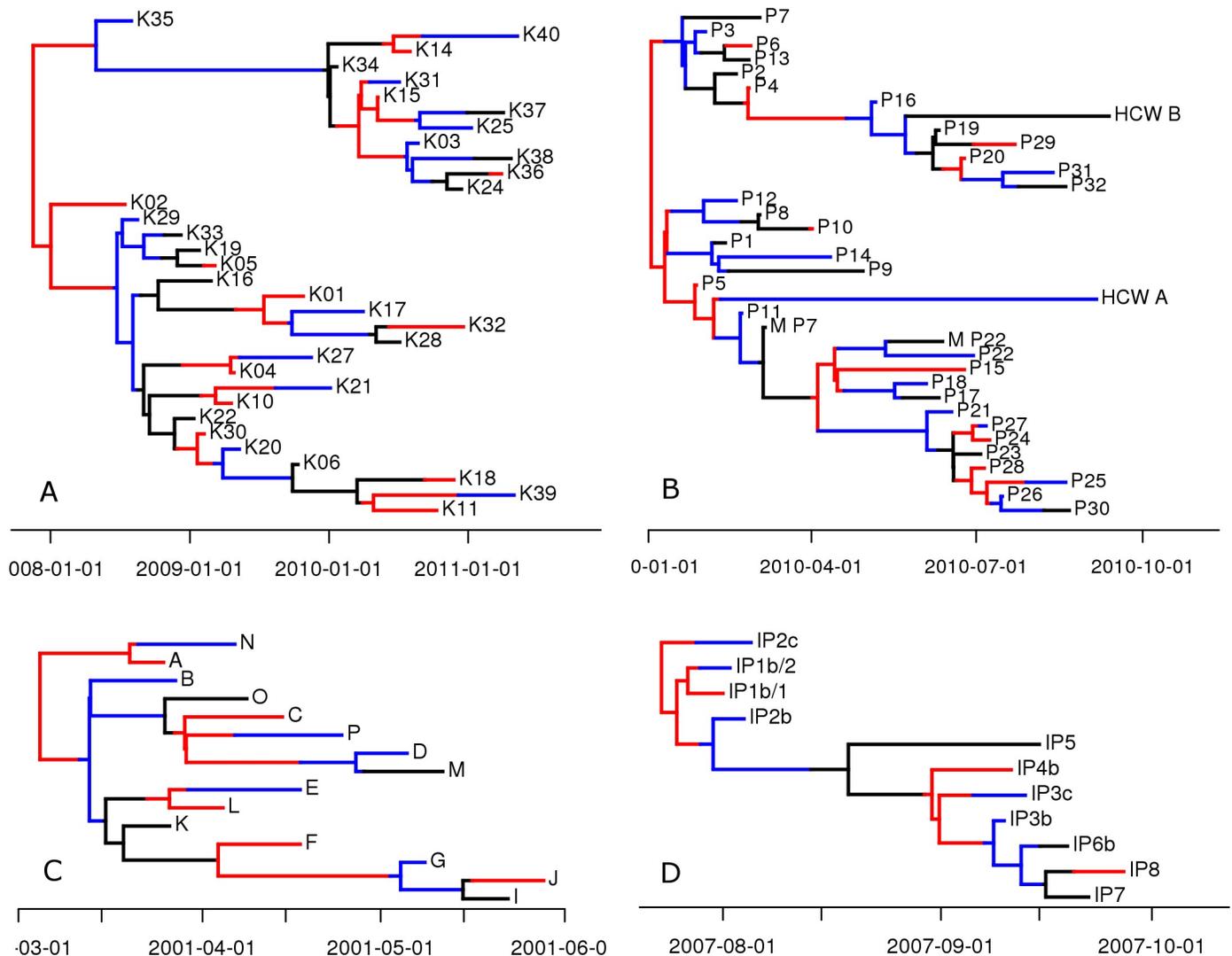


Fig 3. Consensus MPC transmission and phylogenetic trees for four of the five analysed datasets. Each tree is one posterior sample matching the MPC tree topology. Colours are used to indicate the hosts in the transmission tree: connected branches with identical colour are in the same host, and a change of colour along a branch is a transmission event. (A) Mtb data [7]; (B) MRSA data [25]; (C) FMD2001 data [9, 26]; (D) FMD2007 data [9, 27].

<https://doi.org/10.1371/journal.pcbi.1005495.g003>

tree. Importantly, this results in a more recent dating of root of the tree: early 2008 (Fig 3A) instead of early 2004 [7].

The MRSA data were analysed with an informative prior for the mean sampling interval m_s and a shape parameter a_s based on data on the intervals between hospitalisation and the first positive sample. The estimated mutation rate is similar to literature estimates [34, 35], but the posterior median m_s of 31 days is considerably higher than the prior mean of 15 days (Table 4). This may be explained by the two health-care workers (HCW_A and HCW_B), which have very long posterior sampling intervals that were not part of the data informing the prior (Edmonds' consensus tree, Fig 2B). In contrast with the original analysis, we now identify a transmission tree rather than only a phylogenetic tree, resulting in the observation that the two health-care workers may not have infected any patient in spite of their long infection-to-sampling interval. Almost all transmission events with low support (<20%) involved unsequenced

Table 5. Summary statistics for H7N7 dataset.

	Sequenced gene			All genes
	HA	NA	PB2	
MCMC sampling				
Continuous samples ^a				
M	2023	1844	1418	1363
m_G	1216	853	1262	1239
m_S	190	192	171	171
r	166	120	133	248
t_{inf} (range)	894; 8122	829; 6870	1004; 6485	1392; 10605
phylogenetic tree ^d	2898	2077	1960	1334
Infectors^b				
between chains	92.5%	95.4%	95.4%	93.3%
autocorrelation	95.9%	97.1%	95.4%	95.4%
Parameter inference^c				
$\log_{10}(\mu)$	-4.44 (-4.55; -4.35)	-4.44 (-4.56; -4.33)	-4.57 (-4.67; -4.46)	-4.51 (-4.58; -4.45)
m_G	4.0 (3.5; 4.6)	3.6 (3.1; 4.3)	4.0 (3.5; 4.6)	4.6 (4.0; 5.4)
m_S	7.8 (6.8; 8.9)	7.8 (6.9; 8.9)	7.8 (6.9; 9.0)	8.5 (7.4; 9.8)
r	0.58 (0.16; 1.5)	0.45 (0.060; 1.6)	0.33 (0.060; 1.1)	0.67 (0.20; 1.6)
Phylogenetic inference				
Tree parsimony scores ^c	102 (101; 103)	83 (82; 85)	100 (99; 101)	313 (312; 315)
#SNPs in data	90	73	94	257

Results are based on five MCMC chains of 25,000 samples each, with $a_S = 10$; $\mu_S = 7$; $\sigma_S = 0.5$; $a_G = 3$; $\mu_G = 5$; $\sigma_G = 1$; $a_r = b_r = 1$. SNP = single nucleotide polymorphism.

^a effective sample sizes

^b fraction of Fisher's exact tests with $P > 0.05$

^c medians and 95% credible intervals

^d path-distance approximate ESS [23]

<https://doi.org/10.1371/journal.pcbi.1005495.t005>

hosts. Two of them were identified as possible infector (P5 and P7), in the initial stage of the outbreak, when only few samples were sequenced. This indicates that some unsequenced hosts may have played a role in transmission, but that it is not clear which. Finally, a major difference between our results and those in the original publication is the shape of the phylogenetic tree and the dating of the tree root: around 1st January (Fig 3B) instead of 1st September the year before [25].

Analysis of the FMD2001 and FMD2007 datasets resulted in posterior sampling intervals with means of 14 and 10 days, respectively, close to the 8.5 days estimated from epidemic data [36]. Generation intervals were about the same (Table 4). Both datasets contained more SNPs than the Mtb and MRSA data, with unique sequences for each host and higher mutation rates, similar to published rates in FMD outbreak clusters [37]. This resulted in equal Edmonds' and MPC consensus transmission trees, and much higher support for most infectors (Figs 2C, 2D, 3C and 3D). Our transmission tree is almost identical to the one from Ypma et al [11], who used a closely related method but did not allow for transmission after sampling. When comparing to the analysis of these data by Morelli et al [9], the topologies of the phylogenetic trees (Fig 3C and 3D) match the topologies of the genetic networks (Fig S18 in [9]), but the transmission trees are quite different. The main differences are that in the FMD2001 outbreak, they identify farm B as the infector of C, E, K, L, O, and P; and in the FMD2007 outbreak, they have IP4b infecting IP3b, IP3c, IP6b, IP7, and IP8. Differences are likely the result of their use of the

spatial data [9]. Use of additional data is expected to improve inference, although their estimates of infection-to-sampling intervals (about 30 days) were unrealistically long.

The H7N7 dataset was analysed with the sequences of the three genes HA, NA, and PB2 separately, and combined; with informative priors for both the mean sampling and mean

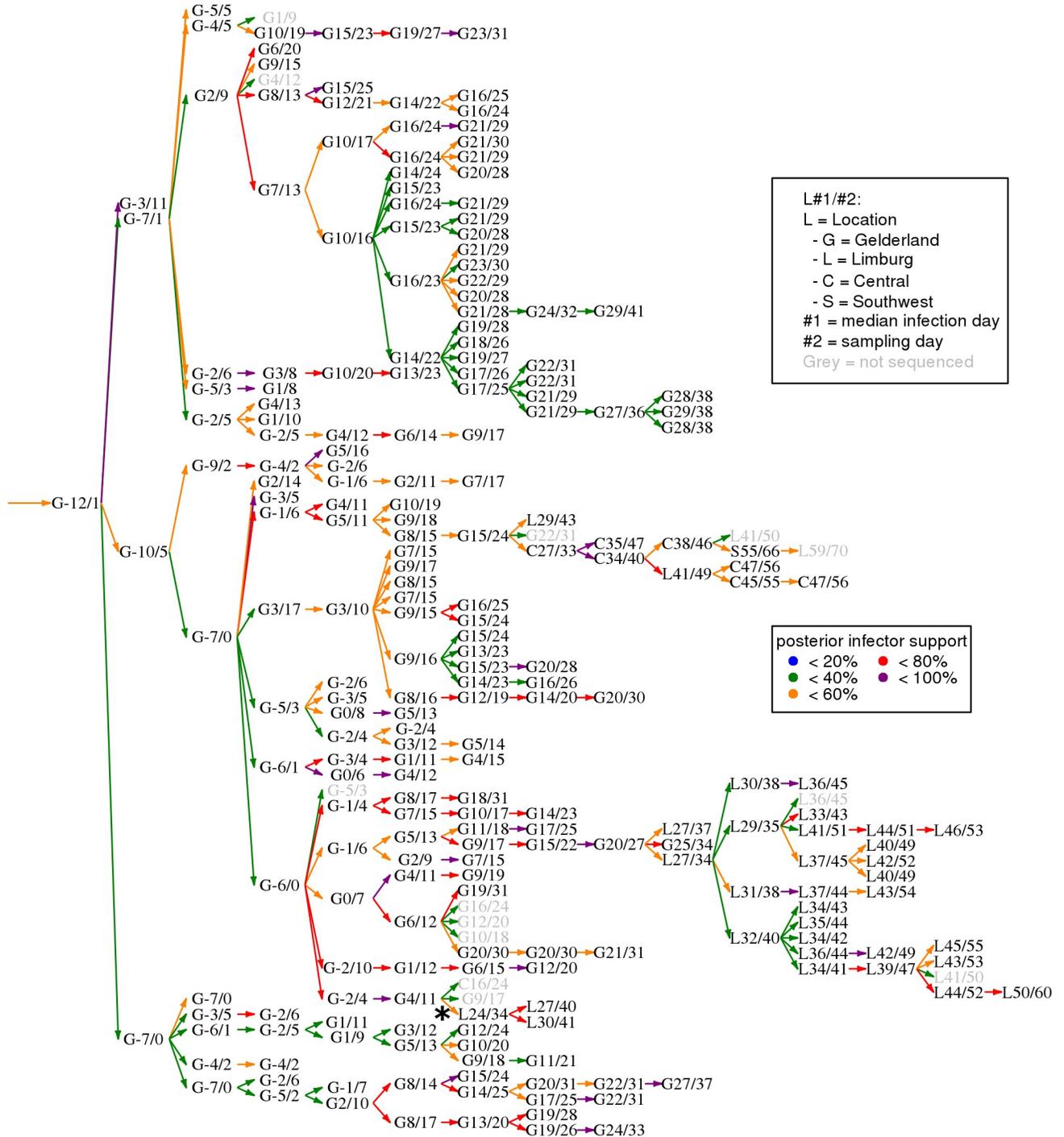


Fig 4. Consensus Edmonds' transmission tree for the H7N7 dataset [12, 28, 30]. Infected premises are (not uniquely) coded by location (as in [12]), median posterior infection day, and sampling day. Coloured arrows indicate most likely infectors, with colours indicating the posterior support for that infector. The asterisk (*) indicates the exceptionally long generation time at the start of a small Limburg cluster.

<https://doi.org/10.1371/journal.pcbi.1005495.g004>

generation intervals. Five parallel chains were run, and mixing was generally good (Table 5); it took about 7 hours on a 2.6GHz CPU to obtain 25,000 unthinned samples in a single chain. Analysis of the three genes combined resulted in a posterior median m_s of 8.4 days, slightly longer than the 7 days on which the informative prior was based [38], and longer than in the separate analyses. The mean generation time was slightly shorter than the prior mean: 4.6 days with all genes. We also calculated the parsimony scores of the posterior sampled trees, defined as the minimum numbers of mutations on the trees that can explain the sequence data [39], and compared these with the numbers of SNPs in the data (Table 5). It appeared that with the genes separately analysed, parsimony scores were 6–13% higher than the numbers of SNPs, indicating some homoplasy in the phylogenetic trees (this was not seen with any of the other datasets). The analysis of all genes together resulted in parsimony scores of 18% higher than the number of SNPs. The estimated mutation rates are among the highest estimates for Avian Influenza Virus, as already noted before in earlier analyses of the same data [28, 40]. Fig 4 shows the Edmonds' consensus tree in generations of infected premises, indicating locations and inferred infection days (full details in S1 Results). Without the use of location data, there is a large Limburg cluster, a Central cluster including two sampled Limburg cases, and a small Limburg cluster of three cases with an exceptionally long generation time (asterisk in Fig 4). A closer look at the sequences makes clear that the first of these cases (L24/34) has 3 SNPs different from assigned infector G4/11, and 4 SNPs different from cases in the large Limburg cluster. Using geographic data as in earlier analyses [12, 30] will probably place these cases within that cluster.

Discussion

We developed a new method to reconstruct outbreaks of infectious diseases with pathogen sequence data from all cases in an outbreak. Our aim was to have an easily accessible and widely applicable method. For ease of access, we developed efficient MCMC updating steps which we implemented in a new R package, *phybreak*. We tested the method on newly simulated data, previously published simulated data, and published datasets. Our model is fast: 25,000 iterations took roughly 30 minutes with the Mtb and MRSA datasets of about 30 hosts, and 7 hours with the full three-genes H7N7 dataset in 241 hosts. Two chains with 50,000 posterior samples proved sufficient (measured by ESS) for tree inference (infectors and infection times) and most model parameters with most simulated and published data. The package contains functions to enter the data, to run the MCMC chain, and to analyse the output, e.g. by making consensus trees and plotting these (as in Figs 2 and 3).

Analysis of simulated datasets showed that the sampling times play an important role in transmission tree reconstruction. Firstly, the use of prior information on the sampling interval distribution (shape parameter as well as mean) greatly improves mixing of the MCMC chain (Tables 1 and 3). Secondly, the use of (correct) prior information on the sampling interval distribution can significantly improve infection time estimation as well as transmission tree reconstruction (Table 3). Thirdly, the extent to which sampling times are correlated with infection times determines how well the method is capable of reconstructing transmission trees, which appears from the fact that outbreaks are less well reconstructed with wider sampling interval distributions (Table 1 vs S1 Results) and the low support for the posterior infectors in the Mtb analysis, where sampling intervals were much longer than generation intervals. Therefore, it is advisable to use prior information on sampling intervals in the analysis (if available), and also to base conclusions not only on the summary transmission tree, but also on the posterior support of links in that tree.

We tested the method on five published datasets, with outbreaks of viral and bacterial infections, and in diverse settings of open and closed populations, in human and veterinary context.

The method performed well on all datasets in terms of MCMC chain mixing and tree reconstruction. With naive priors on mean sampling intervals and mutation rates, we obtained estimates that were all very accurate when compared to literature, and the inferred transmission trees seemed as good, or even better when considering estimated infection times. With two datasets (MRSA and H7N7) we included some prior information on sampling and/or generation intervals, which mainly affected the inferred infection times, but not so much the transmission trees. It is possible that not all cases have been observed in these outbreaks, especially in the Mtb and MRSA outbreaks, an assumption nevertheless made by our model. If not too many cases are missing, the analyses of simulations show that this does not disturb identification of infector-host pairs that are in the data. It will only limitedly affect identification of transmission clusters, because if a host's true infector is not in the data, the true infector's infector is often selected as the most likely infector. Only some of the infection times may have been estimated too early.

For wide applicability, we kept the underlying model simple without putting prior constraints on the order of unobserved events such as infection and coalescence times. Four submodels with only one or two parameters each were used for sampling, transmission, within-host pathogen dynamics, and nucleotide substitution. The sampling model, a gamma distribution for the interval between infection and sampling, has a direct link to inferred infection times, and is the model for which it is most likely that prior information is available from epidemiological data in the same or other outbreaks. We used simulated data to study the effect of uninformative or incorrect prior information on shape parameter a_S and mean m_S . It appears that an incorrect a_S or an incorrect informative prior for m_S does reduce accuracy of inferred infection times. However, consensus trees are hardly affected, at least if the number of SNPs is in the order of the number of hosts as we saw in the actual datasets (Table 1 and Table 2 *Slow Clock*). Only the precision of consensus trees is reduced, i.e. there are fewer inferred infectors with high support. Results with the *Fast Clock* simulations did show a significant reduction in consensus tree accuracy. In that case, there are so many SNPs that the phylogenetic tree topology and times of coalescent nodes are almost fixed; then, too much variance in sampling intervals (low a_S) results in many incorrect placements of infection events on that tree. Possibly, with so many SNPs it could be more efficient to first make the phylogenetic tree, and then use that tree to infer transmission events [7, 10], but it is questionable whether genome-wide mutation rates are ever so high that this becomes a real issue [41].

The submodel for transmission is relevant for transmission tree inference in describing the times between subsequent infection events. Transmission is modelled as a homogeneous branching process, implicitly assuming that there was a small outbreak in a large population, with a reproduction number (mean number of secondary cases per primary case) of 1. If all, or almost all, infectors are in the data, the generation interval distribution reflects the course of infectiousness, separating the cases in time along the tree. This interpretation may be obscured with many unobserved cases, as in the absence of the actual infector, the method often identifies an earlier ancestor in the transmission tree as infector (Table 2). Apart from not taking heterogeneity across hosts into account (an extension we wish to leave for future development, see below), the current model also neglects the possibility that susceptibles can have contact with several infecteds in a smaller population or more structured contact network. That could be modelled by a force of infection, which would more realistically describe contraction of the generation interval during the peak of the outbreak, and provide estimates for relevant quantities such as reproduction ratios [6]. However, it requires information about uninfected susceptibles in the same population and a more complicated transmission model, which is a significant disadvantage when it comes to general applicability, one of our primary aims. More importantly, for transmission tree inference it does not seem to be a problem: the analyses of

the published simulations were almost as accurate as in the original publication [12], and these simulations were in very small populations with almost all hosts infected.

The role of the within-host model is to integrate over all possible phylogenetic mini-trees and mutation events within the hosts, and through that, to obtain a posterior distribution of all transmission trees consistent with the (genetic) data. For this, we used a coalescence model based on a linearly growing within-host population, combined with a Jukes-Cantor substitution model. These models contain each only one parameter, but we think that—as long as only few mutations occur in each host, as in our own simulations, the published *Slow Clock* simulations, and most datasets—for most applications more complex models are not needed for the following reasons. First, the gross structure of the phylogenetic tree topology and branch lengths result from transmission and sampling models, and only the finer within-host details are determined by the within-host model. With only few mutations within each host, precise mini-tree inference is not possible, and for our aim of inferring transmission trees, unnecessary. Second, and confirming this imprecise mini-tree inference, most tree proposal steps include simulation of the within-host phylogenetic mini-trees, resulting in good mixing of transmission and phylogenetic tree topologies. The fact that proposing from the prior distribution works so well indicates that the sequence data do not contain much information on within-host branch lengths. Third, if there are few SNPs, the posterior contains almost only phylogenetic trees with fewest mutations (maximum parsimony). It is therefore not likely that tree inference will improve with more general substitution models. Fourth, inference of transmission trees and infection times appears not to be biased if the underlying simulation model was more realistic (Table 2). If data do contain many SNPs, as in the *Fast Clock* simulations, more detailed and realistic models for within-host pathogen growth and nucleotide substitution do probably improve inference, especially on the phylogenetic tree. Nonetheless, even then our method was still capable of correctly inferring the infection times and transmission trees with almost the same accuracy as in the original publication.

With two exceptions, the parsimony scores of posterior tree samples were always equal to the number of SNPs in the datasets (the minimum possible). The first exception is the set of *Fast Clock* published simulations, which had so many SNPs that many of the same mutations had occurred in parallel. The second exception is the H7N7 dataset. In that case, the analyses of the three genes separately resulted in parsimony scores with 6–12 (6%–13%) more mutations than the number of SNPs, whereas the analysis of all genes together resulted in a parsimony score of 313 (median) to explain only 257 SNPs, a surplus of 56 mutations (18%). The results for separate genes could indicate positive selection, confirming the analysis by Bataille et al [28], who even identified specific mutations that had occurred multiple times. The even higher discrepancy for the combined analysis is suggestive of reassortment events, also recognised by Bataille et al [28].

The proposed method and implementation opens perspectives for further extending the methodology to reconstruct phylogenetic and transmission trees from pathogen sequence data. One possible set of extensions arises from changes to the models embedded in our method, to include additional aspects of outbreak dynamics. For instance, the generation time distribution (infectiousness curve) could be made dependent on the sampling interval, which may be relevant for the MRSA outbreak analysis in which the two health-care workers may have transmitted the bacterium until late after infection. This dependence is implicit in methods in which transmission is modelled more mechanistically (e.g. [11, 12, 16]), but we chose not to do that to keep the model more generic. Another important extension would be to relax the assumption of a complete bottleneck at transmission; the bottleneck may be larger in reality [42, 43] and it has previously been relaxed by looking at transmission pairs [44] or modelling it as separate transmission events [18], but not yet in a timed transmission tree. In our

model, this would mean that a host can carry multiple phylogenetic mini-trees, rooted at the same infection time to the same infector. A third extension would be to include the possibility of reassortment of genes within a host, primarily motivated by the results of the H7N7 analysis. This may be done by modelling the coalescent process within hosts, the phylogenetic mini-trees, differently for different genes, but constrained by a single transmission tree. Finally, it would be possible to allow for multiple index cases, which may play a role in open populations with possible re-introductions (as in the MRSA setting), or when only a subset of a large epidemic is analysed (the FMD2001 dataset). This is implemented in models using genetic models based on pairwise genetic distances [8, 16] and with a model assuming coalescence at transmission [45], but is considered a major challenge with a within-host coalescent model [46]. Multiple index cases could also reflect unobserved hosts in the outbreak itself, recently addressed by Didelot et al [24] in their two-step approach of first inferring a phylogenetic and then a transmission tree.

A second type of extension would stem from incorporating additional data. An example is the use of data that make particular transmissions more or less likely, such as contact tracing data, or censoring times for infection times per host or transmission times between sets of hosts, motivated by the MRSA dataset in which admission and discharge days are known for each patient. Sampling of infection times and infectors could be constrained by these additional data (as in [12, 30]) and could then become less dependent on the sampling times and sampling interval distribution, as in the current implementation. Another example is the use of spatial data in combination with a spatial transmission kernel, so that the likelihood of infectors includes a distance-dependency, a possible extension motivated by the FMD and H7N7 analyses (as in [9, 30]). A third example is the use of host characteristics to model infectivity as a function of covariates. With the MRSA data, it would then be possible to test for increased infectivity of the health-care workers, or to test for differences in transmissibility in the three wards. In general, the use of additional host data would make dealing with hosts for which a sequence is not available less problematic: the method currently can include these hosts, but without additional data their role is unclear and they are often placed at the end of transmission chains in consensus trees (Fig 2B, Fig 3).

Methods

Data

We developed our model for fully observed outbreaks of size n hosts. Data consist of the sampling times \mathbf{S} and DNA sequences \mathbf{G} , which means that for each host i we know the time of sampling or diagnosis S_i and the sequence G_i associated with the sampling time. Hosts without available sequence are given a sequence with noninformative nucleotides (only 'n').

We illustrate the method with the following five datasets from earlier publications (all in [S1 Data](#)):

1. Tuberculosis (*Mycobacterium tuberculosis*, Mtb). This dataset was analysed by Didelot et al [7]. It consists of 33 Mtb cases in a population of drug users (approximate population size 400), with samples collected in a 2.5 years time frame. The 4.4 Mbp long sequences contained 20 SNPs. Analysis of this dataset tests the performance of this method in an outbreak with relatively few cases in a large population.
2. Methicillin-resistant *Staphylococcus aureus* (MRSA). This is the dataset from Nubel et al [25], with 36 MRSA cases in a neonatal ICU sampled within a time period of 7 months. Sampling dates were available for all cases, but sequences only for 28 cases, revealing 26 SNPs in the non-repetitive core genome of 2.7 Mbp. Analysis of this dataset tests for the

performance of this method in an outbreak in a small population, including cases without sequence.

3. Foot-and-mouth disease (FMD2001). This is the dataset from Cottam et al [26] also analysed by several others [9, 11], with 15 infected premises within a time period of 2 months. Sequences were available for all cases, with 85 SNPs among 8196 nucleotides. Analysis of this dataset and the next tests for the performance of this method in a small completely sampled outbreak in a large population and allows comparison of the estimated transmission tree to earlier results.
4. Foot-and-mouth disease (FMD2007). This is the dataset from Cottam et al [27], also analysed by Morelli [9], with 11 infected premises within a time period of 2 months. Sequences were available for all cases, with 27 SNPs among 8176 nucleotides
5. H7N7 avian influenza (H7N7). This dataset has been analysed by several authors [12, 28–30], and consists of 241 poultry farms in a time period of about 2.5 months. Sequences of the HA, NA, and PB2 genes were available on GISAID for 228 farms, with associated sampling dates. The total number of SNPs was 257 in 5541 nucleotides. For the 13 unsampled farms we used the culling date minus 2 days as the observation day (the mean sampling-to-culling interval was 2.4 days in the 228 sampled farms). We analysed the data for the three genes separately, and combined (concatenated). To inform a prior distribution for the interval from infection to sampling, we used estimated infection times from Boender et al [38]. Analysis of this dataset tests for the performance of this method in a large outbreak, including cases without sequence.

The model and likelihood

The model describes the spread of an infectious pathogen in a population through contact transmission, the dynamics of the pathogen within the infected hosts, and mutation in the DNA or RNA of that pathogen. Furthermore, it describes how these dynamics are observed through sampling of pathogens in infected hosts. We infer the transmission tree and parameters describing the relevant processes by a Bayesian analysis, using Markov-Chain Monte Carlo (MCMC) to obtain samples from the posterior distributions of model parameters and transmission trees (infectors and infection times of all cases). We first introduce the models and likelihood functions; then we explain how we update the transmission trees and parameters in the MCMC chain.

The posterior distribution is given by

$$\Pr(\mathbf{I}, \mathbf{M}, P, \boldsymbol{\theta} | \mathbf{S}, \mathbf{G}) \propto \Pr(\mathbf{S}, \mathbf{G} | \mathbf{I}, \mathbf{M}, P, \boldsymbol{\theta}) \cdot \Pr(\mathbf{I}, \mathbf{M}, P, \boldsymbol{\theta}). \quad (1)$$

Eq (1) is the probability for the unobserved infection times \mathbf{I} , infectors \mathbf{M} , phylogenetic tree P , and model parameters $\boldsymbol{\theta}$, given the data (sampling times and sequences). The posterior probability can be split up into separate likelihood terms representing the four processes, times a prior probability for the parameters (see [S1 Methods](#)):

$$\Pr(\mathbf{I}, \mathbf{M}, P, \boldsymbol{\theta} | \mathbf{S}, \mathbf{G}) \propto \Pr(\mathbf{G} | P, \boldsymbol{\theta}) \cdot \Pr(P | \mathbf{S}, \mathbf{I}, \mathbf{M}, \boldsymbol{\theta}) \cdot \Pr(\mathbf{S} | \mathbf{I}, \boldsymbol{\theta}) \cdot \Pr(\mathbf{I}, \mathbf{M} | \boldsymbol{\theta}) \cdot \Pr(\boldsymbol{\theta}). \quad (2)$$

We now introduce the four models, the associated likelihoods, and prior distributions for associated parameters.

Transmission. We assume that the outbreak started with a single case. Each case produced secondary cases at random generation intervals after their own infection (Gamma distribution with shape a_G and mean m_G). We consider that all untimed transmission tree

topologies are equally likely, so that the probability of the transmission tree only depends on its timing. The outbreak is described by the vectors \mathbf{I} and \mathbf{M} with infection times I_i and infectors M_i for all numbered cases i . The infector of the index case is 0. The likelihood is the product of probability densities ($d_{\Gamma(a_G, m_G)}(\cdot)$) of all generation times in the outbreak:

$$\Pr(\mathbf{I}, \mathbf{M} | a_G, m_G) = \prod_{i: M_i > 0} d_{\Gamma(a_G, m_G)}(I_i - I_{M_i}). \tag{3}$$

Sampling. We assume that all cases are observed and sampled once at random times after they were infected, according to a Gamma distribution with shape a_S and mean m_S . Transmission and sampling are independent, so transmission can take place after sampling, and a case can be sampled earlier than its infector. The likelihood is the product of probability densities of all sampling intervals in the outbreak:

$$\Pr(\mathbf{S} | \mathbf{I}, a_S, m_S) = \prod_i d_{\Gamma(a_S, m_S)}(S_i - I_i). \tag{4}$$

Within-host dynamics. The main function of the within-host model is to allow for a stochastic coalescent process within the host. Each host i harbours its own phylogenetic mini-tree P_i , with the tips being the transmission and sampling events, and the root being the time of infection, before the first coalescent node. Thus, samples are assumed to be clonal lineages. In its simplest form, a phylogenetic mini-tree consists of a single branch from infection time to sampling time. The likelihood is the product of all likelihoods per host:

$$\Pr(P | \mathbf{S}, \mathbf{I}, \mathbf{M}, r) = \prod_i \Pr(P_i | S_i, \mathbf{I}, \mathbf{M}, r), \tag{5}$$

in which r is the parameter describing the within-host dynamics (see below). The likelihoods per mini-tree are dependent on all infection times and infectors, because these determine the transmission times with host i as infector.

Going backwards in time, coalescence between any pair of lineages within a host takes place at rate $1/w(\tau, r)$, where $w(\tau, r) = r\tau$ denotes the linearly increasing within-host pathogen population size at (forward) time τ since infection of the host. With this particular function coalescent nodes tend to be close to the transmission events if r is small, whereas they tend to be soon after infection of the infector if r is large. This function also naturally results in only one lineage at the time of infection (complete transmission bottleneck), as the coalescence rate goes to infinity near the time of infection.

In the complete phylogenetic tree P , three types of nodes x are distinguished (Fig 1A): nodes $x = 1 \dots n$ are the sampling nodes of the corresponding hosts $i = 1 \dots n$, i.e. the tips of the tree at which sampling took place; nodes $x = n+1 \dots 2n-1$ are the coalescent nodes; nodes $x = 2n \dots 3n-1$ are the transmission nodes, i.e. the points in the tree at which a lineage goes from one host to the next. By h_x we identify the host in which node x resides; for transmission nodes it identifies the primary host (infector). The mini-tree P_i is the set of nodes within host i , and τ_x is the time of node x since infection of host h_x , so τ_i is the time of sampling. Let $L_i(\tau)$ denote the number of lineages in host i at time τ since infection:

$$L_i(\tau) = 1 + \sum_{x \in P_i \cap n < x < 2n} u(\tau - \tau_x) - \sum_{x \in P_i \cap x \geq 2n} u(\tau - \tau_x) - u(\tau - \tau_i), \tag{6}$$

in which $u(\tau)$ is the heaviside step function, i.e. $u(\tau) = 0$ if $\tau < 0$, and $u(\tau) = 1$ if $\tau \geq 0$. In other words, $L_i(0) = 1$ by definition because of the complete transmission bottleneck, and then it

increases by 1 at each coalescent node and decreases by 1 at each transmission event and at sampling. The likelihood for each mini-tree can then be written as

$$\Pr(P_i|S_i, \mathbf{I}, \mathbf{M}, r) = \exp\left(-\int_0^\infty \binom{L_i(\tau)}{2} \frac{1}{w(\tau, r)} d\tau\right) \prod_{x|P_i \cap n < x < 2n} \frac{1}{w(\tau_x, r)}, \quad (7)$$

with $\binom{0}{2} \equiv \binom{1}{2} \equiv 0$. The first term is the probability to have no coalescent events during the intervals in which there are two or more lineages, the second term is the product of coalescent rates at the coalescent nodes.

Mutation. We use a single fixed mutation rate μ for all sites, with mutation resulting in any of the four nucleotides with equal probability (Jukes-Cantor). This parameterisation means that the effective rate of nucleotide change is 0.75μ . Given the phylogenetic tree, this results in the likelihood:

$$\Pr(\mathbf{G}|P, \mu) = \prod_{loci} \sum_{\{A,C,G,T\}^{3n-1}} \prod_x \left(\frac{1}{4} - \frac{1}{4}\exp(-\mu(t_x - t_{v_x}))\right)^{I_{mut}(1-N)} \cdot \left(\frac{1}{4} + \frac{3}{4}\exp(-\mu(t_x - t_{v_x}))\right)^{(1-I_{mut})(1-N)}. \quad (8)$$

Here, we multiply over all coalescent and transmission nodes x , which occur at time t_x and have parent node v_x ; 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 noninformative nucleotide ('n' in the sequence). The likelihood is calculated using Felsenstein's pruning algorithm [47].

Prior distributions. Here we describe our general choice of prior distributions, not the particular parameterization in our analyses (Section *Evaluating the method*). We chose fixed values for a_G and a_S , the shape parameters of generation and sampling intervals. For their means m_G and m_S , we used prior distributions with means μ_G and μ_S and standard deviations σ_G and σ_S , which are translated into Gamma-distributed priors for rate parameters $b_G = a_G/m_G$ and $b_S = a_S/m_S$, distributed as $\Gamma(a_{0,G}, b_{0,G})$ and $\Gamma(a_{0,S}, b_{0,S})$ (see [S1 Methods](#)). For the slope r of the within-host growth model, we chose a Gamma-distributed prior with shape and rate a_r and b_r . We chose $\log(\mu)$ to have a uniform (improper) prior distribution, equivalent to $\Pr(\mu) \propto 1/\mu$.

Inference method

We use Bayesian statistics to infer transmission trees and estimate the model parameters from the data, and MCMC methods to obtain samples from the posterior distribution. The procedure is implemented as a package in R (*phybreak*), which can be downloaded from GitHub (www.github.com/donkeyshot/phybreak) and is available on CRAN (cran.r-project.org/web/packages/phybreak/index.html). The package also contains functions to simulate data, and to summarize the MCMC output.

The main novelty of our method lies in the proposal steps for the phylogenetic and transmission trees, used to generate the MCMC chain. It makes use of the hierarchical tree perspective, in which the phylogenetic tree is described as a collection of phylogenetic mini-trees (one for each host), connected through the transmission tree. Most proposals are done by taking one host, changing its position in the transmission tree, and simulating the phylogenetic mini-trees in the hosts involved in that change. In a second type of proposal, the transmission tree is changed while keeping the phylogenetic tree intact. A third type of proposal only resimulates the within-host mini-tree topology, keeping the transmission tree and coalescent times intact.

Initialization of the MCMC chain requires initial values for the six model parameters (a_G , m_G , a_S , m_S , r , and μ). The transmission tree is initialized by generating an infection time for

each host (sampling day minus random sampling interval). The first infected host is the index case, and for the remaining hosts an infector is randomly chosen, weighed by the density of the generation time distribution. Finally, the phylogenetic mini-trees in each host are simulated according to the coalescent model and combined with one another to create a complete phylogenetic tree.

Each MCMC iteration cycle starts with updates of the transmission and phylogenetic trees, followed by updates of the model parameters. To start with the latter, the parameters m_S and m_G are directly sampled from their posterior distribution given the current infection times and transmission tree (Gibbs update). This is done by sampling the rate parameters b_S and b_G , which were given conjugate prior distributions (see above). If $T_S = \sum S_i - I_i$ is the sum of n sampling intervals in the tree, $a_{0,S}$ and $b_{0,S}$ are the shape and rate of the prior distribution for b_S , then a new posterior value is drawn as

$$b_S \sim \Gamma(\text{shape} = a_{0,S} + a_S n, \text{rate} = b_{0,S} + T_S), \tag{9}$$

from which m_S is calculated as a_S/b_S . Posterior values for m_G are drawn from a similar distribution, with $T_G = \sum I_i - I_{M_i}$ the sum of $n-1$ generation intervals. The parameters r and μ are updated by Metropolis-Hastings sampling; proposals r' and μ' are generated from lognormal distributions $LN(r, \sigma_r)$ and $LN(\mu, \sigma_\mu)$, i.e. with current values as mean. The standard deviations are calculated based on the expected variance of the target distributions, given the outbreak size for σ_r , and number of SNPs for σ_μ (see [S1 Methods](#)).

Updating the phylogenetic and transmission trees. The phylogenetic and transmission trees, described by the unobserved variables $\mathbf{Z} = \{\mathbf{I}, \mathbf{M}, \mathbf{P}\}$, are updated by proposing a new tree with proposal density $H(\mathbf{Z}'|\mathbf{Z}, \mathbf{S}, \boldsymbol{\theta})$, and accepting with Metropolis-Hastings probability (using Eq (1)) α ,

$$\alpha = \min\left(1, \frac{\Pr(\mathbf{S}, \mathbf{G}|\mathbf{Z}', \boldsymbol{\theta}) \cdot \Pr(\mathbf{Z}', \boldsymbol{\theta}) \cdot H(\mathbf{Z}|\mathbf{Z}', \mathbf{S}, \boldsymbol{\theta})}{\Pr(\mathbf{S}, \mathbf{G}|\mathbf{Z}, \boldsymbol{\theta}) \cdot \Pr(\mathbf{Z}, \boldsymbol{\theta}) \cdot H(\mathbf{Z}'|\mathbf{Z}, \mathbf{S}, \boldsymbol{\theta})}\right). \tag{10}$$

Per MCMC iteration cycle, n proposals are done with each host as a focal host once, in random order. Each proposal starts by taking a focal host i , drawing a sampling interval $T \sim \Gamma(\frac{2}{3}a_S, m_S)$ from a Gamma distribution with shape parameter $\frac{2}{3}a_S$ and mean m_S , and calculating a preliminary proposal for the infection time $I_i' = S_i - T$. The $\frac{2}{3}$ is chosen to have a slightly higher variance for the proposal than in the likelihood, to improve mixing. Based on this preliminary proposal, the topology of the transmission tree is changed (see below), and in most cases the phylogenetic tree as well (80% probability). However, we also allowed for proposal steps without changing the phylogenetic tree (20% probability); this greatly improves mixing of the MCMC chain if there are many SNPs, which more or less fixes the phylogenetic tree topology. The 80%–20% distribution for the two types of proposal was not optimized but chosen such that mixing of the phylogenetic tree is only limitedly less efficient than without the second type of proposal (keeping the phylogenetic tree fixed). It is possible to include a third type of proposal, in which only the topology of the phylogenetic mini-tree of focal host i is resimulated, to improve phylogenetic tree mixing if there are many SNPs. In the default setting in the *phy-break* package, 80% of proposals are of the first type and 20% of the second type, but the user is free to change these percentages. We used these settings in all our analyses, except for the *Fast clock* simulations, where we used a 75%–20%–5% distribution.

Proposals for changes in transmission and phylogenetic trees. Here we describe how changes in the transmission and phylogenetic trees are proposed for six different situations, based on the preliminary proposal for the infection time I_i' and on whether the index case is

involved. Fig 5 shows the proposed changes. More detail on the proposal distribution and calculation of acceptance probability is given in the S1 Methods.

- A. The focal host i is index case, and the preliminary I_i' is before the first transmission event. In that case, the infection time of host i becomes I_i' , and no topological changes are made in the transmission tree (Fig 5A).
- B. The focal host i is index case, and the preliminary I_i' is after the first transmission event, but before host i 's second transmission event, if there is any. In that case, the infection time of host i becomes I_i' , and host i 's first infectee becomes index case, transmitting to i (Fig 5B).
- C. The focal host i is index case, and the preliminary I_i' is after host i 's second transmission event, if there is any. In that case, the infection times of host i and its first infectee are switched, and host i 's first infectee becomes index case. They may or may not exchange infectees, with 50% probability (Fig 5C).
- D. The focal host i is not index case, and the preliminary I_i' is before infection of the index case. In that case, the infection time of host i becomes I_i' , and host i becomes index case, transmitting to the original index case (Fig 5D).
- E. The focal host i is not index case, and the preliminary I_i' is after infection of the index case, but before host i 's first transmission event. In that case, the infection time of host i becomes I_i' , and a new infector is proposed from all hosts infected before I_i' (Fig 5E).
- F. The focal host i is not index case, and the preliminary I_i' is after host i 's first transmission event. In that case, the infection times of host i and its first infectee are switched, as well as their position in the transmission tree. They may or may not exchange infectees, with 50% probability (Fig 5F).

Each change in the transmission tree is followed by proposing new phylogenetic mini-trees for all hosts involved, i.e. if their infection time was changed or transmission nodes were added or removed (grey hosts in Fig 5).

Proposals for changes in the transmission tree only. Here we describe how changes in the transmission tree are proposed without changing the phylogenetic tree, based on the preliminary I_i' and on whether the index case is involved. Fig 6 shows the proposed changes. More detail on the proposal distribution and calculation of acceptance probability is given in the S1 Methods.

- G. The focal host i is the index case. If the preliminary I_i' is before the first coalescence node, the infection time of host i becomes I_i' , and no changes are made in the transmission and phylogenetic trees. If the preliminary I_i' is after the first coalescence node, the proposal is rejected.
- H. The focal host i is not the index case, and the preliminary I_i' is after the most recent common ancestor (MRCA) of the samples of host i and his infector j , which is a coalescent node in infector j . In that case, the infection time of host i becomes I_i' , and infectees may move from host i to infector j or vice versa (Fig 6A).
- I. The focal host i is not the index case, but his infector j is the index case, and the preliminary I_i' is before the MRCA of the samples of host i and his infector j . In that case, an infection time I_j' is proposed for the infector j . If I_j' is after the MRCA, the infection time of the infector j becomes I_j' , and the infection time of host i becomes the original infection time of his infector j . Infectees may move from host i to infector j or vice versa (Fig 6B). If I_j' is before the MRCA, the proposal is rejected.

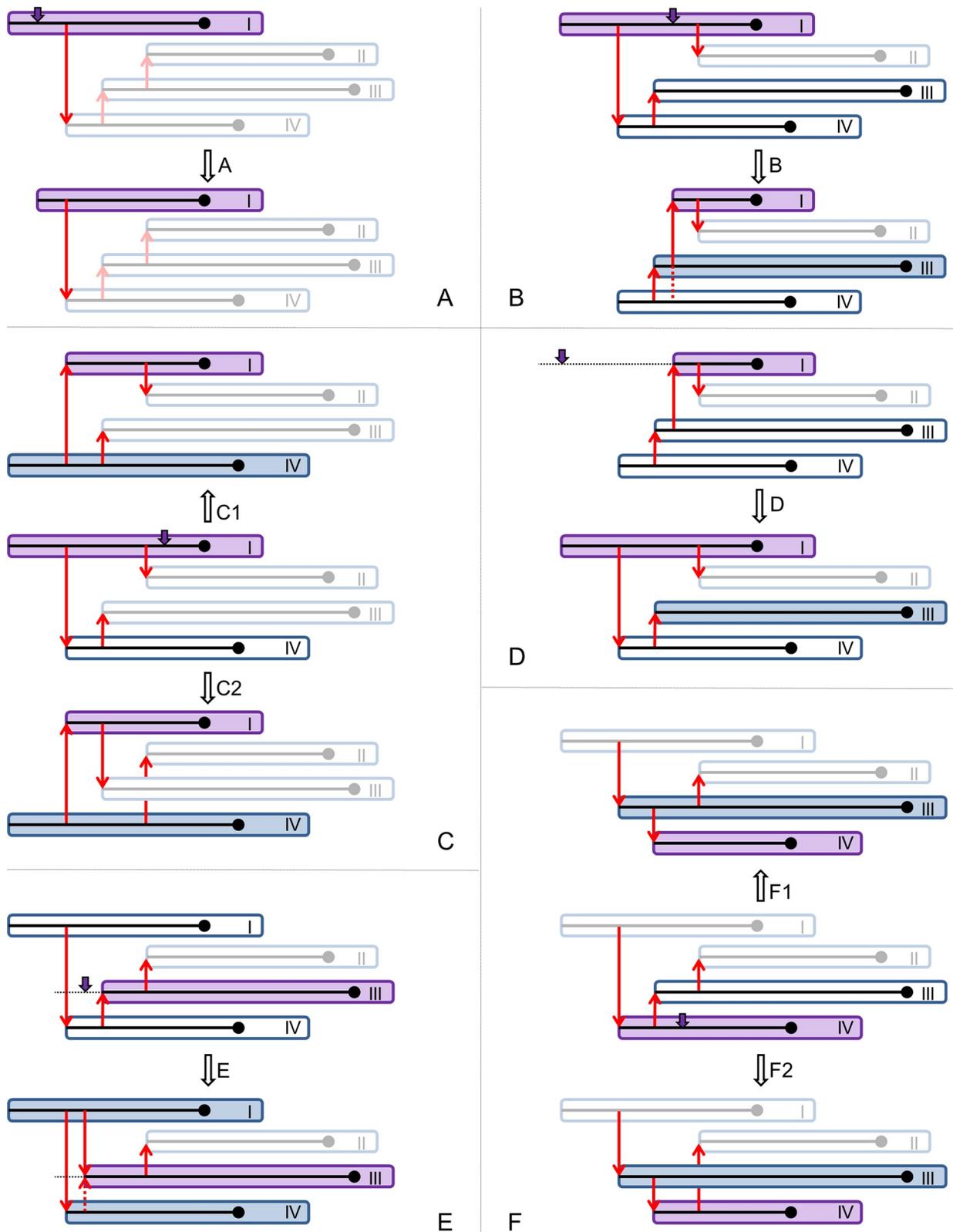


Fig 5. Graphics depicting proposal steps A-F for new transmission and phylogenetic trees. In panels A, B, D, and E, the initial situation is at the top, and the proposal below. In panels C and F, the initial situation is in the middle, and two alternative proposal above and below. Every panel

shows an outbreak with four hosts, with red arrows indicating transmission: the purple host is the focal host, with the purple arrow indicating the proposal for the new infection time I_i' ; filled hosts have a new phylogenetic mini-tree proposed; greyed-out hosts do not play a role in the proposal. (A) the focal host is the index case, and I_i' is before the first transmission event; (B) the focal host is the index case, and I_i' is after the first, but before the second secondary case; (C) the focal host is the index case and I_i' is after his second secondary case; (D) the focal host is not the index case and I_i' is before infection of the index case; (E) the focal host is not the index case and I_i' is before his first secondary case; (F) the focal host is not the index case and I_i' is after his first secondary case.

<https://doi.org/10.1371/journal.pcbi.1005495.g005>

- J. The focal host i is not the index case, and neither is his infector j , and the preliminary I_i' is before the MRCA of the samples of host i and his infector j , but after the MRCA of the samples of host i and infector j 's infector. In that case, an infection time I_j' is proposed for the infector j . If I_j' is after the MRCA, the infection time of the infector j becomes I_j' , and the infection time of host i becomes I_i' . Infectees may move between host i , infector j , and infector j 's infector (Fig 6C). If I_j' is before the MRCA of host i and infector j , or I_i' is before the MRCA of host i and infector j 's infector, the proposal is rejected.

Proposal for changing the phylogenetic mini-tree only

In a single proposal path K, only a new topology of the phylogenetic mini-tree of focal host i is proposed; the coalescent times are kept unchanged.

Evaluating the method

We took three approaches to evaluate our method: analysis of newly simulated data, analysis of published simulated data [12], and analysis of published observed data. When not specified, the following parameter settings and priors were used: shape parameters for sampling and generation interval distributions $a_S = a_G = 3$, uninformative priors for mean sampling and generation intervals with $\mu_S = \mu_G = 1$ and $\sigma_S = \sigma_G = \infty$, and a prior for within-host growth parameter r with $a_r = b_r = 3$. The prior for $\log(\mu)$ (mutation rate) is always uniform.

Analyses were done by two MCMC chains, in each taking 25,000 samples (25,000 MCMC cycles). Burn-ins were different: 5000 MCMC cycles for the newly simulated data, 25,000 for the published simulated data [12], and 5000 for the observed data. With the H7N7 data, five MCMC chains were run, with a burn-in of 5000 samples, followed by 25,000 samples. Burn-in lengths of simulated data were based on visual inspection of convergence for two datasets, and then choosing a burn-in period at least 10 times longer than necessary for all the other simulations, followed by comparing ESS and infector sampling in the parallel chains. The *Slow clock* published simulations had not all converged in the uninformative analysis (S1 Results). For the published data, all chains were inspected visually to confirm convergence.

Analysis of newly simulated data. Four outbreak scenarios were simulated, each replicated 25 times: outbreak sizes of 20 and 50 cases, each with $a_G = a_S = 3$, resulting in overlapping generations and cases sampled earlier than their infector, or $a_G = a_S = 10$, resulting in more discrete generations and cases mostly sampled in order of infection. Further, the mean generation and sampling intervals were $m_G = m_S = 1$ year, the mutation rate $\mu = 10^{-4}$ per year in a DNA sequence with 10^4 sites resulting in a genome-wide mutation rate of 1 per year and a number of SNPs in the same order of magnitude as the outbreak size. For the within-host model we used $r = 1$ per year.

The transmission trees were simulated assuming populations of size 35 or 86 individuals and $R_0 = 1.5$, corresponding to expected final outbreak sizes of about 20 and 50 [22], respectively. Simulations started with one infected individual. All individuals were assumed to be equally infectious, resulting in a Poisson-distributed number of contacts at times since

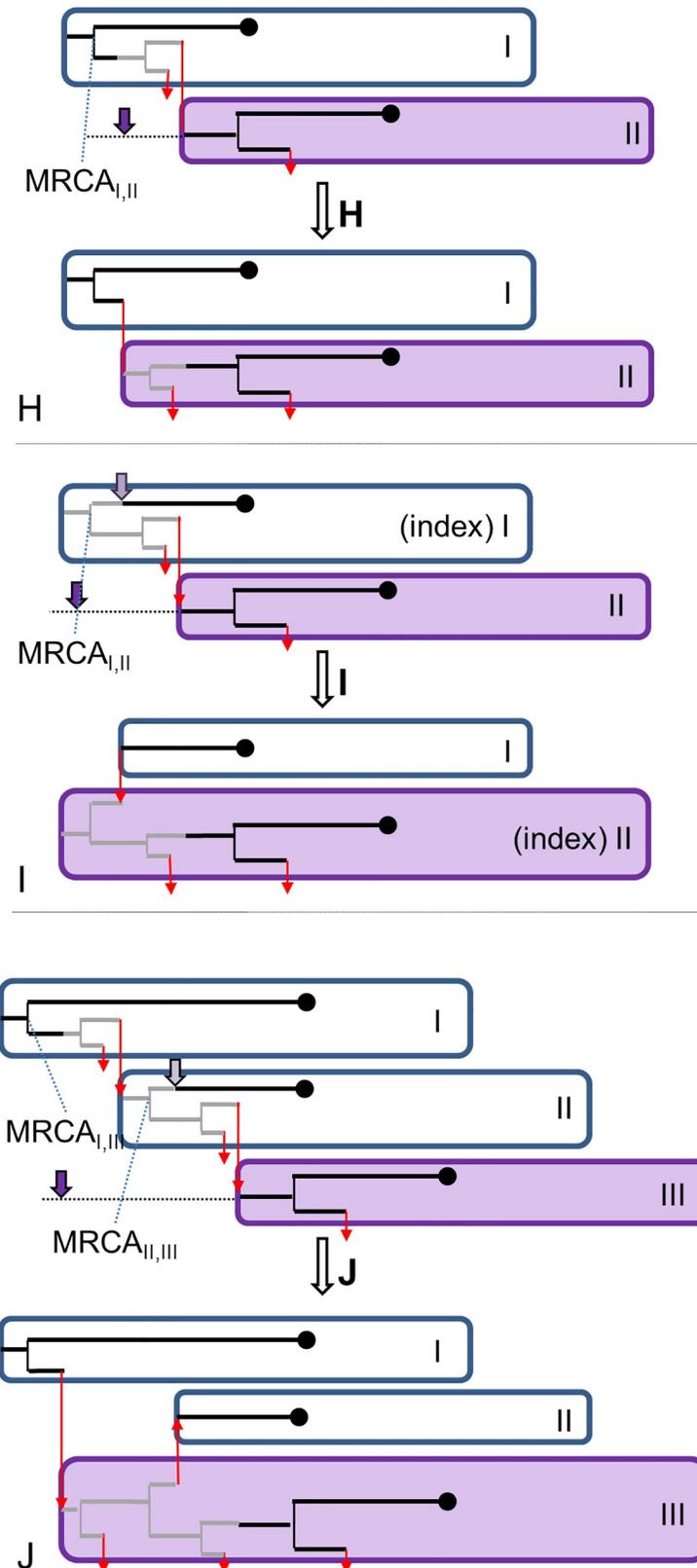


Fig 6. Graphics depicting proposal steps H-J for new transmission trees, keeping the phylogenetic tree unchanged. In all panels, the initial situation is at the top, and the proposal below. Every panel shows

part of an outbreak, with red arrows indicating transmission to depicted or undepicted hosts. Only in panel B host I must be the index case. The purple host is the focal host, with the dark purple arrow indicating the proposal for the new infection time I_i^j ; the light purple arrow in panels B and C indicate the proposal for the new infection time I_i^j of the focal host's infector. The grey parts of the phylogenetic tree are moved between the hosts. (H) the focal host is not the index case, and I_i^j is after MRCA_{I,II} of the focal host and his infector; (I) the focal host is not the index case, and I_i^j is before MRCA_{I,II} of the focal host and his infector (the index case), and I_i^j is after MRCA_{I,II}; (J) the focal host is not the index case, and I_i^j is before MRCA_{I,III} of the focal host and his infector, but after the MRCA_{I,III} of the focal host and his infector's infector; also, I_i^j is after MRCA_{I,III}.

<https://doi.org/10.1371/journal.pcbi.1005495.g006>

infection drawn from the generation time distribution; these contacts were made with randomly selected individuals and resulted in transmissions if that individual had not been infected before. Simulations were repeated until 25 outbreaks were obtained of the desired size.

Given the infection times, sampling times were drawn, and phylogenetic mini-trees were simulated for each host. These were combined into one phylogenetic tree on which random mutation events were placed according to a Poisson process with rate 1. Each mutation event was randomly assigned to one site, and generated one of the four nucleotides with equal probabilities (reducing the effective mutation rate by 25%). By giving the root an arbitrary sequence, the sampled sequences were obtained by following the paths from root to sample and changing the nucleotides at the mutation events.

The simulated data (sampling times and sequences) were analysed with four sets of parameter settings:

- Reference: $a_G = 3$, all other parameters at simulation value (except for μ);
- Informative Correct: a_S at simulation value, informative prior for m_S with $\mu_S = 1$ and $\sigma_S = 0.1$;
- Uninformative: a_S at simulation value;
- Informative Wrong: a_S at simulation value, informative prior for m_S with $\mu_S = 2$ and $\sigma_S = 0.1$.

In addition, the data were analysed with the *Outbreaker* [8] and *TransPhylo* [7, 24] packages in R, with the correct generation and sampling intervals, in *Outbreaker* discretized in steps of 0.1 year. *Outbreaker* analyses consisted of a 200,000 iterations burn-in, followed by 2000 samples with thinning interval of 500. *TransPhylo* analyses consisted of 500,000 iterations burn-in, followed by 500,000 samples without thinning; the required phylogenetic trees were maximum clade credibility (mcc) trees obtained in BEAST v2 [19], assuming a constant coalescent population model. Both the *Outbreaker* and *TransPhylo* outputs were summarized by Edmonds' consensus transmission trees (see below).

Analysis of published simulated data. We used two sets of 25 simulated outbreaks, identified as *Fast clock* and *Slow clock* in the original paper [12], in which full details on the simulations can be found. Briefly summarizing some characteristics, 50 hosts were placed on a grid and a spatial transmission model was run, with exponential transmission kernel. Outbreaks with fewer than 45 cases were discarded. An SEIR (susceptible–exposed–infectious–removed) transmission model was used, with fixed latent period of 2 days and normally distributed infectious period (mean(sd) of 10(1) days). Sampling occurred at the time of removal. Phylogenetic mini-trees were simulated using a logistic within-host growth model $w(\tau) = 0.1(1 + e^6)/(1 + e^{6-1.5\tau})$, starting at $w(0) = 0.1$, then growing to $w(4) = 20.2$ and going to $w(\infty) = 40.4$. Sequences were generated with a 14,000 base pair genome and a mutation rate of 10^{-5} per site per day (*Slow clock*) or $5 \cdot 10^{-4}$ per site per day (*Fast clock*), using a Hasegawa-Kishino-Yano (HKY) substitution model. The *Slow Clock* resulted in a mean number of mutations of 0.14 per

day, or 0.98 per mean generation time of 7 days (latent period plus half infectious period), equivalent to the rate used in the new simulations; the *Fast Clock* was 50 times as fast.

The simulated data (sampling times and sequences, not locations and removal times) were analysed with three levels of prior knowledge on the sampling interval distribution:

- Naive: default settings;
- Uninformative: $a_S = 144$ (coefficient of variation of 0.083, as in the simulation);
- Informative: $a_S = 144$, an informative prior for m_S ($\mu_S = 12$, $\sigma_S = 1$).

Analysis of published datasets. The published Mtb, FMD2001, and FMD2007 datasets were analysed with default settings. The MRSA data contained information on times between hospital entry and first positive sample for 32 patients. Because of their mean and standard deviation of 20 days, we analysed these data with different prior information on the sampling interval only: $a_S = 1$, $\mu_S = 15$, $\sigma_S = 5$. For the H7N7 outbreak data, infection times of the flocks had been estimated [38], from which the mean and standard deviation of the sampling interval was calculated (7.0 and 2.2 days). We used this to inform the sampling intervals with: $a_S = 10$, $\mu_S = 7$, $\sigma_S = 0.5$. Because transmission after culling is not possible, we also used a weak informative prior for the mean generation interval: $\mu_G = 5$, $\sigma_G = 1$.

Performance and outcome measures. The aim of the method is to reconstruct outbreaks in terms of infection times of all hosts and the transmission tree. This requires good mixing of the MCMC chain, especially of infection times and infectors, and a useful method to summarize all sampled transmission trees into a consensus tree.

To test for good mixing, we used effective sample sizes (ESS, calculated with the coda package in R) to evaluate mixing of the parameters and infection times. There are no strict thresholds, but in BEAST, an ESS < 100 is considered too low, whereas an ESS > 200 is considered sufficient [48]. Phylogenetic tree topology mixing was evaluated by the approximate topological ESS [23], available through the *rwty* package in R. Mixing of the transmission tree topology (infector per host) was evaluated as follows. To test for 200 independently sampled infectors per host, the chains were thinned by 250, giving 100 sampled infectors per chain. Then two Fisher's exact tests were done for each host. The first test was to compare the posterior frequency distributions of infectors between the two chains, with a two-row contingency table, entry $\{i, j\}$ counting how often infector j occurred in chain i (100 infectors per chain in total). The second test was to assess independency of subsequent samples within the chains, i.e. absence of autocorrelation, with a contingency table in which entry $\{i, j\}$ counts how often infector i was followed by infector j in the chains (198 pairs of infectors). We used the proportion of successful tests (i.e. $P > 0.05$) as a measure of mixing, expecting 95% successful tests with good mixing.

Two methods were used to make consensus transmission tree topologies (who infected whom), both based on the frequencies of infectors for each host among all posterior trees. The support of host j being the infector of host i is defined as the proportion of posterior trees in which host i infected host j . The first consensus tree is the maximum parent credibility (MPC) tree [12], which is the tree among all posterior trees that has the highest product of infector supports. The second consensus tree is the tree constructed using an adaptation of Edmonds' algorithm, which starts by taking the infector with highest support for each host, and resolves cycles if there are any [21]. Whereas in the original algorithm the sum of weights between nodes is minimized conditional on the absence of cycles, we maximize the sum of supports. Because the algorithm requires prior choice of an index case, we adapted it by repeating the algorithm for all supported index cases, and selecting the tree with highest sum of posterior supports.

Posterior infection times were summarized either outside the context of a consensus tree, i.e. based on all MCMC samples, or for a particular consensus tree, i.e. for each host based only on those samples in which the infector was the consensus infector. The latter is to improve consistency between topology and infection times, although even then consistency is not guaranteed. For plotting transmission trees only, we used the Edmonds' consensus tree; for plotting transmission and phylogenetic trees together, we used the MPC consensus tree, which comes with a consistent phylogenetic tree because it is one of the sampled trees.

Supporting information

S1 Results. Tables with additional results on simulated data.

(DOCX)

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

(DOCX)

S1 Data. Sequence data and sampling times of analysed actual datasets.

(XLSX)

Acknowledgments

We wish to thank Matthew Hall for sharing his simulated data [12], authors of the publications [25–28] for publicly sharing their outbreak data, and Guus Koch for his permission to publish the H7N7 sequence data with this publication.

Author Contributions

Conceptualization: DK JAB XD CC JW.

Data curation: DK.

Formal analysis: DK.

Investigation: DK.

Methodology: DK JAB XD CC JW.

Software: DK JAB.

Supervision: DK JW.

Validation: DK.

Visualization: DK.

Writing – original draft: DK.

Writing – review & editing: DK JAB XD CC JW.

References

1. Gilchrist CA, Turner SD, Riley MF, Petri WA, Jr., Hewlett EL. Whole-genome sequencing in outbreak analysis. *Clin Microbiol Rev.* 2015; 28(3):541–63. PubMed Central PMCID: PMC4399107. <https://doi.org/10.1128/CMR.00075-13> PMID: 25876885
2. Koser CU, Ellington MJ, Peacock SJ. Whole-genome sequencing to control antimicrobial resistance. *Trends Genet.* 2014; 30(9):401–7. PubMed Central PMCID: PMC4156311. <https://doi.org/10.1016/j.tig.2014.07.003> PMID: 25096945

3. Pybus OG, Rambaut A. Evolutionary analysis of the dynamics of viral infectious disease. *Nat Rev Genet.* 2009; 10(8):540–50. <https://doi.org/10.1038/nrg2583> PMID: 19564871
4. Volz EM, Koelle K, Bedford T. Viral phylodynamics. *PLoS Comput Biol.* 2013; 9(3):e1002947. PubMed Central PMCID: PMC3605911. <https://doi.org/10.1371/journal.pcbi.1002947> PMID: 23555203
5. Kenah E. Semiparametric Relative-risk Regression for Infectious Disease Transmission Data. *J Am Stat Assoc.* 2015; 110(509):313–25. PubMed Central PMCID: PMC4489164. <https://doi.org/10.1080/01621459.2014.896807> PMID: 26146425
6. Kenah E, Britton T, Halloran ME, Longini IM Jr. Molecular Infectious Disease Epidemiology: Survival Analysis and Algorithms Linking Phylogenies to Transmission Trees. *PLoS Comput Biol.* 2016; 12(4):e1004869. PubMed Central PMCID: PMC4829193. <https://doi.org/10.1371/journal.pcbi.1004869> PMID: 27070316
7. Didelot X, Gardy J, Colijn C. Bayesian inference of infectious disease transmission from whole-genome sequence data. *Mol Biol Evol.* 2014; 31(7):1869–79. PubMed Central PMCID: PMC4069612. <https://doi.org/10.1093/molbev/msu121> PMID: 24714079
8. Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N. Bayesian reconstruction of disease outbreaks by combining epidemiologic and genomic data. *PLoS Comput Biol.* 2014; 10(1):e1003457. PubMed Central PMCID: PMC3900386. <https://doi.org/10.1371/journal.pcbi.1003457> PMID: 24465202
9. Morelli MJ, Thebaud G, Chadoeuf J, King DP, Haydon DT, Soubeyrand S. A Bayesian inference framework to reconstruct transmission trees using epidemiological and genetic data. *PLoS Comput Biol.* 2012; 8(11):e1002768. PubMed Central PMCID: PMC3499255. <https://doi.org/10.1371/journal.pcbi.1002768> PMID: 23166481
10. Numminen E, Chewapreecha C, Siren J, Turner C, Turner P, Bentley SD, et al. Two-phase importance sampling for inference about transmission trees. *Proc Biol Sci.* 2014; 281(1794):20141324. PubMed Central PMCID: PMC4211445. <https://doi.org/10.1098/rspb.2014.1324> PMID: 25253455
11. Ypma RJ, van Ballegooijen WM, Wallinga J. Relating phylogenetic trees to transmission trees of infectious disease outbreaks. *Genetics.* 2013; 195(3):1055–62. PubMed Central PMCID: PMC3813836. <https://doi.org/10.1534/genetics.113.154856> PMID: 24037268
12. Hall M, Woolhouse M, Rambaut A. Epidemic Reconstruction in a Phylogenetics Framework: Transmission Trees as Partitions of the Node Set. *PLoS Comput Biol.* 2015; 11(12):e1004613. PubMed Central PMCID: PMC4701012. <https://doi.org/10.1371/journal.pcbi.1004613> PMID: 26717515
13. Kanamori H, Parobek CM, Weber DJ, van Duin D, Rutala WA, Cairns BA, et al. Next-Generation Sequencing and Comparative Analysis of Sequential Outbreaks Caused by Multidrug-Resistant *Acinetobacter baumannii* at a Large Academic Burn Center. *Antimicrob Agents Chemother.* 2016; 60(3):1249–57. PubMed Central PMCID: PMC4775949.
14. Onori R, Gaiarsa S, Comandatore F, Pongolini S, Brisse S, Colombo A, et al. Tracking Nosocomial *Klebsiella pneumoniae* Infections and Outbreaks by Whole-Genome Analysis: Small-Scale Italian Scenario within a Single Hospital. *J Clin Microbiol.* 2015; 53(9):2861–8. PubMed Central PMCID: PMC4540926. <https://doi.org/10.1128/JCM.00545-15> PMID: 26135860
15. Stoesser N, Giess A, Batty EM, Sheppard AE, Walker AS, Wilson DJ, et al. Genome sequencing of an extended series of NDM-producing *Klebsiella pneumoniae* isolates from neonatal infections in a Nepali hospital characterizes the extent of community- versus hospital-associated transmission in an endemic setting. *Antimicrob Agents Chemother.* 2014; 58(12):7347–57. PubMed Central PMCID: PMC4249533. <https://doi.org/10.1128/AAC.03900-14> PMID: 25267672
16. Worby CJ, O'Neill PD, Kypraios T, Robotham JV, De Angelis D, Cartwright EJ, et al. Reconstructing transmission trees for communicable diseases using densely sampled genetic data. *Ann Appl Stat.* 2016; 10(1):395–417. PubMed Central PMCID: PMC4817375. PMID: 27042253
17. Lau MS, Marion G, Streftaris G, Gibson G. A Systematic Bayesian Integration of Epidemiological and Genetic Data. *PLoS Comput Biol.* 2015; 11(11):e1004633. PubMed Central PMCID: PMC4658172. <https://doi.org/10.1371/journal.pcbi.1004633> PMID: 26599399
18. De Maio N, Wu C-H, Wilson DJ. SCOTTI: efficient reconstruction of transmission within outbreaks with the structured coalescent. *PLoS Comput Biol.* 2016; 12.
19. Drummond AJ, Bouckaert RR. Bayesian evolutionary analysis with BEAST 2. Cambridge: Cambridge University Press; 2015. xii, 249 pages p.
20. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol Biol Evol.* 2012; 29(8):1969–73. PubMed Central PMCID: PMC3408070. <https://doi.org/10.1093/molbev/mss075> PMID: 22367748
21. Gibbons A. Algorithmic graph theory. Cambridge: Cambridge University Press; 1985. xii, 259 p. p.

22. Diekmann O, Heesterbeek H, Britton T. Mathematical tools for understanding infectious disease dynamics. Princeton, N.J.: Princeton University Press; 2013. xiv, 502 pages p.
23. Lanfear R, Hua X, Warren DL. Estimating the Effective Sample Size of Tree Topologies from Bayesian Phylogenetic Analyses. *Genome Biol Evol.* 2016; 8(8):2319–32. PubMed Central PMCID: PMC5010905. <https://doi.org/10.1093/gbe/evw171> PMID: 27435794
24. Didelot X, Fraser C, Gardy J, Colijn C. Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks. *Mol Biol Evol.* 2017;accepted.
25. Nubel U, Nachtnebel M, Falkenhorst G, Benzler J, Hecht J, Kube M, et al. MRSA transmission on a neonatal intensive care unit: epidemiological and genome-based phylogenetic analyses. *PLoS One.* 2013; 8(1):e54898. PubMed Central PMCID: PMC3561456. <https://doi.org/10.1371/journal.pone.0054898> PMID: 23382995
26. Cottam EM, Thebaud G, Wadsworth J, Gloster J, Mansley L, Paton DJ, et al. Integrating genetic and epidemiological data to determine transmission pathways of foot-and-mouth disease virus. *Proc Biol Sci.* 2008; 275(1637):887–95. PubMed Central PMCID: PMC2599933. <https://doi.org/10.1098/rspb.2007.1442> PMID: 18230598
27. Cottam EM, Wadsworth J, Shaw AE, Rowlands RJ, Goatley L, Maan S, et al. Transmission pathways of foot-and-mouth disease virus in the United Kingdom in 2007. *PLoS Pathog.* 2008; 4(4):e1000050. PubMed Central PMCID: PMC2277462. <https://doi.org/10.1371/journal.ppat.1000050> PMID: 18421380
28. Bataille A, van der Meer F, Stegeman A, Koch G. Evolutionary analysis of inter-farm transmission dynamics in a highly pathogenic avian influenza epidemic. *PLoS Pathog.* 2011; 7(6):e1002094. PubMed Central PMCID: PMC3121798. <https://doi.org/10.1371/journal.ppat.1002094> PMID: 21731491
29. Ypma RJ, Bataille AM, Stegeman A, Koch G, Wallinga J, van Ballegooijen WM. Unravelling transmission trees of infectious diseases by combining genetic and epidemiological data. *Proc Biol Sci.* 2012; 279(1728):444–50. PubMed Central PMCID: PMC3234549. <https://doi.org/10.1098/rspb.2011.0913> PMID: 21733899
30. Ypma RJ, Jonges M, Bataille A, Stegeman A, Koch G, van Boven M, et al. Genetic data provide evidence for wind-mediated transmission of highly pathogenic avian influenza. *J Infect Dis.* 2013; 207(5):730–5. <https://doi.org/10.1093/infdis/jis757> PMID: 23230058
31. Soetens LC, Boshuizen HC, Korthals Altes H. Contribution of seasonality in transmission of Mycobacterium tuberculosis to seasonality in tuberculosis disease: a simulation study. *Am J Epidemiol.* 2013; 178(8):1281–8. <https://doi.org/10.1093/aje/kwt114> PMID: 23880353
32. Ford CB, Lin PL, Chase MR, Shah RR, Iartchouk O, Galagan J, et al. Use of whole genome sequencing to estimate the mutation rate of Mycobacterium tuberculosis during latent infection. *Nat Genet.* 2011; 43(5):482–6. PubMed Central PMCID: PMC3101871. <https://doi.org/10.1038/ng.811> PMID: 21516081
33. Walker TM, Ip CL, Harrell RH, Evans JT, Kapatai G, Dedicoat MJ, et al. Whole-genome sequencing to delineate Mycobacterium tuberculosis outbreaks: a retrospective observational study. *Lancet Infect Dis.* 2013; 13(2):137–46. PubMed Central PMCID: PMC3556524. [https://doi.org/10.1016/S1473-3099\(12\)70277-3](https://doi.org/10.1016/S1473-3099(12)70277-3) PMID: 23158499
34. Nubel U, Dordel J, Kurt K, Strommenger B, Westh H, Shukla SK, et al. A timescale for evolution, population expansion, and spatial spread of an emerging clone of methicillin-resistant Staphylococcus aureus. *PLoS Pathog.* 2010; 6(4):e1000855. PubMed Central PMCID: PMC2851736. <https://doi.org/10.1371/journal.ppat.1000855> PMID: 20386717
35. Young BC, Golubchik T, Batty EM, Fung R, Lerner-Svensson H, Votintseva AA, et al. Evolutionary dynamics of Staphylococcus aureus during progression from carriage to disease. *Proc Natl Acad Sci U S A.* 2012; 109(12):4550–5. PubMed Central PMCID: PMC3311376. <https://doi.org/10.1073/pnas.1113219109> PMID: 22393007
36. Chis Ster I, Dodd PJ, Ferguson NM. Within-farm transmission dynamics of foot and mouth disease as revealed by the 2001 epidemic in Great Britain. *Epidemics.* 2012; 4(3):158–69. <https://doi.org/10.1016/j.epidem.2012.07.002> PMID: 22939313
37. Pedersen CE, Frandsen P, Wekesa SN, Heller R, Sangula AK, Wadsworth J, et al. Time Clustered Sampling Can Inflate the Inferred Substitution Rate in Foot-And-Mouth Disease Virus Analyses. *PLoS One.* 2015; 10(12):e0143605. PubMed Central PMCID: PMC4667911. <https://doi.org/10.1371/journal.pone.0143605> PMID: 26630483
38. Boender GJ, Hagenaars TJ, Bouma A, Nodelijk G, Elbers AR, de Jong MC, et al. Risk maps for the spread of highly pathogenic avian influenza in poultry. *PLoS Comput Biol.* 2007; 3(4):e71. PubMed Central PMCID: PMC1853123. <https://doi.org/10.1371/journal.pcbi.0030071> PMID: 17447838

39. Fitch WM. Toward defining the course of evolution: minimum change for a specific tree topology. *Syst Zool.* 1971; 20:406–16.
40. Chen R, Holmes EC. Avian influenza virus exhibits rapid evolutionary dynamics. *Mol Biol Evol.* 2006; 23(12):2336–41. <https://doi.org/10.1093/molbev/msl102> PMID: 16945980
41. Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. Measurably evolving pathogens in the genomic era. *Trends Ecol Evol.* 2015; 30(6):306–13. PubMed Central PMCID: PMC4457702. <https://doi.org/10.1016/j.tree.2015.03.009> PMID: 25887947
42. Shaw GM, Hunter E. HIV transmission. *Cold Spring Harb Perspect Med.* 2012; 2(11). PubMed Central PMCID: PMC3543106.
43. Varble A, Albrecht RA, Backes S, Crumiller M, Bouvier NM, Sachs D, et al. Influenza A virus transmission bottlenecks are defined by infection route and recipient host. *Cell Host Microbe.* 2014; 16(5):691–700. PubMed Central PMCID: PMC4272616. <https://doi.org/10.1016/j.chom.2014.09.020> PMID: 25456074
44. Worby CJ, Chang HH, Hanage WP, Lipsitch M. The distribution of pairwise genetic distances: a tool for investigating disease transmission. *Genetics.* 2014; 198(4):1395–404. PubMed Central PMCID: PMC4256759. <https://doi.org/10.1534/genetics.114.171538> PMID: 25313129
45. 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. *Proc Biol Sci.* 2014; 281(1782):20133251. PubMed Central PMCID: PMC3973266.
46. Baele G, Suchard MA, Rambaut A, Lemey P. Emerging concepts of data integration in pathogen phylogenetics. *Syst Biol.* 2016.
47. Felsenstein J. Evolutionary trees from DNA sequences: a maximum likelihood approach. *J Mol Evol.* 1981; 17(6):368–76. PMID: 7288891
48. Drummond AJ, Ho SY, Phillips MJ, Rambaut A. Relaxed phylogenetics and dating with confidence. *PLoS Biol.* 2006; 4(5):e88. PubMed Central PMCID: PMC1395354. <https://doi.org/10.1371/journal.pbio.0040088> PMID: 16683862