

Beyond the SNP Threshold: Identifying Outbreak Clusters Using Inferred Transmissions

James Stimson,^{*,1} Jennifer Gardy,^{2,3} Barun Mathema,⁴ Valeriu Crudu,⁵ Ted Cohen,⁶ and Caroline Colijn^{1,7}

¹Department of Mathematics, Imperial College London, London, UK

²British Columbia Centre for Disease Control, Communicable Disease Prevention and Control Services, Vancouver, Canada.

³School of Population and Public Health, University of British Columbia, Vancouver, Canada

⁴Department of Epidemiology, Columbia University Mailman School of Public Health, New York, USA

⁵Phthisiopneumology Institute, Chisinau, Republic of Moldova

⁶Yale University School of Public Health, New Haven

⁷Department of Mathematics, Simon Fraser University, Vancouver, Canada

*Corresponding author: E-mail: james.stimson16@imperial.ac.uk.

Associate editor: Thomas Leitner

Abstract

Whole-genome sequencing (WGS) is increasingly used to aid the understanding of pathogen transmission. A first step in analyzing WGS data is usually to define “transmission clusters,” sets of cases that are potentially linked by direct transmission. This is often done by including two cases in the same cluster if they are separated by fewer single-nucleotide polymorphisms (SNPs) than a specified threshold. However, there is little agreement as to what an appropriate threshold should be. We propose a probabilistic alternative, suggesting that the key inferential target for transmission clusters is the number of transmissions separating cases. We characterize this by combining the number of SNP differences and the length of time over which those differences have accumulated, using information about case timing, molecular clock, and transmission processes. Our framework has the advantage of allowing for variable mutation rates across the genome and can incorporate other epidemiological data. We use two tuberculosis studies to illustrate the impact of our approach: with British Columbia data by using spatial divisions; with Republic of Moldova data by incorporating antibiotic resistance. Simulation results indicate that our transmission-based method is better in identifying direct transmissions than a SNP threshold, with dissimilarity between clusterings of on average 0.27 bits compared with 0.37 bits for the SNP-threshold method and 0.84 bits for randomly permuted data. These results show that it is likely to outperform the SNP-threshold method where clock rates are variable and sample collection times are spread out. We implement the method in the R package *transcluster*.

Key words: SNP, transmission clusters, whole-genome sequencing, public health.

Introduction

Whole-genome sequencing (WGS) of pathogens has become an essential tool for improving understanding of how infectious diseases spread between hosts, particularly in the case of tuberculosis (TB) (Hatherell et al. 2016). The phylogeny derived from pathogen genomic data helps us to infer likely transmission events. Typically, samples are taken from patients in the field, the date and other epidemiological data are recorded, and the pathogen’s genome is sequenced. A first step is typically to assign cases to clusters; for infectious diseases, a cluster is a group of closely related infections that is usually interpreted as resulting from recent transmission (Poon 2016). These clusters are chosen primarily with the aim of making meaningful subdivisions of the data, with the added benefit of making the amount of data fed into attempts to reconstruct outbreaks and to transmission inference models more tractable. However the assignment method is often somewhat ad hoc.

The simplest way to determine sequence relatedness is to count the number of single-nucleotide polymorphisms (SNPs) that differ between two sequences. The SNP-threshold approach places two individuals in the same putative transmission cluster if there are fewer than a threshold number of SNPs between their sequenced pathogen genomes. Many existing methods to identify outbreak clusters rely on SNP thresholds, as surveyed recently (Hatherell et al. 2016) in the case of TB. Similar methods are also used for other pathogens (Dallman et al. 2015; Octavia et al. 2015). However, there is little agreement in the literature as to what such a threshold should be— see table 1 for TB SNP thresholds used in some recent studies. The contexts in which these thresholds are applied differ from study to study, so these numbers are not always strictly comparable, but they do indicate the wide range of values that can reasonably be adopted when determining whether or not cases are closely related. By itself, the number of SNP differences between genomes does not directly imply a probability of recent

Table 1. SNP Thresholds Used in Recent TB Studies.

Authors	Lower SNP Threshold	Upper SNP Threshold
Bryant, Harris, et al. (2013b)	≤ 6 (relapse)	>1,000 (re-infection)
Clark et al. (2013)	< 50	>50
Guerra-Assunção et al. (2015)	≤ 10 (relapse)	>100 (re-infection)
Lee et al. (2015)	< 2	(not specified)
Roetzer et al. (2013)	≤ 3	(not specified)
Walker et al. (2013)	≤ 5	>12
Yang et al. (2017)	≤ 12	(not specified)

The lower threshold indicates the number of SNPs below which cases are positively identified as belonging to the same cluster. Where different, the upper threshold indicates the number of SNPs above which cases are identified as clearly not belonging together. Unless otherwise stated, intermediate values are indeterminate.

transmission. This is implicitly recognized in some sources. For example, we have from Walker et al. (2013): “We predicted that the maximum number of genetic changes at 3 years would be five SNPs and at 10 years would be ten SNPs.” Indeed, other studies directly question the use of SNP thresholds, such as Guerra-Assunção et al. (2015), Bergholz et al. (2014) in the context of food-borne pathogens, and Azarian et al. (2016) in an analysis of the spread of methicillin-resistant *Staphylococcus aureus* (MRSA). Nevertheless, the use of a single SNP threshold is often employed in practice; for example the 12 SNP threshold, used for inferring likely transmission between a pair of TB cases by Public Health England (Walker et al. 2014) amongst others, is perhaps the most common in TB.

The appropriate SNP cut-off for inferring transmission is likely to depend critically on the context. There are many sources of uncertainty. Nucleotide mutation rates vary between pathogens, can vary at different stages of infection, and are subject to the effects of selection pressure. Culture processes (e.g. liquid vs. solid culture, single colony picks vs. sweeps) may affect the diversity in samples that are sent for sequencing. Furthermore, the process of producing finalized SNP data from patient-derived biological samples is a multi-stage procedure where there are choices to be made—including how stringently quality filtering is applied to raw genomic data—which will in general result in different SNP differences being reported. As such, it is important that during every step of the pipeline from sampling from patients and processing the data, to building the models and drawing conclusions from them, that we are aware of sources of uncertainty and attempt to propagate this uncertainty to any conclusions. It is also important that as WGS is rolled out widely as a tool in infectious disease, we recalibrate SNP-based methods to accommodate changes in both sequencing technologies and in the bioinformatics pipelines used to call variant SNPs. Clustering methods that use variant SNP calls exclusively will be most sensitive to such changes.

The fundamental logic behind SNP cut-offs is that it takes time to accrue genetic variation; even in organisms where the molecular clock is variable, it seems uncontroversial to assume that two isolates that differ by only a few SNPs are more likely to be a result of recent transmission than isolates that are 50 SNPs apart. However, the rate at which polymorphisms occur varies not only between organisms (Kuo and

Ochman 2009), but also across a genome; it is affected by selection pressure and by horizontal gene transfer (HGT) (Novichkov et al. 2004), though this is not an issue for TB and there are methods to remove recombination and HGT prior to using SNP cut-offs. As per Barrick and Lenski (2013), it is also important to distinguish between the *mutation rate*, the rate at which spontaneous mutations occur, and the *substitution rate*, the rate of accumulation of changes in a lineage; this depends on both the mutation rate and the effects of selection and drift. Here, when we refer to the clock rate, we mean the substitution rate, as we use the rate to interpret variants measured with sequencing technologies.

This distinction is particularly important for diseases like TB, where selection pressure due to antibiotics can be substantial. Whilst the background SNP accumulation rate for *Mycobacterium tuberculosis* (Mtb) has been estimated at 0.5 SNPs/genome/year (Walker et al. 2013), selection pressure and antibiotic resistance can influence this rate considerably. For example, in Eldholm et al. (2014) we see the observation that “After exclusion of transient mutations in the patient isolates, 4.3 mutations were acquired per year. . . or 2.3 mutations per year when excluding resistance mutations.” The size of the population of bacteria within a host could also affect the number of SNPs observed between that host and those they infect. Unexplained larger variation is also encountered as documented in Korhonen et al. (2016), though high SNP numbers could be a result of re-infection or mixed infection rather than in-host evolution. Where we know that selection or high substitution rates are likely to be present and detected, a higher rate is therefore likely to be appropriate for clustering, and this will affect the relationship between SNPs and transmission events.

It should be noted that there are other approaches to clustering, based on molecular (but not WGS) data and including time and geographical data. For example, Kammerer et al. (2013) apply three different statistical tools to spoligo-type data and mycobacterial interspersed repetitive units (MIRU) data, together with date and location of cases, to show that these tools can successfully identify TB outbreaks. Donker et al. (2016) use variable number tandem repeat (VNTR) data for MRSA cases, with time and location data, to identify clusters based on a hierarchical clustering method (Ypma et al. 2013). There are also software packages such as *vimes* (Jombart and Cori 2017), which provides tools allowing users to integrate different types of data and detect outbreaks. These approaches treat each of the underlying variables as independent inputs, without explicitly modeling the connection between time and the accumulation of genetic differences.

We take a slightly different tack here, and jointly use the sample time and genetic distance, together with a model of SNP acquisition over time and transmission events over time, to base putative transmission clusters on the probability that cases are separated by a threshold number of transmission events. This is motivated by a belief that the number of transmission events between two cases is a natural and intuitive measure of how “clustered” they are in the sense of transmission (and how likely they are to be part of the same

outbreak). This cannot usually be measured directly and must be inferred from other data. However, it is reasonable to assume that appropriate incorporation of the time over which the accumulation of SNPs occurs, as well as the likely time between transmission events, give a more accurate and nuanced measure of the likelihood that cases are linked by a small number of transmission events. We develop a probabilistic approach which permits variation in the SNP accumulation process, allows for faster SNP accumulation for sites under selection and allows for variation in the speed with which individuals infect their contacts. We aim to provide a principled alternative to SNP cut-offs for clustering pathogen genomes into putative transmission clusters.

New Approaches

Two samples are usually considered to be in the same transmission cluster if the number of SNPs between them is less than or equal to a fixed cut-off, or threshold. This is a quick way to explore relatedness among a group of isolates and gain an approximate understanding of the extent of recent (low-distance) transmission, but it is coarse and embeds a number of strong assumptions.

Our proposed probabilistic transmission approach, in contrast, is based on sample pairs being clustered together if we estimate that there were fewer than a threshold number of transmission events between them, with a given probability. It uses the same genetic (SNP) distance information as the SNP-threshold method, but in addition makes use of the sample times, knowledge of the SNP accumulation, and transmission processes. The essential inputs to our method are: the number of SNP differences between sample pairs, the sample dates, the assumed clock rate, and the assumed transmission rate.

In addition, our method can readily be extended to incorporate other factors: we show in Materials and Methods how this can be done for spatial data, and for antibiotic resistance. Building these inputs into our model allows us to create a more nuanced and principled way of identifying transmission clusters, and allows us to apply the method consistently in varied settings for example, in those where drug resistance is suspected to be a factor.

We start by establishing probability distributions for the total length of time (h years) along both lines of descent from the most recent common ancestor (MRCA) of a pair of samples; this depends on the clock process, and helps define the distance between the two samples (there is $h/2$ years of elapsed time from the MRCA to the earlier sampling date). We then compute the probability that at least a threshold number of transmissions took place between the two sampled cases over this time, where the probability distribution for the number of transmissions k is $P(k|h)$ (see table 6 for a summary of the symbols used and their units). This approach gives the flexibility to incorporate sample time information and other data. The method uses sample dates and aligned sequence data (variant calls) together with models of the clock and transmission processes. For a pair of samples, we

use the SNP distance N , the time difference between their sampling dates (δ years) and the clock process to write down the probability distribution $\mathcal{L}(h|N, \delta)$ for when the MRCA of the two sequences existed; this must be before the first sampled case. Integrating over this unknown time, we can find the probability that a certain number of transmissions separate the two cases:

$$P(k|N, \delta) = \int_{h=0}^{\infty} \mathcal{L}(h|N, \delta) P(k|h) dh$$

This is equation (9), developed in more detail in Materials and Methods. To incorporate spatial data, a weighting w is applied to the probabilities to reflect that spatial distance can affect estimates of the number of intermediate transmissions between two sampled individuals; we express this in equation (19):

$$P(k|N, \delta, w) = w \int_{h=0}^{\infty} \mathcal{L}(h|N, \delta) P(k|h) dh$$

Results

We illustrate how the transmission method compares to the SNP-threshold method for a simple toy example. We define the “T cut-off” as the cut-off level for the transmission method, using equation (10); the samples are clustered together where the implied number of transmissions k is less than or equal to T with a probability of 80%, given some clock rate λ and transmission rate β . We see that the transmission method clusters the cases together in a different order to the SNP-threshold method as the cut-off level is incremented. Cases A and B are the closest in SNP distance, but the time elapsed between their sampling dates increases their distance by the transmission distance function relative to cases C and D, which are sampled at the same time as each other. So when we take timing into account, the clustering is altered (also illustrated in fig. 1).

We model the number of intermediate transmissions between two sampled hosts given the total time over which SNPs have likely accumulated. Altering the transmission rate β (by which we mean the rate at which intermediate cases occur in the total time elapsed *between the MRCA of two sampled hosts and the sampling events*; see Materials and Methods) alters the absolute transmission cut-off level at which the clusters change— in this example, increasing β to 3.0 transmissions/year gives the same clusters as in figure 1 but at levels 9, 10, and 12 transmissions rather than 7, 8, and 9 transmissions, respectively. This has no impact on the order of the clustering as the level of the cut-off changes (table 2).

In contrast, altering the clock rate does have a material impact on the way clustering occurs as we increase the transmission threshold. For $\lambda = 0.5$ SNPs/genome/year, cases A and B are closest under the transmission method, just as they are with the SNP-threshold method, and so the clustering is the same for both methods. At $\lambda = 1.5$ SNPs/genome/year,

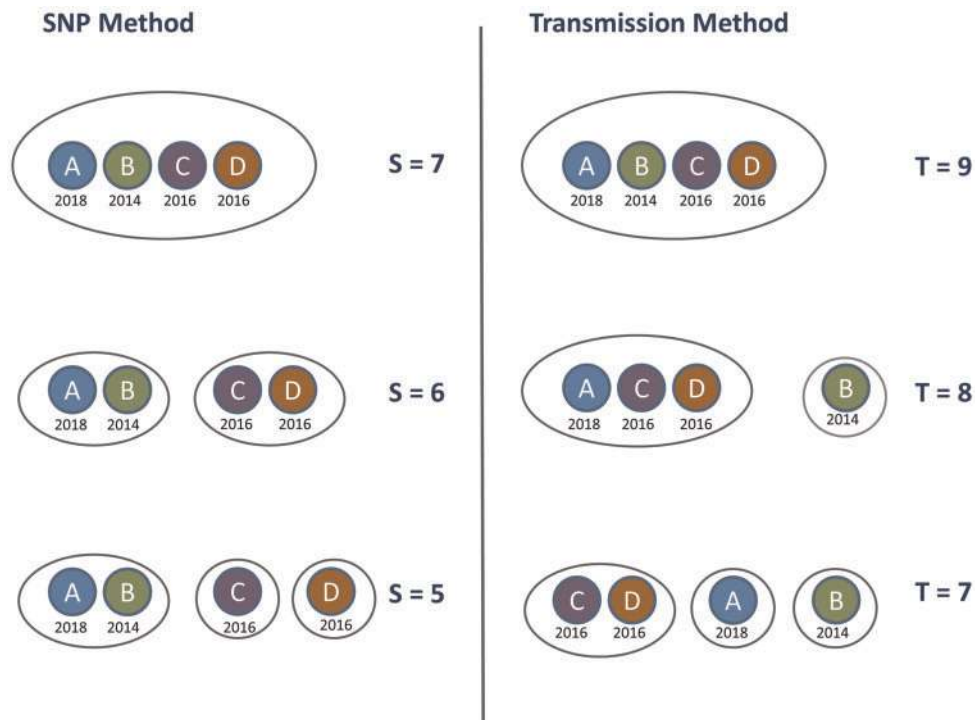


Fig. 1. Clustering on the toy example data set provided in table 2. The left-hand panel shows the clusters obtained by applying the SNP-threshold method with three different thresholds, with the cut-off level denoted by S ; samples are clustered together where the SNP distance is less than or equal to S . The right-hand panel shows the clustering obtained by applying the transmission method, using equation (10), with the cut-off level denoted by T ; samples are clustered together where the implied number of transmissions k is less than or equal to T with a probability of 80%, with clock rate $\lambda = 1.5$ SNPs/genome/year and $\beta = 2.3$ transmissions/year.

Table 2. Model Inputs for Toy Example Data Set.

Label	Sample Date	SNP Dist. to A	SNP Dist. to B	SNP Dist. to C	SNP Dist. to D
A	1/1/2018	0	5	7	7
B	1/1/2014	5	0	8	8
C	1/1/2016	7	8	0	6
D	1/1/2016	7	8	6	0

cases C and D are closest under the transmission method, and the clustering evolves as shown in figure 1.

British Columbia Data

We analyze a data set from British Columbia, comparing the SNP-threshold method to the transmission method in equation (10). The data set comprises 52 samples collected from 51 patients over a 14-year period, and has been prefiltered with the result that all samples are relatively close—within 25 SNPs. Consequently, using the SNP-threshold method with the threshold set to 13 SNPs or higher, all samples are placed in one cluster. When the threshold is 9 SNPs, we obtain a large 42-case cluster, a secondary 8-case cluster and some outliers. As we reduce the threshold further down to 3, the large cluster breaks up but the 8-case cluster persists. We illustrate this in figure 2.

We used the transmission method, using equation (10) with $\beta = 2.0$ transmissions/year and two different average clock rates: $\lambda = 0.5$ and 1.5 SNPs/genome/year (1.5 is larger than

the typical rate for TB but within other outbreak estimates [Bryant, Schürch, et al. 2013]). When λ is low, we can obtain the same clustering as with the SNP cut-off. When λ is high, we have one cluster which contains all the samples for $T > 11$. As with the SNP-threshold method, at $T = 11$ we have a large 42-case cluster, a secondary 8-case cluster, and some outliers. But as we move to $T = 10$, the secondary cluster loses a member, whilst the main cluster stays at size 42. This is because one of the members of the secondary group is very close by the SNP distance to another member of that group, but was sampled more than 10 years before. As with our simple toy example, timing alters the effective distance between samples because the distance takes into account the clock rate and the transmission rate, and so the timing information can affect the clustering.

Furthermore, the probabilistic nature of the approach means that we can see how strongly we predict cases to be linked; in figure 3 we use thicker edges to denote a higher probability of being linked by relatively few transmissions. In addition, we show the effect of incorporating spatial proximity, using equation (19); we assign each of the cases into one of six numbered regions. Including a spatial weighting, reflecting the barrier to the infection moving between different regions, and leaving all other parameters unaltered, changes the clustering that is obtained.

Sensitivity to Clock Rate

An implicit assumption of the SNP-threshold method is that each SNP contributes equally toward the SNP distance. This

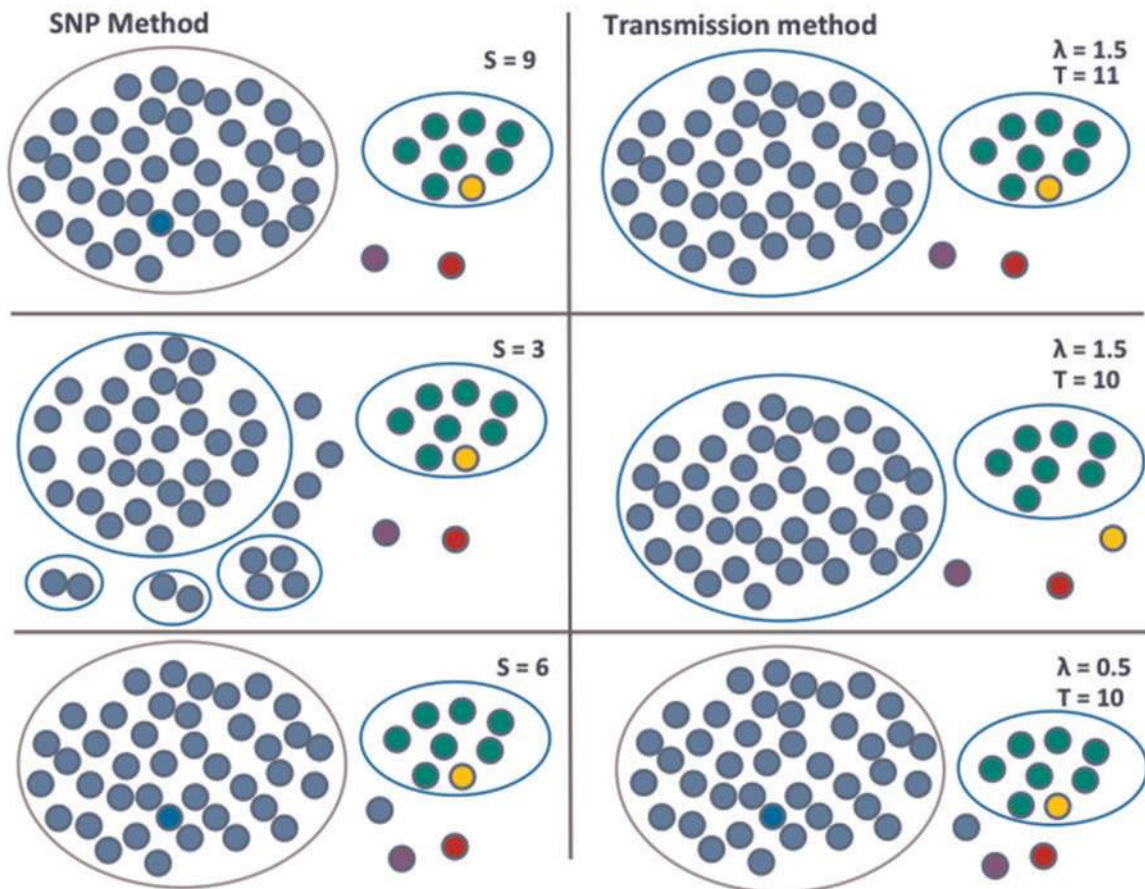


Fig. 2. Clustering on the British Columbia data set. The left-hand panel shows the clusters obtained by applying the SNP-threshold method with three different thresholds, with the cut-off level denoted by S . The largest cluster breaks up as the level is lowered whilst the size 8 cluster remains intact. The right-hand panel shows the clustering obtained by applying the transmission method, using equation (10); samples are clustered together where the implied number of transmissions k is less than or equal to T with a probability of 80%. As shown in the top two thirds, with clock rate $\lambda = 1.5$ SNPs/genome/year and $\beta = 2.0$ transmissions/year, the size 8 cluster loses a member whilst the largest cluster stays the same as the level is lowered. When λ is low, h is larger, so the MRCA of a cluster gets pushed back further in time. In this case, the value of δ between two cases has a limited impact on the estimated number of transmissions; the SNP difference is dominant, and we recover the same clustering that is obtained with the SNP-threshold method. This is shown in the lower third, where the clock rate $\lambda = 0.5$ SNPs/genome/year and $\beta = 1.2$ transmissions/year.

implies that the clock rate or substitution process is constant across the set of isolates and across the genome. When the same threshold is used in different settings and across different pathogen subtypes, the implicit assumption is that the same substitution process holds in these settings. In our new transmission method, the effective distance between any two samples is inversely proportional to the assumed mean clock rate. A lower clock rate means that more time is needed in order for a fixed number of SNPs to be generated; this gives room for more potential intermediate transmission events. A higher clock rate means that the fixed time between samples has a greater effect on the clustering, as the time between samples places a greater constraint on the range of possible heights h ; the fixed time “uses up” more of the time available than it would under a low clock rate (because there is less total estimated time available, a higher portion of it is in the time period δ). We show this in table 3: the transmission clustering method approaches the same results as the SNP clustering method as the assumed clock rate is reduced. We use the variation of information dissimilarity measure given

by *clue* (Meilă 2007) to compare the clusters produced by the two methods.

Moldova Data

This data set comprises 422 samples collected over a period of less than 2 years. For this data—with any reasonable choices of parameters and a fixed substitution rate for all sites—as shown in table 4, our new transmission method does not differ from the SNP-threshold method. This can be explained by two factors that work together: the small distance in time between any two samples and the large SNP differences between cases. There is not enough variation in the timing information relative to the SNP distances for an appreciable difference to emerge between the two clustering methods.

Use of Drug Resistance-Confering SNPs

We can, however, explore the role of drug resistance-conferring SNPs on the clustering. Information on the location of resistance-conferring sites for TB was obtained using

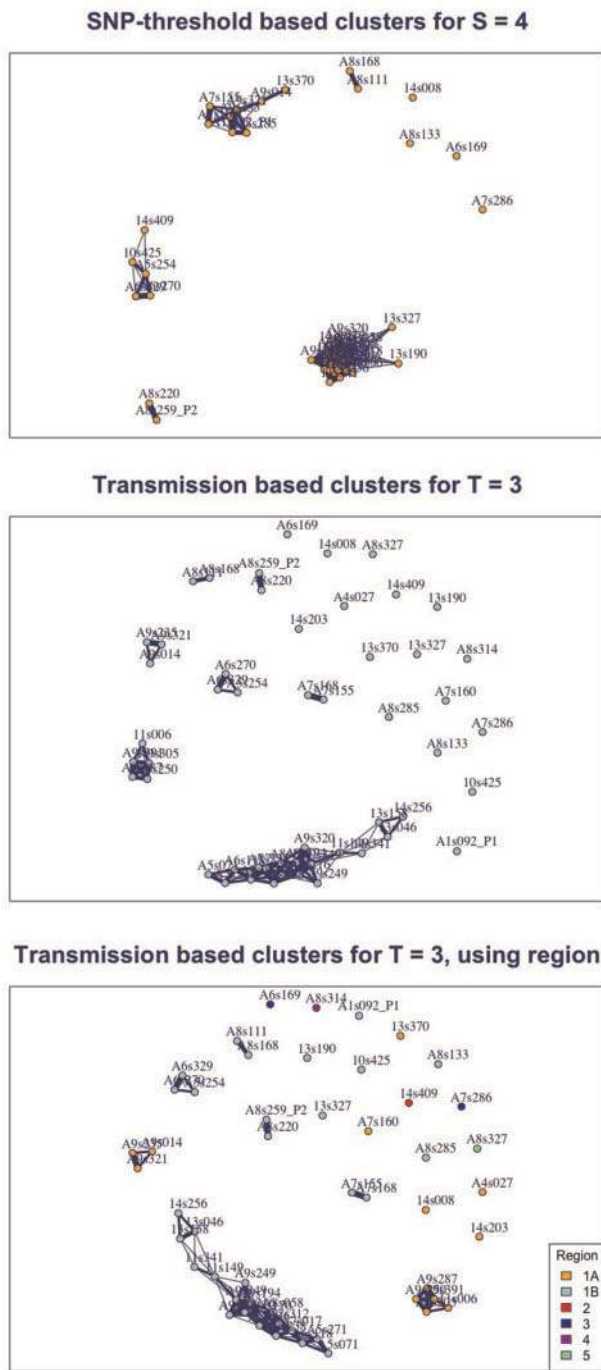


FIG. 3. Three views of the same British Columbia TB data illustrating the contrasting effect of implementing the SNP-threshold and transmission methods and showing estimates of how close individual cases are to each other. In the top figure, edges between nodes indicate that cases are within 4 SNPs of each other. In the lower figures, edges indicate that cases are 80% likely to be within 3 transmission events of each other, given a clock rate $\lambda = 1.5$ SNPs/genome/year and $\beta = 2.0$ transmissions/year. The middle figure is based on equation (10), and the bottom figure uses equation (19), with weighting $w = 20\%$ where two cases are assigned to differing regions. The thicker the edges, the closer the cases are: for the SNP-based clusters the thickest edges correspond to no SNP difference, the thinnest to a distance 4 SNPs; for the transmission-based clusters the thickest edges correspond to one likely transmission event, the thinnest to 3.

PhyResSE (Feuerriegel et al. 2015) and a resistance-conferring SNP distance matrix was computed for the Moldova data by filtering against this information. Selection is likely to lead to resistance-conferring SNPs arising more quickly than other SNPs: for example, one TB study (Eldholm et al. 2014) gives a mutation rate of 4.3 SNPs per genome per year when they are included, in contrast to the 0.5 SNPs per genome per year that is typically estimated for TB (Walker et al. 2013). Resistance acquisition may further increase the rate of acquisition of additional SNPs through multiple resistance, compensatory mutations, or other mechanisms. For this analysis we used a clock rate for the drug-resistant sites, as in equation (18), which is five times higher than for the sites which are not resistance conferring.

Overall, resistance-conferring SNPs in the Moldova data set form only 0.6% of the total number of SNPs. However, restricting to those sample pairs where the SNP distance is less than or equal to 20, they form 8% of the total. If a high proportion of the SNPs between two cases are resistance-conferring SNPs, then this effectively shortens the distance between the cases, making them more likely to be joined together in a transmission cluster. For several sample pairs in this data set, the proportion of resistance-conferring SNPs that differ between the two samples is approaching 35%, whilst for some other pairs there are none at all. For this reason we see a difference when we take resistance into consideration, as seen in table 4 and figure 4. The largest cluster is not shown in detail in the figure and is more robust with respect to the effect of resistance-conferring SNPs than the smaller clusters.

Simulated Data

To explore the performance of the clustering methods in a setting where the “ground truth” is known, we simulate data and compare the SNP-threshold and transmission (as given by eq. 10) methods. The “true” clusters are generated from simulated transmission networks produced by *TransPhylo* (Didelot et al. 2017).

We consider clustering cases based on direct transmission, so that two cases are joined in a cluster if one infected the other, and we compare clusters generated by the SNP-threshold method with those generated by the transmission method. In order to compare the appropriate set of clusters, we find the best match that the method achieves against the true cluster over an appropriately wide range of threshold levels. Then we use the variation of information dissimilarity measure given by *clue* (Meilă 2007) to compare the results of the two methods to the true clusters. We also compare randomly permuted simulated data to the simulated clusters to provide a yardstick of accuracy. This is achieved by fixing the number of clusters to be the number of the true clusters, and then randomly allocating each sample case to one of those clusters. The results in table 5 show that the transmission method is consistently better than the SNP-threshold method in identifying direct transmissions within an outbreak. Both methods perform significantly better than the randomly generated data.

Table 3. Effect of Varying the Clock Rate Using the British Columbia Data.

SNP Threshold	Largest Cluster	λ	β	Trans. Threshold	Largest Cluster	Dissimilarity between SNP and Trans. Methods
S = 12	50	0.5	1.2	T = 22	50	0.000
S = 9, 10, 11	42	0.5	1.2	T = 16	42	0.000
S = 6, 7, 8	41	0.5	1.2	T = 10	41	0.000
S = 5	37	0.5	1.2	T = 8	37	0.000
S = 4	31	0.5	1.2	T = 7	31	0.139
S = 3	29	0.5	1.2	T = 6	29	0.052
S = 2	28	0.5	1.2	T = 3	28	0.113
S = 12	50	1.0	1.2	T = 11	50	0.000
S = 9, 10, 11	42	1.0	1.2	T = 8	42	0.000
S = 6, 7, 8	41	1.0	1.2	T = 5	41	0.058
S = 5	37	1.0	1.2	T = 4	37	0.113
S = 4	31	1.0	1.2	T = 3	29	0.374
S = 3	29	1.0	1.2	T = 3	29	0.174
S = 2	28	1.0	1.2	T = 1	28	0.113
S = 12	50	1.5	1.2	T = 7	52	0.163
S = 9, 10, 11	42	1.5	1.2	T = 6	42	0.000
S = 6, 7, 8	41	1.5	1.2	T = 4	41	0.058
S = 5	37	1.5	1.2	T = 2	35	0.289
S = 4	31	1.5	1.2	T = 2	35	0.400
S = 3	29	1.5	1.2	T = 1	29	0.240
S = 2	28	1.5	1.2	T = 1	29	0.290
S = 12	50	2.0	1.2	T = 6	52	0.163
S = 9, 10, 11	42	2.0	1.2	T = 4	42	0.058
S = 6, 7, 8	41	2.0	1.2	T = 4	42	0.149
S = 5	37	2.0	1.2	T = 2	37	0.412
S = 4	31	2.0	1.2	T = 1	29	0.447
S = 3	29	2.0	1.2	T = 1	29	0.230
S = 2	28	2.0	1.2	T = 1	29	0.240

NOTE.—This table shows how the clock rate affects the transmission method, using equation (10), keeping the transmission rate β constant. For the SNP-threshold method, samples are clustered together where the SNP distance is less than or equal to S . For the transmission method, samples are clustered together where the implied number of transmissions k is less than or equal to T with a probability of 80%. For a clock rate of 0.5 SNPs/genome/year, the transmission method matches all SNP-based clusters for thresholds of 5 SNPs and above. As the clock rate increases, the transmission clustering diverges further from the SNP clustering. As we vary the other parameters, the choice of β is effectively a scale factor and does not affect the pattern of clustering. We use the variation of information dissimilarity measure given by *clue* (Meilä 2007) to compare the results of the two methods.

Identifying direct transmissions is not the aim of either the SNP cut-off or transmission clustering method; rather, both aim to simply group cases into sets of isolates for onward, more intensive (model specific, Bayesian for example) outbreak reconstructions. Testing the ability of SNP versus transmission-based methods to accomplish this using simulated data would require an appropriate simulation setup, which in turn would have a lot of flexibility (and could no doubt be tweaked to ensure that the transmission method performs well, or that the SNP cut-off does). For example, one approach would be to simulate the introduction of new cases whose SNP distance is 25 from existing cases in an existing outbreak. The SNP-threshold method with a threshold of 12 SNPs will always correctly place such new introductions in a new cluster, and will group their descending infections correctly until one or more of them is more than 12 SNPs away from other sampled cases in the cluster. Conversely, if new introductions were only 12 SNPs from existing cases the

Table 4. Comparison of the Methods with the Moldova data, Taking Drug Resistance into Account.

SNP Threshold	λ	β	Trans. Threshold	Dissimilarity SNP and Trans.	Dissimilarity SNP and Trans. with Resistance
S = 12	1.0	1.2	T = 11	0.000	0.132
S = 10	1.0	1.2	T = 9	0.000	0.164
S = 8	1.0	1.2	T = 7	0.000	0.128
S = 7	1.0	1.2	T = 6	0.000	0.245
S = 6	1.0	1.2	T = 5	0.000	0.090
S = 5	1.0	1.2	T = 4	0.000	0.047
S = 4	1.0	1.2	T = 3	0.000	0.023
S = 11	2.0	1.2	T = 4	0.000	0.233
S = 9	2.0	1.2	T = 3	0.000	0.128
S = 7	2.0	1.2	T = 2	0.000	0.308
S = 4	2.0	1.2	T = 1	0.000	0.023

NOTE.—This table shows how the clock rate affects the transmission method, using equation (10), keeping the transmission rate β constant, and the effect of including resistance using equation (18). For $\lambda = 1.0$ and 2.0 SNPs/genome/year, the pattern of clusters is identical using the SNP-threshold and transmission methods across a range of thresholds, but differs when resistance-conferring SNPs are taken into account. For the SNP-threshold method, samples are clustered together where the SNP distance is less than or equal to S . For the transmission method, samples are clustered together where the implied number of transmissions k is less than or equal to T with a probability of 80%. We use the variation of information dissimilarity measure given by *clue* (Meilä 2007) to compare the results of the methods.

SNP-threshold method would misclassify them as linked to existing clusters. In the transmission method, we can compute the probability that a newly introduced case that is 25 SNPs from existing cases will fall within a certain number of transmission events. This gives us the probability that we would infer an incorrect link to an existing cluster. With $\lambda = 1.2$ SNPs/genome/year and $\beta = 1.5$ transmissions/year, the probability that there are more than 10 transmissions for cases with 25 SNPs apart is 99.9%. This falls to 98.3% for more than 15 transmissions. Accordingly, the simulation approach for introducing new clusters will greatly affect the performance of both the SNP-threshold and transmission methods, and so we have not chosen to perform extensive simulations to compare the methods.

Discussion

We have demonstrated how our approach can be consistently applied in different contexts, with timing information, with spatial data in the case of British Columbia, and with resistance data in the case of Moldova. This is an advance on what is possible with the fixed SNP-threshold approach, where there is no general way to adjust thresholds to take this context-specific information into account.

A fixed number of SNPs can arise from different number of transmissions depending on other factors, including the timing of transmission, selection for resistance, the substitution process, location, and factors we have not explicitly modeled (social contacts, host risk factors, pathogen factors). We have seen that sampled cases which are relatively close in genetic distance can nevertheless be separated by large distances in time. In this scenario, a simple SNP cut-off may place samples too close together for outbreak clustering purposes. In contrast, our new method is robust with respect to outlying cases

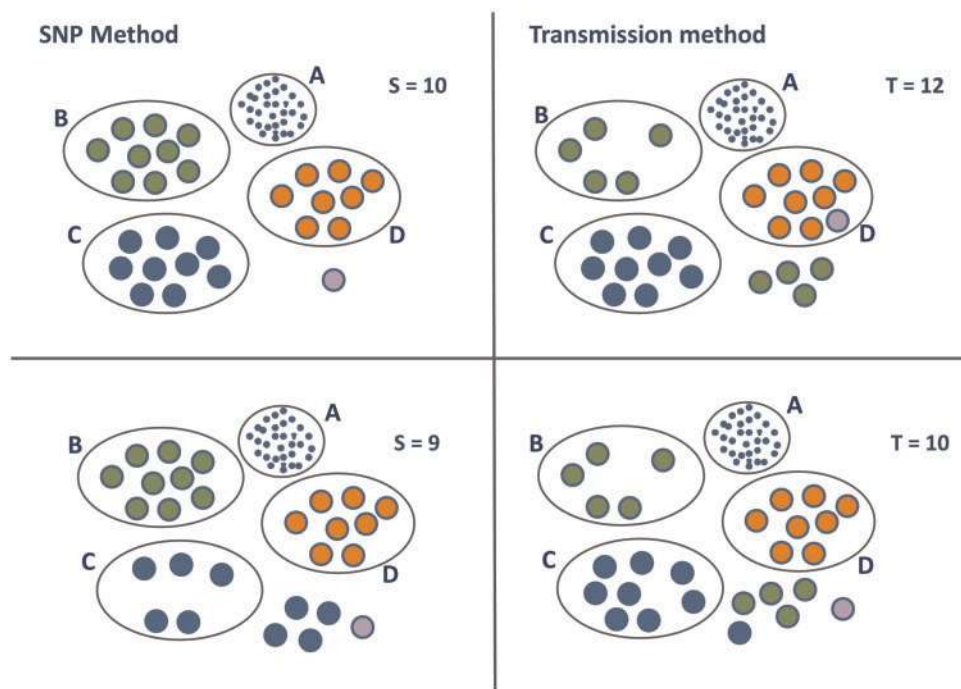


FIG. 4. Clusters in the Moldova data set illustrating the effect of accounting for resistance-conferring SNPs in the transmission method, using equation (18). Clusters B, C, and D are the second to fourth largest clusters in the Moldova data using the SNP-threshold method. The largest cluster A, with 93 members for $S = 10$ is shown for completeness. Isolated cases are shown with no enclosing oval. Colors are chosen to enable identification of the same cases in the four different scenarios. The left-hand panel shows the clusters obtained by applying the SNP-threshold method with two different thresholds, with the cut-off level denoted by S ; samples are clustered together where the SNP distance is less than or equal to S . The right-hand panel shows the clustering obtained by applying the transmission method, using equation (10), with the cut-off level denoted by T ; samples are clustered together where the implied number of transmissions k is less than or equal to T with a probability of 80%, with clock rate $\lambda = 1.5$ SNPs/genome/year and $\beta = 2.0$ transmissions/year.

which have been sampled at very different times compared with the majority of cases. These cases can make inference of timed phylogenetic trees challenging because the low genetic variation is hard to reconcile with the large time distance. Furthermore, true transmission clusters need not be clades in phylogenetic trees, because one cluster could descend from another but be separated by a long time or a large genetic distance (due to sampling effects). Accordingly, the clusters obtained by our method do not necessarily correspond to phylogenetic clades. We briefly discuss the application of our method to timed phylogenetic trees in the [supplementary data, Supplementary Material](#) online, with an example cluster which is not a clade.

Our probabilistic transmission method has certain advantages. It is relatively simple, requiring only the implementation of fast-running algorithms to estimate the time distributions; the heavy machinery to run large simulation methodologies (like MCMC) is not required. The amount of information required for the model is limited and consists of as little as the SNP distances, the timing data and a knowledge about the substitution and transmission processes. Nevertheless it has the flexibility to be able to handle SNPs under selection, SNPs with a different substitution process and variability in the substitution and transmission processes, and it has the scope for extensions to include more epidemiological data. Even in data sets where there is not much timing information to work with, we have seen that the integration of

information on resistance-conferring sites can be used within our framework to fine tune the clustering. Using two distinct processes—transmission, and the accumulation of measurable genetic variation—to define clusters carries the advantage that these processes may be estimable from data. This enables transmission clusters to be formed based on focused discussion and estimation of measurable processes rather than based on fixed cut-offs, and it allows ready adaptation for new pipelines that detect variation.

There are some limitations. Prior knowledge of the substitution and transmission processes is required, and there is some uncertainty in choosing appropriate values. However, the model is typically robust with respect to changes in these variables; in particular, varying the transmission rate does not have a material impact on the clustering because a rescaling of the cut-off will compensate. The choice of a time-varying transmission function $\beta(t)$ is, however, likely to have an impact on results. In particular we would expect a low probability of very quick transmission—as the pathogen numbers are building up in a new host—to have a significant impact, compared with the use of a constant transmission rate, as would a fast rate early diminishing to a much lower rate later. Note also that the parameter t in our model represents the total time since infection to both the sample dates: so we are not modeling the variation of transmission rates in calendar time.

In some diseases, such as TB, there is considerable variation in the latency period, during which the mutation rate may be

Table 5. Dissimilarity Measure Comparing Both the SNP-Threshold and Transmission Methods against Simulated Data, Averaged over the Full Set of Simulations.

Clock Rate	Transmission Method	SNP-Threshold Method	Random
0.5	0.250	0.366	0.871
1.0	0.261	0.366	0.845
1.5	0.265	0.366	0.795
2.0	0.282	0.366	0.842
2.5	0.293	0.366	0.845
3.0	0.303	0.366	0.848

NOTE.—We use the variation of information dissimilarity measure given by *clue* (Meilă 2007) to compare the results of the methods. Lower numbers indicate sets of clusters that are more similar to the true clusters. An outbreak was simulated 100 times with 10 sampled cases from a total of between 20 and 30 cases, depending on the simulation. The measure is obtained by comparing clusters from a range of thresholds to the known clusters, and picking the one with the lowest score. The averages are 0.27 bits for the transmission method, 0.366 for the SNP-threshold method, and 0.84 for the randomly permuted data. The clock rate is the rate used by the transmission method only to relate the number of SNPs to the time distribution, and thus does not affect the SNP-threshold method results—this is why the SNP-threshold method has the same dissimilarity compared with the simulated data whatever the clock rate. The table shows how the dissimilarity varies as the clock rate varies, for a fixed $\beta = 3.0$ transmissions/year, as compared with simulated samples connected by direct transmission. The random column shows the dissimilarity obtained for randomly allocated simulated clusters.

lower than it is during active disease. This variability can be incorporated into the negative binomial model as expressed in equation (14). We do not explicitly model within-host diversity, though this is relevant to identifying direct transmission events (Didelot et al. 2014, 2017; Worby et al. 2014; Hall et al. 2015, 2016). Cases of direct transmission will be clustered together with high probability in our method despite slight inaccuracy in the timing due to both branches of the pair's two-case tree spending time in the same host. Pairs of cases for which the clustering decision is ambiguous are likely to have several intermediate cases between them, with a larger tree height, and so the contribution of in-host diversity in either sampled case will be small. In-host diversity in unsampled cases would not affect our estimates unless it contributed to changes in the molecular clock rate.

WGS data has been noted to be helpful in ruling out transmission but insufficient, on its own, to resolve transmission events (Casali et al. 2016; Campbell et al. 2018). If the primary use of WGS data is only to refute transmission, one might ask why clustering matters. We would argue that the transmissions that are not refuted by WGS are then presumably considered to be possible recent, or direct, or clustered transmissions. Even if the primary use of WGS data is to refute direct transmission, there is a trade-off between the strength of that refutation and the possibility of mistakenly refuting genuine recent transmission events. This is more likely, using SNP cut-offs, where selection (say for antibiotic resistance) has led to higher SNP differences than expected. In addition, in practice WGS data are not only used to refute direct transmission, but to produce clusters that inform onward analyzes, reports on the extent of recent transmission, outbreak analysis and reconstruction and even public health policy; see (Guthrie et al. 2018) for one example.

We have accommodated the possibility of low substitution rates in latency with a non-Poisson model for the clock process, λ , in equation (5) (though we have not implemented this) and to some extent with the option of a nonconstant transmission rate. However, we have not modeled the possibility of a direct relationship between low SNP accumulation and low probability transmission. If this relationship exists—for example if latent cases both do not transmit and do not accumulate SNPs (Colangeli et al. 2014)—then low SNP differences could correspond to fewer intermediate hosts despite long elapsed times. This is an implicit assumption of a SNP-only method; although it may be correct it is a strong assumption, and there is evidence that mutation rates in latency are not reduced compared with active disease (Ford et al. 2011; Lillebaek et al. 2016).

We have not used the probability of sampling in forming our clusters, in contrast to other tools including the vimes package (Jombart and Cori 2017). For example, if it is known that surveillance is strong, then it would be less likely for 10 intermediate cases to be unsampled than for 5 intermediate cases, and this could be built in to a clustering method. Our rationale for not taking this into account is to provide a clustering approach that is as parallel as possible to the SNP cut-offs currently in widespread use while taking additional information on timing, molecular evolution, and transmission into account. It is often the case that the true sampling rate is not known and may change over time, and—particularly for TB in high-resource settings—cases can be missed because they are hard to identify (perhaps being at higher risk of TB due to homelessness or other factors, as in Casali et al. [2016]). In many settings the sampling probability may be uncertain. We have taken the approach of defining the clusters themselves without explicit reference to the sampling probability, with the view that the clusters are central inputs to other analyzes which will take sampling into account (as is done for example in TransPhylo [Didelot et al. 2017]). However, in our approach, changes in the sampling probability would likely be apparent in changes in the temporal and genetic distance between cases over time.

We have also not modeled changes in the transmission process over time in a community (e.g. due to depletion of susceptible individuals, improved infection control, etc.). As with including sampling, this may best be done in a more nuanced analysis after the initial clustering rather than as part of the clustering itself, but in principle, changes to the transmission function over calendar time could be incorporated into the mathematics behind equation (8). However, this would raise interpretation challenges because of the fact that our transmission process reflects the rate of the pathogen moving between hosts where it is known that there is an infected host at the “end” of the chain (since each pair consists of two sampled hosts, whose pathogen was sequenced and who were therefore certainly infected). We do not model the number of contacts over which transmission could have occurred.

The choice of a particular SNP cut-off also takes no account of the inevitable uncertainties involved in the gathering

and processing of raw read data, and does not allow for the modeling of this uncertainty. Different bioinformatics pipelines—and different parameters used within those pipelines—can have a substantial effect on the number of SNP differences reported between cases. It is usual for SNP differences to be taken as given and, although sometimes details are provided—see for example [Katz et al. \(2013\)](#)—it is important to recognize that there can be considerable variation between SNPs reported using different pipelines and parameters. For example, the level of quality scores and read depth cut-offs used will generally have a high impact, as will the precise way in which hypervariable sites and repeat regions are handled (or excluded). As technology improves we may begin to capture variation in repeat regions, or types of variation (e.g. insertions/deletions) that are currently masked, and in that new pipeline 12 SNPs may not carry the interpretation it does today. The model could easily incorporate more genomic information, resulting in a more sophisticated version of the distance function. In particular, large-scale genomic features can readily help to establish that cases belong to separate and therefore distantly related lineages. As variation-calling pipelines evolve, our method could be used to relate each pipeline to number of transmissions or to estimated divergence time; this would form an approach to compare bioinformatics pipelines and data sources, and to curate their use in defining distances between isolates.

TB has distinct phylogeographic lineages which have been reported to have different mutation rates, with lineage 2 (the East Asian and Beijing lineage) having higher mutation rates than lineage 4 (Euro-American) ([Ford et al. 2013](#)). Our approach could unify clustering despite such differences, as the same transmission and probability settings could be used under different SNP accumulation rates. This would provide a consistent approach to clustering in areas where multiple lineages cocirculate, and allow comparison of TB clustering patterns in different settings. The same would be true for adapting to differing natural histories across different pathogen lineages or subpopulations: the choice of β could reflect transmission differences while the other settings remained the same.

The long-term aim of changing how cases are assigned to clusters is to improve the way that WGS and epidemiological data are used and to best capture clusters that correspond to transmissions of an infectious disease. We have found that basing clusters on the number of transmission events, with a probabilistic cut-off, is feasible, can integrate timing and other data, and compares favorably to clustering based on SNP cut-offs.

Materials and Methods

Data

In this article we focus on TB, but our approach is applicable to other pathogens for which WGS can be carried out and where it is appropriate to use SNPs to compare closely related isolates (naturally, parameters will vary). TB provides a convenient model as it avoids the complications associated with HGT, it is an important pathogen worldwide, it has very

diverse epidemiological settings and WGS tools are increasingly used for public health purposes.

British Columbia

The British Columbia Centre for Disease Control (BCCDC)'s Public Health Laboratory (BCPHL) receives all *Mtb* cultures for the province and performs routine MIRU-VNTR genotyping on all *Mtb* isolates. *Mtb* isolates belonging to MIRU-VNTR cluster MClust-012 were revived from archived stocks, DNA extracted, and sequenced using 125 bp paired-end reads on the Illumina HiSeqX platform at the Michael Smith Genome Sciences Centre (Vancouver, British Columbia). The resulting fastq files were analyzed using a pipeline developed by Oxford University and Public Health England. Reads were aligned to the *Mtb* H37Rv reference genome (GenBank ID: NC000962.2), with an average of 92% of the reference genome covered. Single-nucleotide variants (SNVs) were identified across all mapped nonrepetitive sites. Fastq files for all genomes are available at NCBI under BioProject PRJNA413593.

Republic of Moldova

Sample Collection and Epidemiological Data. The study population included patients diagnosed with culture-positive TB at the municipal hospital from October 2013 to December 2014 in the Republic of Moldova. All epidemiological and laboratory data from TB patients are routinely entered into a country-wide web-based TB electronic medical record (EMR) database. Epidemiological data including age, sex, previous TB history, results of chest radiograph, history of incarceration, and place of residence were collected. Laboratory data, including mycobacterial smear grade, culture, and drug-susceptibility testing to first and second line anti-TB agents, were extracted from the EMR. As part of this study, all *M. tuberculosis* patient isolates were subcultured and frozen for genomic analysis.

Variant Calling and Phylogenetic Analysis. DNA was extracted from *M. tuberculosis* grown on Lowenstein-Jensen slants as described previously. Paired-end (250 bp) sequences were generated on the Illumina MiSeq platform. Raw fastq reads were filtered for length and trimmed for low-quality trailing base pairs using Trim Galore, aligned to the H37Rv NC000962.3 reference genome using BWA, with duplicate reads removed using PicardTools. The mpileup function in samtools was used for single-isolate variant calling. Isolates with a high proportion of apparent mixed or heterozygous SNP calls (i.e. those with >25% reads supporting the reference allele) were excluded from analysis. SNPs within 15 bp of insertions or deletions (indels) or with variant quality scores < 100 were excluded. SNPs in or within 50 bp of hypervariable PPE/PE gene families, repeat regions, and mobile elements were excluded ([Eldholm et al. 2015](#)). A phylogenetic tree was constructed in RAxML (GTR-gamma for nucleotide substitution and correcting for SNP ascertainment bias) and annotated with DST results and drug-resistance-associated variants from Mykrobe Predictor ([Bradley et al. 2015](#)). Representative strains from other studies in the region,

including L4 (LAM, Haarlem, Ural) and L2 (Casali et al. 2014; Merker et al. 2015), were also included. Percy256 (Lineage 7) was included as an outgroup. Fastq files for all genomes are available at NCBI under Accession number SRP156366.

Clustering Approach

The overall approach is to use the SNPs and case timing to derive a distribution for the time to the MRCA of each pair of samples, condition on that time to write the probability that the samples were separated by some number of transmission events, and then integrate out the unknown time to the pairs' MRCA. The first step makes use of the molecular clock process and depends on the clock rate and on the number of SNPs under a form of selection (like antibiotic resistance). The second step using information about the transmission process and the natural history of the pathogen.

For each sample, we start with the date on which the sample was taken and the aligned nucleotide sequence for the set of variable sites in our set of samples. For any two samples S_1 and S_2 , we have the SNP distance $N = N(S_1, S_2)$ which is equal to the Hamming distance between their respective nucleotide sequences. We also have the sampling time difference $\delta = \delta(S_1, S_2)$. Without loss of generality, we can assume that S_1 is sampled either at the same time as, or before, S_2 . What we do not know a priori, and therefore we have to estimate, is the total amount of time h over which the SNPs have accumulated (on both branches in total) since the date of the MCRA of S_1 and S_2 . We also refer to h as the "height."

Given the time h , we can use a transmission process to estimate the probability that there are more than some threshold number of transmission events T in a total time h ; we integrate over the unknown h . This transmission process need not be homogeneous (table 6).

We make various assumptions in setting up the model. Both substitutions and transmissions occur according to (possibly nonhomogeneous) Poisson processes over time. Unless it is otherwise stated, the population from which the samples are drawn is homogeneous, so transmission is random and equally likely between hosts irrespective of factors such as location of abode, individual lifestyle, etc. We do not assume that all infected cases are reported and sequenced. However, where we do have sequence data, we assume that it is correct and complete. Reported cases may be sampled more than once. We do not explicitly model reporting and sampling rates. If these change over time, then this would be reflected in the time and genetic distance between nearby cases and consequently in the estimated number of intermediate transmissions between reported cases. Once infected, we assume that a patient becomes infectious immediately, either with a constant probability of infection per unit of time, or in a process yielding a gamma-distributed time to the next infection. This "natural history" model is assumed not to change with calendar time, such that the course of infectiousness proceeds in the same manner from infection to infecting others independent of the calendar time of infection. Our approach is intended to group

Table 6. Symbols Used in the Model and Their Meaning.

Symbol	Meaning	Units
λ	Clock rate	SNPs/genome/year
β	Transmission rate	Transmissions/year
h	Total time over which SNPs occur over both lines of descent before the first sampling date	Years
δ	Time between two sample dates	Years
N	No. of SNP differences between two cases	
N_δ	No. of SNP differences occurring after first sample	
k	Number of transmissions	
S	Cut-off threshold for SNP-threshold method	
T	Cut-off threshold for transmission method	

sequences into clusters, and does not model reported cases for whom there is no sequence data.

Noting that δ is fixed by the sampling times in the data, we estimate the distribution of the time $h/2$ over which the SNPs have had to accumulate before the sample date of S_1 and S_2 . This is equivalent to estimating the date of the MRCA of S_1 and S_2 . Because both branches are free to evolve over this time, $h/2 + h/2 = h$ is the effective overall time between the MRCA and S_1 , and $\delta + h$ is therefore the total evolutionary time separating the two cases.

Estimate of the Height Where Sample Dates Are the Same

The simplest model for the number of SNPs per unit time is a Poisson process with a constant rate λ ; we can also accommodate overdispersion, reflecting a more variable SNP accumulation process suitable for pathogens whose substitutions are not as clock-like (see below). The standard Poisson distribution with parameter λh gives the probability density of the number of SNPs on a given time interval h :

$$P(N|h) = \frac{e^{-\lambda h} (\lambda h)^N}{N!} \quad (1)$$

However, we are interested in the likelihood of the time h as a function of the specified number of SNPs.

We know by standard theory that the arrival time density—that is, the time density until the next SNP—can be modeled by the exponential density function $\lambda e^{-\lambda h}$. Furthermore, the waiting time until the N th SNP is also a Poisson process, as the arrivals are assumed to be independent and identically Poisson distributed. It can be shown (for example in Chapter 2 of [Gallagher 2013]) by repeated convolution of densities that the distribution of the N th arrival time A_N is given by

$$\text{PDF}_{A_N} = \frac{\lambda^N h^{N-1} e^{-\lambda h}}{(N-1)!} \quad (2)$$

for $N > 0$. This is the Erlang distribution, with mean $= N/\lambda$, as expected.

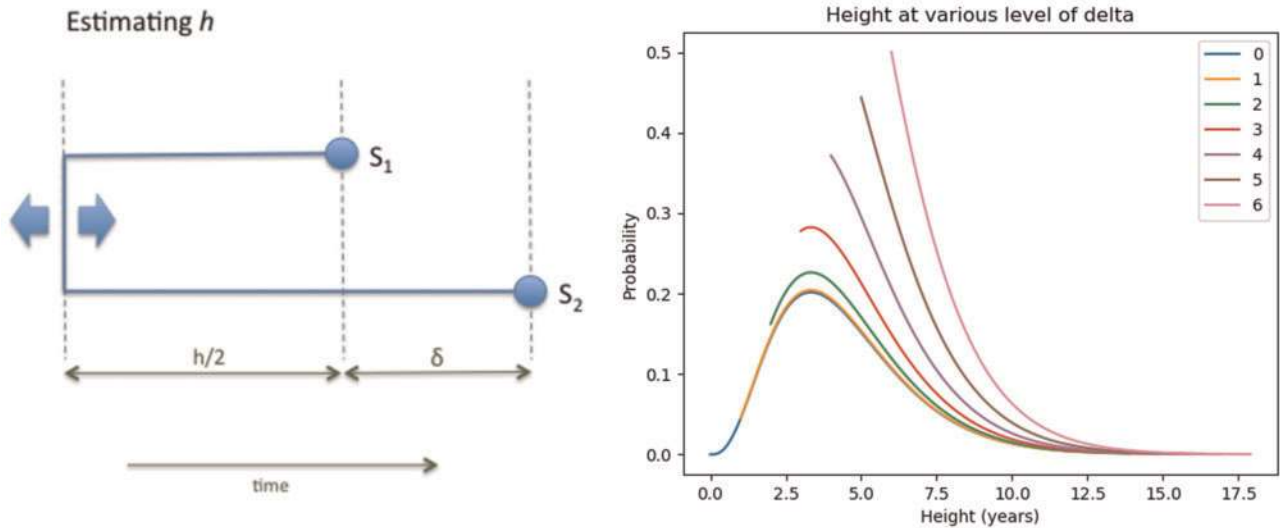


FIG. 5. On the left is a schematic illustration of the notation. h is the total time in years over which SNPs accumulate between two cases before the first sample is taken, whereas the total time over which SNPs can occur is $h + \delta$ years. δ is known and fixed; h is unknown. On the right is a plot of $h + \delta$, where h is given by equation (4), for values of δ ranging from 0 through 6 years, with $N = 3$, $\lambda = 0.9$ SNPs/genome/year, and $\beta = 1.2$ transmissions/year. Since $h + \delta > \delta$, the lines corresponding to higher values of δ begin above 0.

We know that exactly N SNPs have already occurred on a time interval of uncertain length h , and we are interested in the likelihood of h given the data N . Since we already have N SNPs and are waiting for the $(N + 1)$ th, this is given by the arrival time density for the $(N + 1)$ th SNP; by replacing N with $N + 1$ in the above and interpreting it as a function of h , we have:

$$\mathcal{L}(h|N) = \frac{\lambda^{N+1} h^N e^{-\lambda h}}{N!} \quad (3)$$

Note that when $N = 0$, this reduces to $\lambda e^{-\lambda h}$.

Alternatively, we can generalize the arrival time density to a gamma distribution, where the extra parameter allows us to fix the mean but change the variance. This allows us to be more flexible with respect to dispersion than with using the exponential distribution. The gamma density, with two parameters a and b , is

$$f(h; a, b) = \frac{b^a h^{a-1} e^{-bh}}{\Gamma(a)} \quad (4)$$

The mean is a/b and the variance a/b^2 . Note that we can recover the Poisson model result by setting $a = 1$ and $b = \lambda$ (Cameron and Trivedi 2013). In this case, the arrival time density for the $(N + 1)$ th SNP is given by

$$f(h; a(N + 1), b) = \frac{b^{a(N+1)} h^{a(N+1)-1} e^{-bh}}{\Gamma(a(N + 1))} \quad (5)$$

by standard properties of the gamma distribution.

Estimate of the Height Where Sampling Times Differ

In this case, we account for the fact that some of the SNPs may have occurred in the fixed time interval of length δ between the two sample dates. Again, we begin with the

simple model in which the number of SNPs occurring in this time is given by a Poisson distribution, in this case with parameter $\lambda\delta$. We write $N = N_h + N_\delta$, where N_δ is Poisson distributed with parameter $\lambda\delta$.

The number of SNPs N_δ accumulated on the fixed interval of length δ is somewhere between 0 and N inclusive; $0 \leq N_\delta \leq N$. Unconstrained, N_δ is Poisson distributed with parameter $\lambda\delta$. Conditioning on the probability that N_δ does not exceed N gives us the probability density

$$P(N_\delta|\delta) = \left(\frac{e^{-\lambda\delta} (\lambda\delta)^{N_\delta}}{N_\delta!} \right) / \sum_{i=0}^N e^{-\lambda\delta} (\lambda\delta)^i / i!$$

and writing $F(N_\delta)$ for $\sum_{i=0}^N e^{-\lambda\delta} (\lambda\delta)^i / i!$

$$P(N_\delta|\delta) = \left(\frac{e^{-\lambda\delta} (\lambda\delta)^{N_\delta}}{N_\delta!} \right) / F(N_\delta) \quad (6)$$

To obtain the expression for $\mathcal{L}(h|N, \delta)$, we sum over all the possible values of N_δ , giving

$$\mathcal{L}(h|N, \delta) = \sum_{i=0}^N \mathcal{L}(h|i, \delta) P((N - i)|\delta)$$

Substituting into our earlier expression,

$$\mathcal{L}(h|N, \delta) = \sum_{i=0}^N \left(\frac{\lambda^{i+1} h^i e^{-\lambda h}}{i!} \right) \left(\frac{e^{-\lambda\delta} (\lambda\delta)^{(N-i)}}{(N-i)!} \right) / F(N)$$

$$\mathcal{L}(h|N, \delta) = \frac{e^{-\lambda(h+\delta)} \lambda^{N+1}}{F(N)} \sum_{i=0}^N \frac{h^i \delta^{N-i}}{i!(N-i)!} \quad (7)$$

An example plot for the equation above is shown in figure 5.

Modeling Transmissions

We connect SNP distances to transmissions using a model for the number of transmissions likely to have occurred over a given total time period, *conditional* on the two cases being infected at or before the sampling times. This means that unlike a transmission rate in a population-level epidemic model, which typically describes the rate of transmission per unit time given contact between a susceptible and an infectious individual, our transmission process is better described in terms of the rate at which a pathogen lineage will jump to a new host. This is, of course, distinct from the rate at which new transmissions occur in a community and the per-contact rate of transmission of infection between two individuals. We first assume for simplicity that β is a constant function, and that it is a Poisson process; we allow a more general model later. The amount of time over which transmissions can occur between our two cases is $h + \delta$, and the expected number of transmissions is $\beta(h + \delta)$. The number of transmission events k is therefore given by

$$P(k|h, \delta) = \frac{\beta^k (h + \delta)^k e^{-\beta(h+\delta)}}{k!} \tag{8}$$

Integrating over h , we have

$$P(k|N, \delta) = \int_{h=0}^{\infty} \mathcal{L}(h|N, \delta) P(k|h) dh \tag{9}$$

$$= \frac{e^{-\delta(\lambda+\beta)} \lambda^{N+1} \beta^k}{k! F(N)} \int_{h=0}^{\infty} e^{-h(\lambda+\beta)} (h + \delta)^k \sum_{i=0}^N h^i \left(\frac{\delta^{N-i}}{i!(N-i)!} \right) dh \tag{10}$$

This equation expresses the key relationship that allows us to translate raw SNP differences and sample time differences into transmission probability distributions—examples are shown in figure 6. As the sample time between cases increases, it can be seen that this factor makes an increasingly important contribution, relative to the SNP distance, to the distance between cases.

Unless stated otherwise, equation (10) is used to generate the data presented in the Results Section.

Time Varying Transmissions

In our context, a transmission event should be understood as an event in which a pathogen is transferred to a new host, ultimately causing a secondary case in that host. Although there may be undetected transmission events in which the secondary cases never develop disease, our data are on sampled cases with active disease, and the time between successive transmissions should approximately reflect the serial interval between cases with active disease. We allow the number of transmissions $\beta = \beta(t)$ to be a function of time since infection, allowing for a variable risk of infecting others during the course of infection. Once a host is infected, the details of the

natural history of the pathogen affect the generation time—the exact form of the function $\beta(t)$ allows us to incorporate the varying rates of progression from infection to active disease, and then on to transmission. This illustrates that our framework has the flexibility to include more detailed and accurate modeling of the underlying disease dynamics. As stated in Didelot et al. (2017), the generation time distribution can take any form (Fine 2003; Wallinga and Lipsitch 2007), but gamma distributions are often used, as for example in Conlan et al. (2010). We apply the gamma distribution with parameters shape α and scale θ , so that:

$$\beta(t) = \beta(t; \alpha, \theta) = \frac{1}{\Gamma(\alpha)\theta^\alpha} t^{\alpha-1} e^{-t/\theta} \tag{11}$$

and the mean value is $\alpha\theta * dt$ in a given time interval dt .

Putting this together with our Poisson model for the number of SNPs on a time interval (eq. 1, we obtain:

$$P(N|\lambda, \beta(t; \alpha, \theta)) = \int_{t=0}^{\infty} \frac{e^{-\lambda t} (\lambda t)^N}{N!} \left(\frac{1}{\Gamma(\alpha)\theta^\alpha} t^{\alpha-1} e^{-t/\theta} \right) dt \tag{12}$$

$$= \binom{N + \alpha - 1}{N} \left(1 - \frac{\theta\lambda}{\theta\lambda + 1} \right)^\alpha \left(\frac{\theta\lambda}{\theta\lambda + 1} \right)^N \tag{13}$$

This is a negative binomial (denoted NB) distribution for the number of SNPs for one transmission generation,

$$P(N|\lambda, \beta(t; \alpha, \theta)) \sim \text{NB} \left(\alpha; \frac{\theta\lambda}{\theta\lambda + 1} \right)$$

Assuming transmission events are independent of each other, it then follows (by standard properties of the negative binomial) that the probability of N given k transmissions is also distributed as a negative binomial, with

$$P(N|k, \lambda, \beta(t; \alpha, \theta)) \sim \text{NB} \left(k\alpha; \frac{\theta\lambda}{\theta\lambda + 1} \right)$$

$$P(N|k, \lambda, \beta(t; \alpha, \theta)) = \binom{N + k\alpha - 1}{N} \left(1 - \frac{\theta\lambda}{\theta\lambda + 1} \right)^{k\alpha} \left(\frac{\theta\lambda}{\theta\lambda + 1} \right)^N \tag{14}$$

Modeling Resistance-Conferring SNPs

Suppose that we know that there are resistance-conferring SNPs in our sample population, or perhaps other SNPs at sites known to be under selection or simply to have a different rate of substitution. Let us assume they account for a certain fixed proportion of the observed SNP differences. Given N SNPs, assume that m are not resistance conferring and n are, so $N = m + n$. Their respective mutation rates are given by λ_m and λ_n , where $\lambda_n > \lambda_m$. Assuming independence, on a given time interval of length h we have

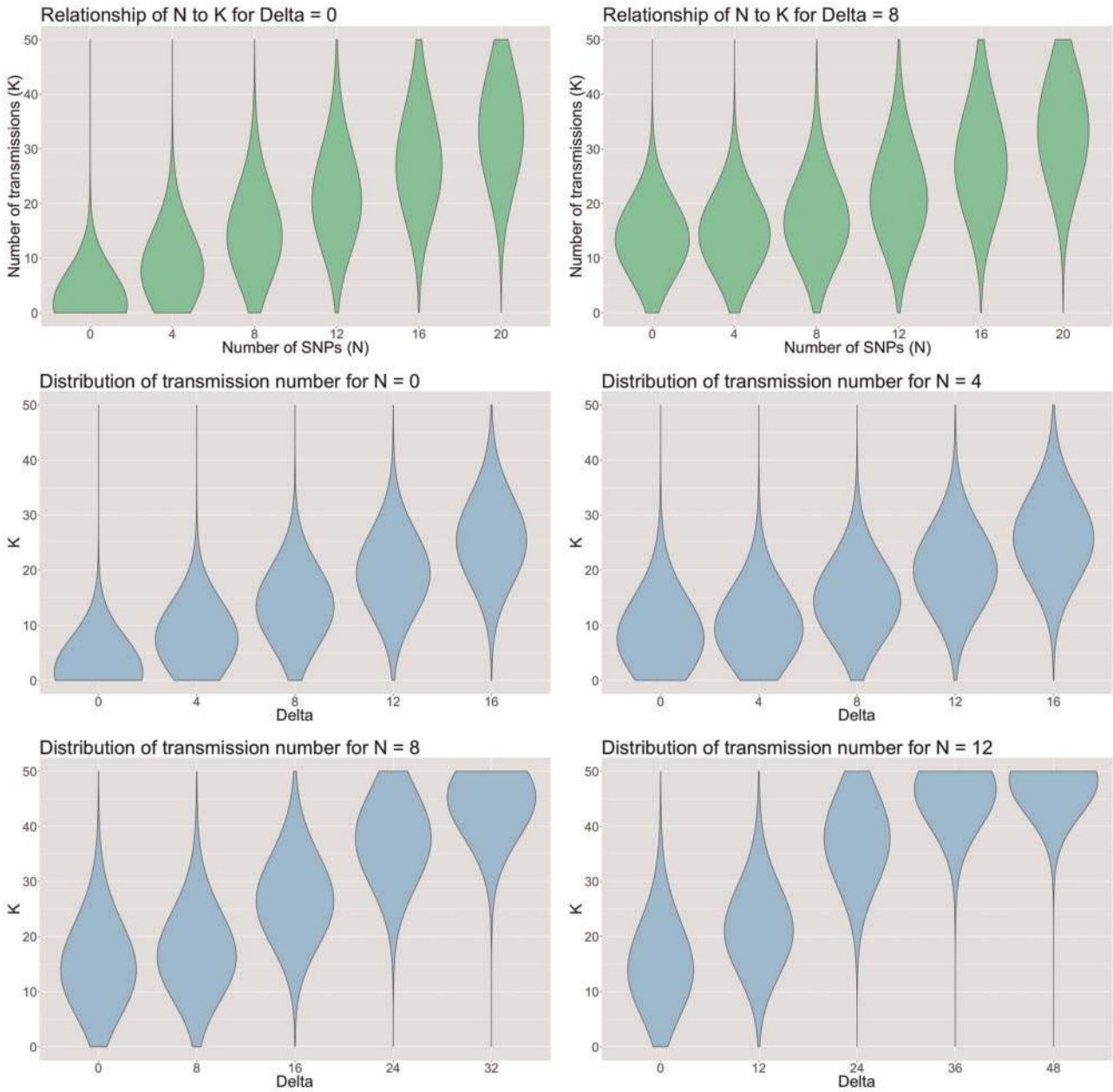


FIG. 6. Probability density for the number of transmissions, given by equation (10), with clock rate $\lambda = 0.9$ SNPs/genome/year and $\beta = 1.5$ transmissions/year. The upper two panels show the densities for delta values of 0 and 4 years for a range of SNP distance between 0 and 20. The lower four panels show the densities for 0, 4, 8 and 12 SNP distance respectively, for $\delta = 0, 4, 8, 12, 16$ years.

$$\begin{aligned}
 P(N|h) &= P(m|h, \lambda_m)P(n|h, \lambda_n) \\
 &= \left(\frac{\lambda_m^m h^m e^{-\lambda_m h}}{m!} \right) \left(\frac{\lambda_n^n h^n e^{-\lambda_n h}}{n!} \right) \\
 &= \left(\frac{\lambda_m^m \lambda_n^n}{m!n!} \right) h^m e^{-\lambda_m h} h^n e^{-\lambda_n h} \\
 &= \Lambda_{mn} h^{m+n} e^{-(\lambda_m + \lambda_n)h}
 \end{aligned}
 \tag{15}$$

$$\Lambda_{mn} = \frac{\lambda_m^m \lambda_n^n}{m!n!}
 \tag{16}$$

Compare this to equation (1), which we can recover by setting $m = N, n = 0, \lambda = \lambda_m,$ and $\lambda_n = 0.$

We have a Poisson process which is the sum of two independent Poisson processes with $\lambda = \lambda_m + \lambda_n.$ As before, we can derive expressions for $\mathcal{L}(h|N = n + m)$ and $P(k|N = n + m),$ so

$$\mathcal{L}(h|m, n) = \Lambda'_{mn} h^{m+n} e^{-(\lambda_m + \lambda_n)h}
 \tag{17}$$

where

where

$$\Lambda'_{mn} = \frac{\lambda_m^{m+1} \lambda_n^{n+1}}{m!n!}$$

A way to illustrate the effect of including resistance-conferring SNPs is to consider the expected value of h . Recall that under equation (1), the mean is given by N/λ . Thinking of our resistance and nonresistance-conferring SNP processes, they have means respectively of n/λ_n and m/λ_m . Thus the combined process has mean $n/\lambda_n + m/\lambda_m$, and we can write the rate parameter λ^* of the combined process as

$$\lambda^* = \frac{n + m}{n/\lambda_n + m/\lambda_m} \quad (18)$$

Note that for large λ_m , λ^* tends to $\lambda_n * (n + m)/n$. The larger the value of λ_m as compared with λ_n , the smaller the contribution that the resistance-conferring SNPs make to the value of h —accordingly, four SNPs likely to have arisen due to inappropriate treatment or another selection process should not contribute as strongly toward separating two cases into different transmission clusters as four “neutral” SNPs. Ideally, the value of λ_m should be estimated from data. Once resistance SNPs have occurred in an individual, they are likely to be transmitted onwards when the individual infects others. These secondary cases share the resistance SNPs with each other ($n = 0$ in these pairs) and they are likely to be placed in the same cluster. Between each secondary case and the infecting case, $n > 0$; our method allows the resistance SNPs to “count for” less time than other SNPs, and the index case is likely clustered with the onward cases.

Spatial Proximity and Other Individual Data

Other factors that affect the likelihood of transmission, such as spatial proximity or other covariate data including contact tracing, demographics or other host factors, can be built into the model.

To incorporate spatial proximity, we assign each of the cases into one of a number of regions R_i where i is the region index. For the British Columbia data set, there are six regions defined, as shown in figure 3. For any pair of cases, a probability weighting is assigned which is equal to 1 in the case that both cases belong to the same region, and a value below 1 for cases which belong to different regions. This weighting w is then applied to the probability of obtaining k transmissions given N SNPs, giving us a modified version of equation (10)

$$P(k|N, \delta, w) = w \int_{h=0}^{\infty} \mathcal{L}(h|N, \delta) P(k|h) dh \quad (19)$$

Simulations

We generate simulated outbreaks and compare the SNP-threshold and transmission methods on them with a technique that measures the similarity of clusters using an information-theoretic approach (Meilă 2007). Outbreaks are simulated using *TransPhylo* (Didelot et al. 2017), which generates a dated transmission network for each simulation, containing both sampled and unsampled cases. From these,

and for all the cases, phylogenetic trees are extracted using *phyloTop* (Kendall et al. 2016). Sequences are then generated with *phangorn* (Schliep 2011) and output as *fasta* format files. For the sampled, and therefore “known,” cases we generate sets of clusters using the SNP-threshold and transmission methods for a range of cut-off levels. We also generate the “true” clustering of the sampled cases implied by the simulated *TransPhylo* transmission networks.

Software Availability

The methods presented here are available as R functions in the *transcluster* package, available at <https://github.com/JamesStimson/transcluster>.

Supplementary Material

Supplementary data are available at *Molecular Biology and Evolution* online.

Acknowledgments

C.C. and J.S. are supported by The Engineering and Physical Sciences Research Council (EPSRC) grant EP/K026003/1 and C.C. is additionally supported by EPSRC grant EP/N014529/1. T.C. is supported in part by U54GM088558 from the National Institute of General Medical Sciences. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institute of General Medical Sciences or the National Institutes of Health. We would like to thank Tyler S. Brown for contributions to the bioinformatics and analysis of the data from Moldova.

References

- Azarian T, Maraqa NF, Cook RL, Johnson JA, Bailey C, Wheeler S, Nolan D, Rathore MH, Morris JG Jr, Salemi M. 2016. Genomic epidemiology of methicillin-resistant *Staphylococcus aureus* in a neonatal intensive care unit. *PLoS One* 11(10):e0164397.
- Barrick JE, Lenski RE. 2013. Genome dynamics during experimental evolution. *Nat Rev Genet.* 14(12):827.
- Bergholz TM, Switt AIM, Wiedmann M. 2014. Omics approaches in food safety: fulfilling the promise? *Trends Microbiol.* 22(5):275–281.
- Bradley P, Gordon NC, Walker TM, Dunn L, Heys S, Huang B, Earle S, Pankhurst LJ, Anson L, De Cesare M, et al. 2015. Rapid antibiotic-resistance predictions from genome sequence data for *Staphylococcus aureus* and *Mycobacterium tuberculosis*. *Nat Commun.* 6:10063.
- Bryant JM, Harris SR, Parkhill J, Dawson R, Diacon AH, van Helden P, Pym A, Mahayiddin AA, Chuchottaworn C, Sanne IM, et al. 2013. Whole-genome sequencing to establish relapse or re-infection with *Mycobacterium tuberculosis*: a retrospective observational study. *Lancet Respir Med.* 1(10):786–792.
- Bryant JM, Schürch AC, van Deutekom H, Harris SR, de Beer JL, de Jager V, Kremer K, van Hijum SA, Siezen RJ, Borgdorff M, et al. 2013. Inferring patient to patient transmission of *Mycobacterium tuberculosis* from whole genome sequencing data. *BMC Microbiol.* 13(1):110.
- Cameron AC, Trivedi PK. (2013). Regression analysis of count data. Vol. 53. Cambridge: Cambridge University Press.
- Campbell F, Strang C, Ferguson N, Cori A, Jombart T. 2018. When are pathogen genome sequences informative of transmission events? *PLoS Pathog.* 14(2):e1006885.
- Casali N, Broda A, Harris SR, Parkhill J, Brown T, Drobniewski F. 2016. Whole genome sequence analysis of a large isoniazid-resistant tuberculosis outbreak in London: a retrospective observational study. *PLoS Med.* 13(10):e1002137.

- Casali N, Nikolayevskyy V, Balabanova Y, Harris SR, Ignatyeva O, Kontsevaya I, Corander J, Bryant J, Parkhill J, Nejentsev S, et al. 2014. Evolution and transmission of drug-resistant tuberculosis in a Russian population. *Nat Genet.* 46(3):279.
- Clark TG, Mallard K, Coll F, Preston M, Assefa S, Harris D, Ogbwang S, Mumbowa F, Kirenga B, O'Sullivan DM, et al. 2013. Elucidating emergence and transmission of multidrug-resistant tuberculosis in treatment experienced patients by whole genome sequencing. *PLoS One* 8(12):e83012.
- Colangeli R, Arcus VL, Cursons RT, Ruthe A, Karalus N, Coley K, Manning SD, Kim S, Marchiano E, Alland D. 2014. Whole genome sequencing of *Mycobacterium tuberculosis* reveals slow growth and low mutation rates during latent infections in humans. *PLoS One* 9(3):e91024.
- Conlan AJ, Rohani P, Lloyd AL, Keeling M, Grenfell BT. 2010. Resolving the impact of waiting time distributions on the persistence of measles. *J R Soc Interface* 7(45):623–640.
- Dallman TJ, Ashton PM, Byrne L, Perry NT, Petrovska L, Ellis R, Allison L, Hanson M, Holmes A, Gunn GJ, et al. 2015. Applying phylogenomics to understand the emergence of Shiga-toxin-producing *Escherichia coli* O157: H7 strains causing severe human disease in the UK. *Microb Genom* 1(3):e000029.
- Didelot X, Fraser C, Gardy J, Colijn C. 2017. Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks. *Mol Biol Evol.* 34(4):997–1007.
- Didelot X, Gardy J, Colijn C. 2014. Bayesian inference of infectious disease transmission from whole-genome sequence data. *Mol Biol Evol.* 31(7):1869–1879.
- Donker T, Bosch T, Ypma R, Haenen A, van Ballegooijen W, Heck M, Schouls L, Wallinga J, Grundmann H. 2016. Monitoring the spread of methicillin-resistant *Staphylococcus aureus* in The Netherlands from a reference laboratory perspective. *J Hosp Infect* 93(4):366–374.
- Eldholm V, Monteserin J, Rieux A, Lopez B, Sobkowiak B, Ritacco V, Balloux F. 2015. Four decades of transmission of a multidrug-resistant *Mycobacterium tuberculosis* outbreak strain. *Nat Commun.* 6:7119.
- Eldholm V, Norheim G, von der Lippe B, Kinander W, Dahle UR, Caugant DA, Mannsåker T, Mengshoel AT, Dyrhol-Riise AM, Balloux F. 2014. Evolution of extensively drug-resistant *Mycobacterium tuberculosis* from a susceptible ancestor in a single patient. *Genome Biol.* 15(11):490.
- Feuerriegel S, Schleusener V, Beckert P, Kohl TA, Miotto P, Cirillo DM, Cabibbe AM, Niemann S, Fellenberg K. 2015. PhyResSE: a web tool delineating *Mycobacterium tuberculosis* antibiotic resistance and lineage from whole-genome sequencing data. *J Clin Microbiol.* 53(6):1908–1914.
- Fine PE. 2003. The interval between successive cases of an infectious disease. *Am J Epidemiol.* 158(11):1039–1047.
- Ford CB, Lin PL, Chase MR, Shah RR, Iartchouk O, Galagan J, Mohaideen N, Ioeberger TR, Sacchetti JC, Lipsitch M, et al. 2011. Use of whole genome sequencing to estimate the mutation rate of *Mycobacterium tuberculosis* during latent infection. *Nat Genet.* 43(5):482–486.
- Ford CB, Shah RR, Maeda MK, Gagneux S, Murray MB, Cohen T, Johnston JC, Gardy J, Lipsitch M, Fortune SM. 2013. *Mycobacterium tuberculosis* mutation rate estimates from different lineages predict substantial differences in the emergence of drug-resistant tuberculosis. *Nat Genet.* 45(7):784–790.
- Gallagher RG. (2013). Stochastic processes: theory for applications. Cambridge (United Kingdom): Cambridge University Press.
- Guerra-Assunção JA, Houben R, Crampin AC, Mzembe T, Mallard K, Coll F, Khan P, Banda L, Chiwaya A, Pereira R, et al. 2015. Relapse or reinfection with tuberculosis: a whole genome sequencing approach in a large population-based cohort with high HIV prevalence and active follow-up. *J Infect Dis.* 211(7):1154–1163.
- Guthrie JL, Delli Pizzi A, Roth D, Kong C, Jorgensen D, Rodrigues M, Tang P, Cook VJ, Johnston J, Gardy JL. 2018. Genotyping and whole-genome sequencing to identify tuberculosis transmission to pediatric patients in British Columbia, Canada, 2005–2014. *J Infect Dis.* 40:1–9.
- Hall M, Woolhouse M, Rambaut A. 2015. Epidemic reconstruction in a phylogenetics framework: transmission trees as partitions of the node set. *PLoS Comput Biol.* 11(12):e1004613.
- Hall MD, Woolhouse MEJ, Rambaut A. 2016. Using genomics data to reconstruct transmission trees during disease outbreaks. *Rev Sci Tech.* 35(1):287–296.
- Hatherell H-A, Colijn C, Stagg HR, Jackson C, Winter JR, Abubakar I. 2016. Interpreting whole genome sequencing for investigating tuberculosis transmission: a systematic review. *BMC Med.* 14(1):21.
- Jombart T, Cori A. (2017). vimes. <http://www.repidemicsconsortium.org/vimes/articles/distributions.html>; last accessed May 5, 2018.
- Kammerer JS, Shang N, Althomsons SP, Haddad MB, Grant J, Navin TR. 2013. Using statistical methods and genotyping to detect tuberculosis outbreaks. *Int J Health Geogr.* 12(1):15.
- Katz LS, Petkau A, Beaulaurier J, Tyler S, Antonova ES, Turnsek MA, Guo Y, Wang S, Paxinos EE, Orata F, et al. 2013. Evolutionary dynamics of *Vibrio cholerae* O1 following a single-source introduction to Haiti. *MBio* 4(4):e00398–e00313.
- Kendall M, Boyd M, Colijn C. (2016). phyloTop. <https://CRAN.R-project.org/package=phyloTop>; last accessed February 23, 2018.
- Korhonen V, Smit P, Haanperä M, Casali N, Ruutu P, Vasankari T, Soini H. 2016. Whole genome analysis of *Mycobacterium tuberculosis* isolates from recurrent episodes of tuberculosis, Finland, 1995–2013. *Clin Microbiol Infect.* 22(6):549–554.
- Kuo C-H, Ochman H. 2009. Inferring clocks when lacking rocks: the variable rates of molecular evolution in bacteria. *Biol Direct* 4(1):35.
- Lee RS, Radomski N, Proulx J-F, Manry J, McIntosh F, Desjardins F, Soualhine H, Domenech P, Reed MB, Menzies D, et al. 2015. Reemergence and amplification of tuberculosis in the Canadian Arctic. *J Infect Dis.* 211(12):1905–1914.
- Lillebaek T, Norman A, Rasmussen EM, Marvig RL, Folkvardsen DB, Andersen ÅB, Jelsbak L. 2016. Substantial molecular evolution and mutation rates in prolonged latent *Mycobacterium tuberculosis* infection in humans. *Int J Med Microbiol.* 306(7):580–585.
- Meilä M. 2007. Comparing clusterings – an information based distance. *J Multivar Anal.* 98(5):873–895.
- Merker M, Blin C, Mona S, Duforet-Frebourg N, Lecher S, Willery E, Blum MG, Rüschi-Gerdes S, Mokrousov I, Aleksic E, et al. 2015. Evolutionary history and global spread of the *Mycobacterium tuberculosis* Beijing lineage. *Nat Genet.* 47(3):242.
- Novichkov PS, Omelchenko MV, Gelfand MS, Mironov AA, Wolf YI, Koonin EV. 2004. Genome-wide molecular clock and horizontal gene transfer in bacterial evolution. *J Bacteriol.* 186(19):6575–6585.
- Octavia S, Wang Q, Tanaka MM, Kaur S, Sintchenko V, Lan R. 2015. Delineating community outbreaks of *Salmonella enterica* serovar Typhimurium by use of whole-genome sequencing: insights into genomic variability within an outbreak. *J Clin Microbiol.* 53(4):1063–1071.
- Poon AF. 2016. Impacts and shortcomings of genetic clustering methods for infectious disease outbreaks. *Virus Evol.* 2(2):vew031.
- Roetzer A, Diel R, Kohl TA, Rückert C, Nübel U, Blom J, Wirth T, Jaenicke S, Schuback S, Rüschi-Gerdes S, et al. 2013. Whole genome sequencing versus traditional genotyping for investigation of a *Mycobacterium tuberculosis* outbreak: a longitudinal molecular epidemiological study. *PLoS Med.* 10(2):e1001387.
- Schliep KP. 2011. phangorn: phylogenetic analysis in R. *Bioinformatics* 27(4):592.
- Walker TM, Ip CL, Harrell RH, Evans JT, Kapatai G, Dedicoat MJ, Eyre DW, Wilson DJ, Hawkey PM, Crook DW, et al. 2013. Whole-genome sequencing to delineate *Mycobacterium tuberculosis* outbreaks: a retrospective observational study. *Lancet Infect Dis.* 13(2):137–146.
- Walker TM, Lalor MK, Broda A, Ortega LS, Morgan M, Parker L, Churchill S, Bennett K, Golubchik T, Giess AP, et al. 2014. Assessment of *Mycobacterium tuberculosis* transmission in Oxfordshire, UK, 2007–12, with whole pathogen genome sequences: an observational study. *Lancet Respir Med.* 2(4):285–292.
- Wallinga J, Lipsitch M. 2007. How generation intervals shape the relationship between growth rates and reproductive numbers. *Proc R Soc Lond B Biol Sci.* 274(1609):599–604.

- Worby CJ, Lipsitch M, Hanage WP. 2014. Within-host bacterial diversity hinders accurate reconstruction of transmission networks from genomic distance data. *PLoS Comput Biol.* 10(3):e1003549.
- Yang C, Luo T, Shen X, Wu J, Gan M, Xu P, Wu Z, Lin S, Tian J, Liu Q, et al. 2017. Transmission of multidrug-resistant *Mycobacterium tuberculosis* in Shanghai, China: a retrospective observational study using whole-genome sequencing and epidemiological investigation. *Lancet Infect Dis.* 17(3):275–284.
- Ypma RJ, Donker T, van Ballegooijen WM, Wallinga J. 2013. Finding evidence for local transmission of contagious disease in molecular epidemiological datasets. *PLoS One* 8(7):e69875.