



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

The utility of whole-genome sequencing to identify likely transmission pairs for pathogens with slow and variable evolution

Citation for published version:

Wood, AJ, Benton, CH, Delahay, RJ, Marion, G, Palkopoulou, E, Pooley, CM, Smith, GC & Kao, RR 2024, 'The utility of whole-genome sequencing to identify likely transmission pairs for pathogens with slow and variable evolution', *Epidemics*, vol. 48, 100787, pp. 1-10. <https://doi.org/10.1016/j.epidem.2024.100787>

Digital Object Identifier (DOI):

[10.1016/j.epidem.2024.100787](https://doi.org/10.1016/j.epidem.2024.100787)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Publisher's PDF, also known as Version of record

Published In:

Epidemics

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.





The utility of whole-genome sequencing to identify likely transmission pairs for pathogens with slow and variable evolution

A.J. Wood ^a, C.H. Benton ^b, R.J. Delahay ^b, G. Marion ^c, E. Palkopoulou ^b, C.M. Pooley ^c, G.C. Smith ^b, R.R. Kao ^{a,d,*}

^a Roslin Institute, University of Edinburgh, United Kingdom

^b Animal & Plant Health Agency, United Kingdom

^c Biomathematics and Statistics Scotland, United Kingdom

^d Royal (Dick) School of Veterinary Studies, University of Edinburgh, United Kingdom

ARTICLE INFO

Keywords:

Mycobacterium Bovis
Bovine tuberculosis
Whole genome sequencing
Precision epidemiology

ABSTRACT

Pathogen whole-genome sequencing (WGS) has been used to track the transmission of infectious diseases in extraordinary detail, especially for pathogens that undergo fast and steady evolution, as is the case with many RNA viruses. However, for other pathogens evolution is less predictable, making interpretation of these data to inform our understanding of their epidemiology more challenging and the value of densely collected pathogen genome data uncertain. Here, we assess the utility of WGS for one such pathogen, in the “who-infected-whom” identification problem. We study samples from hosts (130 cattle, 111 badgers) with confirmed infection of *M. bovis* (causing bovine Tuberculosis), which has an estimated clock rate as slow as ~0.1–1 variations per year. For each potential pathway between hosts, we calculate the relative likelihood that such a transmission event occurred. This is informed by an epidemiological model of transmission, and host life history data. By including WGS data, we shrink the number of plausible pathways significantly, relative to those deemed likely on the basis of life history data alone. Despite our uncertainty relating to the evolution of *M. bovis*, the WGS data are therefore a valuable adjunct to epidemiological investigations, especially for wildlife species whose life history data are sparse.

1. Introduction

Pathogen whole-genome sequencing (WGS) is an invaluable tool for understanding transmission patterns in infectious diseases (Pybus and Rambaut, 2009; Kao et al., 2014; Ladner et al., 2019; Guinat et al., 2021). Underpinning such use of pathogen WGS data is the relationship between the rates at which evolutionary processes occur (the “molecular clock” (Bromham and Penny, 2003)), and the timescale of epidemiological processes of interest (e.g. the generation time between infection events). If many events are observed and occur at a consistent rate, this allows for precise data to inform our understanding of the epidemiology. However, for many bacterial pathogens this evolution can be highly variable, and slow enough that the timescale for a single variation is comparable to or longer than the generation time of the disease itself (Menardo et al., 2019; Sanjuán et al., 2010). Here, it is less clear how useful WGS data may be for epidemiological studies (Biek et al., 2015).

An important example of a slowly evolving pathogen for which considerable WGS data has been gathered is *Mycobacterium bovis* (*M.*

bovis), a member of the *Mycobacterium tuberculosis complex* (MTBC) and the causative agent of bovine Tuberculosis (bTB) in many mammal species (O’Reilly and Daborn, 1995). The *M. bovis* genome comprises $\sim 4.3 \times 10^7$ base pairs (Garnier et al., 2003), and has an estimated evolutionary rate of between 0.1–1 variations per year (Crispell et al., 2017, 2019; Trewby et al., 2016; Akhmetova et al., 2023; Canini et al., 2023). However many estimates from the literature are non-overlapping, consistent with a highly variable evolutionary rate.

Bovine Tuberculosis is a chronic zoonotic disease with considerable animal health and economic impacts in many countries including in the United Kingdom and Ireland O’Brien et al. (2023), Butler et al. (2010), Barnes et al. (2023). Control of bTB is a complex challenge in part due to a low sensitivity in the standard skin test for disease in cattle (De la Rua-Domenech et al., 2006), and in several countries the presence of additional susceptible wildlife species capable of cross-species transmission (e.g. European badgers (*Meles meles*) in the UK and Ireland, wild boars (*Sus scrofa*) in continental Europe, white-tailed

* Corresponding author at: Roslin Institute, University of Edinburgh, United Kingdom.
E-mail address: rowland.kao@ed.ac.uk (R.R. Kao).

<https://doi.org/10.1016/j.epidem.2024.100787>

Received 7 May 2024; Received in revised form 3 July 2024; Accepted 14 August 2024

Available online 23 August 2024

1755-4365/© 2024 Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

deer (*Odocoileus virginianus*) and elk (*Cervus canadensis*) in the USA and Australian brushtail possums (*Trichosurus vulpecula*) in New Zealand).

Studies using *M. bovis* WGS data so far mainly focus on inference of phylogenetic trees and the evolutionary history of the pathogen. This has been used to characterise clusters of transmission (van Tonder et al., 2021; Reis et al., 2021; Perea et al., 2021; Rodrigues et al., 2021; Rossi et al., 2023), and infer rates of between-species transmission from when the sample species changes between nodes in the phylogenetic tree (Biek et al., 2012; Santos et al., 2022; Crispell et al., 2019; van Tonder et al., 2021; Price-Carter et al., 2018; Rossi et al., 2020, 2022; Pereira et al., 2023; Canini et al., 2023). This is a first step towards building more complete transmission trees, and inference of *who-infected-whom*. However, for slower-evolving pathogens different host samples may have very similar or even identical sequences entirely, making it more difficult to resolve specific infector–infectee pairs (Campbell et al., 2018). There are a family of approaches for transmission tree reconstruction that combine evolutionary and epidemiological processes (Duault et al., 2022), including Bayesian inference methods which have been applied to MTBC pathogens (Didelot et al., 2017; Xu et al., 2019; Sobkowiak et al., 2020; Xu et al., 2020). Reconstructions on simulated data (where the transmission tree is known) indicate that tree accuracy is sensitive to sampling density (Sobkowiak et al., 2023), and for *M. bovis* in particular, suffers when wildlife species are undersampled (Duault et al., 2023). These are key challenges especially for multi-species disease systems such as bTB, where understanding the role of different species is central to effective disease control (Godfray et al., 2013).

Towards this, then, here we assess the utility of *M. bovis* WGS data in informing *who-infected-whom*. We study WGS samples taken from cows and badgers with confirmed *M. bovis* infection in the vicinity of a long-term study of infection in badgers (Woodchester Park) (Delahay et al., 2013). These data are detailed but low-density, and have significant sample biases which make a full transmission tree inference unsuitable. Instead, then, starting with *all* possible pathways, we assess how the sample data can be used to directly *rule out* pathways, leaving only those that are supported by the sample data. Our interest is in how WGS data can be used to further discriminate between different pathways, by either corroborating existing data, or contradicting them (ruling out pathways that otherwise appear plausible, or boosting pathways where the host data do not support a transmission, but their WGS data do).

For each potential transmission pathway between sampled hosts, using an underlying model for *M. bovis* evolution and transmission we calculate the *relative likelihood* that such a transmission occurred (defined here as the size of space of model outcomes that are consistent with all sample data on those hosts). This is informed by the host sampling times, life history data (time of birth and death, location and movement history), and rates associated with the epidemiology of bTB. In a second set of calculations we include the number of single nucleotide polymorphisms (SNPs) identified in the pathogen sequences relative to a reference sequence. Any differences found between the two approaches over different pathways indicate the additional discriminatory power of the WGS data. We hypothesise that, despite the slow and variable evolution that *M. bovis* exhibits, WGS will be a powerful complementary source of information for inferring transmission pathways, when used alongside existing life history data.

2. Results

We used cattle WGS samples where hosts had a known birth and death date, and movement history between holdings, and badger WGS samples where hosts had a known sample location (their lifespan is estimated from their capture history, see Section 4). We also removed samples from before the year 2000 due to the unreliability of cattle tracing data before this period, as well as repeat samples from the same host (taking only the first sample). This left *M. bovis* samples from 111 badgers and 130 cattle, spanning between 2000 and 2020

(spatial and temporal distribution shown in Supplementary Material A, Fig S1). Most of the cattle samples were from 2013–2020. We note this represents a low density of sampling relative to all known cases in badgers and cattle in the vicinity of the study area and to all likely infections; recorded and not recorded (see Supplementary Material A, Fig S2 for comparative plots including all cases including those not sequenced).

The 241 samples gives a maximum of $241 \times (241 - 1) = 57840$ potential transmission pathways to assess (each host pair, each direction). We removed those deemed impossible due to the hosts not being alive at the same time, and cattle pairs that according to the cattle movement data had not been in the same holding at the same time. This left 12844 possible direct pathways. For pairs of cows that did not necessarily interact directly, but where each shared a holding with one or more cows that may have acted as an intermediary host, we calculate the relative likelihood of transmission via those intermediaries. Of the 130 cattle (16670 potential cattle–cattle pathways), we found 56 pathways with direct interaction, and 247 pathways via potential intermediate hosts. Longer transmission chains with additional intermediaries were not considered but would be considered as having diminishing added value with increased chain length (see Supplementary Material B, Table S1 for a summary of transmission pathways considered and not considered in this work).

For all transmission pathways remaining, we then calculated the relative likelihood of such a transmission having occurred, via a numerical integration. This calculation is outlined in Section 4.3, with full integral expressions (which are not analytically solvable) presented in Supplementary Material C. Briefly here, we model a 2-body susceptible–exposed–infectious process, where substitutions occur on the *M. bovis* genome via a Poisson process. We define here the *relative likelihood* of transmission as the size of the state space of model trajectories under our model and a given choice of model parameters, that would be consistent with the observed data and genetic distance between samples.

2.1. Relative likelihood over all transmission pathways

Fig. 1 illustrates relative transmission likelihoods across all hosts in the data, with hosts ordered by time of sampling. For transmission between cattle pairs via an intermediary host, these pathways are illustrated in Supplementary Material, Fig S6. Line thickness indicates the relative transmission likelihood. For each host pair H_1, H_2 , the line is for the pathway between those hosts with the highest relative likelihood ($H_1 \rightarrow H_2/H_2 \rightarrow H_1$). For the variability in relative likelihood in these datasets, many pathways result in line thickness that is too fine to be resolvable. Over both sets of relative likelihoods there is a clear temporal pattern, with higher relative likelihoods for host pairs sampled closer together in time (see Supplementary Material A, Fig S1 for the temporal distribution of samples). Higher relative likelihood pathways are mainly same-species (see Supplementary Material D, Figs S4, S5 for a breakdown by species pair).

The lines between hosts representing the relative likelihoods are characteristically sparser in Fig. 1B than in Fig. 1A, indicating a larger spread. A broader distribution of relative likelihoods (on a log scale) indicates stronger discrimination between different pathways. The standard deviation of the relative likelihoods over all pathways quantifies the degree to which the inclusion of WGS data has further discriminated between different pathways. Excluding those with likelihood zero, the standard deviation of relative likelihoods evaluated on the basis of life history data is 26.4 (powers of 10), and with life history and WGS data combined is 46.7 (see Supplementary Material D, Fig S7 for full distributions).

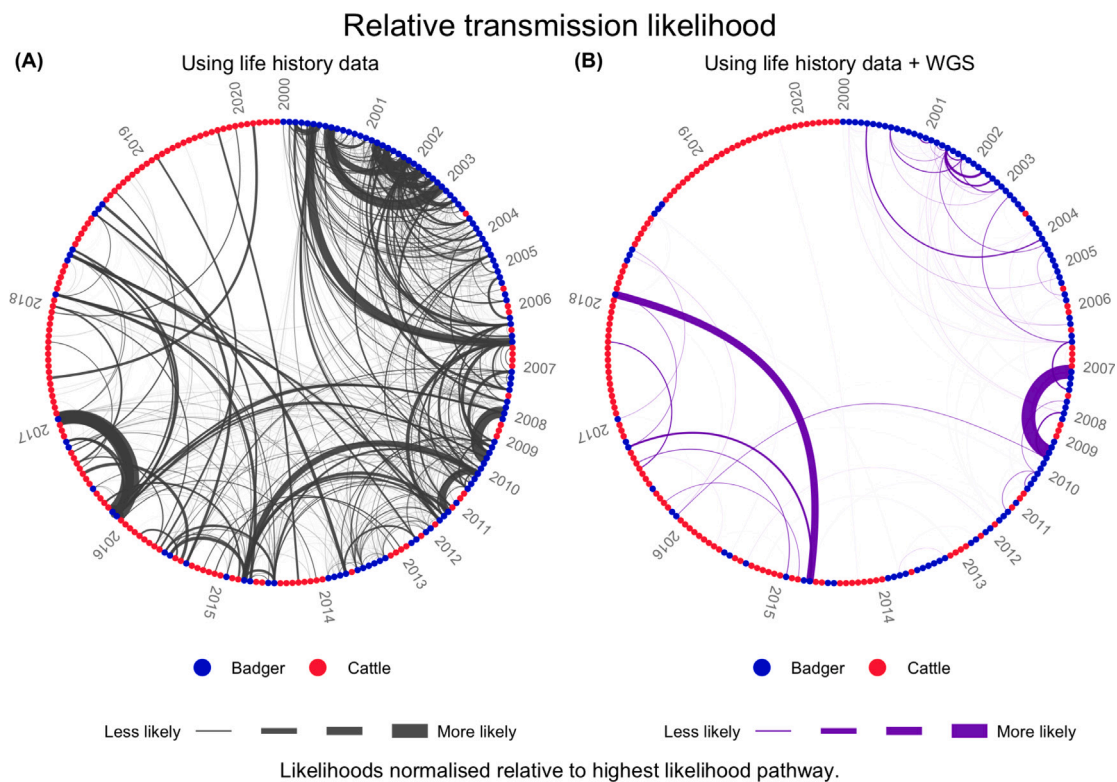


Fig. 1. Relative direct transmission likelihoods between all hosts, evaluated on the basis of (A) life history data, (B) life history data and WGS data. The width of the connecting line is the relative likelihood. For each pair of hosts H_1, H_2 (each host being a point, the colour denoting the species), the relative likelihood is evaluated in both directions ($\mathcal{L}_{H_1 \rightarrow H_2}, \mathcal{L}_{H_2 \rightarrow H_1}$), and we show here the pathway with the higher relative likelihood. Samples are time-ordered in a clockwise arrangement, with the earliest samples at “twelve o’clock”. The relative likelihoods in each set are scaled relative to the pathway with the highest relative likelihood; on this scale many pathways result in line thickness too fine to be resolvable. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Table 1
A broad categorisation of transmission pathways, considering the relative likelihood evaluated on the basis of life history data alone, and that evaluated on the inclusion of WGS data.

Label	Description	Examples
(i)	The pathway is relatively <i>likely</i> on the basis of both life history data and WGS data. The transmission is supported by both data sources.	Supp. Mat., S3
(ii)	The pathway is relatively <i>unlikely</i> on the basis of life history data and WGS data. The transmission pathway can be reasonably ruled out.	Supp. Mat., Table S4
(iii)	The pathway is relatively <i>likely</i> on the basis of life history data but WGS is <i>contradictory</i> to this. The transmission can be reasonably ruled out on the basis of WGS data.	Supp. Mat., Table S5
(iv)	The pathway is relatively <i>unlikely</i> on the basis of life history data but WGS is <i>contradictory</i> to this. Despite the life history data being unresponsive, the WGS signature merits further investigation (e.g. seeking a missing intermediary, or considering a transmission scenario outwith prior assumptions of the epidemiology of bTB).	Supp. Mat., Table S6

2.2. Comparing the discriminatory power of WGS over different species

Fig. 2A shows the distribution of relative likelihoods over all transmission pathways, on the inclusion of WGS data, broken down by pathway type. The specific value of the relative likelihood for all transmission pathways decreases with the inclusion of WGS data. The degree of y-axis dispersion then illustrates that the WGS data further discriminates between transmission pathways with similar relative likelihoods, when evaluated on the basis of life history data alone.

To quantify this discriminatory power, in Fig. 2B we have compared the ranks of transmission pathways in order of relative likelihood. This allows us to classify different transmission pathways in to four broad categories (see Table 1), to indicate whether the transmission pathway is supported, and also whether the WGS and life history data corroborate. Pathways that deviate from the diagonal in Fig. 2B fall in to categories (iii) and (iv), where WGS contradicts the life history data ((iii): deviating bottom-right, (iv) deviating top-left).

The two cattle–cattle pathways with the highest relative likelihood from life history data are between hosts labelled C.054→C.034 (rank 65 out of 13091 pathways, including pathways via an intermediary) and C.104→C.080 (rank 66). These pairs each spent extended time together in the same holding. Taking the WGS data, the genetic distance between the samples for C.054 and C.034 was 4 (C.054 had 1 unique SNP relative to a reference sample not seen in C.034, which in turn had 3 unique SNPs not seen in C.054), and the relative likelihood with these data remains high compared to other pathways (rank 22, category (i) in Table 1). Conversely, the genetic distance between the samples for C.104 and C.080 was 33 (with 18 and 15 unique SNPs respectively). Such a divergence between infector and infectee sequences is extremely improbable given the time between observation and possible infection (see Eq. (10)) ruling it out as a potential transmission pathway (rank 4575, category (iii) in Table 1). The relative likelihoods of these pathways are illustrated in Supplementary Information E, Fig S8 A,

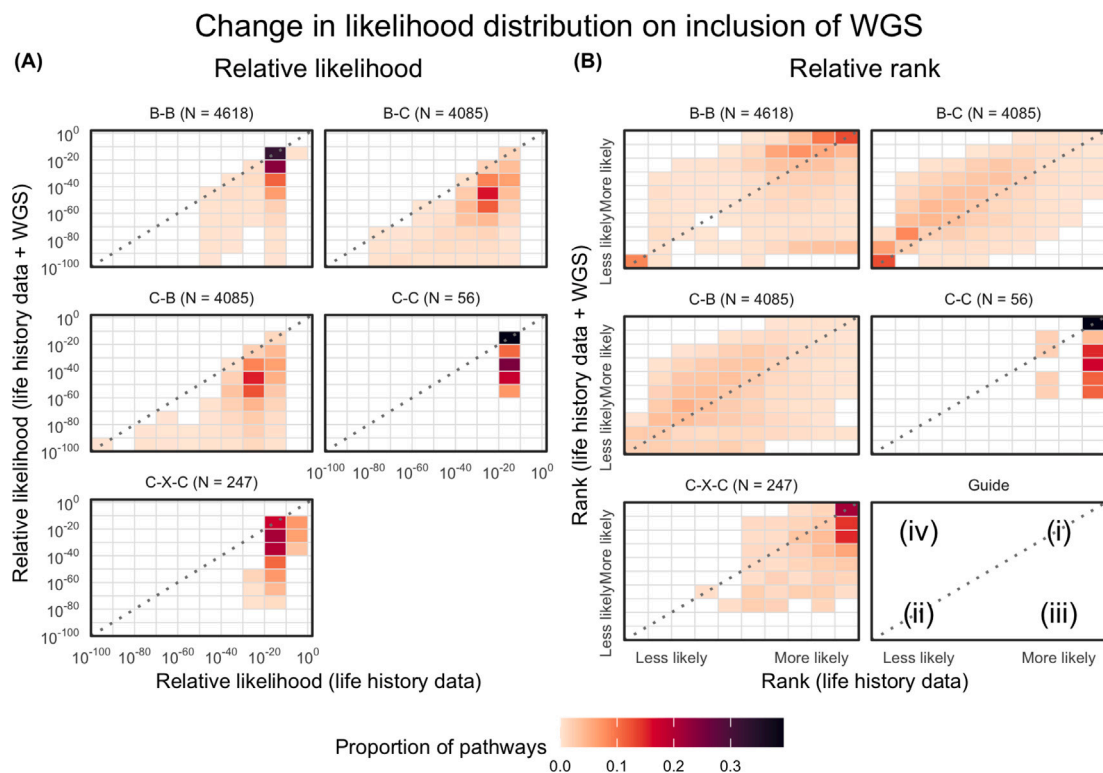


Fig. 2. (A) Density plots showing the distribution of relative likelihoods for individual transmission pathways, on the basis of life history data (x -axis), and life history data and WGS data (y -axis). Darker squares indicate a higher proportion of pathways. Each subplot covers a different species pair, noting significantly different total numbers in each. We exclude pathways deemed impossible from life history data alone. (B) The rank order of all pathways from highest to lowest relative likelihood, again on the basis of life history data (x -axis), and the inclusion of WGS data (y -axis). Points on the diagonal indicate pathways where the additional WGS data broadly corroborates the relative likelihood evaluated on the basis of host life history data alone; these are either deemed to now be relatively less likely (bottom right), or more likely (top left). The position of each of these pathways places them into broad categories (i)–(iv) in Table 1.

amongst similar pathways. We also include in Supplementary Material E, Fig S8B, C, a set of pathways infecting two specific hosts: a badger labelled B.104 (the infectee in the transmission pathway with highest relative likelihood on the basis of life history data), and cow C.085 (the cow with the most interactions with other sampled cattle, direct and via-intermediary). This mimics an epidemiological investigation (asking who infected a particular host). For both, the inclusion of WGS provides a significant degree of discrimination, with the majority of pathways fitting category (iii) in Table 1. For C.085 the WGS data identify the most likely source of infection amongst sampled hosts to be a badger, B.088. This badger was sampled ~ 1 km from a holding that the cow had spent time in, and the WGS samples had a genetic distance of 2 SNPs.

Further examples of specific transmission pathways and relative likelihoods are given in Supplementary Material, Tables S3–S6, with details of the host life histories, and genetic distances amongst *M. bovis* sequences.

Sensitivity to parameter choices

The rank ordering of transmission pathways here are from median relative likelihoods over many parameter draws. The broad rank ordering of pathways remains consistent when looking at individual parameter draws (distributions given in Supplementary Material, Table S2). The median “range of ranks” for pathways under different draws was 362 (with 11077 total having likelihood greater than zero, see Supplementary Material, Fig S10 A). In other words, more likely pathways are consistently more likely, and less likely pathways consistently less likely. However, this variability emphasises the need to consider these relative likelihoods with their broad confidence intervals in mind (as illustrated in e.g. Supplementary Material, Fig S8).

We also assess the sensitivity of the relative (log) likelihood to changing each of the variables in the model via a *Boruta* analysis. This is a standard machine learning method applied to multivariable data that assesses how “important” different variables are in explaining an outcome (the relative likelihood). This is based on how much accuracy a statistical model *loses* if a variable was effectively removed; more important variables see greater losses. This showed (Supplementary Material, Fig S10B) that the relative likelihood is initially most sensitive to changes in the infection rate β_0 , and to the rate of transition to infectiousness σ and time to test-sensitivity Ω to a lesser degree. However on the inclusion of WGS data, the substitution rate μ is overwhelmingly the most important variable; the relative likelihood is most sensitive to changes in μ as opposed to other variables.

3. Discussion

Pathogen whole genome sequencing data, coupled with an understanding of that pathogen’s evolution, can be used to estimate a time to a most recent common ancestor. This in turn may inform the relative likelihood that one host infected another. Used in conjunction with other host data, WGS provides clear utility for inferring transmission pathways when pathogen evolution is steady and predictable, and a reasonable number of substitutions are expected to occur between infection and onward transmission. The aim of this work was to assess the utility of WGS in a far more challenging context applicable to many important pathogens: when evolution is sufficiently slow that the substitution rate is itself difficult to estimate. This applies to wide range of bacterial pathogens, including mycobacteria such as *M. bovis*, which evolves at an estimated rate as low as one substitution on the genome every 10 years. In this instance, the number of variations

between isolates from different hosts can be zero. The “who-infected-whom” identification problem is especially pressing for *M. bovis* control in multi-host systems where infection may be transmitted amongst susceptible domestic and wild species.

We described a method for calculating the *relative likelihood* of transmission between specific hosts in a sample of cattle and badgers with confirmed *M. bovis* infection. We find that a relatively large number of plausible transmission pathways can be identified on the basis of life history data alone (when the hosts were alive, and where). However, the inclusion of data on genetic distances between *M. bovis* isolates permitted the number of plausible pathways to be reduced considerably, ruling out those where the genetic distance between samples would be highly improbable based on the evolutionary characteristics of the pathogen. This illustrates the utility of WGS in identifying a subset of more likely transmission pathways, despite uncertainty over the evolutionary rate of the pathogen (using a wide range of $\mu \in U(0.1, 0.9)$ substitutions/genome/year).

The manner in which WGS was “useful” in our analyses depended on the specific transmission pathway, as categorised in Table 1. First, WGS may *corroborate* the relative likelihood evaluated on the basis of life histories, either as (i) more likely or (ii) less likely. Second, in some instances (iii) the pathway is deemed likely given life history data, but the genetic distance between pathogen samples is highly improbable, ruling out that pathway. These are important for “who-infected-whom” investigations in further pruning pathways that appear plausible based on all other data. The final category (iv) is where an unlikely pathway emerges as comparatively *more* likely on being informed by WGS. An example of this may be hosts that were never in close proximity to one another, but the genetic distance is consistent with a direct transmission having occurred. These pathways merit further investigation. For example, there may have been an unobserved intermediary(s) in a longer transmission chain connecting the two hosts, in which case the life history data can be used to identify where intermediaries may have been. Examples include additional cattle in the movement network, or wild animals at intermediate physical locations. We could also interpret these results as a challenge to the assumptions we have made about the epidemiology of bTB; that is, a mechanism of transmission not captured by our model assumptions (e.g. fomites, rare long-distance badger movements (Gaughran et al., 2018)).

The transmission relative likelihoods are evaluated directly by integrating over the different epidemiological stages (exposure, infectiousness, onward infection), bounded by known delays in test sensitivity and life history data of the hosts. The epidemiology of bTB and evolution of *M. bovis* are embedded within these formulae as functional expressions, which may be generalised to more sophisticated epidemiology (e.g. a more elaborate compartmental model) and models of pathogen evolution. The method presented here accounts for transmission between cattle via a possible single intermediary that has not been sampled, and this could be generalised to longer chains of transmission. An estimate of the potential number of intermediary hosts could also be important in defining the probability of a wildlife host (in areas where wildlife are not sampled) being involved, given information on the testing frequency in cattle.

Method limitations and assumptions

This approach is by construction highly sensitive to the underlying data and model assumptions. The main limitation is how accurately the model and data capture the true movements and interactions of the hosts, and the epidemiology of bTB. Owing to the slow progression of disease and variability in disease severity amongst hosts, many of the parameters used in this work are highly uncertain.

On the evolution of the *M. bovis* bacterium, we model substitutions as occurring randomly and independently at some prescribed rate (Eq. (9)). While the range of substitution rates reflects that estimated in the literature (Crispell et al., 2017, 2019; Trewby et al., 2016; Akhmetova et al., 2023; Canini et al., 2023), these estimates are often

non-overlapping and suggest that the “true” substitution rate may be non-regular in time, and also vary depend on the pathogenesis of *M. bovis* and the immunological response to infection in specific hosts, as well as the specific lineage of *M. bovis*. Inevitably, the collection of more WGS data will allow us to build a more precise understanding of the evolutionary process and these systematic variations, and in turn improve the value of those data in such epidemiological investigations.

A limitation specific to the cattle data is that while we identify “same holding, same time” instances for cattle pairs, holdings themselves may be geographically fragmented meaning some cattle on the same holding may not have had any opportunity to interact. In addition we may miss important interactions such as those “over-the-fence” between adjacent holdings, or transmission “between” two hosts that had never interacted directly, via fomites (such as a common piece of equipment, or contamination of shared grazing land).

Our model choice (a two-body susceptible-exposed-infected model) is relatively simple and, before considering lifespan and location, places few constraints as to the types of transmission permitted. The same approach used here could be applied with a more complex disease model that considered e.g. details of the environment and natural boundaries, social structures and roaming behaviour. Such an approach would mainly rule out more pathways on top of those done so here, rather than returning an entirely different set of plausible/improbable pathways.

Our model uncertainty is greatest in the badgers, as these wild animals do not have precisely tracked locations. The badger data we use here are exceptional, as the social structure and hence the resident territories of badgers in Woodchester Park are well documented. We estimated a distance-dependency on the infection rate from the difference between within and between-social group infection rates (Eq. (7)). The main effect of changing this distance dependency is to systematically suppress/amplify pathways at higher distances (see Supplementary Material F, Fig S9), but any distance-dependency on infection probability between two hosts will vary based on factors including but not limited to hosts’ tendency to move and the influence of the distribution of suitable habitat and physical boundaries to movement, and the nature of their interactions with others (e.g. badgers biting one another in territorial disputes).

The infection rates used in this work can be refined further, however sensitivity analysis from Section 2 does indicate that the relative likelihood is overwhelmingly more sensitive to a tuning of the substitution rate. Thus though we would expect some changes to individual rankings, the overall conclusions as to the discriminatory value of WGS even with such low substitution rates, are unlikely to change.

Finally, we do not consider historical negative tests from either species, which may inform lower bounds on infection time. We chose not to include this due to the poor sensitivity of the diagnostic tests for bTB (De la Rua-Domenech et al., 2006; Nuñez-García et al., 2018); as all sampled hosts in these data recorded a positive test eventually, previous tests may in fact have been false negatives. We were then not confident in including a hard cutoff at, say, the most recent negative test, however these data could be used in a more careful manner in the future, if considering test sensitivity and overall disease prevalence.

Sample coverage

For the data used here, the sample coverage relative to all likely infections within the study area in the period studied is very low. We analysed 130 cattle between 2000–2020, which is only a small subset of the 5852 cattle in total found with *M. bovis* infection (either a skin test reactor, or visible lesions found at slaughter) within a 10 km radius of the Woodchester Park study area (per the SAM test database). Similarly, over 715 badgers were found with confirmed infection (either by culture or serological testing) during this period in Woodchester Park alone, inclusive of the 111 sequences studied (Supplementary Material, Fig S2). With the compounding effect of not all infections being detected by the tests available, this makes it

unlikely that many “true” transmission pairs are fully characterised by these data, particularly between-species due to strong spatial and temporal biases in host sampling (see Supplementary Material, Fig S1). For these reasons we hesitate to make strong claims about “direction of transmission” from this analysis (as compared with analyses that build the broader phylogenetic tree (Crispell et al., 2019; Canini et al., 2023; van Tonder et al., 2021)). Rather, this work indicates that WGS is highly effective for *ruling out* transmission pathways across sampled hosts, including instances when transmission appeared likely on the basis of other information (Supplementary Material, Fig S8B). This approach provides a simple method for assessing when an infection source has been missed (i.e. when for a specific host, no plausible infection source is found). A simple application would be to infer a wildlife source for a herd breakdown, where no suitable cattle source was found after dense cattle sampling.

Specific utility of WGS for disease transmission in wildlife

Knowledge of cattle movements alone allowed us to eliminate the vast majority of potential transmission pathways amongst cattle (notwithstanding limitations of these data as described above), though WGS was effective at discriminating between remaining pathways (e.g., Supplementary Material, Figs S5, S6, S8C). Conversely, without knowing such detailed movement data for individual wild badgers (even though the data used were exceptional for a wildlife species), we had to consider all pathways involving hosts with overlapping lifespans. For these pathways WGS was especially effective and often contradicted the relative likelihood evaluated on the basis of only life history data (e.g. Fig. 2B). This indicates that WGS has particular value for interactions involving wildlife species despite the variability in pathogen evolution, as other data on movements and propensity to transmit disease are likely to be even more limited. This is especially relevant to *M. bovis*, where a principle challenge in disease control in several countries is understanding the role of the wildlife reservoir in sustaining infection in domesticated species. In a separate case of cattle infection in the Low Risk Area of England where badger culling had taken place, there was sufficient WGS data on wildlife to potentially determine the field in which cattle grazed and were infected, as we had more detailed wildlife data and could potentially exclude one of the two areas where they had grazed. Such detailed analysis could only be determined from robust WGS data in both species.

4. Data and methods

4.1. Data

This study uses life history data on badgers and cattle located in and around Woodchester Park, a Site of Special Scientific Interest in south-west England. This lies within the bTB “high-risk” zone, which covers all of south-west England and some counties beyond. In this zone cattle testing is generally mandated every 6 months (Animal and Plant Health Agency, 2020).

WGS samples

The WGS dataset contains *M. bovis* sequences with metadata on host species, sample date, and for a subset of the hosts a grid-reference location. WGS samples here are from between 2000 and 2020. In this period the primary method of bacterial characterisation in GB cattle was not WGS, but rather spoligotyping and variable number tandem repeat typing (both of which sample only short regions of the genome) (Animal and Plant Health Agency, 2021, 2022), with surveillance in wildlife beyond Woodchester Park extremely limited. The sequences dataset used here had been aligned to the reference genome AF2122/97 (Malone et al., 2017).

Woodchester park dataset

The badger population in Woodchester Park has been intensively studied through a regular programme of capture-test-release since the 1970’s (Delahay et al., 2005, 2013; McDonald et al., 2018). The study archives the life histories of captured badgers, including bTB test results, social group affiliation and capture locations. Clinical samples are collected from each captured badger for the purposes of serological testing and microbiological culture to determine their infection status. All captured animals are re-released unharmed following examination and sampling, regardless of their bTB status. Tests are conducted via culture or serological testing.

Exact times of birth and death of the badgers are not known. Most badgers are captured more than once, implying that almost all badgers are represented in the study data. Some badgers are first captured as cubs, whereas others may be found dead (often after a road traffic collision).

Cattle movements

The Cattle Tracing System (CTS, version dated 1 November 2020) archives the movements of all individual cattle in Great Britain, and the grid-reference locations of holdings (Department for Environment, Food & Rural Affairs, 2023). From these data we associate with each WGS sample a time of birth and death, and movement history of the host cow, and the movement history of all cattle that had shared a holding with that cow at some point in its lifetime.

Cattle bTB test results

The SAM database (version containing entries up to 2021) contains data on cattle bTB test results, and herd breakdowns (Animal and Plant Health Agency, 2023). We extract from SAM the eartags of all cattle that had registered a positive bTB test, regardless of whether they were sequenced.

4.2. Data preparation

For each host, we require in addition to their WGS sample:

- time of sampling;
- time of birth and death (estimated for badgers based on observation times);
- location (observation location for badgers, full movement history for cattle).

Badgers

For each badger isolate, the corresponding isolate ID is linked to a badger in the Woodchester Park dataset. Captured badgers are released after examination and sampling regardless of infection status, so it is possible that they remain alive for several years after the sampling event. For badgers first captured as a cub, and later recorded as dead, we fix these dates as the times of birth and death respectively. For badgers first observed as an adult, or never observed dead, we use an estimated lifespan. For all badgers the birth date can be estimated as being between January–February, as per the modal birth time reported by Neal and Cheeseman (Cheeseman, 1996). For badgers first observed as a cub, we take the birth date to be February 1st of the year of first capture. For badgers first observed as an adult, we instead take the birth date as February 1st of the previous. For the death date, we either take the date of observation if the badger was found dead, or six months after the date of last capture if never found dead. The badger’s location is taken as the location of the badger at the time the sample was taken.

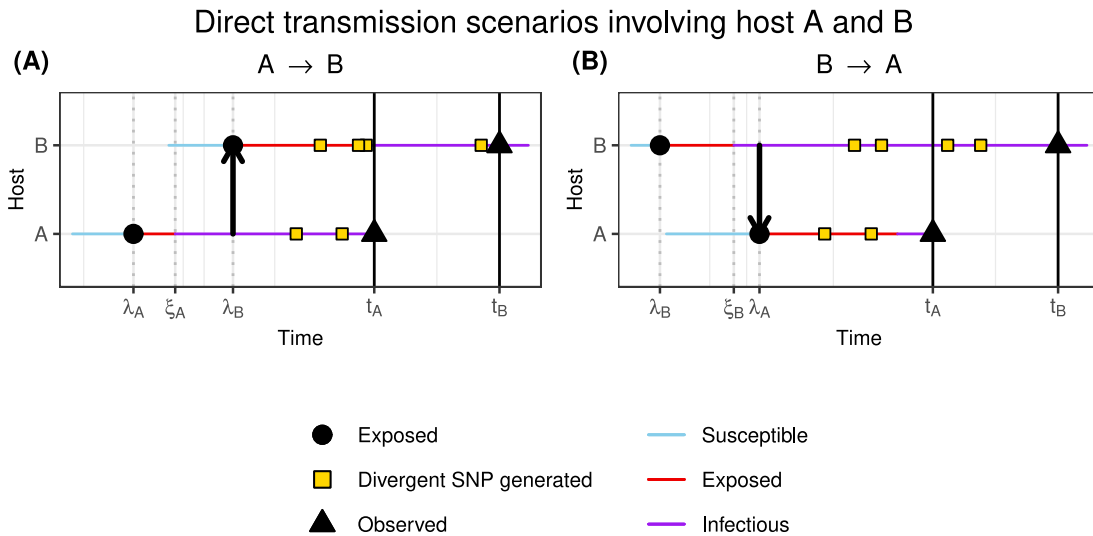


Fig. 3. For two hosts A and B that are sampled at times t_A and t_B respectively, two transmission pathways that would be consistent with an observed genetic distance of 6 SNPs, with A having 2 SNPs not seen in B, and B with 4 SNPs not seen in A, both with respect to a reference genome. Hosts are born susceptible (blue), exposed to infection at times denoted λ (red), and become infectious at ξ (purple). The generation of SNPs on the pathogen genome inside each host are denoted with yellow squares. (A): A exposes B to infection. (B): B exposes A to infection. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Cattle

For each cattle isolate, the corresponding isolate ID is linked to a cow in the CTS dataset (table tblAnimal), usually via an eartag. For hosts without eartags, we instead use the geographic location and time of the sample to identify from SAM potential cattle reactors. If a unique match is found, that eartag is used.

Cattle reactors are generally slaughtered very shortly after detection, thus the sample time is usually very shortly before the time of death.

4.3. Evaluation of relative transmission likelihoods

The relative likelihood that a transmission occurred from one given host to another is expressed as an integral over the ensemble of scenarios where one host infects another, constrained by:

- the *epidemiology* of bTB and the timescales in the epidemiological processes;
- the known *lifespan* of the badgers and cattle observed;
- the known *locations* of the badgers and cattle relative to one another;
- the *observation times* when samples were taken;
- the *evolutionary* process of *M. bovis* in the generation of substitutions on the genome.

Here we describe how each of these model assumptions and data inform the relative likelihood. The exact integral expressions for the relative likelihood are presented in Supplementary Material, Section C. Parameter ranges are summarised in Supplementary Material, Table S2.

4.3.1. Epidemiology of *M. bovis*

We model bTB as a three-stage SEI (Susceptible, Exposed, Infectious) process (Fig. 3). (I)nfectious hosts expose (S)usceptible hosts at rate β , denoting the time of exposure of host X as λ_X . (E)xposed individuals transition to the infectious stage at rate σ , with the time of transition for host X denoted ξ_X . The times λ and ξ are not generally observable.

Consider a system with two hosts A and B where at time λ_A A is exposed, with B susceptible. In an isolated system, A first transitions to infectiousness at some later time $\xi_A \geq \lambda_A$, after which A exposes B at some later time $\lambda_B \geq \xi_A$ (Fig. 3A).

The time ξ_A depends on λ_A ; if A is exposed at λ_A , the probability that A is not yet infectious by time ξ is

$$P(\text{A not yet \{I\} at } \xi_A \geq \lambda_A \mid \text{exposed at } \lambda_A) = e^{-\sigma(\xi_A - \lambda_A)} \tag{1}$$

thus the distribution of transition times ξ_A is

$$P(\text{A becomes \{I\} between } (\xi_A, \xi_A + d\xi_A) \mid \lambda_A) = -\partial_{\xi_A} (e^{-\sigma(\xi_A - \lambda_A)}) d\xi_A = d\xi_A \sigma e^{-\sigma(\xi_A - \lambda_A)}. \tag{2}$$

In a similar way, the time λ_B is conditioned on ξ_A , and is distributed as:

$$P(\text{B becomes \{E\} between } (\lambda_B, \lambda_B + d\lambda_B) \mid \xi_A) = -\partial_{\lambda_B} (e^{-\beta(\lambda_B - \xi_A)}) d\lambda_B = d\lambda_B \beta e^{-\beta(\lambda_B - \xi_A)}. \tag{3}$$

To evaluate the relative likelihood of transmission from one host to another, we integrate over all possible times of exposure ξ and transition λ within the system, constrained by the lifespans of the animals and times of observation.

4.3.2. Test sensitivity

We assume that hosts are not immediately test-sensitive after exposure; that is, if tested, they would report a false negative. Rather, they only become test-sensitive after a period Ω . In our model, we take this time to be independent of the time of transition to *infectiousness*. We incorporate this with a term W_X

$$W_X(t_X, \lambda_X) = P(\text{test at } t_X \text{ returns positive result} \mid \text{exposed at } \lambda_X) = \theta(t_X - [\lambda_X + \Omega]) = \begin{cases} 1 & t_X \geq \lambda_X + \Omega \\ 0 & \text{otherwise} \end{cases}. \tag{4}$$

θ is a Heaviside step function; $\theta(x) = 1$ if $x \geq 0$, and 0 otherwise. This excludes scenarios where a host tests positive very shortly after being exposed to infection.

The time to test-sensitivity is another highly uncertain quantity (see e.g. Conlan et al. (2012), where the period was estimated as either within (0, 7.7) or (24, 517) days, dependent on model choice). We use the univariate range $\Omega \in U(0, 84)$ d based on a challenge study taking a skin test 12 weeks post-exposure (Dean et al., 2005).

4.3.3. Lifespan of hosts

The relative likelihood of transmission between a pair of hosts also depends on the overlap of their lifespans. If a host X was born at $t_X^{(B)}$ and died at $t_X^{(D)}$, we exclude scenarios occurring in the range (t_1, t_2) that are beyond the lifetime of X, using a Heaviside step term

$$\begin{aligned} V_X(t_1, t_2) &= P(\text{host X is alive at both } t_1 \text{ and } t_2) \\ &= \theta(t_1 - t_X^{(B)}) \theta(t_X^{(D)} - t_2) \\ &= \begin{cases} 1 & t_X^{(B)} \leq t_1, t_2 \leq t_X^{(D)} \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (5)$$

4.3.4. Physical distance between hosts

Direct host-to-host transmission occurs from close contact, thus the relative transmission likelihood should depend on the separation between hosts (affecting how likely they are to interact). We capture this with a distance-dependence on the infection rate β . To account for the different movement patterns of wild badgers and domesticated cattle, we treat the interaction distance differently depending on the host species.

Cattle–cattle

For cattle–cattle transmission, we take the infection rate to be non-zero only when they were in the same herd at the same time. We identify interactions from their movement histories (CTS table `vMovementDirect`). The within-herd infection rate in those periods is taken to be that at zero separation. For all other times, the infection rate is zero:

$$\beta_{C \rightarrow C} = \begin{cases} \beta_0 & \text{when in the same herd together} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

The relative likelihood of direct transmission between two cows that were never in the same herd at the same time is therefore zero. For β_0 we use a gamma-distributed distribution based on the distribution obtained in Brooks-Pollock et al. (2014), with a posterior estimate of the rate of within-holding onward transmission between cattle as $\beta_0 \sim 0.61 \text{ y}^{-1}$ (95% CI: [0.0503, 1.54]).

Also from cattle movement histories we identify *indirect* interactions between two cows. These are when the first cow interacted with one or more cows, that *later* went on to interact with the second, and could therefore have acted as an intermediary in a three-cow transmission chain. For these interactions with intermediaries, the transmission rate is also taken to be that at zero distance. We multiply each relative likelihood by the number of intermediaries found in CTS, to estimate the relative likelihood of transmission via any one of these intermediaries.

Badger–badger

Medium-to-high density badger populations in the UK are typically organised into social groups with relatively well-defined territories (Rogers et al., 1997; Vicente et al., 2007). This was the case for Woodchester Park, where the social groups of individual badgers were closely monitored and group territories were well-defined during the period of study. The majority of movements and interactions between badgers have been observed to be within the territories of their own social group (Rogers et al., 1998). Between-social group movements are less frequent but serve important ecological purposes such as reducing inbreeding (Dugdale et al., 2008; Benton et al., 2018). The amount of between-group activity likely varies by badger depending on factors including age and sex (Böhm et al., 2008; Macdonald et al., 2008).

We take the effective rate of infection $\beta_{B \rightarrow B}$ between two badgers to fall exponentially with the separation d between them:

$$\beta_{B \rightarrow B} = \beta_0 e^{-\alpha d} \quad (7)$$

where α characterises how steep the fall with distance is. The social structure suggests that the majority of badger–badger transmission of *M. bovis* will be between badgers of the same social group, but it

is more difficult to quantify the relative frequency of transmission pathways between badgers in adjacent social groups. Smith et al. (2001) note that *M. bovis* infection originating specifically from a bite wound may be indicative of a between-group transmission, where the bite is inflicted by an infectious badger in a neighbouring social group, as a form of territorial aggression. From their study, of 33 infections 10 were from bite wounds, suggesting that up to 30% of infections may have been from inter-group pathways. If we estimate that a given social group will have approximately 5 social groups that could be considered adjacent, the proportion of infections within a social group from each neighbouring groups is 0.06.

Using the location data of the badgers at Woodchester Park, the mean distance between a social group and its 5 nearest neighbouring groups is 0.50 km. We then use a value $\alpha = 4.92 \text{ km}^{-1}$ such that the infection rate at 0.50 km is 0.06/0.70 times that at zero distance.

Badger–cattle and cattle–badger

The cattle–badger and badger–cattle transmission rates also depend on distance, however the distance varies due to the movement of cattle between holdings:

$$\beta_{B \rightarrow C, C \rightarrow B} = \beta_0 e^{-\alpha d(t)} \quad (8)$$

thus the infection rate changes in a step-like fashion whenever a cow moves holding, based on the physical distance between each holding and where the badger was observed. The same distance-dependency α is used as in Eq. (7), reflecting the estimated roaming of the badgers. The physical locations of the holdings (used to determine a physical distance to the badger sampling location) are derived from CTS, table `tblLocation`.

4.3.5. Evolutionary process of *M. bovis*

We assume that within-host single nucleotide polymorphisms (SNPs) of *M. bovis* are generated randomly and independently under Poisson dynamics, at a mean rate μ . Under Poisson dynamics, over an interval Δt , the probability that *exactly* m SNPs are generated is

$$\begin{aligned} Q(m, \Delta t) &= P(\text{exactly } m \text{ SNPs are generated in interval } \Delta t) \\ &= \frac{(\mu \cdot \Delta t)^m}{m!} e^{-\mu \Delta t} \end{aligned} \quad (9)$$

In an interval Δt the expected number of SNPs to have been generated is $\langle m \rangle = \mu \Delta t$, with standard deviation $\sqrt{\langle m^2 \rangle - \langle m \rangle^2} = \sqrt{\mu \Delta t}$. *M. bovis* exhibits a slow rate of molecular evolution; the substitution rate μ is estimated across different studies to be within the range of $\mu \approx [0.1-0.9]$ substitutions, per genome, per year (Crispell et al., 2017, 2019; Trewby et al., 2016; Akhmetova et al., 2023; Canini et al., 2023).

From the aggregation of all WGS samples, we identify the most common nucleotide per site, using this sequence of nucleotides as a reference. Then, for each pair of hosts A and B we count (i) the number of divergent SNPs (differences between A's sequence and the reference) observed in A but *not* in B (denoted m_A), and; (ii) the number of divergent SNPs observed in B but *not* in A (m_B). Note that the sum $m_A + m_B$ is the overall genetic distance between A and B, independent of the reference used. In a scenario where A infected B or vice versa, the infection event is the point of divergence for those two *M. bovis* strains, so any SNPs unique to one host only must have been generated *between* infection and observation (Fig. 3).

Following Section 4.3.1, suppose the infection $A \rightarrow B$ occurred at time λ_B . A and B are later sampled and their strains sequenced at times t_A, t_B respectively. A's sample has m_A SNPs not seen in B's sample, and B's has m_B SNPs not seen in A's. By Eq. (9), the probability of this is

$$\begin{aligned} P(m_A, m_B \mid \lambda_A) &= Q(m_A, t_A - \lambda_B) \cdot Q(m_B, t_B - \lambda_B) \\ &= \frac{(\mu \cdot [t_A - \lambda_B])^{m_A}}{m_A!} e^{-\mu \cdot [t_A - \lambda_B]} \cdot \frac{(\mu \cdot [t_B - \lambda_B])^{m_B}}{m_B!} e^{-\mu \cdot [t_B - \lambda_B]} \end{aligned} \quad (10)$$

This incorporates the *different* sampling times of hosts, and the evolutionary processes that would happen *independently* in the two hosts post-infection.

Before incorporating WGS, the space of possible scenarios for transmission from one host to another is constrained by the host lifespans, and when those hosts interacted. The addition of WGS then constrains us to only scenarios where a specific number of SNPs are generated between infection and observation. We therefore expect the relative likelihoods for all pathways to fall compared to the equivalent calculation not including WGS data (which in effect sets the expression in Eq. (9) to 1).

4.3.6. Calculation

To evaluate the relative transmission likelihood between two hosts we numerically integrate over all possible times for the epidemiological events (Eqs. (2), (3)), constrained by the delay to test sensitivity (Eq. (4)) and the host lifespans (Eq. (5)), with the infection rate dependent on the host distance (Eqs. (6)–(8)). The relative likelihoods with WGS information are then further weighted by the probability of the observed genetic distance occurring (Eq. (9)). Full expressions are presented in Supplementary Material, Section C.

For each transmission pathway, we evaluate the relative likelihood from 300 samples of the parameter distributions (Supplementary Material, Table S2). For the 250 host pairs with the highest median relative likelihood after this, we do a further 1000 evaluations, after which we record the median relative likelihood, and [0.05, 0.95] confidence interval.

Analysis code was written in R (version 4.3.1) (R Core Team, 2022). To evaluate relative likelihoods we used multivariable numerical integration, with a nested application of the `integrate` function (stats package (R Core Team, 2022)). The implementation in R was verified against a symbolic implementation in *Mathematica* (version 12.3.1.0 (Wolfram Research, Inc., 2023)), using the `NIntegrate` function.

CRediT authorship contribution statement

A.J. Wood: Writing – review & editing, Writing – original draft, Methodology, Investigation, Formal analysis. **C.H. Benton:** Writing – review & editing, Methodology. **R.J. Delahay:** Writing – review & editing, Methodology. **G. Marion:** Writing – review & editing, Methodology. **E. Palkopoulou:** Writing – review & editing, Methodology. **C.M. Pooley:** Writing – review & editing, Methodology. **G.C. Smith:** Writing – review & editing, Methodology. **R.R. Kao:** Writing – review & editing, Writing – original draft, Supervision, Methodology, Conceptualization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

The cattle movement and testing data (CTS, SAM), the Woodchester Park study data and the WGS data are not publicly available, and are provided to the authors under a data sharing agreement with the Animal and Plant Health Agency, who can be contacted via enquiries@apha.gov.uk.

Acknowledgements

We thank Gianluigi Rossi for helpful discussions on relative transmission likelihoods, and providing comments on the manuscript.

Code availability

Analysis code is available at <https://git.ecdf.ed.ac.uk/awood310/mBovis-WGS-transmission-likelihoods>.

Funding statement

This work has been supported by the UKRI, United Kingdom grants BB/W007290/1 and BB/W007711/1 “Developing better modelling inference tools to inform disease control for bovine Tuberculosis using epidemiological and pathogen genetic information”, the Roslin Institute Strategic Programme grant BBS/E/RL/230002D, and also supported by the GB bovine TB research budget, project code SE3337, held and administered centrally by UK Government’s Department for Environment, Food and Rural Affairs (Defra) on behalf of England, Scotland and Wales. WGS samples used in study were funded by Defra, United Kingdom (projects APHATBSB4030 and APHATBOR1089). GM and CP were supported by the Scottish Government’s Rural and Environment Science and Analytical Services Division (RESAS), United Kingdom.

Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.epidem.2024.100787>.

References

- Akhmetova, A., Guerrero, J., McAdam, P., Salvador, L.C., Crispell, J., Lavery, J., Presho, E., Kao, R.R., Biek, R., Menzies, F., et al., 2023. Genomic epidemiology of *Mycobacterium bovis* infection in sympatric badger and cattle populations in Northern Ireland. *Microbial Genom.* 9 (5).
- Animal and Plant Health Agency, 2020. Guidance — Bovine TB testing intervals. <https://www.gov.uk/guidance/bovine-tb-testing-intervals>. (Online; Accessed 17 November 2023).
- Animal and Plant Health Agency, 2021. APHA briefing note 13/21; introduction of the operational use of whole genome sequencing (WGS) for bTB purposes. <http://apha.defra.gov.uk/documents/ov/Briefing-Note-1321.pdf>. (Online; Accessed 9 November 2023).
- Animal and Plant Health Agency, 2022. Bovine tuberculosis in Great Britain in 2022; explanatory supplement to the annual reports. <https://www.gov.uk/government/publications/bovine-tb-epidemiology-and-surveillance-in-great-britain-2022>. (Online; Accessed 9 November 2023).
- Animal and Plant Health Agency, 2023. SAM. <https://www.data.gov.uk/dataset/a553b50e-60ab-45f7-9b6c-90d0c21ada80/sam>. (Online; Accessed 10 November 2023).
- Barnes, A.P., Moxey, A., Brocklehurst, S., Barratt, A., McKendrick, I.J., Innocent, G., Ahmadi, B.V., 2023. The consequential costs of bovine tuberculosis (bTB) breakdowns in England and Wales. *Prevent. Vet. Med.* 211, 105808. <http://dx.doi.org/10.1016/j.prevetmed.2022.105808>.
- Benton, C.H., Delahay, R.J., Smith, F.A., Robertson, A., McDonald, R.A., Young, A.J., Burke, T.A., Hodgson, D., 2018. Inbreeding intensifies sex-and age-dependent disease in a wild mammal. *J. Anim. Ecol.* 87 (6), 1500–1511.
- Biek, R., O’Hare, A., Wright, D., Mallon, T., McCormick, C., Orton, R.J., McDowell, S., Trewby, H., Skuce, R.A., Kao, R.R., 2012. Whole genome sequencing reveals local transmission patterns of *Mycobacterium bovis* in sympatric cattle and badger populations. *PLoS Pathog.* 8 (11), e1003008.
- Biek, R., Pybus, O.G., Lloyd-Smith, J.O., Didelot, X., 2015. Measurably evolving pathogens in the genomic era. *Trends Ecol. Evolut.* 30 (6), 306–313.
- Böhm, M., Palphramand, K.L., Newton-Cross, G., Hutchings, M.R., White, P.C., 2008. Dynamic interactions among badgers: implications for sociality and disease transmission. *J. Anim. Ecol.* 77 (4), 735–745.
- Bromham, L., Penny, D., 2003. The modern molecular clock. *Nature Rev. Genet.* 4 (3), 216–224.
- Brooks-Pollock, E., Roberts, G.O., Keeling, M.J., 2014. A dynamic model of bovine tuberculosis spread and control in Great Britain. *Nature* 511 (7508), 228–231.
- Butler, A.J., Lobley, M., Winter, M., 2010. Economic impact assessment of bovine tuberculosis in the south west of England.
- 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.
- Canini, L., Modenesi, G., Courcoul, A., Boschirolli, M.-L., Durand, B., Michelet, L., 2023. Deciphering the role of host species for two *Mycobacterium bovis* genotypes from the European 3 clonal complex circulation within a cattle-badger-wild boar multihost system. *MicrobiologyOpen* 12 (1), e1331.

- Cheeseman, C., 1996. Badgers. T & AD Poyser Natural History.
- Conlan, A.J., McKinley, T.J., Karolemeas, K., Pollock, E.B., Goodchild, A.V., Mitchell, A.P., Birch, C.P., Clifton-Hadley, R.S., Wood, J.L., 2012. Estimating the hidden burden of bovine tuberculosis in Great Britain. Public Library of Science San Francisco, USA.
- Crispell, J., Benton, C.H., Balaz, D., De Maio, N., Ahkmetova, A., Allen, A., Biek, R., Presho, E.L., Dale, J., Hewinson, G., et al., 2019. Combining genomics and epidemiology to analyse bi-directional transmission of *Mycobacterium bovis* in a multi-host system. *Elife* 8, e45833.
- Crispell, J., Zadoks, R.N., Harris, S.R., Paterson, B., Collins, D.M., de Lisle, G.W., Livingstone, P., Neill, M.A., Biek, R., Lycett, S.J., et al., 2017. Using whole genome sequencing to investigate transmission in a multi-host system: bovine tuberculosis in New Zealand. *BMC Genomics* 18, 1–12.
- De la Rua-Domenech, R., Goodchild, A., Vordermeier, H., Hewinson, R., Christiansen, K., Clifton-Hadley, R., 2006. Ante mortem diagnosis of tuberculosis in cattle: a review of the tuberculin tests, γ -interferon assay and other ancillary diagnostic techniques. *Res. Veterinary Sci.* 81 (2), 190–210.
- Dean, G.S., Rhodes, S.G., Coad, M., Whelan, A.O., Cockle, P.J., Clifford, D.J., Hewinson, R.G., Vordermeier, H.M., 2005. Minimum infective dose of *Mycobacterium bovis* in cattle. *Infect. Immunity* 73 (10), 6467–6471.
- Delahay, R., Smith, G., Ward, A., Cheeseman, C., 2005. Options for the management of bovine tuberculosis transmission from badger (*Meles meles*) to cattle: evidence from a long-term study. *Mammal Study* 30, S73–S81. [http://dx.doi.org/10.3106/1348-6160\(2005\)30\[S73:OFTMOB\]2.0.CO;2](http://dx.doi.org/10.3106/1348-6160(2005)30[S73:OFTMOB]2.0.CO;2).
- Delahay, R.J., Walker, N., Smith, G., Wilkinson, D., Clifton-Hadley, R., Cheeseman, C., Tomlinson, A., Chambers, M., 2013. Long-term temporal trends and estimated transmission rates for *Mycobacterium bovis* infection in an undisturbed high-density badger (*Meles meles*) population. *Epidemiol. Infect.* 141 (7), 1445–1456.
- Department for Environment, Food & Rural Affairs, 2023. CTS online. <https://secure.services.defra.gov.uk/wps/portal/ctso>. (Online; Accessed 10 November 2023).
- Didelot, X., Fraser, C., Gardy, J., Colijn, C., 2017. Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks. *Molecular Biol. Evolution* 34 (4), 997–1007.
- Duault, H., Durand, B., Canini, L., 2022. Methods combining genomic and epidemiological data in the reconstruction of transmission trees: a systematic review. *Pathogens* 11 (2), 252.
- Duault, H., Durand, B., Canini, L., 2023. Outbreak reconstruction with a slowly evolving multi-host pathogen: a comparative study of three existing methods on *Mycobacterium bovis* outbreaks. *bioRxiv*.
- Dugdale, H.L., Macdonald, D.W., Pope, L.C., Johnson, P.J., Burke, T., 2008. Reproductive skew and relatedness in social groups of European badgers, *meles meles*. *Molecular Ecol.* 17 (7), 1815–1827.
- Garnier, T., Eiglmeier, K., Camus, J.-C., Medina, N., Mansoor, H., Pryor, M., Duthoy, S., Grondin, S., Lacroix, C., Monsempé, C., et al., 2003. The complete genome sequence of *Mycobacterium bovis*. *Proc. Natl. Acad. Sci.* 100 (13), 7877–7882.
- Gaughran, A., Kelly, D.J., MacWhite, T., Mullen, E., Maher, P., Good, M., Marples, N.M., 2018. Super-ranging: a new ranging strategy in European badgers. *PLoS One* 13 (2), e0191818.
- Godfray, H.C.J., Donnelly, C.A., Kao, R.R., Macdonald, D.W., McDonald, R.A., Petrokofsky, G., Wood, J.L., Woodroffe, R., Young, D.B., McLean, A.R., 2013. A restatement of the natural science evidence base relevant to the control of bovine tuberculosis in Great Britain. *Proc. R. Soc. B: Biol. Sci.* 280 (1768), 20131634.
- Guinat, C., Vergne, T., Kocher, A., Chakraborty, D., Paul, M.C., Ducatez, M., Stadler, T., 2021. What can phylodynamics bring to animal health research? *Trends Ecol. Evolut.* 36 (9), 837–847.
- Kao, R.R., Haydon, D.T., Lycett, S.J., Murcia, P.R., 2014. Supersize me: how whole-genome sequencing and big data are transforming epidemiology. *Trends Microbiol.* 22 (5), 282–291.
- Ladner, J.T., Grubaugh, N.D., Pybus, O.G., Andersen, K.G., 2019. Precision epidemiology for infectious disease control. *Nat. Med.* 25 (2), 206–211.
- Macdonald, D.W., Newman, C., Buesching, C.D., Johnson, P.J., 2008. Male-biased movement in a high-density population of the Eurasian badger (*meles meles*). *J. Mammal.* 89 (5), 1077–1086.
- Malone, K.M., Farrell, D., Stuber, T.P., Schubert, O.T., Aebersold, R., Robbe-Austerman, S., Gordon, S.V., 2017. Updated reference genome sequence and annotation of *Mycobacterium bovis* AF2122/97. *Genome Announcements* 5 (14), 10–1128.
- McDonald, J.L., Robertson, A., Silk, M.J., 2018. Wildlife disease ecology from the individual to the population: Insights from a long-term study of a naturally infected European badger population. *J. Anim. Ecol.* 87 (1), 101–112.
- Menardo, F., Duchêne, S., Brites, D., Gagneux, S., 2019. The molecular clock of *Mycobacterium tuberculosis*. *PLoS Pathogens* 15 (9), e1008067.
- Núñez-García, J., Downs, S.H., Parry, J.E., Abernethy, D.A., Broughan, J.M., Cameron, A.R., Cook, A.J., De La Rua-domenech, R., Goodchild, A.V., Gunn, J., et al., 2018. Meta-analyses of the sensitivity and specificity of ante-mortem and post-mortem diagnostic tests for bovine tuberculosis in the UK and Ireland. *Prevent. Vet. Med.* 153, 94–107.
- O'Brien, D.J., Kao, R.R., Little, R.A., Enticott, G., Riley, S.J., 2023. The road not traveled: Bovine Tuberculosis in England, Wales, and Michigan, USA. *One Health Cases*.
- O'Reilly, L.M., Daborn, C., 1995. The epidemiology of *Mycobacterium bovis* infections in animals and man: a review. *Tubercle Lung Dis.* 76, 1–46.
- Perea, C., Ciaravino, G., Stuber, T., Thacker, T.C., Robbe-Austerman, S., Allepuz, A., Val, B.P.d., 2021. Whole-genome SNP analysis identifies putative *Mycobacterium bovis* transmission clusters in livestock and wildlife in Catalonia, Spain. *Microorganisms* 9 (8), 1629.
- Pereira, A.C., Reis, A.C., Cunha, M.V., 2023. Genomic epidemiology sheds light on the emergence and spread of *Mycobacterium bovis* Eu2 Clonal Complex in Portugal. *Emerg. Microbes Infect.* 12 (2), 2253340.
- Price-Carter, M., Brauning, R., De Lisle, G.W., Livingstone, P., Neill, M., Sinclair, J., Paterson, B., Atkinson, G., Knowles, G., Crews, K., et al., 2018. Whole genome sequencing for determining the source of *Mycobacterium bovis* infections in livestock herds and wildlife in New Zealand. *Front. Veterinary Sci.* 5, 272.
- Pybus, O.G., Rambaut, A., 2009. Evolutionary analysis of the dynamics of viral infectious disease. *Nature Rev. Genet.* 10 (8), 540–550.
- R Core Team, 2022. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, URL <https://www.R-project.org/>.
- Reis, A.C., Salvador, L.C., Robbe-Austerman, S., Tenreiro, R., Botelho, A., Albuquerque, T., Cunha, M.V., 2021. Whole genome sequencing refines knowledge on the population structure of *Mycobacterium bovis* from a multi-host tuberculosis system. *Microorganisms* 9 (8), 1585.
- Rodrigues, R.d., Ribeiro Araújo, F., Rivera Dávila, A.M., Etges, R.N., Parkhill, J., van Tonder, A.J., 2021. Genomic and temporal analyses of *Mycobacterium bovis* in southern Brazil. *Microbial Genom.* 7 (5), 000569.
- Rogers, L., Cheeseman, C., Mallinson, P., Clifton-Hadley, R., 1997. The demography of a high-density badger (*Meles meles*) population in the west of England. *J. Zool.* 242 (4), 705–728.
- Rogers, L., Delahay, R., Cheeseman, C., Langton, S., Smith, G., Clifton-Hadley, R., 1998. Movement of badgers (*Meles meles*) in a high-density population: individual, population and disease effects. *Proc. R. Soc. London. Series B* 265 (1403), 1269–1276.
- Rossi, G., Crispell, J., Balaz, D., Lycett, S.J., Benton, C.H., Delahay, R.J., Kao, R.R., 2020. Identifying likely transmissions in *Mycobacterium bovis* infected populations of cattle and badgers using the Kolmogorov Forward Equations. *Sci. Rep.* 10 (1), 21980.
- Rossi, G., Crispell, J., Brough, T., Lycett, S.J., White, P.C., Allen, A., Ellis, R.J., Gordon, S.V., Harwood, R., Palkopoulou, E., et al., 2022. Phylodynamic analysis of an emergent *Mycobacterium bovis* outbreak in an area with no previously known wildlife infections. *J. Appl. Ecol.* 59 (1), 210–222.
- Rossi, G., Shih, B.B.-J., Egbe, N.F., Motta, P., Duchatel, F., Kelly, R.F., Ndip, L., Sander, M., Tanya, V.N., Lycett, S.J., et al., 2023. Unraveling the epidemiology of *Mycobacterium bovis* using whole-genome sequencing combined with environmental and demographic data. *Front. Veterinary Sci.* 10, 1086001.
- Sanjuán, R., Nebot, M.R., Chirico, N., Mansky, L.M., Belshaw, R., 2010. Viral mutation rates. *J. Virol.* 84 (19), 9733–9748.
- Santos, N., Colino, E.F., Arnal, M.C., de Luco, D.F., Sevilla, I., Garrido, J.M., Fonseca, E., Valente, A.M., Balseiro, A., Queirós, J., et al., 2022. Complementary roles of wild boar and red deer to animal tuberculosis maintenance in multi-host communities. *Epidemics* 41, 100633.
- Smith, G., Cheeseman, C., Wilkinson, D., Clifton-Hadley, R., 2001. A model of bovine tuberculosis in the badger *Meles meles*: the inclusion of cattle and the use of a live test. *J. Appl. Ecol.* 520–535.
- Sobkowiak, B., Banda, L., Mzembe, T., Crampin, A.C., Glynn, J.R., Clark, T.G., 2020. Bayesian reconstruction of *Mycobacterium tuberculosis* transmission networks in a high incidence area over two decades in Malawi reveals associated risk factors and genomic variants. *Microbial Genom.* 6 (4), e000361.
- Sobkowiak, B., Romanowski, K., Sekirov, I., Gardy, J.L., Johnston, J.C., 2023. Comparing *Mycobacterium tuberculosis* transmission reconstruction models from whole genome sequence data. *Epidemiol. Infect.* 151, e105.
- Trewby, H., Wright, D., Breadon, E.L., Lycett, S.J., Mallon, T.R., McCormick, C., Johnson, P., Orton, R.J., Allen, A.R., Galbraith, J., et al., 2016. Use of bacterial whole-genome sequencing to investigate local persistence and spread in bovine tuberculosis. *Epidemics* 14, 26–35.
- van Tonder, A.J., Thornton, M.J., Conlan, A.J., Jolley, K.A., Goolding, L., Mitchell, A.P., Dale, J., Palkopoulou, E., Hogarth, P.J., Hewinson, R.G., et al., 2021. Inferring *Mycobacterium bovis* transmission between cattle and badgers using isolates from the randomised badger culling trial. *PLoS Pathog.* 17 (11), e1010075.
- Vicente, J., Delahay, R., Walker, N., Cheeseman, C., 2007. Social organization and movement influence the incidence of bovine tuberculosis in an undisturbed high-density badger *Meles meles* population. *J. Anim. Ecol.* 76 (2), 348–360.
- Wolfram Research, Inc., 2023. *Mathematica*, version 12.3.1.0. Champaign, IL, URL <https://www.wolfram.com/mathematica>.
- Xu, Y., Cancino-Muñoz, I., Torres-Puente, M., Villamayor, L.M., Borrás, R., Borrás-Máñez, M., Bosque, M., Camarena, J.J., Colomer-Roig, E., Colomina, J., et al., 2019. High-resolution mapping of tuberculosis transmission: Whole genome sequencing and phylogenetic modelling of a cohort from Valencia Region, Spain. *PLoS Med.* 16 (10), e1002961.
- Xu, Y., Stockdale, J.E., Naidu, V., Hatherell, H., Stimson, J., Stagg, H.R., Abubakar, I., Colijn, C., 2020. Transmission analysis of a large tuberculosis outbreak in London: a mathematical modelling study using genomic data. *Microbial Genom.* 6 (11), e000450.