1

2

# Quantifying transmission of emerging zoonoses:

# Using mathematical models to maximize the value of surveillance data

5

6

7

8  Monique R. Ambrose[1], Adam J. Kucharski[2,3], Pierre Formenty[4], Jean-Jacques Muyembe-
9  Tamfum[5], Anne W. Rimoin[6], and James O. Lloyd-Smith[1,3*]

10

11

12

13  [1] Department of Ecology and Evolutionary Biology, University of California, Los Angeles

14  [2] London School of Hygiene and Tropical Medicine

15  [3] Fogarty International Center, National Institutes of Health, Bethesda, MD

16  [4] World Health Organization, Geneva, Switzerland

17  [5] Democratic Republic of the Congo National Institute for Biomedical Research

18  [6] Department of Epidemiology, University of California, Los Angeles

19

20

21  * Corresponding author

22  E-mail: jlloydsmith@ucla.edu (JOL-S)

1

## Abstract

23

24      Understanding and quantifying the transmission of zoonotic pathogens is essential for

25  directing public health responses, especially for pathogens capable of transmission between

26  humans. However, determining a pathogen's transmission dynamics is complicated by

27  challenges often encountered in zoonotic disease surveillance, including unobserved sources of

28  transmission (both human and zoonotic), limited spatial information, and unknown scope of

29  surveillance. In this work, we present a model-based inference method that addresses these

30  challenges for subcritical zoonotic pathogens using a spatial model with two levels of mixing.

31  After demonstrating the robustness of the method using simulation studies, we apply the new

32  method to a dataset of human monkeypox cases detected during an active surveillance program

33  from 1982-1986 in the Democratic Republic of the Congo (DRC). Our results provide estimates

34  of the reproductive number and spillover rate of monkeypox during this surveillance period and

35  suggest that most human-to-human transmission events occur over distances of 30km or less.

36  Taking advantage of contact-tracing data available for a subset of monkeypox cases, we find that

37  around 80% of contact-traced links could be correctly recovered from transmission trees inferred

38  using only date and location. Our results highlight the importance of identifying the appropriate

39  spatial scale of transmission, and show how even imperfect spatiotemporal data can be

40  incorporated into models to obtain reliable estimates of human-to-human transmission patterns.

## Author Summary

42      Surveillance datasets are often the only sources of information about the ecology and

43  epidemiology of zoonotic infectious diseases. Methods that can extract as much information as

44  possible from these datasets therefore provide a key advantage for informing our understanding

2

45 of the disease dynamics and improving our ability to choose the optimal intervention strategy.

46 We developed and tested a likelihood-based inference method based on a mechanistic model of

47 the spillover and human-to-human transmission processes. We first used simulated datasets to

48 explore which information about the disease dynamics of a subcritical zoonotic pathogen could

49 be successfully extracted from a line-list surveillance dataset with non-localized spatial

50 information and unknown geographic coverage. We then applied the method to a dataset of

51 human monkeypox cases detected during an active surveillance program in the Democratic

52 Republic of the Congo between 1982 and 1986 to obtain estimates of the reproductive number,

53 spillover rate, and spatial dispersal of monkeypox in humans.

## Introduction

55     Many recent infectious disease threats have been caused by pathogens with zoonotic

56 origins, including Ebola, pandemic H1N1 influenza, and SARS- and MERS- Coronaviruses, and

57 zoonotic pathogens are expected to be a primary source of future emerging infectious diseases

58 [1–8]. By definition, zoonotic pathogens can transmit from animals to humans; those also

59 capable of human-to-human transmission are of particular public health concern [5,9]. Infectious

60 disease surveillance serves a crucial role for detecting and gathering information on zoonotic

61 pathogens: data obtained through surveillance are often the primary resource available for

62 informing public health management decisions [10]. Developing methods that improve our

63 ability to infer information about a pathogen's transmission dynamics from available

64 surveillance data is therefore an essential frontier for understanding and ultimately combating

65 these pathogens [11,12].

66      For zoonoses, three epidemiological measures are crucial for summarizing transmission

67      dynamics and informing risk assessments. The first of these is the spillover rate, which indicates

68      how frequently the pathogen is transmitted from the animal reservoir into humans and helps

69      inform the total expected disease incidence [13]. The second measure describes the pathogen's

70      potential for further spread once in the human population and is commonly assessed using the

71      reproductive number ($R$), which gives the average number of secondary human cases caused by

72      an infectious individual [14,15]. Values of $R$ greater than one indicate that the pathogen is

73      capable of sustained (i.e. 'supercritical') transmission in humans. Pathogens with subcritical

74      transmission ($R$ less than one but greater than zero) can cause limited chains of transmission in

75      humans after a zoonotic introduction, and they pose a risk of acquiring ability for supercritical

76      transmission via evolutionary or environmental change [2,5,16]. The third epidemiological

77      measure is the distance over which human-to-human transmission occurs, which informs how

78      the disease will spread spatially and the risk of it being introduced into new populations.

79      Combined, these three measures can help evaluate the current public health threat posed by the

80      pathogen, the risk of future emergence, and the most effective approaches for disease

81      management.

82      Estimating epidemiological measures is a challenging task in any pathogen system, and

83      the unique properties of zoonotic diseases can exacerbate these difficulties. Infectious disease

84      surveillance often records temporal information and certain aspects of spatial information about

85      human cases, but the underlying transmission events are seldom observed. In a zoonotic system,

86      this means that an observed human infection could have been caused by a previous human case

87      or by zoonotic spillover. Without intensive contact tracing, or sequence data in the case of fast-

88      evolving pathogens, quantifying the relative contribution of zoonotic versus human-to-human

89  transmission is a major challenge; identifying the source of infection for specific individuals is

90  an even bigger one.

91      Epidemiological analyses are often hindered by data truncation and unknown

92  denominators [17,18]. In many disease surveillance systems, the total set of localities under

93  surveillance (i.e. those that would appear in the dataset if a case occurred there) can be separated

94  into 'observed localities,' which appear in the dataset because they reported one or more cases,

95  and 'silent localities,' which have no cases during the surveillance period and therefore do not

96  appear in the dataset. This form of truncation, where localities with zero cases are absent from

97  the dataset, obscures the true scope of the surveillance effort. Without knowledge of the total

98  number of localities under observation (the 'unknown denominator'), accurately estimating the

99  spillover rate and probability of human-to-human transmission between localities is not

100  straightforward. Simply disregarding these silent localities in the analysis is the functional

101  equivalent of selectively removing zeros from the dataset and can lead to problematic inference

102  biases.

103      Complicating inference efforts further is the fact that surveillance datasets often report

104  the geographic location of cases only at a coarse resolution, obscuring information about a

105  transmission process that occurs on a much finer scale [19–21]. Precise spatial information is

106  often absent from historic datasets and data collected in remote or low-resource areas, replaced

107  by the names of the locality and broader administrative units where the case occurred. For

108  example, only the village name and the region and country to which the village belongs may be

109  recorded in a dataset. Furthermore, linking a village name to spatial coordinates is often

110  impossible when maps of the region do not exist or only unofficial local names are used.

111  Although collecting exact spatial coordinates has become more practical in contemporary disease

5

5
bioRxiv preprint doi: https://doi.org/10.1101/677021; this version posted June 19, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

112 surveillance, privacy and confidentiality concerns can arise in both human and agricultural

113 contexts when data contains high-resolution spatial information [19,20,22–25], leading to data

114 being reported in a non-localized manner. Methods that can use this inexact spatial information

115 are especially needed for zoonotic diseases, where any additional information about the

116 proximity of human cases to one another can improve the power to distinguish between human-

117 to-human transmission and zoonotic spillover.

118     Despite these challenges, a series of research efforts have expanded our ability to

119 estimate the transmission properties of zoonotic pathogens from case onset data. A key set of

120 methods revolve around inferring $R$ from the sizes of case clusters (a cluster is defined as a group

121 of cases that occur in close spatiotemporal proximity to one another) or from the proportion of

122 observed cases that were infected by zoonotic spillover [16,26–30]. However, these approaches

123 either require detailed case investigations to determine whether a case was infected by a zoonotic

124 or human source or assume that each cluster is caused by one single spillover event followed by

125 human-to-human transmission. A likelihood-based approach for estimating $R$ for human-to-

126 human transmission using only symptom onset dates of cases was introduced by Wallinga and

127 Teunis [31]. This method was extended to apply to zoonotic systems by Lo Iacono et al. [32], but

128 the extension requires that chains of exclusively human-to-human transmission can be identified,

129 and is thus not applicable to many zoonotic surveillance systems where human and zoonotic

130 transmissions are intermixed. A different approach was taken by White and Pagano [33], who

131 introduced a different likelihood-based method that compares the observed number of cases on

132 each day with the expected number, as calculated using the number and timing of previous cases.

133 Though the White and Pagano approach was only applicable to human-to-human transmission, it

134 was expanded by Kucharski et al. [34] to work in zoonotic spillover systems in scenarios where a

6

135    control measure, implemented at a known point in time, causes an abrupt reduction in spillover.

136    A related approach that requires knowledge of the human and animal reservoir population sizes

137    was also explored in Lo Iacono et al. [35]. Crucially, however, none of these methods

138    incorporate information about the spatial location of cases to improve inference power or to

139    estimate patterns of spatial spread. Spatial data is a powerful tool in transmission inference in

140    single-species studies (e.g. [36–39]), but has largely been excluded from analyses of zoonotic

141    transmission, which often implicitly assume homogenous mixing across the study area or that

142    human-to-human transmission can only occur within a locality. One recent exception to this is

143    the analysis by Cauchemez et al. [40], which includes transmission at several spatial levels.

144        In this work, we present model-based inference methods that allow us to infer $R$, the

145    spillover rate, and properties of spatial spread among humans from surveillance datasets with

146    non-localized spatial information and an unknown total number of surveilled localities. Our

147    approach builds on methods introduced by White and Pagano [33] and Kucharski et al. [34], but

148    allows continuous spillover throughout the surveillance period and makes use of available spatial

149    information on case location. While the method could be readily adjusted to incorporate more

150    precise geographic information should it be available, in this study we focus on the more

151    challenging scenario in which only the names of the locality and broader administrative units

152    where a case occurred are known. To make use of this form of non-localized spatial data, our

153    model considers two scales of spatial mixing and transmission (Fig 1A), reminiscent of the

154    'epidemics with two levels of mixing' structure utilized in Ball et al. [41] and Demiris and

155    O'Neill [42]. The first mixing level is the locality in which the case occurred, such as a village,

156    conceptualized as a group of individuals geographically separated from other localities. We

157    assume that individuals within the same locality have more frequent contact with one another

7

158 than with individuals from other localities, and therefore that infection is more likely to be

159 transmitted within a locality. However, the total number of localities under surveillance is

160 unknown because only localities with one or more cases appear in the dataset (the 'unknown

161 denominator' problem discussed above). We refer to the second spatial level as the 'broader

162 contact zone.' It describes a collection of localities that all occur within the same administrative

163 unit and likely share some amount of human movement. When multiple types of administrative

164 units of different sizes are reported in the dataset (e.g., districts, regions, provinces, etc.), the

165 ideal choice for broader contact zone is the smallest administrative unit that contains inter-

166 locality human-to-human transmission events. If this scale is not known *a priori*, inferring the

167 appropriate scale of administrative unit is necessary.

168

169 **Fig 1. Model schematic. A.** The schematic illustrates the spatial scales considered in the model

170 and the types of transmission that occurs at different scales. Human cases are represented in

171 black if they were infected by zoonotic spillover, blue if they were infected by within-locality

172 human-to-human transmission, and orange if infected by between-locality human-to-human

173 transmission. Individuals who are not infected are colored gray and do not appear in the

174 surveillance dataset. Similarly, if zero individuals in a locality are infected, that 'silent locality'

175 does not appear in the dataset (represented by the gray locality in the broader contact zone). **B.**

176 The possible sources of human infection, which in aggregate determine the number of new

177 infections on day $t$, locality $v$. The number of cases arising from spillover and human-to-human

178 transmissions follow Poisson distributions with means $\lambda_Z$ and $\lambda_{\{s,w\},\{t,v\}}$, respectively.

179

180    We tested the method against a variety of datasets simulated using different

181    epidemiological parameters, offspring distributions for human-to-human transmission, and

182    spatial transmission kernels. To assess the performance of the method, we compared the

183    estimated and true values for epidemiological measures such as the reproductive number and

184    spillover rate, and also examined how well the method was able to estimate the probable

185    transmission source of each case. When silent localities were not accounted for, substantial

186    biases arose in zoonotic spillover rate estimates. However, a modified method that accounts for

187    these silent localities was successful in a wide range of circumstances. We therefore applied this

188    'corrected-denominator method' to a dataset on human monkeypox cases from an active

189    surveillance effort conducted in the Democratic Republic of the Congo (formerly Zaire) in the

190    1980s [43] (Fig 2). Gaining insights to the disease dynamics of human monkeypox is particularly

191    relevant given the recent increase in monkeypox incidence and outbreaks and the growing list of

192    countries and regions reporting human monkeypox cases [44–51]. Using the high-coverage

193    1980s surveillance dataset to quantify the pathogen's transmission dynamics will improve our

194    understanding of what drives its spread and lays the groundwork to assess what has changed over

195    the past decades to give rise to observed increases. With the 1980s monkeypox surveillance

196    dataset, we repeated the analyses using four different assumptions about the appropriate spatial

197    scale to represent the 'broader contact zone' over which human-to-human transmissions take

198    place and selected the preferred option using the deviance information criterion (DIC) method

199    for model comparison. In the monkeypox dataset, contact-tracing data are available for a subset

200    of the cases, providing a rare opportunity to compare inferred transmission sources with those

201    suggested by epidemiological investigation. In addition, some localities were associated with

202    known GPS coordinates, enabling us to estimate the spatial transmission kernel in greater detail.

203    As such, our monkeypox analysis yielded estimates of $R$ and the spillover rate during the 1980s

204    surveillance period, as well as insights into the spatial scale of human transmission of

205    monkeypox.

206

207    **Fig 2. Map and time-series showing locations and dates of human monkeypox cases.** The

208    size of points on the map indicate the number of cases and the color of points corresponds to the

209    region in which the cases occurred. Dark lines indicate region boundaries while light lines

210    indicate the official boundaries for districts (though in the monkeypox surveillance dataset these

211    are sometimes further divided into administrative subregions).

212

213    **Results**

214    **Overview of the approach**

215        We first validated the inference framework using a simulation study, then applied the

216    validated method to a dataset on human monkeypox cases to estimate key epidemiological

217    parameters and the spatial scale of transmission. To generate simulated test datasets and perform

218    parameter inference, we used a mathematical model of the zoonotic pathogen's transmission into

219    and among humans. The model tracks the number of human cases that occur in each locality on

220    each day; infections can arise from spillover from the zoonotic reservoir or from human-to-

221    human transmission (Fig 1B). Three key parameters govern the behavior of the system. The

222    spillover rate ($\lambda_z$) describes the average number of human cases caused by animal-to-human

223    transmission ('primary cases') in each locality per day. The reproductive number of the pathogen

10

224    (*R*) determines the average number of ('secondary') cases caused by each infected human. And

225    the spatial dispersal of the pathogen is controlled by the fraction of cases arising from human-to-

226    human transmission that occur in the same locality as the source case ($\sigma$) and the rules governing

227    inter-locality transmission events. Two spatial scales of transmission are included in the model:

228    within the locality of the case and between localities in the same broader contact zone. Using this

229    model (described further in Methods 4.1) and values for the three parameters, the likelihood of

230    observing $N_{t,v}$ cases on each day *t* and locality *v* can be calculated. Markov chain Monte Carlo

231    (MCMC) methods were used to infer posterior parameter distributions for a given dataset of

232    cases.

233    **Robustness of model-based inference method**

234    **Basic method (assumes the total number of localities under surveillance is known).** To

235    assess the accuracy and precision of our method's estimates of spillover and transmission

236    parameters, we simulated datasets with known parameter values and compared these true values

237    with the inferred values. We investigated a range of *R* and $\lambda_z$ values in the neighborhood of

238    values previously estimated for monkeypox [16,52], with *R* ranging from 0.2 to 0.6 and $\lambda_z$

239    ranging from 0.0001 to 0.0007 expected spillover events per locality per day ($\lambda_z$ values

240    correspond to 59 to 415 expected spillover events in the five year simulation period, across all

241    localities). Transmission events between humans had a probability $\sigma$=0.75 of occurring within a

242    locality and otherwise were equally likely between any localities in the same broader contact

243    zone. We were interested in seeing how well the inference methods are able to use the spatial-

244    temporal arrangement of cases to estimate the true parameter values.

245    Across 125 simulations (25 simulations for each of five parameter sets), estimated values

246    clustered around the true parameter values. The true value for $R$ was included in the 95%

247    credible interval (CI) 119 times (95.2%) and for $\lambda_z$ was included 121 times (96.8%) (Fig 3A). On

248    average, the posterior mean estimate of $R$ differed from the true value by 8.6%; the analogous

249    percent errors for $\lambda_z$ and σ estimates were 6.3% and 7.0%, respectively (S1 Table).

250

251    **Fig 3. Comparison of true and inferred parameter values in simulation study**. Within each

252    color, large points outlined in black indicate the true parameter set and smaller points indicate the

253    inferred parameter values from simulated datasets (lines show the 95% credible interval).

254    Inferences were performed **A)** when the true number of localities under surveillance was known,

255    **B)** when the true number was unknown and it was assumed that the number of observed

256    localities was the total number of localities, and **C)** when the true number of localities was

257    unknown and the corrected-denominator method was used to control for the locality observation

258    process.

259

260    However, this method assumes that the true number of localities under surveillance is

261    known. In real-world situations, 'silent' localities that experience zero cases often do not appear

262    in the dataset, resulting in an unknown true number of surveilled localities. We investigated

263    possible biases in parameter estimates that could arise from assuming that the number of

264    localities that reported one or more cases represents the total number of localities under

265    surveillance. To do so, we used the same set of simulated datasets as described above, but

266    removed knowledge about the number of silent localities. In these datasets, silent localities make

267    up between 21% and 85% of all localities under surveillance, with the proportion driven

268    primarily by the spillover rate. Estimates for the reproductive number $R$ were not strongly

269    impacted (95.2% of the 95% CIs contained the true value with an average percent error of 8.4%),

270    but the spillover rate $\lambda_z$ was consistently overestimated (Fig 3B). The true value for $\lambda_z$ was

271    contained in none of the simulations' 95% CIs and the posterior mean had an average percent

272    error of 153% (S1 Table).

273         To further investigate the effect of this data truncation (whereby localities with zero cases

274    do not appear in the dataset), we performed inference assuming that the observed localities

275    represented all, 1/2, or 1/5 of the total localities under surveillance. While this assumption had a

276    relatively small impact on the estimated $R$, it greatly impacted the inferred $\lambda_z$ (which is measured

277    as the number of spillover events *per locality* per day and is therefore strongly affected by

278    changes in the assumed number of localities) (S1 Fig). Assuming that a larger fraction of

279    surveilled localities appear in the dataset resulted in substantially higher estimated spillover

280    rates.

281    **Corrected-denominator method (conditions on the locality observation process).** Because

282    the total number of localities assumed to be under surveillance has a substantial impact on

283    parameter estimates, we developed a modified version of the likelihood function that accounts

284    for localities that were under surveillance but never observed in the dataset. This approach

285    calculates the likelihood of the observed dataset conditional on the fact that only localities with

286    one or more cases are included (details on the modified likelihood function can be found in

287    Methods and S1 Text).

13

288    We tested the performance of the corrected-denominator method against simulated

289    datasets, looking at the same parameter sets as in the first section. The inferred parameter values

290    cluster well with their corresponding true values (Fig 3C): mean percent error in $R$ estimates was

291    8.4% and in $\lambda_z$ estimates was 14.0%. Across the 125 simulations, the true parameter value was

292    included in the 95% CI 116 times (92.8%) for $R$ and 117 times (93.6%) for $\lambda_z$ (S1 Table).

293    Because an estimate of the true number of localities under surveillance would help

294    determine the size of the population that could be detected for a given system, we assessed how

295    well we could approximate this value. Given the number of localities with one or more cases and

296    the mean parameter estimates, it is possible to calculate the expected total number of localities

297    under surveillance (see S1 Text). Estimates of the true number of localities calculated for the

298    simulated datasets center on the correct value (S2 Fig). The magnitude of estimate error is driven

299    by the spillover rate, which largely determines the proportion of localities that are observed by

300    surveillance. The mean percent error across simulations with spillover rate of 0.0001, 0.00036,

301    and 0.0007 were 25.4%, 7.9%, and 2.4%, respectively, while simulations with spillover rates of

302    0.004 and above almost always recorded at least one case in each locality during the five year

303    surveillance period and therefore tended to estimate the exact true number of localities.

304    **Inferring the sources of transmission events.** We investigated how well sampled transmission

305    trees recovered the source of individual cases as well as higher-order measures, such as the

306    fraction of cases originating from zoonotic, within-locality, and between-locality transmission.

307    We tested our method using 125 simulated datasets, with 25 datasets simulated for each of five

308    sets of true parameter values (these are the same datasets as discussed above, simulated with $R$

309    between 0.2 and 0.6 and spillover rate between 0.0001 and 0.0007). Two hundred plausible

310    transmission trees were sampled for each simulated dataset.

14

311    When comparing the overall fraction of cases attributed to each source type (zoonotic

312    versus within-locality versus between-locality transmission), the sampled transmission trees

313    closely match the true transmission patterns (Fig 4). On average, the difference between the true

314    fraction of cases caused by zoonotic spillover and the fraction inferred in a tree was 0.022

315    (standard deviation 0.018), the difference for within-locality transmission was 0.006 (standard

316    deviation 0.005), and the difference for between-locality transmission was 0.022 (standard

317    deviation 0.018).

318

319    **Fig 4. Comparison of the true and inferred fraction of transmissions from each source type.**

320    For each of five parameter sets, 25 datasets were simulated and 200 transmission trees were

321    sampled for each of these simulated datasets. **A.** Stacked bars show the true fraction of

322    transmissions from zoonotic (bottom bar, medium-darkness), within-locality (middle bar, light

323    color), and between-locality (top bar, darkest color). Points on the bars indicate the inferred

324    values. If the fraction of transmissions for each source is perfectly inferred, points will lie exactly

325    on the transition between bar colors. **B.** Box plots summarize the error in the inferred fraction of

326    cases originating from each source type. The error size is small across all parameter sets,

327    especially for within-locality human-to-human transmission. The upper whisker was calculated

328    as $\min(\max(x), Q_3+1.5*IQR)$ and the lower whisker was calculated as $\max(\min(x), Q_1-1.5*IQR)$.

329

330    The success at recovering individual transmission links was high overall but varied

331    slightly depending on the true parameters underlying the simulation (S3 Fig). On average,

332    sampled transmission trees inferred 85.9% of all sources correctly. Better performance was

333     observed for lower spillover rates and lower $R$, presumably due to the fewer opportunities for

334     misattribution of cases. Some transmission links were more likely to be captured than others: on

335     average 90.9% and 90.1% of sampled trees correctly inferred links with zoonotic and within-

336     locality sources, respectively, but only 36.8% of trees correctly identified the source of between-

337     locality transmission events.

338     **Epidemiological insights into monkeypox**

339     **Applying the corrected-denominator method to 1980s monkeypox surveillance data.**

340     Between 1982 and 1986, the active monkeypox surveillance program in the Democratic

341     Republic of the Congo detected 331 human cases in 171 localities [43]. For each human case, we

342     know the name of the locality as well as the district or administrative subregion (henceforth

343     referred to simply as 'district') and region to which it belongs. However, the total number of

344     localities that would have been detected by surveillance had they experienced a case is unknown.

345     We therefore used the corrected-denominator method to generate estimates under four different

346     assumptions about which administrative unit most suitably represents the broader contact zone.

347     The country-level, region-level, and district-level models correspond to progressively smaller

348     choices of broader contact zones, while the locality-level model assumes that all instances of

349     human-to-human transmission occur within a locality. We anticipate that assuming an

350     excessively large broader contact zone could result in overestimating $R$ and underestimating $\lambda_z$ if

351     too many spurious human-to-human transmission events are inferred from pairs of cases that just

352     happen to occur within a generation-time interval of one another, while assuming an

353     inappropriately small broader contact zone could result in the opposite parameter biases if the

354     model is unable to detect actual incidents of human-to-human transmission because the cases

355     occur in different (assumed) broader contact zones.

16

356    In the monkeypox analysis, the size of the administrative unit used as the broader contact

357    zone has a strong effect on the resulting parameter estimates (Fig 5A). When larger

358    administrative units are assumed to represent the broader contact zone, a given pair of cases is

359    more likely to belong to the same broader contact zone, giving the model more opportunities to

360    infer inter-locality human-to-human transmission events and resulting in larger estimated

361    reproductive number $R$ and a smaller spillover rate $\lambda_z$. Mean values of the posterior distribution

362    of $R$ range from 0.29 when transmission is assumed to occur only within localities to 0.52 when

363    transmission is assumed to occur among all localities in the country (Table 1).

364

365    **Fig 5. Assumptions about the broader contact zone and the total number of localities under**

366    **surveillance affect parameter estimates for the monkeypox dataset.** Estimates and 95% CIs

367    for the reproductive number ($R$) and the spillover rate ($\lambda_z$) of the monkeypox dataset are shown

368    for each of the four choices of spatial scale for the broader contact zone (locality = green, district

369    = blue, region = purple, country = red). **A.** Inference performed using the corrected-denominator

370    method that accounts for silent localities. Light background dots are draws from the posterior,

371    larger dots designate the mean value, and bars indicate the 95% CI. **B.** Inference performed

372    assuming that the fraction of localities under surveillance with one or more monkeypox cases ($p$)

373    is 1/5, 1/2, or 1. For each assumption about the total number of localities, parameter estimates

374    fall roughly along the line $R = 1 - \frac{V*T*\lambda_z}{N}$ (indicated by grey lines), where $V$ is the true number

375    of localities under surveillance, $T$ is the duration of surveillance, and $N$ is to total number of

376    cases. The position of estimates along this line depends on the spatial model used. Note that the

377    slope of each line is proportional to -1/$p$ because $V$ = (number of observed localities) / $p$. Dots

17

378     represent the mean posterior estimates and bars indicate the 95% CI. The four darker dots show

379     the mean estimates from panel **A.**

380

381     **Table 1. District model performs best for the monkeypox dataset in DIC model**
382     **comparisons.**

| Approach for dealing with silent localities | Model | ΔDIC | mean $R$ | mean $\lambda_z$ | mean σ |
|---|---|---|---|---|---|
| Corrected-denominator method | Locality | 23.11 | 0.290 | 0.000387 | 1 |
| | **District** | **0.0** | **0.381** | **0.000309** | **0.696** |
| | Region | 5.88 | 0.418 | 0.000271 | 0.622 |
| | Country | 5.82 | 0.522 | 0.000188 | 0.464 |
| Assume all surveilled localities were observed | Locality | 21.98 | 0.272 | 0.000785 | 1 |
| | **District** | **0.0** | **0.372** | **0.000676** | **0.717** |
| | Region | 6.25 | 0.413 | 0.000633 | 0.656 |
| | Country | 10.92 | 0.479 | 0.000564 | 0.568 |
| Assume 1/2 of surveilled localities were observed | Locality | 17.06 | 0.290 | 0.000382 | 1 |
| | **District** | **0.0** | **0.381** | **0.000334** | **0.756** |
| | Region | 3.12 | 0.424 | 0.000311 | 0.684 |
| | Country | 6.79 | 0.488 | 0.000276 | 0.598 |
| Assume 1/5 of surveilled localities were observed | Locality | 15.05 | 0.310 | 0.000148 | 1 |
| | **District** | **0.0** | **0.395** | **0.000130** | **0.777** |
| | Region | 2.01 | 0.439 | 0.000121 | 0.704 |
| | Country | 5.34 | 0.500 | 0.000108 | 0.622 |

383     Parameter inference for the monkeypox dataset was performed using four different approaches

384     for dealing with the silent locality problem: the corrected-denominator method (which conditions

385     on the observation process for localities under surveillance) and three assumptions about the

386     fraction of localities under surveillance that were observed. For each of these approaches,

387     inference was repeated using four choices for the broader contact zone and the DIC was

388     calculated. Parameter estimates and ΔDIC values are shown. The model with lowest ΔDIC is

389     preferred and is shown in bold text.

18

390

391    We used the mean parameter estimates obtained using each of the four broader contact

392    zone assumptions to generate estimates of the expected total number of localities under

393    surveillance. While only 171 localities were observed in the dataset, estimates of the total

394    number of surveilled localities ranged from 337 (using the locality-level model) to 408 (using the

395    country-level model). The district-level and region-level models generated similar estimates of

396    351 and 366 total localities, respectively.

397    **Insights into how underlying assumptions drive monkeypox estimates**. We investigated how

398    different assumptions about the true number of localities and the spatial scale of human-to-

399    human transmission would affect the parameter estimates for the monkeypox system. To explore

400    how the presence of silent localities affects results, we repeated the analysis using the basic

401    method (which does not account for silent localities) under the assumption that the localities

402    observed in the monkeypox dataset represent all, 1/2, and 1/5 of the total number of localities

403    that were under surveillance. Furthermore, for each of these assumptions about the total number

404    of localities under surveillance, we repeated the analysis using the four different choices of

405    broader contact zone to determine how the assumed spatial scales of inter-locality transmission

406    impacted inference results.

407    Both the choice of broader contact zone and the assumed total number of localities have a

408    large impact on estimates of $R$ and $\lambda_z$ (Fig 5B). As noted above, models assuming smaller

409    broader contact zones allow fewer opportunities for human-to-human transmissions to be

410    inferred, and these models estimate substantially lower $R$ values and correspondingly higher

411    spillover rates. In contrast, assuming that a smaller fraction of surveilled localities were observed

19

412   leads to slightly higher estimates of $R$ and substantially lower estimates of $\lambda_z$ because the

413   presence of many silent localities drives the estimate of the number of spillover events *per*

414   *locality* per day lower. Estimates of $R$ are most strongly affected by the choice of broader contact

415   zone, while estimates of $\lambda_z$ are most strongly impacted by assumed fraction of localities

416   observed. For all assumptions of broader contact zone and total number of localities, the means

417   of the parameters' posterior distributions fall along the line

418   $$R = 1 - \frac{V*T*\lambda_z}{N} \quad , \tag{1}$$

419   where $V$ is the true number of localities under surveillance, $T$ is the number of days over which

420   surveillance occurred, and $N$ is to total number of cases in the monkeypox dataset. This

421   relationship arises because the expected number of total cases is equal to the expected number of

422   spillover events ($V * T * \lambda_Z$) multiplied by the total number of human cases expected to occur

423   from each spillover event ($1 / (1 - R)$ for $0<R<1$). Each assumption about the total number of

424   localities under surveillance corresponds to a separate line along which parameter estimates fall

425   (Fig 5B). The position of the parameter estimates along this line depends on the spatio-temporal

426   distribution of the $N$ cases and the assumed spatial scale of human-to-human transmission.

427   **District-level broader contact zone preferred in model comparisons**. To assess which broader

428   contact zone assumption is most appropriate for the monkeypox system, we used the deviance

429   information criterion (DIC) to perform model comparisons for the corrected-denominator

430   method as well as for each assumption about the number of surveilled localities. For the

431   corrected-denominator method, the district-level model had the best DIC score, followed by the

432   region and country-level models (Table 1). The locality-level model received a much larger DIC

433   value, indicating that the data strongly support models that allow transmission between localities.

434    Similarly, for each of the three assumptions about the true number of surveilled localities, the

435    district-scale model performed best in DIC model comparisons (Table 1).

436    **Inferring the sources and distances of transmission events**. We used the district-level

437    corrected-denominator method to sample 20,000 transmission trees for the monkeypox dataset.

438    The sampled transmission trees attributed an average of 60.8% (standard deviation of 2.2%) of

439    cases to zoonotic spillover, 28.5% (standard deviation of 0.9%) of cases to within-locality

440    human-to-human transmission, and 10.7% (standard deviation of 2.1%) of cases to between-

441    locality human-to-human transmission. For comparison, the results using the three other broader

442    contact zone assumptions are shown in S4A Fig. Each model's trees include a similar number of

443    within-locality human-to-human transmission events, but increasing the spatial scale of the

444    broader contact zone increases the number of inferred between-locality transmission events.

445         To characterize the distance range over which inter-locality transmission occurs, we

446    focused on links in the sampled transmission trees that occurred between cases with known GPS

447    coordinates (280 out of 331 monkeypox cases had recorded GPS coordinates). The number of

448    transmission events in each sampled tree that occurred over a certain distance was then compared

449    to the number of transmission events expected to occur over each distance if transmission

450    between all localities in a broader contact zone was equally likely (see Methods 4.3 for how this

451    'null distribution' was calculated).

452         For all models allowing inter-locality transmission, more transmission events were

453    inferred to occur across ≤ 30 kilometers than expected based on the null distribution (Fig 6, S4B

454    Fig). For each inferred transmission tree, a binomial test was used to examine whether more

455    transmissions were inferred to occur over ≤ 30 kilometers than expected based on the null

21

456 distribution of transmission distances. Out of 20,000 sampled trees for each model, p-values of

457 less than 0.1 were obtained in 93% of the district, 72% of the region, and 81% of the country-

458 level models' trees. The median p-values for these three models were 0.007, 0.030, and 0.012,

459 respectively (S5 Fig shows the full distributions of p-values obtained across all sampled trees).

460

461 **Fig 6. Distance of inferred inter-locality human-to-human transmission events.** Shaded bars

462 show the difference between the mean proportion of inter-locality human-to-human

463 transmissions inferred to occur over a given distance by the district model and the proportion

464 expected based on the spatial distribution of localities (the 'null expectation'). Error bars show

465 the standard deviation among all inferred transmission trees.

466

467 **Comparison of sampled transmission trees with contact-tracing data**. Contact-tracing, where

468 the contacts of a case were recorded and follow-up investigations determined whether or not the

469 contacts had become infected, was done for a subset of monkeypox cases. Instances where a

470 contact developed an infection are presumed to be instances of human-to-human transmission.

471 For each of these epidemiologically contact-traced links, we looked at how frequently the

472 sampled transmission trees for each model captured the transmission link.

473   Of the 53 case pairs linked through contact tracing, an average of 79.5% (standard

474 deviation of 4.2%) were recovered in each of the district model's sampled transmission trees (Fig

475 7A). The highest success was seen for pairs of epidemiologically-linked cases whose dates of

476 symptom onset were between 7 and 25 days apart (Fig 7B). Although it is generally believed that

477 the generation interval for human-to-human transmission of monkeypox is between 7 and 23

478  days [43,53], several case pairs that occurred more than 23 days apart were epidemiologically

479  linked through contact-tracing. It is possible that these links, which were often missed in the

480  sampled transmission trees, are not true instances of human-to-human transmission. Cases that

481  occurred in different localities were also less likely to be linked in a sampled transmission tree,

482  though even for these inter-locality pairs, the district-level model tended to perform better than

483  the other three models (S6 Fig). The four models had similar success at recovering within-

484  locality links. In all models, when a link was incorrectly inferred, it frequently was inferred to

485  originate from zoonotic spillover instead. Although the district model had the highest success at

486  recovering contact-traced links, the sampled trees from all models recovered an average of >76%

487  of contact pairs.

488

489  **Fig 7. Comparison of epidemiologically contact-traced links with sampled transmission**

490  **trees. A.** Circles (left axis) show the fraction of sampled trees that infer the epidemiologically-

491  traced source. Open circles represent inter-locality links while closed circles represent intra-

492  locality links. Crosses (right axis) indicate the probability that a link is instead inferred to have a

493  zoonotic source. Results are shown for the model assuming the district-level broader contact

494  zone. Links are sorted from lowest to highest success. **B.** The fraction of sampled transmission

495  trees that recover a contact-traced link is influenced by the number of days between the symptom

496  onset of source and recipient cases. Circles (left axis) show how often a given link was inferred

497  as a function of the generation interval while the gray curve (right axis) shows the probability

498  density for the generation interval assumed by the model.

499

23

500        Comparison of the transmission tree generated using only contact-tracing data with the

501        trees created using the district-level and locality-level models highlights how much our

502        perception of the transmission dynamics depends on assumptions about spatial spread (Fig 8).

503        Most of the within-locality transmission links detected through epidemiological contact-tracing

504        appear in the locality-level model's tree, though the locality-level tree suggests substantially

505        more human-to-human transmission events than captured in the contact-tracing tree. However,

506        the locality-level tree misses all inter-locality links. The district-level model's tree captures most

507        of the links indicated by the locality-level tree, and also suggests that inter-locality transmission

508        is occurring, though it has low power to determine exactly which case pairs are linked through

509        inter-locality transmission.

510

511        **Fig 8. Comparison of monkeypox transmission trees created from contact-tracing, the**

512        **locality-level model, and the district-level model**. Points represent cases and edges indicate

513        inferred transmission links between cases. Edge thickness corresponds to the frequency with

514        which a given transmission link was inferred while edge color indicates whether a pair of linked

515        cases occurred within the same (blue) or different (red) localities. The darkness of a point's fill

516        indicates how frequently the case was inferred to have a zoonotic source, so transmission links

517        often go from black points (cases caused by zoonotic spillover) to white points (cases infected by

518        a human source).

519

520    **Sensitivity analyses**

521        We conducted a variety of sensitivity analysis tests using simulated datasets to assess

522    how robust the method was over a range of parameter values and assumption violations (full

523    descriptions are provided in S1 Text). The method continued to perform well even at very high

524    spillover rates (S7 Fig, S2 Table) and when the offspring distribution used in simulations

525    differed from the one assumed in the inference (S8 Fig, S3 Table). In some situations, assuming

526    a larger broader contact zone than the one used for simulations could lead to an overestimation of

527    $R$ and an underestimation of $\lambda_z$ (S4 and S5 Tables). This outcome is consistent with what was

528    observed in the monkeypox analysis where assuming a larger spatial scale for the broader contact

529    zone corresponded to a higher estimate of $R$ and a smaller estimate of the spillover rate (Fig 5).

530    When simulations were run with highly structured, non-homogeneous spillover, substantial

531    biases in the inference results occurred because this spillover process gives rise to clusters of

532    primary cases that the model mistakes as arising from human-to-human transmission (S9 Fig).

533    **Discussion**

534    **Principal findings**

535        In this work, we developed and tested a method to infer fundamental epidemiological

536    parameters and transmission patterns for zoonotic pathogens from epidemiological surveillance

537    data with aggregated spatial information. When tested against simulated datasets, the method

538    successfully recovered estimates of $R$ and spillover rate close to the true values and also inferred

539    the fraction of cases resulting from zoonotic, within-locality, and between-locality sources with a

540    high degree of accuracy. The 'unknown denominator problem' that occurs when the total number

25

541    of localities under surveillance is unknown can cause large biases in parameter estimates, so we

542    modified the inference method to account for this observational process and enable unbiased

543    estimation in the presence of this common data gap.

544         We applied the method to a rich surveillance dataset of human monkeypox in the Congo

545    basin from the 1980s and found that human-to-human transmission of monkeypox between

546    localities plays an important role in the pathogen's spread. Of the four assumptions we tested for

547    the spatial scale of the broader contact zone, the district-level model was best supported by DIC

548    model comparisons and validation with contact-tracing. In addition, the signal of elevated inter-

549    locality transmission occurring over $\leq$ 30 kilometers suggests that most inter-locality

550    transmissions occur in a relatively small neighborhood, consistent with the limited transportation

551    infrastructure in the DRC. This further corroborates that the district-level model, which is the

552    smallest spatial aggregation scale that still permits inter-locality transmission, is likely the most

553    appropriate choice for capturing inter-locality transmission patterns of human monkeypox.

554         The district-level model estimates a reproductive number for human monkeypox of 0.38

555    (0.31-0.45 95% CI). This value is slightly higher than previous estimates of $R$ for the 1980s DRC

556    monkeypox dataset, which was estimated as 0.30 (90% CI 0.22-0.40) in Blumberg and Lloyd-

557    Smith [16], as 0.32 (90% CI 0.22-0.40) in Lloyd-Smith et al. [54], and as 0.28 in Jezek et al.

558    [52]. There are several explanations for the higher estimate we obtained. The previous studies

559    may have underestimated the reproductive number, particularly if contact-tracing or cluster

560    formation methods were liable to miss transmissions that occurred between localities. Indeed, the

561    estimate obtained using the locality-level model ($R = 0.29$) closely matches previous estimates. It

562    is also possible that the district-level model may overestimate the amount of human-to-human

563    transmission in the same way that the region- and country-level models picked up a higher signal

26

564 of human-to-human transmission than the district-level model due to their larger broader contact

565 zone sizes. The size of the DRC's districts and administrative subregions used for the district-

566 level model vary in size, but average around fifteen thousand square kilometers, or around one

567 hundred forty kilometers across, encompassing a much greater distance than most human-to-

568 human transmission events likely occur over. We therefore expect that the true value of $R$ is

569 bounded by the estimates of the locality-level and the district-level models.

570       In addition to providing an estimate of monkeypox's reproductive number, the methods

571 give insight into the frequency of spillover and the spatial scale of human-to-human

572 transmission. The district-level model estimates a mean spillover rate of around 0.11 spillover

573 events per locality per year, which corresponds to roughly one spillover event every nine years in

574 each locality. It also estimated that around 70% of human-to-human transmissions occur within a

575 locality. This finding contrasts with the assumption that human-to-human transmission occurs

576 within a locality, which is commonly used to generate transmission clusters, and suggests that

577 estimates generated using that assumption may substantially underestimate the amount of

578 human-to-human transmission occurring in the system. The importance of inter-locality contacts

579 has been reported for the neighboring country of Uganda, where a survey by le Polain de

580 Waroux et al. [55] on rural movement and social contact patterns indicated that 12% of social

581 contacts occurred outside participants' village of residence.

582       Among human monkeypox cases with recorded geographical coordinates, a clear signal

583 emerged of higher rates of human-to-human transmission between localities $\leq 30$ kilometers

584 apart. This pattern seems reasonable given the infrastructure and general difficulty of

585 transportation in the more remote regions of the DRC. It also suggests a similar pattern of

27

586     movement as found in the le Polain de Waroux et al. [55] survey. Their analyses indicate that

587     90% of people who traveled outside their village of residence remained within 12 km.

588     **Spatial scale of transmission and aggregated spatial data**

589         The potential biases introduced when analyzing data reported at a course spatial scale

590     have been explored in a wide range of contexts [56–58], yet the implications of using this type of

591     spatial information to infer the transmission dynamics of an infectious disease is not obvious.

592     When spatial information is only reported at the level of large spatial zones like districts, regions,

593     or countries, no finer-scale information is available to inform which human cases transmitted

594     infection to one another between different localities. Here we explored how the size of these

595     spatial zones would affect inference for the monkeypox system by repeating the analysis using

596     spatial information at the district, region, or country resolution. The large differences in

597     parameter estimates generated under different broader contact zone assumptions in the

598     monkeypox analysis illustrates how sensitive inference results can be to the spatial scale

599     assumed for human-to-human transmission, and suggests that reporting spatial data at too large a

600     scale or ignoring inter-locality transmissions can lead to substantial estimate biases.

601         In the context of monkeypox in the DRC, analysis of simulations using the exact

602     geographic coordinates reported for 80% of localities in the monkeypox surveillance dataset

603     replicated the increasing estimates of $R$ and decreasing estimates of spillover rate as the spatial

604     aggregation scale increased (S4 and S5 Tables). However, the magnitude of the effect in

605     simulated datasets was smaller than in the monkeypox analysis. This could be a result of the

606     particular assumptions about inter-locality transmission patterns used in the simulations, but it

607     also opens the question of whether outside large-scale factors such as seasonality or fluctuations

608    in surveillance effort might induce temporal autocorrelation among unlinked human cases,

609    giving rise to temporal clustering of cases that the model interprets as human-to-human

610    transmission.

611        This analysis serves to emphasize the importance of selecting an appropriate spatial scale

612    and using caution when interpreting results obtained using spatially aggregated data. Many

613    methods implicitly assume a certain scale of spatial transmission, often ignoring the possibility

614    of longer-range transmissions, so careful consideration of whether that scale is appropriate for

615    the system is essential.

616        In general, recording precise spatial locations of cases is vital for increasing the

617    inferential power of modeling analyses. Developing methods that maintain spatial information

618    without risking a breach in confidentiality is a nontrivial challenge, but progress has already been

619    made in generating possible solutions such as geographic masking or the verified neighbor

620    approach [59,60].

621    **Model assumptions and future directions**

622        In this work, we assumed that the spillover rate was homogenous through time and space,

623    but more complex disease dynamics in the reservoir or spatiotemporal heterogeneity in animal-

624    human contacts may cause nontrivial deviations from this assumption in real-world systems. Of

625    particular concern is the possibility that outbreaks in the reservoir could cause periods of

626    amplified local spillover, which could create a clustering pattern of human cases potentially

627    indistinguishable from human-to-human transmission. Without information about disease

628    dynamics in the reservoir, accounting for this heterogeneous spillover will be challenging, but

629 certain types of pathogen dynamics, such as seasonal epidemics or expanding wave-fronts of

630 infection, could be incorporated into the model structure.

631 Similarly, spatially and temporally variable surveillance intensity could potentially mimic

632 the signal of human-to-human transmission clusters and result in overestimates of the

633 reproductive number. Future surveillance programs could help mitigate this challenge by

634 recording a measure of surveillance effort undertaken at each location and time.

635 This work assumes that $R$ is constant across all localities; however, to obtain a full picture

636 of pathogen emergence risk, it may be necessary to consider the heterogeneity in transmission

637 intensity among different human populations, as well as the interplay between where $R$ is highest

638 versus where spillover tends to occur [61]. In some zoonotic systems, for instance, spillover

639 predominantly occurs into remote villages and towns that are in close proximity to forested

640 regions. However, we generally expect these villages to have lower levels of human-to-human

641 transmission than the more densely-packed cities [62–64]. A pathogen may even be incapable of

642 supercritical spread until it reaches such a city. Therefore, to assess the probability a pathogen

643 will successfully emerge and to determine which populations to target with control measures, it

644 may be necessary to establish not only the spillover rate and $R$ across different populations, but

645 also the rate of dispersal of the pathogen between those populations [61].

646 Several assumptions may need to be modified when applying this method to other

647 zoonotic systems. Because we assume that the source of human-to-human transmission events

648 will show symptoms before the recipient, the likelihood function can treat human cases as

649 occurring independently conditional on preceding cases. For zoonotic diseases in which infected

650 individuals frequently transmit the pathogen before showing symptoms (or when asymptomatic

651     cases contribute non-negligibly to transmission), the likelihood expression would need to be

652     modified substantially, and the lack of independence between cases might make a simulation-

653     based inference approach necessary.

654        We assume that sufficiently few infections occur relative to the population size that

655     depletion of susceptible individuals does not affect transmission dynamics. While appropriate

656     when there are few human infections or in the early stages of invasion, this assumption could

657     bias estimates if applied in a system with sufficiently high levels of human infection or where

658     transmission occurs primarily among highly clustered contacts, such as individuals within a

659     household. We also note that in the monkeypox example we are estimating the *effective*

660     reproductive number, which takes into account existing population immunity. If the goal instead

661     were to establish the basic reproductive number (the reproductive number for the pathogen in a

662     fully susceptible human population), accounting for past exposure to the pathogen or other cross-

663     immunizing pathogens or vaccines would be necessary.

664        The current methods assume that all human cases that occur during the surveillance

665     period inside the surveillance area are observed. This assumption is plausible for the analysis of

666     the 1980s monkeypox dataset, given the unusually high resources and experience level of this

667     surveillance effort in the aftermath of the smallpox eradication program and the use of serology

668     to detect missed cases retrospectively [43]. However, most zoonotic surveillance systems operate

669     with limited resources and have a much lower detection rate. Ignoring unobserved cases will lead

670     to underestimation of the spillover rate, while the effect on estimation of $R$ will depend on the

671     nature of the surveillance program. For instance, in the chain-size analyses of Ferguson et al.

672     [28] and Blumberg and Lloyd-Smith [16], $R$ is underestimated when the detection probability of

673     each case is independent of one another or when right-censoring occurs but overestimated when

31

674    a detected case triggers a retrospective investigation that detects all cases in that transmission

675    chain.

**Conclusions**

677    This work expands our ability to assess and quantify important zoonotic pathogen traits

678    from commonly available epidemiological surveillance data, even in the absence of exact spatial

679    information or a complete count of localities under surveillance. We anticipate that these

680    methods will have greatest value in the common circumstance when the source of cases,

681    particularly whether a case came from an animal or human source, cannot be readily established.

682    In such situations, the ability to infer the pathogen's reproductive number, spillover rate, and

683    spatial spread patterns from available surveillance data, will greatly enhance our understanding

684    of the pathogen's behavior and could provide valuable insights to help guide surveillance design

685    and outbreak response.

# Methods

**Model**

688    In broad terms, the model describes the probability of observing a set of symptom onset

689    times and locations of human cases given the timing and location of previous cases and

690    parameters that underlie the transmission process. Human infections can arise from either

691    animal-to-human transmission ('zoonotic spillover') or human-to-human transmission (Fig 1B).

692    Human-to-human contact occurs more frequently within a locality than between localities, but

693    can still occur between localities that belong to the same broader contact zone (Fig 1A).

694     All sources of infection are assumed to generate new cases independently of one another.

695     The number of human cases that become symptomatic on each day in each locality caused by

696     zoonotic spillover is assumed to follow a Poisson distribution with mean $\lambda_z$. For simplicity and

697     because reservoir disease dynamics are rarely well characterized, we assume the Poisson process

698     is homogenous through time and across localities, but this assumption could be modified for a

699     system where more information is available about the reservoir dynamics (e.g., [34]). New

700     infections can also arise from contact with infected humans. The number of new infections that

701     become symptomatic on day $t$ in locality $v$ caused by an infectious individual who became

702     symptomatic on day $s$ in locality $w$ is assumed to be a Poisson-distributed random variable with

703     mean $\lambda_{\{s,w\},\{t,v\}}$.

704     Aggregating cases caused by all sources of infection (both human and zoonotic), the total

705     number of new cases on day $t$ in locality $v$ is a Poisson-distributed random variable with mean

706
$$\mu_{t,v} = \sum_{s=1}^{t-1}\sum_{w=1}^{\mathcal{V}}[N_{s,w} * \lambda_{\{s,w\},\{t,v\}}] + \lambda_z \ , \tag{2}$$

707     where $\mathcal{V}$ is the number of localities under surveillance and $N_{s,w}$ is the number of cases with

708     symptom onset on day $s$ in locality $w$.

709     The mean of the Poisson random variable describing human-to-human transmission,

710     $\lambda_{\{s,w\},\{t,v\}}$, depends on the reproductive number of the pathogen in humans, the generation time

711     distribution, and the coupling between localities:

712
$$\lambda_{\{s,w\}\{t,v\}} = R * g(t-s) * H(v,w) \ , \tag{3}$$

713     where $R$ is the reproductive number of the pathogen; $g(t\text{-}s)$ is the generation time distribution,

714     which gives the probability that a secondary case becomes symptomatic $t\text{-}s$ days after the index

33

715   case shows symptoms; and $H(v,w)$ describes the amount of transmission between localities $v$ and

716   $w$ and takes values between zero (if no transmission can occur between localities $v$ and $w$) and

717   one (if all cases arising from an infected individual in locality $v$ arise in locality $w$). The

718   generation time $g(t\text{-}s)$ is assumed to follow a negative binomial distribution. For this study, we

719   used a mean of 16 days and a dispersion parameter of 728.7 (calculated by fitting a negative

720   binomial distribution to observed generation interval counts for smallpox presented in Fig. 2b of

721   [65]), which is consistent with previous estimates of the generation time for both smallpox and

722   monkeypox [43,53,65,66].

723       The factor that describes the amount of transmission that occurs between localities $v$ and

724   $w$ ($H(v,w)$) could reflect Euclidean distance, travel time, inclusion in different spatial zones, or

725   any other available measurement. To accommodate the imperfect spatial information available

726   for many zoonotic surveillance systems, this study focused on developing methods for the

727   situation when only a locality name and an aggregated spatial zone (such as district or country) is

728   reported for cases, rather than an exact position. We assume that inter-locality transmission

729   occurs only among localities within the same broader contact zone (Fig 1A). Because

730   transmission will be greater within a locality than between localities, a proportion $\sigma$ of secondary

731   cases are assumed to occur in the same locality as the source case and a proportion ($1\text{-}\sigma$) of

732   secondary cases are assumed to occur amongst the outside localities that are within the same

733   broader contact zone as the source case. This outside transmission is assumed to be divided

734   equally among all localities within the index case's broader contact zone:

735   $$H(v,w) = \begin{cases} 0 \;, \; Z_v \neq Z_w \\ \sigma \;, \quad v = w \\ \frac{(1-\sigma)}{(\mathcal{V}_v - 1)} \;, \quad Z_v = Z_w, v \neq w \end{cases} , \qquad (4)$$

34

736     where $Z_v$ indicates the broader contact zone of locality $v$ and $\mathcal{V}_v$ is the total number of localities

737     in the broader contact zone of locality $v$. For a given locality $v$, the sum of $H(v,w)$ across all $w$

738     equals one. To observe the effect of assuming different broader contact zones, the monkeypox

739     case study was repeated under four different assumptions about the spatial scale of human-to-

740     human transmission: locality, district, region, and country-level.


741     **Model inference**

742     **Likelihood function**. Using the model described above, a likelihood function was used to

743     evaluate a parameter set ($\theta = \{R, \lambda_z, \sigma\}$) given the data ($D = N_{t,v}$ cases observed on each day $t$ and

744     locality $v$):


745     $$\mathcal{L}(\theta|D) = \prod_{t=1}^{T} \prod_{v=1}^{V} \frac{e^{\mu_{t,v}} \mu_{t,v}^{N_{t,v}}}{N_{t,v}!} \; , \tag{5}$$


746     where $T$ is the number of days surveillance was conducted and $V$ is the total number of localities

747     under surveillance.


748         While this approach works well when the total number of surveilled localities is known

749     (see Fig 3A), localities often only appear in the dataset if they have reported cases; as a result we

750     may not know the total number of localities under surveillance. Ignoring localities with zero

751     cases can lead to biased parameter estimates (see Fig 3B). We explored several alternative

752     approaches to account for these silent localities; the preferred approach rescales the likelihood

753     function to reflect that localities with zero cases are not included in the data. Several

754     approximations are made in this approach to estimate unknown parameters and improve

755     computational tractability. The details of the derivation for the model are given in S1 Text, and

756     the final likelihood function is:


35

757
$$\mathcal{L}(\theta|D) = \prod_{w=1}^{W} \frac{\prod_{t=1}^{T} \frac{e^{\mu_{t,w}} \mu_{t,w}^{N_{t,w}}}{N_{t,w}!}}{\left(1 - e^{-\lambda_z T - \left(\frac{R\,T\,\lambda_z\,(1-\sigma)(E[V]-1)}{E[V]-1-R(E[V]-2+\sigma)}\right)}\right)} \, ,$$
(6)

758     where $W$ is the number of observed localities (localities with one or more cases) and $E[V]$ is the

759     expected number of localities given the parameter values and the number of observed localities.

760     **Parameter estimation**. Markov chain Monte Carlo (MCMC) was used to obtain the posterior

761     distributions of the model parameters. The fraction of transmissions occurring within a locality

762     ($\sigma$) and the reproductive number ($R$) were given uniform priors on zero to one. The expected

763     number of spillover events per locality per day ($\lambda_z$) was given a uniform prior with a lower bound

764     of zero and an upper bound selected to be far above the converged posterior distribution (ranging

765     from 0.0075 to 1, see S10 Fig for comparison of spillover priors and posterior distributions).

766         The chains were run for 100,000 steps, with a burn-in of 20,000. They satisfied visual

767     inspection for convergence. In addition, the Gelman and Rubin multiple sequence diagnostic was

768     evaluated for three parallel chains from each of the models for the monkeypox dataset [67]. The

769     Gelman-Rubin potential scale reduction values were less than 1.00033 across all models,

770     indicating that the chains have converged close to the target distribution [68].

771     **DIC model comparisons**

772         For the monkeypox dataset, four assumptions about the choice of broader contact zone

773     were compared using the deviance information criterion (DIC). This approach combines a

774     complexity measure, used to capture the effective number of parameters in each model, with a

775     measure of fit in order to perform model comparisons. Models are rewarded for better

776     'goodness-of-fit' to the data and penalized for increasing model complexity. Similarly to the

777     well-known Akaike information criterion (AIC) model comparisons, models with smaller DIC

778     values are preferred. As a rule of thumb, a difference between models' scores of four or more

779     generally indicates that the model with the larger value is 'considerably less' well supported by

780     the empirical evidence [69]. The values necessary to calculate the DIC can be readily obtained

781     from the MCMC output [70].

782     **Transmission tree reconstruction**

783         The origin of cases (zoonotic spillover, intra-locality human-to-human transmission, or

784     inter-locality human-to-human transmission) and the distances of inter-locality human-to-human

785     transmission events (when case localities are known) can be established given a particular

786     transmission tree. To gain estimates of these measures, trees were sampled based on the model

787     and the parameter posterior distributions. From the MCMC output (representing draws from the

788     posterior distribution), $d_1$ sets of parameter estimates were drawn to create $d_1$ transmission-

789     probability matrices ($P$). The entry $P_{ij}$ describes the probability that individual $i$ was infected by

790     individual $j$. The diagonal values of the matrix represent the probability a case originated from

791     zoonotic spillover. For a case $i$ observed to occur on day $t$ in locality $v$, the probability that case $j$

792     was the source of case $i$ ($P_{ij}$) was taken to be the proportion of $\mu_{t,v}$ (the expected total number of

793     cases on that day and locality; defined in equation 2) contributed by case $j$. By sampling $d_2$

794     transmission trees from each of these transmission-probability matrices, we calculated the

795     proportion of cases that resulted from spillover, within-locality transmission, and between-

796     locality transmission in each sampled tree. When testing the method using 125 simulated

797     datasets, 200 sampled transmission trees were generated for each dataset, with $d_1 = 20$ and $d_2$

798     =10. For the monkeypox dataset, 20,000 transmission trees were generated with $d_1 = 200$ and $d_2$

799     =100.

800       For inferred inter-locality human-to-human transmission events in the monkeypox

801    dataset, if the GPS coordinates were known for both localities in a transmission pair, the

802    transmission distance was calculated using the *gdist* function in the R package Imap [71]. The

803    'null distribution,' used for comparing the number of inferred inter-locality transmission events

804    with the number expected to occur if spatial location played no role in transmission, was

805    calculated by pooling all cases for which locality GPS coordinates are known, sampling all inter-

806    locality pairs permitted by the model, and recording the distance between the localities in each

807    pair.

808    **Simulation of test datasets**

809       To test the effectiveness of the methods, datasets with known parameter values were

810    simulated using the model explained above. Simulations were run over 1825 days

811    (approximately 5 years) and 325 surveilled localities. The localities were assumed to be

812    partitioned across thirty districts and six regions, with the distribution of localities across districts

813    and regions similar to that observed for the monkeypox dataset. The generation time interval (the

814    number of days between symptom onset of the source and recipient cases) was assumed to

815    follow a negative binomial distribution with a mean of 16 days and a dispersion parameter of

816    728.7 (as described above), with a maximum generation time interval of 40 days. A number of

817    parameter sets, as well as different underlying model structures, were used for simulations (S6

818    Table). Simulation parameters were chosen to approximate the monkeypox dataset, with $\sigma$ set at

819    0.75, $R$ ranging from 0.2 to 0.6, and $\lambda_z$ ranging from 0.0001 to 0.1. Unless otherwise specified,

820    simulations were performed assuming the district-level model. Details on the models used for

821    sensitivity analyses that use the exact spatial location of cases or allow highly structured and

822    non-homogenous spillover patterns are provided in S1 Text.

823 **Monkeypox data**

824      Data on human monkeypox cases in the Democratic Republic of the Congo (DRC),

825  formerly 'Zaire,' were collected as part of an intensive surveillance program supported by the

826  World Health Organization. During the peak surveillance period, between 1982 and 1986 [72],

827  data on 331 cases of laboratory-confirmed human monkeypox were recorded (see Fig 2, S1

828  Data) [43]. As part of field investigations, mobile teams visited the locality of a monkeypox case

829  to collect information about the case, such as the date of fever and rash onset (for this study, the

830  symptom onset date was taken to be the fever onset date; if the date of onset was not recorded,

831  the rash onset date was used instead), as well as to identify individuals who had had close contact

832  with the case [52,73]. If one of these contacts developed monkeypox within 7 to 21 days of first

833  exposure, the presumptive source case was recorded (S2 Data) [43,73].

834      Between 1982 and 1986, human monkeypox cases were observed in 171 distinct

835  localities, distributed among 30 districts and administrative subregions (simply referred to as

836  'districts') and 6 regions. The total number of localities that could have been detected by

837  surveillance is unknown. Of the 171 observed localities, GPS coordinates are available for 136

838  localities (which corresponds to 280 out of 331 cases). The district, region, and country of a

839  locality were always recorded.

840

# References

841     1.     Morse SS. Factors in the Emergence of Infectious Diseases. Emerg Infect Dis J [Internet].

843          1995;1(1):7. Available from: http://wwwnc.cdc.gov/eid/article/1/1/95-0102

844     2.     Woolhouse MEJ, Gowtage-Sequeria S. Host Range and Emerging and Reemerging

845          Pathogens. Emerg Infect Dis [Internet]. 2005 Dec;11(12):1842–7. Available from:

846          http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3367654/

847     3.     Xu R-H, He J-F, Evans MR, Peng G-W, Field HE, Yu D-W, et al. Epidemiologic Clues to

848          SARS Origin in China. Emerg Infect Dis [Internet]. 2004 Jun;10(6):1030–7. Available

849          from: http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3323155/

850     4.     Jones KE, Patel NG, Levy MA, Storeygard A, Balk D, Gittleman JL, et al. Global trends

851          in emerging infectious diseases. Nature [Internet]. 2008 Feb 21;451(7181):990–3.

852          Available from: http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5960580/

853     5.     Lloyd-Smith JO, George D, Pepin KM, Pitzer VE, Pulliam JRC, Dobson AP, et al.

854          Epidemic dynamics at the human-animal interface. Science [Internet]. 2009 Dec

855          4;326(5958):1362–7. Available from:

856          http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3891603/

857     6.     Gire SK, Goba A, Andersen KG, Sealfon RSG, Park DJ, Kanneh L, et al. Genomic

858          surveillance elucidates Ebola virus origin and transmission during the 2014 outbreak.

859          Science (80- ) [Internet]. 2014 Sep 11;345(6202):1369–72. Available from:

860          http://science.sciencemag.org/content/345/6202/1369.abstract

861     7.     Olival KJ, Hosseini PR, Zambrana-Torrelio C, Ross N, Bogich TL, Daszak P. Host and

862          viral traits predict zoonotic spillover from mammals. Nature [Internet]. 2017 Jun

863      29;546(7660):646–50. Available from: http://dx.doi.org/10.1038/nature22975

864    8.    Memish ZA, Cotten M, Meyer B, Watson SJ, Alsahafi AJ, Al Rabeeah AA, et al. Human

865      Infection with MERS Coronavirus after Exposure to Infected Camels, Saudi Arabia, 2013.

866      Emerg Infect Dis [Internet]. 2014 Jun;20(6):1012–5. Available from:

867      http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4036761/

868    9.    Woolhouse M, Gaunt E. Ecological Origins of Novel Human Pathogens. Crit Rev

869      Microbiol [Internet]. 2007 Jan 1;33(4):231–42. Available from:

870      https://doi.org/10.1080/10408410701647560

871    10.    National Research Council (US) Committee on Achieving Sustainable Global Capacity

872      for Surveillance and Response to Emerging Diseases of Zoonotic Origin; Keusch GT,

873      Pappaioanou M, Gonzalez MC, et al.  editors. Sustaining Global Surveillance and

874      Response to Emerging Zoonotic Diseases. Washingt Natl Acad Press [Internet]. 2009;

875      Available from: https://www.ncbi.nlm.nih.gov/books/NBK215317/

876    11.    Heesterbeek H, Anderson R, Andreasen V, Bansal S, De Angelis D, Dye C, et al.

877      Modeling infectious disease dynamics in the complex landscape of global health. Science

878      [Internet]. 2015 Mar 13;347(6227):aaa4339-aaa4339. Available from:

879      http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4445966/

880    12.    Lloyd-Smith JO, Funk S, McLean AR, Riley S, Wood JLN. Nine challenges in modelling

881      the emergence of novel pathogens. Epidemics [Internet]. 2015;10:35–9. Available from:

882      http://www.sciencedirect.com/science/article/pii/S1755436514000504

883    13.    Plowright RK, Parrish CR, McCallum H, Hudson PJ, Ko AI, Graham AL, et al. Pathways

884      to zoonotic spillover. Nat Rev Microbiol. 2017;15(8):502–10.

885     14.     Anderson RM, May RM. Infectious diseases of humans : dynamics and control. Oxford.

886             New York: Oxford University Press; 1991.

887     15.     Keeling MJ, Rohani P. Modeling Infectious Diseases in Humans and Animals. Princeton

888             Univ. Press; 2008. 366 p.

889     16.     Blumberg S, Lloyd-Smith JO. Inference of <italic>R</italic>$_0$ and Transmission

890             Heterogeneity from the Size Distribution of Stuttering Chains. PLoS Comput Biol

891             [Internet]. 2013;9(5):e1002993. Available from:

892             http://dx.doi.org/10.1371%2Fjournal.pcbi.1002993

893     17.     Kravitz HM. Denominator Difficulties. Wiley StatsRef Stat Ref Online. 2014;

894     18.     Tatem AJ. Mapping the denominator: spatial demography in the measurement of progress.

895             Int Health [Internet]. 2014 Sep 14;6(3):153–5. Available from:

896             http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4161992/

897     19.     Ozonoff A, Jeffery C, Manjourides J, White LF, Pagano M. Effect of spatial resolution on

898             cluster detection: a simulation study. Int J Health Geogr [Internet]. 2007;6(1):52.

899             Available from: https://doi.org/10.1186/1476-072X-6-52

900     20.     Zhang Z, Manjourides J, Cohen T, Hu Y, Jiang Q. Spatial measurement errors in the field

901             of spatial epidemiology. Int J Health Geogr [Internet]. 2016;15(1):21. Available from:

902             https://doi.org/10.1186/s12942-016-0049-5

903     21.     Dudas G, Carvalho LM, Bedford T, Tatem AJ, Baele G, Faria NR, et al. Virus genomes

904             reveal factors that spread and sustained the Ebola epidemic. Nature [Internet]. 2017 Apr

905             20;544(7650):309–15. Available from:

906             http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5712493/

907   22.   Curtis AJ, Mills JW, Leitner M. Spatial confidentiality and GIS: re-engineering mortality

908          locations from published maps about Hurricane Katrina. Int J Health Geogr [Internet].

909          2006;5(1):44. Available from: https://doi.org/10.1186/1476-072X-5-44

910   23.   National Research C. Putting People on the Map: Protecting Confidentiality with Linked

911          Social-Spatial Data. Stern PC, Gutman MP, editors. Washington, DC: National

912          Academies Press; 2007.

913   24.   Gutmann M, Witkowski K, Colyer C, O'Rourke JM, McNally J. Providing Spatial Data

914          for Secondary Analysis: Issues and Current Practices relating to Confidentiality. Popul

915          Res Policy Rev [Internet]. 2008;27(6):639–65. Available from:

916          http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2600804/

917   25.   de Montjoye Y-A, Hidalgo CA, Verleysen M, Blondel VD. Unique in the Crowd: The

918          privacy bounds of human mobility. 2013 Mar 25;3:1376. Available from:

919          http://dx.doi.org/10.1038/srep01376

920   26.   De Serres G, Gay NJ, Farrington CP. Epidemiology of Transmissible Diseases after

921          Elimination. Am J Epidemiol [Internet]. 2000 Jun 1;151(11):1039–48. Available from:

922          http://dx.doi.org/10.1093/oxfordjournals.aje.a010145

923   27.   Jansen VAA, Stollenwerk N, Jensen HJ, Ramsay ME, Edmunds WJ, Rhodes CJ. Measles

924          Outbreaks in a Population with Declining Vaccine Uptake. Science (80- ) [Internet]. 2003

925          Aug 7;301(5634):804 LP-804. Available from:

926          http://science.sciencemag.org/content/301/5634/804.abstract

927   28.   Ferguson NM, Fraser C, Donnelly CA, Ghani AC, Anderson RM. Public Health Risk

928          from the Avian H5N1 Influenza Epidemic. Science (80- ) [Internet]. 2004;304(5673):968–

929       9. Available from: http://www.sciencemag.org/content/304/5673/968.short

930   29.   Cauchemez S, Epperson S, Biggerstaff M, Swerdlow D, Finelli L, Ferguson NM. Using

931        Routine Surveillance Data to Estimate the Epidemic Potential of Emerging Zoonoses:

932        Application to the Emergence of US Swine Origin Influenza A H3N2v Virus. PLoS Med

933        [Internet]. 2013 Mar 5;10(3):e1001399. Available from:

934        http://dx.doi.org/10.1371%2Fjournal.pmed.1001399

935   30.   Blumberg S, Lloyd-Smith JO. Comparing methods for estimating R(0) from the size

936        distribution of subcritical transmission chains. Epidemics [Internet]. 2013 Sep

937        3;5(3):10.1016/j.epidem.2013.05.002. Available from:

938        http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3821076/

939   31.   Wallinga J, Teunis P. Different Epidemic Curves for Severe Acute Respiratory Syndrome

940        Reveal Similar Impacts of Control Measures. Am J Epidemiol  [Internet]. 2004 Sep

941        15;160(6):509–16. Available from:

942        http://aje.oxfordjournals.org/content/160/6/509.abstract

943   32.   Lo Iacono G, Cunningham AA, Fichet-Calvet E, Garry RF, Grant DS, Khan SH, et al.

944        Using Modelling to Disentangle the Relative Contributions of Zoonotic and

945        Anthroponotic Transmission: The Case of Lassa Fever. PLoS Negl Trop Dis [Internet].

946        2015 Jan 8;9(1):e3398. Available from:

947        http://dx.doi.org/10.1371%2Fjournal.pntd.0003398

948   33.   White LF, Pagano M. A likelihood-based method for real-time estimation of the serial

949        interval and reproductive number of an epidemic. Stat Med [Internet]. 2008 Jul

950        20;27(16):2999–3016. Available from:

951   http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3951165/

952 34. Kucharski A, Mills H, Pinsent A, Fraser C, Van Kerkhove M, Donnelly C, et al.

953   Distinguishing Between Reservoir Exposure and Human-to-Human Transmission for

954   Emerging Pathogens Using Case Onset Data. PLoS Curr Outbreaks. 2014;

955 35. Lo Iacono G, Cunningham AA, Fichet-Calvet E, Garry RF, Grant DS, Leach M, et al. A

956   Unified Framework for the Infection Dynamics of Zoonotic Spillover and Spread. Foley J,

957   editor. PLoS Negl Trop Dis [Internet]. 2016 Sep 2;10(9):e0004957. Available from:

958   http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5010258/

959 36. Keeling MJ, Woolhouse MEJ, Shaw DJ, Matthews L, Chase-Topping M, Haydon DT, et

960   al. Dynamics of the 2001 UK Foot and Mouth Epidemic: Stochastic Dispersal in a

961   Heterogeneous Landscape. Science (80- ) [Internet]. 2001;294(5543):813–7. Available

962   from: http://www.sciencemag.org/content/294/5543/813.abstract

963 37. Höhle M, Jørgensen E, O'Neill PD. Inference in disease transmission experiments by

964   using stochastic epidemic models. J R Stat Soc Ser C (Applied Stat [Internet].

965   2005;54(2):349–66. Available from: http:https://doi.org/10.1111/j.1467-

966   9876.2005.00488.x

967 38. Boender GJ, Hagenaars TJ, Bouma A, Nodelijk G, Elbers ARW, de Jong MCM, et al.

968   Risk Maps for the Spread of Highly Pathogenic Avian Influenza in Poultry. PLOS

969   Comput Biol [Internet]. 2007 Apr 20;3(4):e71. Available from:

970   https://doi.org/10.1371/journal.pcbi.0030071

971 39. Ypma RJF, Bataille AMA, Stegeman A, Koch G, Wallinga J, van Ballegooijen WM.

972   Unravelling transmission trees of infectious diseases by combining genetic and

973     epidemiological data. Proc R Soc B Biol Sci [Internet]. 2012;279(1728):444–50.

974     Available from: http://rspb.royalsocietypublishing.org/content/279/1728/444.abstract

975  40.  Cauchemez S, Nouvellet P, Cori A, Jombart T, Garske T, Clapham H, et al. Unraveling

976     the drivers of MERS-CoV transmission. Proc Natl Acad Sci  [Internet]. 2016 Aug

977     9;113(32):9081–6. Available from: http://www.pnas.org/content/113/32/9081.abstract

978  41.  Ball F, Mollison D, Scalia-Tomba G. Epidemics with Two Levels of Mixing. Ann Appl

979     Probab [Internet]. 1997;7(1):46–89. Available from: http://www.jstor.org/stable/2245132

980  42.  DEMIRIS N, O'NEILL PD. Bayesian inference for epidemics with two levels of mixing.

981     Scand J Stat [Internet]. 2005;32(2):265–80. Available from:

982     http:https://doi.org/10.1111/j.1467-9469.2005.00420.x

983  43.  Ježek Z, Fenner F. Human monkeypox [Internet]. S Karger Ag; 1988. Available from:

984     http://books.google.com/books?id=fyupMQEACAAJ

985  44.  Yinka-Ogunleye A, Aruna O, Ogoina D, Aworabhi N, Eteng W, Badaru S, et al.

986     Reemergence of Human Monkeypox in Nigeria, 2017. Emerg Infect Dis [Internet]. 2018

987     Jun;24(6):1149–51. Available from:

988     http://www.ncbi.nlm.nih.gov/pmc/articles/PMC6004876/

989  45.  Centers for Disease Control and Prevention. About Monkeypox [Internet]. 2018 [cited

990     2018 Oct 21]. Available from: https://www.cdc.gov/poxvirus/monkeypox/about.html

991  46.  Vaughan A, Aarons E, Astbury J, Balasegaram S, Beadsworth M, Beck CR, et al. Two

992     cases of monkeypox imported to the United Kingdom, September 2018. Eurosurveillance

993     [Internet]. 2018;23(38). Available from: https://doi.org/10.2807/1560-

994     7917.ES.2018.23.38.1800509%0A

46

995  47.  World Health Organization. Monkeypox - Nigeria [Internet]. Disease Outbreak News.

996       2018 [cited 2018 Oct 21]. Available from: http://www.who.int/csr/don/05-october-2018-

997       monkeypox-nigeria/en/

998  48.  World Health Organization. Monkeypox - Cameroon [Internet]. Disease Outbreak News.

999       2018 [cited 2018 Oct 21]. Available from: http://www.who.int/csr/don/05-june-2018-

1000      monkeypox-cameroon/en/

1001 49.  Rimoin AW, Mulembakani PM, Johnston SC, Smith JOL, Kisalu NK, Kinkela TL, et al.

1002      Major increase in human monkeypox incidence 30 years after smallpox vaccination

1003      campaigns cease in the Democratic Republic of Congo. Proc Natl Acad Sci [Internet].

1004      2010; Available from: http://www.pnas.org/content/early/2010/08/24/1005769107.abstract

1005 50.  Public Health England. Monkeypox case in England [Internet]. GOV.UK. 2018 [cited

1006      2018 Oct 21]. Available from: https://www.gov.uk/government/news/monkeypox-case-in-

1007      england

1008 51.  State of Isreal Ministry of Health. Monkeypox Patient Diagnosed [Internet]. State of Isreal

1009      Ministry of Health: Press Releases. 2018 [cited 2018 Oct 21]. Available from:

1010      https://www.health.gov.il/English/News_and_Events/Spokespersons_Messages/Pages/121

1011      02018_1.aspx

1012 52.  Ježek Z, Grab B, Szczeniowski M V, Paluku KM, Mutombo M. Human monkeypox:

1013      secondary attack rates. Bull World Health Organ [Internet]. 1988;66(4):465–70. Available

1014      from: http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2491159/

1015 53.  Fenner F, Henderson DA, Arita I, Jezek Z, Ladnyi ID. Smallpox and its Eradication.

1016      World Health Organization; 1988.

1017  54.  Lloyd-Smith JO, Schreiber SJ, Kopp PE, Getz WM. Superspreading and the effect of

1018        individual variation on disease emergence. Nature [Internet]. 2005;438(November):355–9.

1019        Available from: http://www.ncbi.nlm.nih.gov/pubmed/16292310

1020  55.  le Polain de Waroux O, Cohuet S, Ndazima D, Kucharski AJ, Juan-Giner A, Flasche S, et

1021        al. Characteristics of human encounters and social mixing patterns relevant to infectious

1022        diseases spread by close contact: a survey in Southwest Uganda. BMC Infect Dis

1023        [Internet]. 2018 Apr 11;18:172. Available from:

1024        http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5896105/

1025  56.  Clark WA V., Avery KL. The effects of data aggregation in statistical analysis. Geogr

1026        Anal. 1976;8:428–38.

1027  57.  Openshaw S, Taylor PJ. A million or so correlation coefficients: three experiments on the

1028        modifiable areal unit problem. Stat Appl Spat Sci. 1979;127–44.

1029  58.  Beale L, Abellan JJ, Hodgson S, Jarup L. Methodologic issues and approaches to spatial

1030        epidemiology. Environ Health Perspect. 2008;116:1105–10.

1031  59.  Richter W. The verified neighbor approach to geoprivacy: An improved method for

1032        geographic masking. J Expo Sci Environ Epidemiol [Internet]. 2017 Sep 20;28:109.

1033        Available from: http://dx.doi.org/10.1038/jes.2017.17

1034  60.  Zandbergen PA. Ensuring Confidentiality of Geocoded Health Data: Assessing

1035        Geographic Masking Strategies for Individual-Level Data. Adv Med [Internet].

1036        2014;2014. Available from: https://doi.org/10.1155/2014/567049

1037  61.  Sebastian J. Schreiber, James O. Lloyd-Smith. Invasion Dynamics in Spatially

1038        Heterogeneous Environments. Am Nat [Internet]. 2009;174(4):490–505. Available from:

1039        http://www.jstor.org/stable/10.1086/605405

1040  62.  Arita I, Wickett J, Fenner F. Impact of Population Density on Immunization Programmes.

1041        J Hyg (Lond) [Internet]. 1986;96(3):459–66. Available from:

1042        http://www.jstor.org/stable/3863139

1043  63.  Grenfell, Bolker. Cities and villages: infection hierarchies in a measles metapopulation.

1044        Ecol Lett [Internet]. 1998 Jul 1;1(1):63–70. Available from:

1045        http://dx.doi.org/10.1046/j.1461-0248.1998.00016.x

1046  64.  Neiderud C-J. How urbanization affects the epidemiology of emerging infectious diseases.

1047        Infect Ecol Epidemiol [Internet]. 2015 Jun 24;5:10.3402/iee.v5.27060. Available from:

1048        http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4481042/

1049  65.  NISHIURA H, EICHNER M. Infectiousness of smallpox relative to disease age: estimates

1050        based on transmission network and incubation period. Epidemiol Infect [Internet]. 2007

1051        Oct 7;135(7):1145–50. Available from:

1052        http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2870668/

1053  66.  Fine PEM. The Interval between Successive Cases of an Infectious Disease. Am J

1054        Epidemiol [Internet]. 2003 Dec 1;158(11):1039–47. Available from:

1055        http://dx.doi.org/10.1093/aje/kwg251

1056  67.  Gelman A, Rubin DB. Inference from Iterative Simulation Using Multiple Sequences. Stat

1057        Sci [Internet]. 1992;7(4):457–72. Available from:

1058        https://projecteuclid.org:443/euclid.ss/1177011136

1059  68.  Brooks SP, Gelman A. General Methods for Monitoring Convergence of Iterative

1060        Simulations. J Comput Graph Stat [Internet]. 1998 Dec 1;7(4):434–55. Available from:

1061    https://www.tandfonline.com/doi/abs/10.1080/10618600.1998.10474787

1062    69.    Burnham KP, Anderson DR. Model Selection and Multimodel Inference: A Practical

1063            Information-Theoretic Approach. 2nd ed. New York: Springer-Verlag New York, Inc;

1064            2002.

1065    70.    Spiegelhalter D, Best N, P Carlin B. Bayesian Deviance, the Effective Number of

1066            Parameters, and the Comparison of Arbitrarily Complex Models. Vol. 64, Journal of

1067            Royal Statistical Society. 1998.

1068    71.    Wallace JR. Imap: Interactive Mapping [Internet]. 2012. p. R package version 1.32.

1069            Available from: https://cran.r-project.org/package=Imap

1070    72.    Anonymous. The current status of human monkeypox: Memorandum from a WHO

1071            meeting. Bull World Health Organ [Internet]. 1984;62(5):703–13. Available from:

1072            https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2536211/pdf/bullwho00094-0031.pdf

1073    73.    Jezek Z, Marennikova SS, Mutumbo M, Nakano JH, Paluku KM, Szczeniowski M.

1074            Human Monkeypox: A Study of 2,510 Contacts of 214 Patients. Vol. 154, The Journal of

1075            infectious diseases. 1986. 551-555 p.

1076

1077

# Supporting information captions

1078

1079

1080 **S1 Text. Additional information on methods.** Supplementary text describing the corrected-

1081 denominator likelihood, the estimation of the total number of localities under surveillance, the

1082 simulation methods, and the sensitivity analyses.

1083

1084 **S1 Fig. Effect of assumed fraction of localities observed on parameter estimates.** The true

1085 parameter values are indicated by a large black dot and while smaller points indicate the inferred

1086 values from 25 simulated datasets (lines show the 95% credible interval). For each dataset,

1087 inference was performed assuming that 1/5, 1/2, and all of the localities under surveillance were

1088 observed. For these simulations, the true percentage of localities observed ranged from 46% to

1089 57%, with a mean of 52%.

1090

1091 **S2 Fig. Estimated number of localities under surveillance** (calculated given the number of

1092 observed localities and the estimated parameter values). Large colored dots indicate the

1093 estimated number of localities under surveillance for each simulated dataset while the smaller

1094 dots show the number of localities observed in the dataset. The true number of localities is

1095 represented by the horizontal dashed line. Each color corresponds to a different parameter set

1096 used for simulations.

1097

51

1098 **S3 Fig. Accuracy of inferred transmission trees at inferring the correct source of cases.** For

1099 each simulated dataset (25 simulations for each of 5 parameter sets), 200 transmission trees were

1100 drawn. Points show the mean fraction of cases inferred correctly in a sampled transmission tree

1101 and bars indicate the standard deviation.

1102

1103 **S4 Fig. Inferred sources of monkeypox cases. A.** The fraction of cases inferred to have

1104 originated from each source using each of the four spatial models (locality-green, district-blue,

1105 region-purple, country-red). **B.** Difference in the proportion of inter-locality human-to-human

1106 transmissions inferred by the models to occur over a given transmission distance versus expected

1107 based on the spatial distribution of localities. The p-values indicate the probability of observing

1108 as many or more transmissions over distances of ≤ 30 kilometers based on the null model (i.e.

1109 assuming distance plays no role in determining which localities are linked by inferred

1110 transmission events). The median p-value of sampled transmission trees is given, and the full

1111 distribution of p-values can be seen in S5 Fig.

1112

1113 **S5 Fig. The distribution of p-values obtained across sampled transmission trees**. P-values

1114 obtained from a binomial test examining whether the number of transmission events inferred to

1115 occur across thirty or fewer kilometers is greater than that expected based on the null

1116 distribution. Each p-value corresponds to a sampled transmission tree.

1117

1118    **S6 Fig. Comparison of epidemiologically contact-traced links with sampled transmission**

1119    **trees.** Circles (left axis) show the fraction of sampled trees that infer the epidemiologically-

1120    traced source. Open circles represent inter-locality links while closed circles represent intra-

1121    locality links. Bars (right axis) indicate the probability that a link is instead inferred to have a

1122    zoonotic source. Results are shown for models that use the country-level (red), region-level

1123    (purple), district-level (blue), and locality-level (green) broader contact zones. Links are sorted

1124    from lowest to highest success in the district model.

1125

1126    **S7 Fig. Effect of increasing spillover rate on parameter estimate success.** Within each color,

1127    large points outlined in black indicate the true parameter set and smaller points indicate the

1128    inferred parameter values from 25 simulated datasets (lines show the 95% credible interval).

1129    Warmer colors correspond with higher spillover rates. Note the log-scale x-axis.

1130

1131    **S8 Fig. Parameter estimate residuals for data simulated using a negative binomial versus**

1132    **Poisson offspring distribution.** Because the inference method assumes a Poisson offspring

1133    distribution, we compared the inference successes for datasets simulated assuming a Poisson

1134    offspring distribution versus datasets simulated assuming a negative binomial offspring

1135    distribution. The residuals in parameter estimates for 25 simulations are shown for **A)** the

1136    reproductive number and **B)** the spillover rate.

1137

1138    **S9 Fig. Strongly heterogeneous spillover causes bias in parameter estimates.** The true

1139    parameter value is indicated by the large dot while smaller points indicate the inferred values

1140    from 25 simulated datasets (lines show the 95% credible interval). Simulations were conducted

1141    to mimic pockets of zoonotic disease moving through the reservoir population. To capture the

1142    idea that, at any given time, only a small subset of localities might be experiencing high levels of

1143    spillover while the rest of the localities experienced no spillover, the simulations assumed that

1144    every 25 days a new set of three localities experienced the full force of spillover for the entire

1145    system. This gave rise to clusters of primary cases, which tend to be misclassified as human-to-

1146    human transmission events by our inference approach, which assumes homogeneous spillover

1147    rates.

1148

1149    **S10 Fig. Comparison of prior and posterior distributions for spillover rate $\lambda_z$.** Black bars

1150    represent posterior distribution while red lines mark limits of the uniform prior distribution. One

1151    representative simulation is shown for each of the nine parameter sets. Notice that the posterior

1152    distribution is always relatively far from upper bound of the prior.

1153

1154    **S1 Table. Comparison of inference method success over the same simulated datasets.**

1155

1156    **S2 Table. Success of the corrected denominator inference method for datasets simulated**

1157    **with increasing spillover rates.**

1158

1159 **S3 Table. Success of the corrected denominator inference method for datasets simulated**

1160 **with different offspring distributions.**

1161

1162 **S4 Table. Comparison of parameter estimates inferred using models of increasing spatial**

1163 **scale – data simulated using the 'nearest five neighbors' inter-locality transmission rule**

1164 where localities take the same GPS coordinates as in the DRC monkeypox surveillance dataset

1165 (true R is 0.36, true spillover rate is 0.00036; mean parameter estimates from inference on 25

1166 simulated datasets).

1167

1168 **S5 Table. Comparison of parameter estimates inferred using models of increasing spatial**

1169 **scale – data simulated assuming inter-locality transmission can occur between any localities**

1170 **located within 30 km of one another**, where localities take the same GPS coordinates as in the

1171 DRC monkeypox surveillance dataset (true R is 0.36, true spillover rate is 0.00036; mean

1172 parameter estimates from inference on 25 simulated datasets).

1173

1174 **S6 Table. Description of datasets simulated.**

1175

1176 **S7 Table. Parameter descriptions.**

1177

1178    **S1 Data. Case records.** For all individuals included in the analyses, records the case

1179    identification number, the locality identification number, the day of surveillance when disease

1180    onset occurred (the first day of fever when known, otherwise the first day of the rash), the names

1181    of the district and region where the case occurred, and masked GPS coordinates of the locality.

1182    The geographic masking technique known as 'donut masking' was used to obscure the exact

1183    location of cases and preserve privacy. For each locality with a recorded location, two random

1184    values were drawn: the first determines the direction and the second determines the distance

1185    from the original point. The new location is within 0.1 degrees from the original point but not

1186    closer than 0.02 degrees.

1187

1188    **S2 Data. Contact-tracing links.** Each row provides the case identification numbers for a pair of

1189    cases that was identified as a probable transmission link through epidemiological contact-tracing.

1190

1191

56

# Fig 1

**A**

**Locality**  **Broader contact zone**  **Surveillance area**



Individuals can be infected by zoonotic spillover or by previous human cases. Most human-to-human transmission occurs between individuals in the same **locality.**

Transmission can also occur between individuals in different localities if the localities are in the same **broader contact zone.** The total number of localities may be unknown if a locality only appears in the dataset when one or more cases occur there.

The **surveillance area** encompasses all localities monitored by surveillance, including the silent localities that do not appear in the dataset.

**B**  **Possible transmission sources of a case observed on day $t$, locality $v$:**



**Zoonotic spillover** causes an expected $\lambda_z$ new cases on day $t$, locality $v$

**Within-locality transmission:** an individual previously observed on day $s$ in locality $v$ causes an expected $\lambda_{\{s,v\},\{t,v\}}$ new cases on day $t$ in the same locality

**Between-locality transmission:** an individual previously observed on day $s$ in locality $w$ causes an expected $\lambda_{\{s,w\},\{t,v\}}$ new cases on day $t$, locality $v$

**Fig 2**

**Fig 3**

# Fig 4

**Fig 5**

**Fig 6**

**Fig 7**

# Fig 8



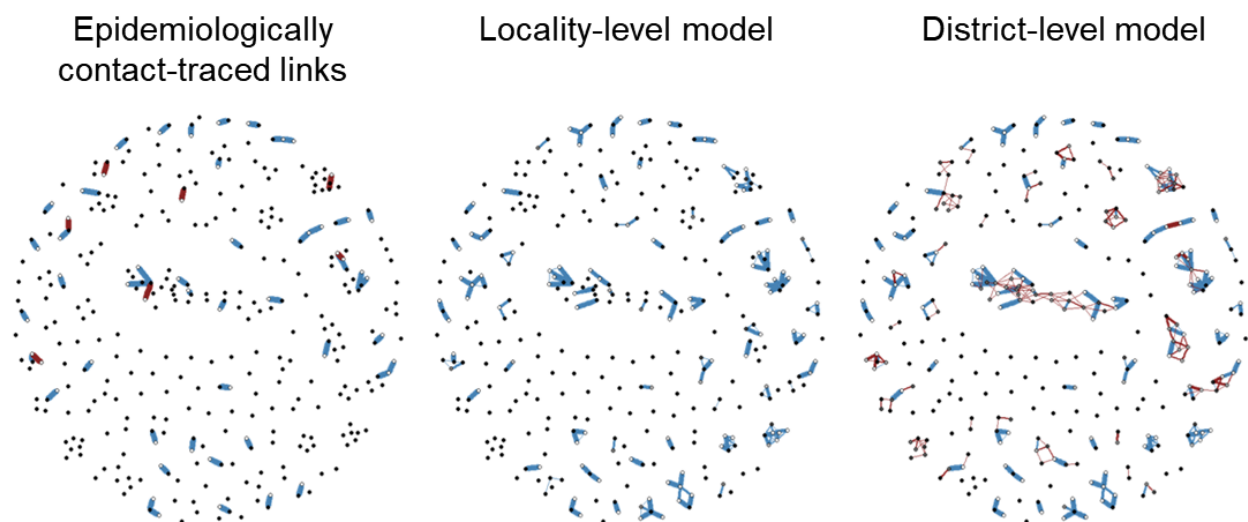Epidemiologically contact-traced links     Locality-level model     District-level model

# S1 Text. Supplementary material on methods

1192 **Corrected denominator method: Derivation for the conditional likelihood**

1193 **function**

1194     The model described in the main text tells us that the number of new human cases on day

1195 $t$ in locality $v$ follows a Poisson distribution with mean

1196
$$\mu_{t,v} = \sum_{s=1}^{t-1}\sum_{w=1}^{\mathcal{V}}[N_{s,w} * \lambda_{\{s,w\},\{t,v\}}] + \lambda_z \ , \tag{1}$$

1197 which represents the sum of the expected numbers of cases caused by spillover and all previous

1198 human cases (S7 Table provides a description of parameters). Based on this model, the

1199 likelihood of a set of parameters ($\theta = \{R, \lambda_z, \sigma\}$) given surveillance data ($D = N_{t,v}$ cases

1200 observed on each day $t$ and locality $v$) is:

1201
$$\mathcal{L}(\theta|D) = \prod_{t=1}^{T}\prod_{v=1}^{V}\frac{e^{\mu_{t,v}}\mu_{t,v}^{N_{t,v}}}{N_{t,v}!} \ . \tag{2}$$

1202     A challenge in applying this likelihood function to surveillance data arises when the total

1203 number of localities under surveillance, $V$, is unknown. Instead, we observe $W$ localities that

1204 have one or more observed cases. If we re-arrange the product functions in the likelihood

1205 function, it becomes more apparent that we are taking the product of the likelihood for each

1206 locality:

1207
$$\mathcal{L}(\theta|D) = \prod_{v=1}^{V}\prod_{t=1}^{T}\frac{e^{\mu_{t,v}}\mu_{t,v}^{N_{t,v}}}{N_{t,v}!} \ . \tag{3}$$

1208 However, because we only observe localities with one or more cases in the surveillance data, we

1209 need that conditioning to be reflected in the likelihood. In other words, we now want to express

1210 the likelihood of a particular time-series of cases in a locality *conditional on that locality having*

65

1211    *one or more cases.* This can be done for each locality by multiplying its component of the

1212    likelihood by the inverse of the probability ($q$) of having one or more cases:

1213    $$\mathcal{L}(\theta|D) = \prod_{w=1}^{W} \frac{\prod_{t=1}^{T} \frac{e^{\mu_{t,w}} \mu_{t,w}^{N_{t,w}}}{N_{t,w}!}}{q} \quad .$$    (4)

1214    It is now necessary to calculate the probability a surveilled locality experiences one or

1215    more cases. This probability is equivalent to one minus the probability of no cases occurring at a

1216    locality during the surveillance period. The following section explains how the probability of

1217    zero cases occurring at a given locality (here denoted $p$) is calculated.

1218    For zero cases to occur in a locality, there must be no zoonotic spillover into that locality

1219    as well as no human-to-human transmission from an outside locality. The zoonotic component is

1220    relatively straightforward to calculate, as it is simply the probability of zero spillover events on

1221    each of the $T$ days (which equals $e^{-\lambda_z T}$). The probability of no transmission from an outside

1222    human source is a bit more complicated and can be broken down by the generation of the outside

1223    case to avoid double-counting. The generation of a case indicates how many human-to-human

1224    transmission events occurred leading to the case. We refer to cases resulting from zoonotic

1225    spillover as primary cases. Individuals infected by primary cases are second generation cases,

1226    individuals infected by second generation cases are third generation cases, etc. For there to be no

1227    cases in a locality, no transmission may have occurred into that locality from outside cases in any

1228    generation:

66

$$P(no\ transmission\ from\ cases\ in\ other\ localities)$$

$$= P(no\ transmission\ from\ primary\ cases)$$

$$* P\left(\begin{array}{c}no\ transmission\ from\ second\ generation\ cases|no\ transmission\ from\ primary\\cases\end{array}\right)$$

$$* P\left(\begin{array}{c}no\ transmission\ from\ third\ generation\ cases|no\ transmission\ from\ primary\\or\ second\ generation\ cases\end{array}\right)$$

$$* \ldots$$

1229    The number of cases caused by a given case (of any generation) in the target locality is

1230    described by a Poisson distribution with expected value equal to $R\frac{(1-\sigma)}{(V_w-1)}$, where $V_w$ is the

1231    number of localities within the target locality's broader contact zone. Because each case

1232    transmits disease independently of one another (conditioned on the previous cases), the

1233    probability that no generation $i$ cases cause infections in the target locality is $e^{-R\frac{(1-\sigma)}{(V_w-1)}n_i}$, where

1234    $n_i$ is the total number of $i^{\text{th}}$ generation cases within the broader contact zone (given knowledge

1235    that none of the cases from previous generations transmitted to the target locality). Incorporating

1236    this information, the probability of observing zero cases in a locality ($p$) becomes:

1237    $$p = e^{-\lambda_z T} * \prod_{i=1}^{\infty} e^{-R\frac{(1-\sigma)}{(V_w-1)}n_i}$$

1238    $$= e^{-\lambda_z T} * e^{-R\frac{(1-\sigma)}{(V_w-1)}\sum_{i=1}^{\infty} n_i}.\qquad\qquad(5)$$

1239    We next need to calculate estimates for the expected values of each of the $n_i$. The

1240    expected number of primary cases in the entire broader contact zone (given that no spillover

1241    events occurred into the target locality) is the expected number of spillover events per locality

1242    ($\lambda_z$) multiplied by the number of localities under consideration ($V_w - 1$), multiplied by the

1243    number of surveillance days ($T$). For subsequent case generations, we can calculate the expected

1244    number of cases in generation $i+1$ as the number of cases caused by the $i^{\text{th}}$ generation in their

1245 own localities plus those caused in the $\mathcal{V}_w - 2$ other possible localities (there are $\mathcal{V}_w - 2$ other

1246 possible localities because the case's current locality and the target locality have already been

1247 counted):

1248
$$\mathbb{E}[n_{i+1}] = n_i \left( R\sigma + R \sum_{v=1}^{\mathcal{V}_w - 2} \frac{(1-\sigma)}{(\mathcal{V}_w - 1)} \right)$$

1249
$$= R \, n_i \frac{(\mathcal{V}_w + \sigma - 2)}{(\mathcal{V}_w - 1)}. \tag{6}$$

1250 If we approximate the values of $n_i$ with $\mathbb{E}[n_i]$, we get

1251
$$\mathbb{E}[n_{i+1}] \approx \lambda_z T(\mathcal{V}_w - 1) \left[ R \, \frac{(\mathcal{V}_w + \sigma - 2)}{(\mathcal{V}_w - 1)} \right]^i. \tag{7}$$

1252 Returning to our estimation of $p$, we can approximate $n_i$ values with $\mathbb{E}[n_i]$ and get

1253
$$p \approx e^{-\lambda_z T} * e^{-R\frac{(1-\sigma)}{(\mathcal{V}_w - 1)} \sum_{i=1}^{\infty} \mathbb{E}[n_i]}$$

1254
$$= e^{-\lambda_z T} * e^{-R\frac{(1-\sigma)}{(\mathcal{V}_w - 1)} \sum_{j=0}^{\infty} \lambda_z T(\mathcal{V}_w - 1) \left[ R \frac{(\mathcal{V}_w + \sigma - 2)}{(\mathcal{V}_w - 1)} \right]^j}$$

1255
$$= e^{-\lambda_z T} * e^{-R\frac{(1-\sigma)}{(\mathcal{V}_w - 1)} * \frac{\lambda_z T(\mathcal{V}_w - 1)}{1 - R\frac{(\mathcal{V}_w + \sigma - 2)}{(\mathcal{V}_w - 1)}}}$$

1256
$$= e^{-\lambda_z T} * e^{\frac{-R \lambda_z T(1-\sigma)(\mathcal{V}_w - 1)}{\mathcal{V}_w - 1 - R(\mathcal{V}_w + \sigma - 2)}}$$

1257
$$= e^{-\lambda_z T - \frac{R \lambda_z T(1-\sigma)(\mathcal{V}_w - 1)}{\mathcal{V}_w - 1 - R(\mathcal{V}_w + \sigma - 2)}}. \tag{8}$$

1258 With some additional algebraic simplification, we can insert this value in the original equation:

1259
$$\mathcal{L}(\theta|D) = \prod_{w=1}^{W} \frac{\prod_{t=1}^{T} \frac{e^{\mu_{t,w}} \mu_{t,w}^{N_{t,w}}}{N_{t,w}!}}{1 - e^{-\lambda_z T - \frac{R \lambda_z T(1-\sigma)(\mathcal{V}_w - 1)}{\mathcal{V}_w - 1 - R(\mathcal{V}_w + \sigma - 2)}}}. \tag{9}$$

68

1260    This expression still includes the parameter $\mathcal{V}_w$, though fortunately the sensitivity of results to the

1261    value of this parameter is relatively low. We therefore approximate $\mathcal{V}_w$ using the expected

1262    number of localities under surveillance in the broader contact zone. This calculation is explained

1263    in the following section.


## Estimating total number of localities under surveillance

1265        We wish to use the estimated parameter values for $R$, $\lambda_z$, and $\sigma$ in conjunction with the

1266    number of observed localities in a broader contact zone ($W_w$) to estimate the total number of

1267    localities under surveillance in that broader contact zone ($V_w$). If we let $q$ be the probability a

1268    locality is observed (has one or more cases during the surveillance period), then we expect $V_w * q$

1269    $\approx W_w$. From the section above, we approximate $q = 1-p$ as:

1270
$$q \approx 1 - e^{-\lambda_z T - \frac{R \lambda_z T (1-\sigma)(\mathcal{V}_w - 1)}{\mathcal{V}_w - 1 - R (\mathcal{V}_w + \sigma - 2)}}. \tag{10}$$

1271    So we estimate $V_w$ as the value that satisfies the equation:

1272
$$0 = \mathcal{V}_w \left( 1 - e^{-\lambda_z T - \frac{R \lambda_z T (1-\sigma)(\mathcal{V}_w - 1)}{\mathcal{V}_w - 1 - R (\mathcal{V}_w + \sigma - 2)}} \right) - W_w. \tag{11}$$


## Simulation methods

### Simulations with exact spatial locations

1275        Although the model assumes that inter-locality transmission with a broader contact zone

1276    is equal between all locality pairs, we expect that the actual amount of shared transmission

1277    between two localities is strongly influenced by the distance between those localities. We

1278    conducted two simulations using localities with set geographic locations and inter-locality

1279    transmissions depending on the spatial relationship of the localities. We took the 178 GPS

1280    records available from monkeypox surveillance in the DRC during the 1980s and simulated

1281    transmission across localities with the same coordinates and the same district and region

1282    boundaries. Two types of inter-locality transmission rules were explored. In the first of these,

1283    inter-locality transmissions were assumed to occur equally into a source locality's five closest

1284    neighbors. In the second set of simulations, inter-locality transmissions from a source locality

1285    were assumed to occur equally among all outside localities within 30 km of the source locality.

1286    **Simulations with highly structured and non-homogeneous spillover patterns**

1287    To illustrate how highly structured and non-homogeneous spillover could bias parameter

1288    estimates, we simulated an extreme case of a zoonotic epidemic traveling through time and

1289    space. We imagined that disease dynamics in the reservoir would occur in a single location for

1290    25 days before moving to a new spot, in an extreme form of a traveling zoonotic epidemic. For

1291    each 25 day period, three localities (selected to be in the same district when possible) would be

1292    selected to experience all of the spillover in the entire system. Aside from this extreme spillover

1293    pattern, the simulation followed the district-level model.

1294    **Sensitivity analyses**

1295    **Sensitivity of parameter inference to elevated or heterogeneous spillover**

1296    To test whether a high rate of spillover would inundate the system with so many cases

1297    that the temporal clustering patterns resulting from human-to-human transmission could be

1298    obscured, we simulated datasets with spillover rates up to 0.1. This value corresponds with an

1299    expected 59,312.5 spillover events during the five year simulation, which corresponds to an

70

1300    average of 36.5 per year in each locality. At this rate of spillover, there is an average of only ten

1301    days between spillover events, a shorter period than the mean generation time for human-to-

1302    human transmission events, which was sixteen days. Across the range of spillover rates tested,

1303    the method did very well at both point estimates and capturing the true parameter values within

1304    the 95% CI (an average of 94.3% of CIs included the true value of $R$ and 94.9% included the true

1305    value of $\lambda_z$; S7 Fig, S2 Table). As the spillover rate increased from 0.0001 to 0.1, estimates of $R$

1306    tended to improve (posterior means closer to true value and smaller CIs). While the absolute

1307    error on estimates of $\lambda_z$ increased as spillover rate increased, the relative error tended to decrease.

1308    As such, it appears that elevated spillover rates, far from obscuring patterns, may actually

1309    correspond with improved estimates, presumably due to the increased inference power resulting

1310    from a larger number of cases.

1311        Spillover is unlikely to occur homogeneously through time and space in real-world

1312    settings. As an illustration of the potential effect this occurrence could have on parameter

1313    estimates, we simulated an extreme case (see 'Simulations with highly structured and non-

1314    homogeneous spillover patterns,' above) where spillover occurs into three localities at a time.

1315    The parameter inference results for this situation were strongly biased (S10 Fig).

1316    **Sensitivity of parameter inference to offspring distribution assumptions**

1317        The model used in this study assumes that the number of new cases caused by an

1318    infectious individual follows a Poisson distribution, but previous work suggests that the offspring

1319    distribution is often better characterized by a negative binomial distribution, which allows for a

1320    greater amount of variation between individuals [1]. We simulated datasets using a negative

1321    binomial offspring distribution (using a dispersion parameter $k$=0.58 in accordance with previous

71

1322    estimates for monkeypox from [1]) and examined how well our inference method, which

1323    assumes a Poisson offspring distribution, estimated the true parameter values. Estimates for these

1324    datasets were only marginally less accurate than estimates for datasets generated with a Poisson

1325    offspring distribution (with an average percent error of 10.9% as opposed to 8.2% for $R$ and of

1326    11.6% as opposed to 10.4% for spillover rate estimates) (S8 Fig, S3 Table). As such, there are

1327    unlikely to be strong biases introduced from a mis-specified offspring distribution for the

1328    monkeypox dataset, though this bias could increase if applied to pathogens with more extreme

1329    transmission variance.

1330    **Sensitivity of parameter inference to broader contact zone assumption**

1331        To examine how assuming different broader contact zones would affect inference results,

1332    we compared parameter estimates obtained under three choices of broader contact zones for data

1333    simulated under two inter-locality transmission rules. We simulated disease spread in a system

1334    where localities were placed in the same arrangement as seen in 178 localities with GPS

1335    coordinates included in the monkeypox surveillance system, district and region arrangement

1336    were the same as in the 1980s surveillance, and human-to-human transmission could occur either

1337    between a locality and its five closest neighbors or between localities located within 30 km of

1338    one another. Inference results again showed increasing estimates of $R$ and decreasing estimates

1339    of spillover rate as the size of the assumed broader contact zone increased (S4 and S5 Table).
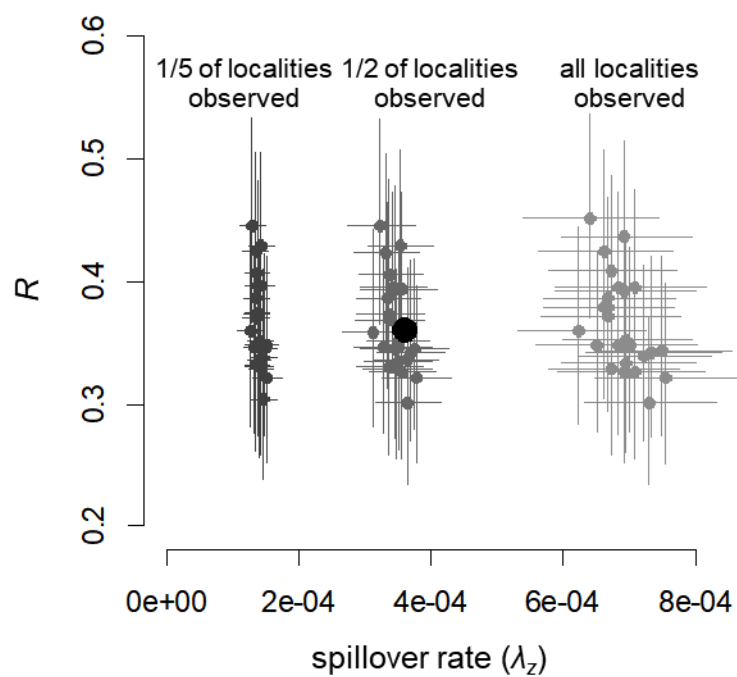
1340

## Supplementary material references

1.    Lloyd-Smith JO, Schreiber SJ, Kopp PE, Getz WM. Superspreading and the effect of individual variation on disease emergence. Nature [Internet]. 2005;438(November):355–9. Available from: http://www.ncbi.nlm.nih.gov/pubmed/16292310
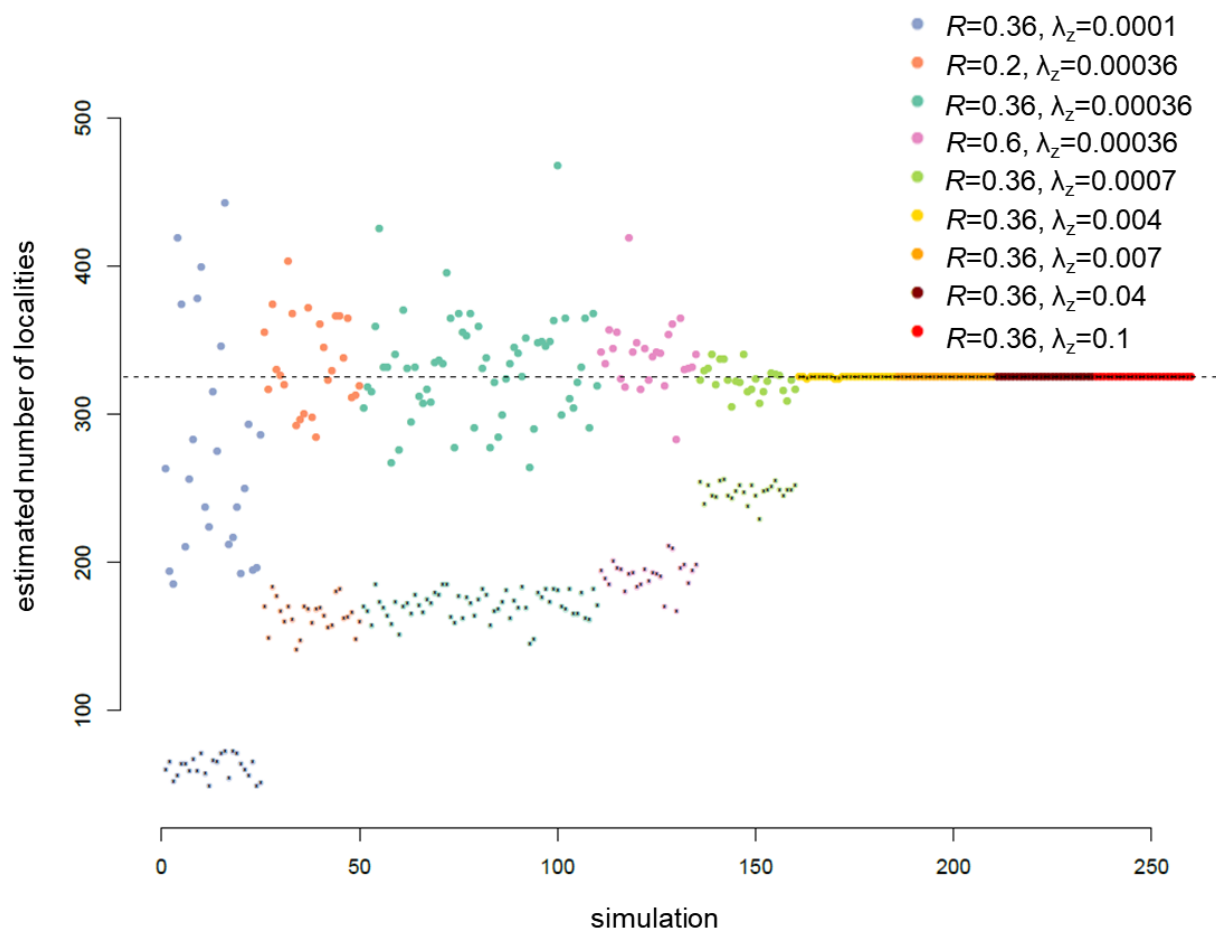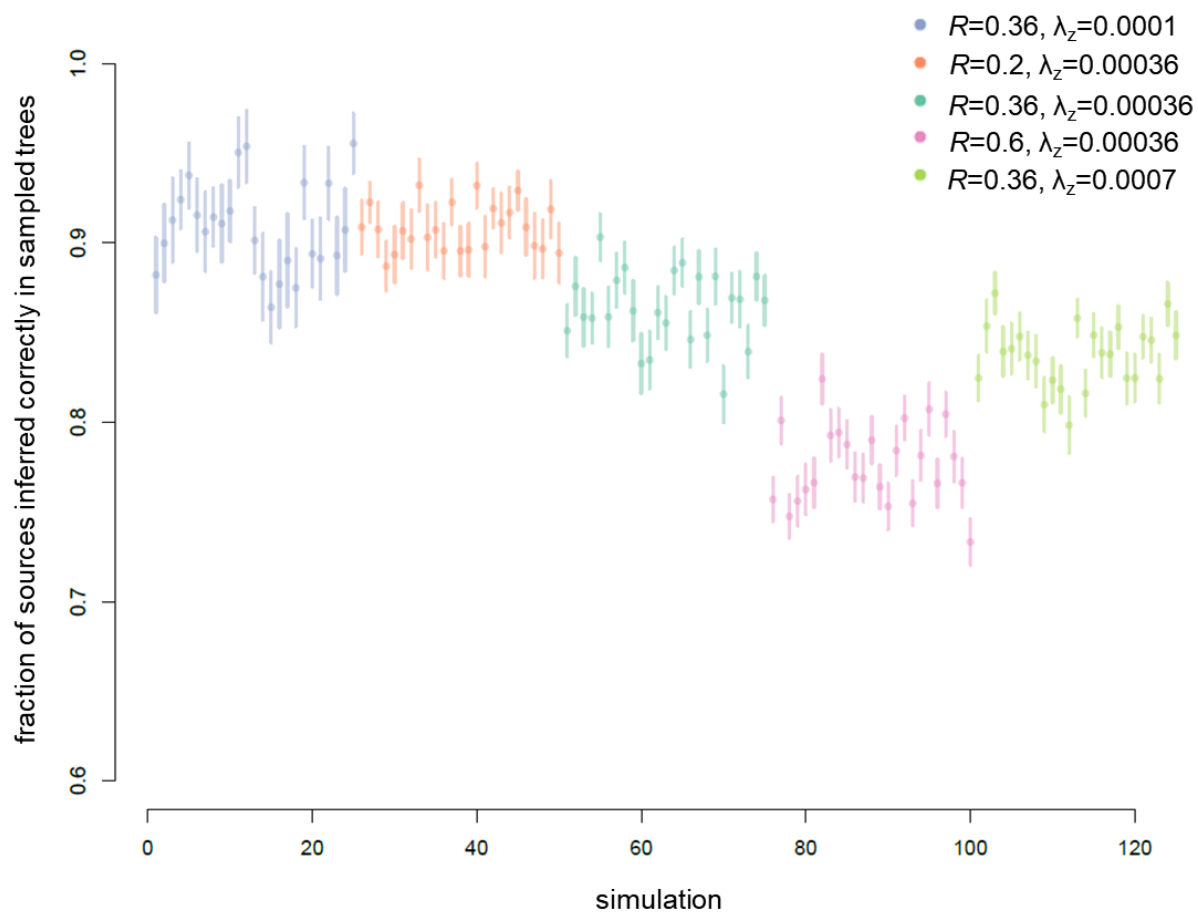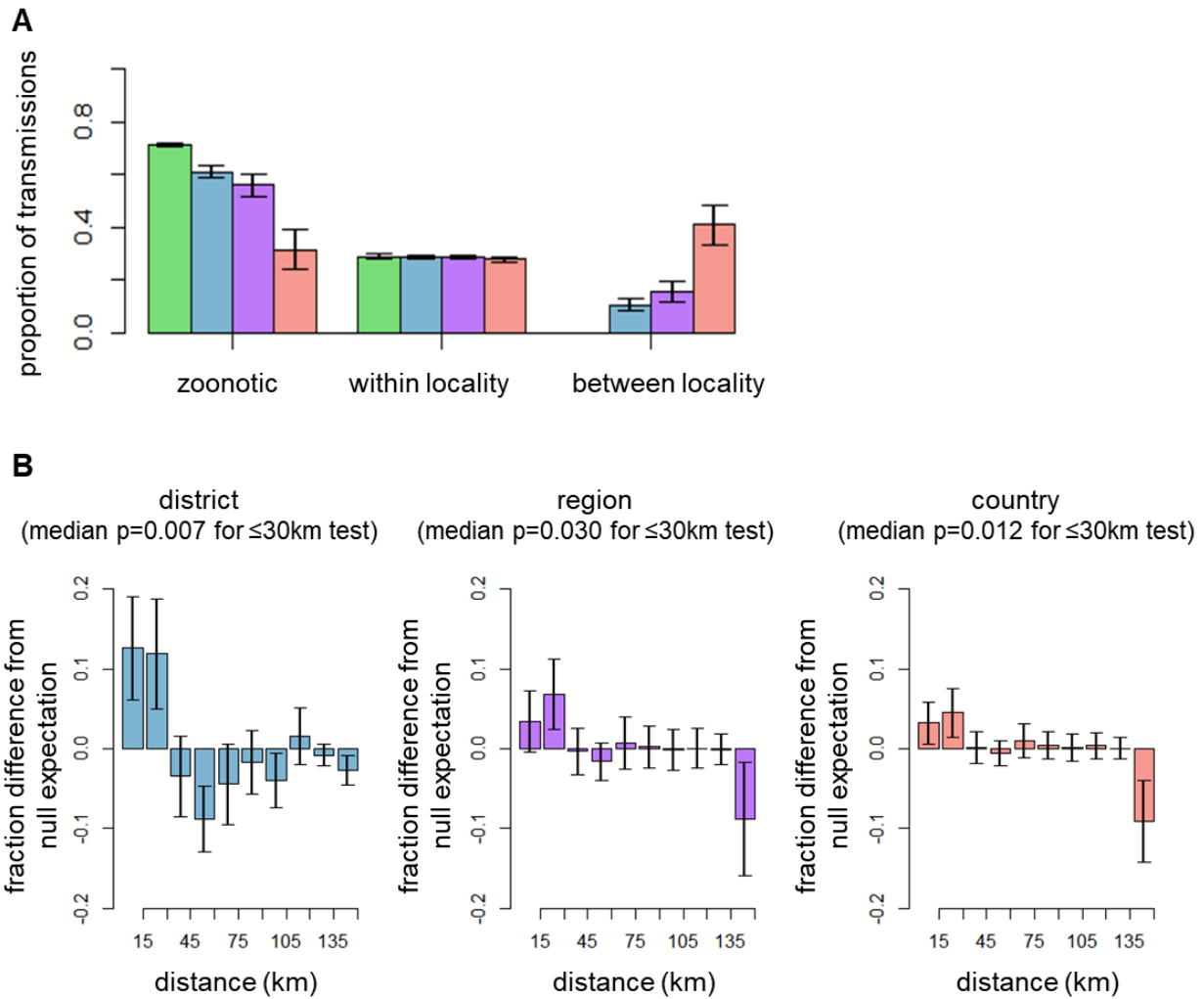
# S1 Fig

## S2 Fig

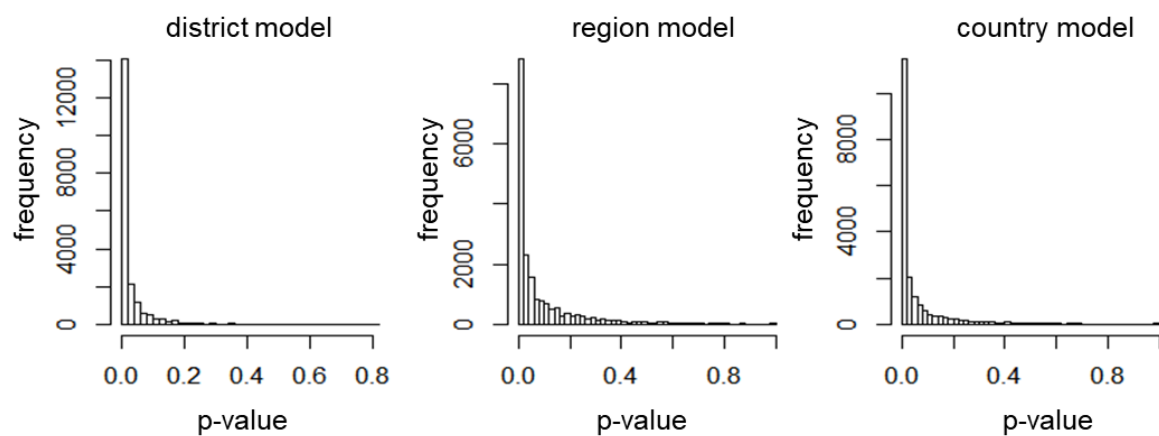# S3 Fig

# S4 Fig

# S5 Fig

# S6 Fig

## S7 Fig

## S8 Fig

## S9 Fig

# S10 Fig

# S1 Table

| Inference Approach | R | | | $\lambda_z$ | | | σ | | |
|---|---|---|---|---|---|---|---|---|---|
| | Fraction of 95%CIs include true value | Average error size | Average percent error | Fraction of 95%CIs include true value | Average error size | Average percent error | Fraction of 95%CIs include true value | Average error size | Average percent error |
| True number of localities known | 95.2% (119/125) | 0.0293 | 8.6% | 96.8% (121/125) | 1.99E-05 | 6.3% | 96.0% (120/125) | 0.0522 | 7.0% |
| Assume all localities are observed | 95.2% (119/125) | 0.0288 | 8.4% | 0.0% (0/125) | 3.30E-04 | 153.0% | 94.4% (118/125) | 0.0575 | 7.7% |
| Corrected denominator method (account for silent localities) | 92.8% (116/125) | 0.0298 | 8.4% | 93.6% (117/125) | 3.59E-05 | 14.0% | 88.0% (110/125) | 0.0665 | 8.9% |

# S2 Table

| True $\lambda_z$ value | $R$ | | | | $\lambda_z$ | | | | $\sigma$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Fraction of 95% CIs include true value | Average error size | Average percent error | Average width of CI | Fraction of 95% CIs include true value | Average error size | Average percent error | Average width of CI | Fraction of 95% CIs include true value | Average error size | Average percent error | Average width of CI |
| 0.0001 | 96% (24/25) | 0.0458 | 12.7% | 0.226 | 96% (24/25) | 3.54 E-05 | 35.4% | 1.75 E-04 | 88% (22/25) | 0.1094 | 14.58% | 0.395 |
| 0.00036 | 92% (23/25) | 0.0279 | 7.8% | 0.137 | 88% (22/25) | 3.68 E-05 | 10.2% | 1.72 E-04 | 84% (21/25) | 0.0587 | 7.83% | 0.239 |
| 0.0007 | 100% (25/25) | 0.0213 | 5.9% | 0.108 | 96% (24/25) | 3.85 E-05 | 5.5% | 2.06 E-04 | 88% (22/25) | 0.0493 | 6.57% | 0.187 |
| 0.004 | 88% (22/25) | 0.0173 | 4.8% | 0.068 | 96% (24/25) | 1.04 E-04 | 2.6% | 4.71 E-04 | 100% (25/25) | 0.0255 | 3.39% | 0.122 |
| 0.007 | 92% (23/25) | 0.0121 | 3.4% | 0.060 | 96% (24/25) | 1.27 E-04 | 1.8% | 7.05 E-04 | 92% (23/25) | 0.0261 | 3.49% | 0.113 |
| 0.04 | 92% (23/25) | 0.0121 | 3.4% | 0.050 | 92% (23/25) | 7.50 E-04 | 1.9% | 3.15 E-03 | 96% (24/25) | 0.0215 | 2.87% | 0.100 |
| 0.1 | 100% (25/25) | 0.0071 | 2.0% | 0.038 | 100% (25/25) | 1.12 E-03 | 1.1% | 5.90 E-03 | 100% (25/25) | 0.0134 | 1.79% | 0.082 |

# S3 Table

| Offspring distribution | $R$ | | | $\lambda_z$ | | | $\sigma$ | | |
|---|---|---|---|---|---|---|---|---|---|
| | Fraction of 95% CIs include true value | Average error size | Average percent error | Fraction of 95% CIs include true value | Average error size | Average percent error | Fraction of 95% CIs include true value | Average error size | Average percent error |
| Poisson | 91.7% (55/60) | 0.0289 | 8.0% | 93.3% (56/60) | 3.76E-05 | 10.4% | 83.3% (50/60) | 0.0649 | 8.7% |
| Negative binomial ($k$=0.58) | 86.7% (52/60) | 0.0393 | 10.9% | 90.0% (54/60) | 4.18E-05 | 11.6% | 91.7% (55/60) | 0.0555 | 7.4% |

# S4 Table

| Model used for inference | mean $R$ | mean $\lambda_z$ |
|---|---|---|
| District | 0.314 | 0.000346 |
| Region | 0.323 | 0.000343 |
| Country | 0.354 | 0.000328 |

# S5 Table

| Model used for inference | mean $R$ | mean $\lambda_z$ |
|---|---|---|
| District | 0.348 | 0.000385 |
| Region | 0.357 | 0.000355 |
| Country | 0.379 | 0.000334 |

# S6 Table

| Inter-locality transmission rule | Offspring distribution | True $R$ | True $\lambda_z$ | # datasets simulated |
|---|---|---|---|---|
| Broader contact zone: district-level | Poisson | 0.36 | 0.00036 | 60 |
| Broader contact zone: district-level | Poisson | 0.2 | 0.00036 | 25 |
| Broader contact zone: district-level | Poisson | 0.6 | 0.00036 | 25 |
| Broader contact zone: district-level | Poisson | 0.36 | 0.0001 | 25 |
| Broader contact zone: district-level | Poisson | 0.36 | 0.0007 | 25 |
| Broader contact zone: district-level | Poisson | 0.36 | 0.004 | 25 |
| Broader contact zone: district-level | Poisson | 0.95 | 0.007 | 25 |
| Broader contact zone: district-level | Poisson | 0.36 | 0.04 | 25 |
| Broader contact zone: district-level | Poisson | 0.36 | 0.1 | 25 |
| Broader contact zone: district-level | NBinom ($k$=0.58) | 0.36 | 0.00036 | 60 |
| Broader contact zone: district-level | Poisson | 0.01 | 0.00036 (intensity heterogeneous through time and space) | 25 |
| Localities have same spatial coordinates as recorded for DRC monkeypox localities, inter-locality transmission with closest 5 neighbors | Poisson | 0.36 | 0.00036 | 25 |
| Localities have same spatial coordinates as recorded for DRC monkeypox localities, inter-locality transmission with neighbors within 30 km | Poisson | 0.36 | 0.00036 | 25 |

# S7 Table

| Symbol | Description |
|--------|-------------|
| $\mu_{t,v}$ | Expected number of cases observed on day $t$, in locality $v$ |
| $N_{t,v}$ | Actual number of cases observed on day $t$, in locality $v$ |
| $N$ | Actual number of cases observed across all localities over the course of surveillance |
| $V$ | Total number of localities under surveillance |
| $V_w$ | Total number of localities under surveillance in the broader contact zone of locality $w$ |
| $W$ | Number of localities with one or more cases (the number of localities that appear in the surveillance dataset) |
| $W_w$ | Number of localities with one or more cases in the broader contact zone of locality $w$ |
| $T$ | Duration of surveillance: number of days surveillance was conducted |
| $\lambda_z$ | Spillover rate: the expected number of spillover events per day in a given locality |
| $\lambda_{\{s,w\},\{t,v\}}$ | The expected number of new infections that become symptomatic on day $t$ in locality $v$ caused by an infectious individual who became symptomatic on day $s$ in locality $w$ |
| $R$ | Reproductive number: the average number of secondary cases caused by an infectious individual |
| $\sigma$ | Within-locality transmission proportion: the fraction of cases arising from human-to-human transmission that occur in the same locality as the source case |