

1 **Misdiagnosis prevents accurate monitoring of transmission and burden for sub-critical**
2 **pathogens: a case study of *Plasmodium knowlesi* malaria**

3 John H. Huber

4

5 Department of Biological Sciences and Eck Institute for Global Health, University of Notre
6 Dame, Notre Dame, IN, USA

7

8 Corresponding Author: huber.j.h@wustl.edu (JHH)

9

10 **ABSTRACT**

11 Maintaining surveillance of emerging infectious diseases presents challenges for monitoring their
12 transmission and burden. Incomplete observation of infections and imperfect diagnosis reduce
13 the observed sizes of transmission chains relative to their true sizes. Previous studies have
14 examined the effect of incomplete observation on estimates of pathogen transmission and
15 burden. However, each study assumed that, if observed, each infection was correctly diagnosed.
16 Here, I leveraged principles from branching process theory to examine how misdiagnosis could
17 contribute to bias in estimates of transmission and burden for emerging infectious diseases.
18 Using the zoonotic *Plasmodium knowlesi* malaria as a case study, I found that, even when
19 assuming complete observation of infections, the number of misdiagnosed cases within a
20 transmission chain for every correctly diagnosed case could range from 0 (0 – 4) when R_0 was
21 0.1 to 86 (0 – 837) when R_0 was 0.9. Data on transmission chain sizes obtained using an
22 imperfect diagnostic could consistently lead to underestimates of R_0 , the basic reproduction
23 number, and simulations revealed that such data on up to 1,000 observed transmission chains
24 was not powered to detect changes in transmission. My results demonstrate that misdiagnosis
25 may hinder effective monitoring of emerging infectious diseases and that sensitivity of
26 diagnostics should be considered in evaluations of surveillance systems.

27 INTRODUCTION

28 For pathogens with sub-critical transmission (i.e., $R_0 < 1$), a robust surveillance system that
29 identifies and correctly diagnoses infections is necessary to monitor changes in pathogen
30 transmission and burden (1). Such pathogen surveillance is important both for measuring
31 progress towards elimination of diseases with immediate public health importance, such as
32 measles (2–4) and malaria (5), and for assessing the future threat of emerging infectious diseases
33 (6), such as avian influenza (7), human monkeypox (1,8), and Middle East respiratory syndrome
34 coronavirus (2,9).

35 Considerable work has been devoted to advance a mathematical framework that
36 leverages the data collected by surveillance systems to obtain estimates of transmission and
37 burden for pathogens with sub-critical dynamics (1,2,4,10,11). These studies have improved our
38 understanding of a wide range of emerging infectious diseases and have critically evaluated the
39 sensitivity of these estimates to the quality of data from the surveillance system. Crucially, each
40 study modeled variation in surveillance quality through variation in the ascertainment fraction
41 (i.e., the proportion of infections that are detected) and assumed that, once detected, all infections
42 were correctly diagnosed. In reality, however, non-specific clinical and biological features are
43 likely to limit the sensitivity of clinical diagnosis, particularly for emerging infectious diseases
44 (12,13). The extent to which misdiagnosis affects estimates of transmission and burden for
45 pathogens with sub-critical dynamics remains largely unaddressed.

46 The zoonotic *Plasmodium knowlesi* malaria offers a natural case study to examine the
47 impact of misdiagnosis on estimates of transmission and burden. Endemic to Southeast Asia
48 (14), *P. knowlesi* is a vector-borne disease with most or all infections in humans caused by
49 spillover transmission from the long- and pig-tailed macaque reservoir (15,16). The extent of

50 transmission between humans is currently unknown (17). Due to morphological similarities with
51 other *Plasmodium* spp., *P. knowlesi* is routinely misdiagnosed by light microscopy (18). A recent
52 systematic review and meta-analysis estimated that the sensitivity of light microscopy for
53 diagnosing *P. knowlesi* infections was less than 1% (19). This high rate of misdiagnosis greatly
54 affects the quality of surveillance data on *P. knowlesi*, potentially biasing estimates of
55 transmission and burden.

56 In this study, I aimed to evaluate the extent to which misdiagnosis of a pathogen affected
57 the ability to monitor its change in transmission and burden. Using *P. knowlesi* as a case study, I
58 leveraged an established framework based upon branching process theory to first quantify the
59 potential magnitude of underestimation of pathogen burden on account of misdiagnosis. Next, I
60 considered how underestimates of pathogen burden could lead to bias in estimates of R_0 , the
61 basic reproduction number. Finally, I quantified the degree to which misdiagnosis reduced the
62 statistical power to detect changes in transmission from surveillance data for emerging infectious
63 diseases, such as *P. knowlesi*.

64

65 **METHODS**

66 **Branching Process Framework of Sub-Critical Transmission**

67 To explore the effects of misdiagnosis on the monitoring of sub-critical transmission (i.e., $R_0 <$
68 1) of *P. knowlesi*, I extended a framework that uses branching process theory to estimate a
69 pathogen's R_0 from its size distribution of stuttering transmission chains. Here, I followed
70 Blumberg and Lloyd-Smith (1,11) and defined a transmission chain as a primary infection (i.e., a
71 spillover infection from a zoonotic reservoir) and all secondary infections arising from that
72 primary infection through at least one generation of pathogen transmission.

73 Assuming that the number of secondary infections caused through one generation of
74 pathogen transmission followed a negative binomial distribution with mean R_0 and dispersion
75 parameter κ , I used the branching framework to calculate summary statistics of the transmission
76 chains. Specifically, I solved for the probability that a transmission chain was truly of size j
77 infections, r_j , and the mean size of transmission chains, μ .

78 Following Blumberg and Lloyd-Smith (11), I considered two models of observation of
79 infections: (i) independent observation and (ii) size-dependent observation. The model of
80 independent observation assumes that each infection is subject to an independent probability of
81 observation and correct diagnosis, p_{ind} , that is equal to the product of the observation
82 probability, p_{det} , and the diagnostic sensitivity, se . By comparison, the model of size-dependent
83 observation assumes that observation of transmission chains occurs through sentinel infections.
84 Each infection within a transmission chain is a sentinel infection with probability, p_{sent} , and, if
85 there is at least sentinel infection within the transmission chain, then all infections within the
86 transmission chain are observed. Diagnosis of each infection occurs independently and is subject
87 to sensitivity, se .

88 I then computed the mean observed transmission chain size, μ^* , as a function of the
89 transmission parameters (R_0 and κ), the observation model (p_{det} or p_{sent}), and the diagnostic
90 accuracy (se). This allowed me to relate the distribution of observed transmission chain sizes to
91 the distribution of true transmission chain sizes and quantify bias in the maximum-likelihood
92 estimates of transmission, $\hat{R}_0 = 1 - \frac{1}{\mu^*}$. If all infections are observed and correctly diagnosed,
93 then $\mu^* = \mu$ and thus $\hat{R}_0 = R_0$. Violations of this assumption, either through incomplete
94 observation or misdiagnosis, introduce bias into transmission estimates. A full description of the
95 branching process framework can be found in the Supplement.

96

97 **Analyses**

98 ***Quantifying the Bounds of Total Burden***

99 To first demonstrate how misdiagnosis, in addition to incomplete observation, may lead to an
100 underestimate of *P. knowlesi* burden, I computed the probability distribution of the true size of a
101 transmission chain conditional upon the observed size of a transmission chain. That is, given that
102 I observed a transmission chain of size \hat{j} , the probability that the transmission chain is truly of
103 size j is equal to

104

$$105 \quad \Pr(j|\hat{j}) = \frac{\Pr(\hat{j}|j) \Pr(j)}{\Pr(\hat{j})}. \quad (1)$$

106

107 In eq. (1), $\Pr(j)$ is the probability that a transmission chain is of size j , r_j , computed using eq.
108 (S1), and $\Pr(\hat{j})$ is the probability that a transmission chain is of observed size j , s'_j , computed
109 using eq. (S2) for the model of independent observation and using the numerator of eq. (S7) for
110 the model of size-dependent observation. For the model of independent observation, the
111 probability of observing a transmission chain of size \hat{j} given that the transmission chain is truly
112 of size j is

113

$$114 \quad \Pr(\hat{j}|j) = \binom{j}{\hat{j}} \cdot p_{ind}^{\hat{j}} \cdot (1 - p_{ind})^{j-\hat{j}}, \quad (2)$$

115

116 where p_{ind} is equal to the product of the probability of detection, p_{det} , and the sensitivity of
117 diagnosis, se . By comparison, for the model of size-dependent observation, this quantity is
118 computed as

119

$$120 \quad \Pr(\hat{j}|j) = (1 - (1 - p_{sent})^j) \cdot \binom{j}{\hat{j}} \cdot se^{\hat{j}} \cdot (1 - se)^{j-\hat{j}}. \quad (3)$$

121

122 Substituting the respective terms into eq. (1), for the model of independent observation, I
123 computed the probability that a transmission chain is of true size j given that it is observed to be
124 of size \hat{j} as

125

$$126 \quad \Pr(j|\hat{j}) = \frac{\binom{j}{\hat{j}} \cdot p_{ind}^{\hat{j}} \cdot (1 - p_{ind})^{j-\hat{j}} r_j}{s'_j}. \quad (4)$$

127

128 For the model of size-dependent observation, I computed this quantity as

129

$$130 \quad \Pr(j|\hat{j}) = \frac{(1 - (1 - p_{sent})^j) \cdot \binom{j}{\hat{j}} \cdot se^{\hat{j}} \cdot (1 - se)^{j-\hat{j}} r_j}{s'_j}. \quad (5)$$

131

132 I used eqs. (4-5) to compute the expected true transmission chain sizes given observed
133 transmission chains of one, two, or three cases while varying the probability of observation, p_{det}
134 or p_{sent} , from 0.1 to 1.0 in increments of 0.1. I sampled the sensitivity of the diagnostic method
135 from the posterior estimate of sensitivity of light microscopy for *P. knowlesi* with mean equal to
136 1.19×10^{-3} (19). For each combination of observed chain size and probability of observation, I

137 calculated the expected true transmission chain size, assuming that the true value of R_0 was
138 equal to 0.1, 0.5, or 0.9. These values of R_0 represent low, medium, and high values of sub-
139 critical transmission and fall within the plausible range of human-to-human transmission of *P.*
140 *knowlesi* (20). The dispersion parameter κ was assumed to be 0.1 in all scenarios, though a
141 supplementary analysis was performed where $\kappa \rightarrow \infty$.

142

143 ***Effect of Misdiagnosis on Estimates of Transmission***

144 On account of incomplete observation and misdiagnosis, the observed burden of *P. knowlesi* may
145 not reflect the true burden. It follows that the mean observed transmission chain size will not
146 equal the true mean transmission chain size, biasing our estimates of R_0 . To explore the extent of
147 this bias in scenarios where R_0 was equal to 0.1, 0.5, or 0.9, I calculated the mean observed
148 transmission chain sizes while varying the probabilities of observation, p_{det} or p_{sent} , from 0.1 to
149 1.0 in increments of 0.1 and while using posterior samples of sensitivity of light microscopy for
150 *P. knowlesi* (19). I then compared the maximum-likelihood estimates of \hat{R}_0 to the true values of
151 R_0 for both the models of independent observation and size-dependent observation. In all
152 scenarios, I assumed that the dispersion parameter κ was 0.1, and a supplementary analysis was
153 performed where $\kappa \rightarrow \infty$.

154

155 ***Effect of Misdiagnosis on Statistical Power to Detect Changes in Transmission***

156 Bias in R_0 estimates on account of misdiagnosis could reduce the statistical power to detect
157 changes in R_0 over time using data on the size of transmission chains. To measure statistical
158 power as a function of the number of observed transmission chains, I followed an approach taken
159 by Blumberg *et al.* (2). I assumed that R_0 was historically equal to 0.1 and then increased to

160 $R_0 + \Delta R_0$, where ΔR_0 was set to 0.1, 0.5, or 0.9. I then simulated N observed transmission chains
161 and estimated \hat{R}_0 while varying N from 1 to 1,000. I then compared the model in which R_0 was
162 estimated to have changed to the null hypothesis that there was no change in transmission (i.e.,
163 $\Delta R_0 = 0$) using the Akaike Information Criterion (AIC) (21). For each number of observed
164 transmission chains N , I repeated this procedure 1,000 times and computed statistical power as
165 the proportion of simulations in which I detected a change in transmission on the basis of AIC.
166 To measure the minimum effect of misdiagnosis on statistical power, I set p_{det} and p_{sent} equal
167 to 1. Because all infections are observed under this assumption and diagnosis is performed
168 independently across infections in both models, the models of independent and size-dependent
169 observation yield identical results. In all scenarios, I assumed that the dispersion parameter κ
170 was 0.1, and a supplementary analysis was performed were $\kappa \rightarrow \infty$.

171

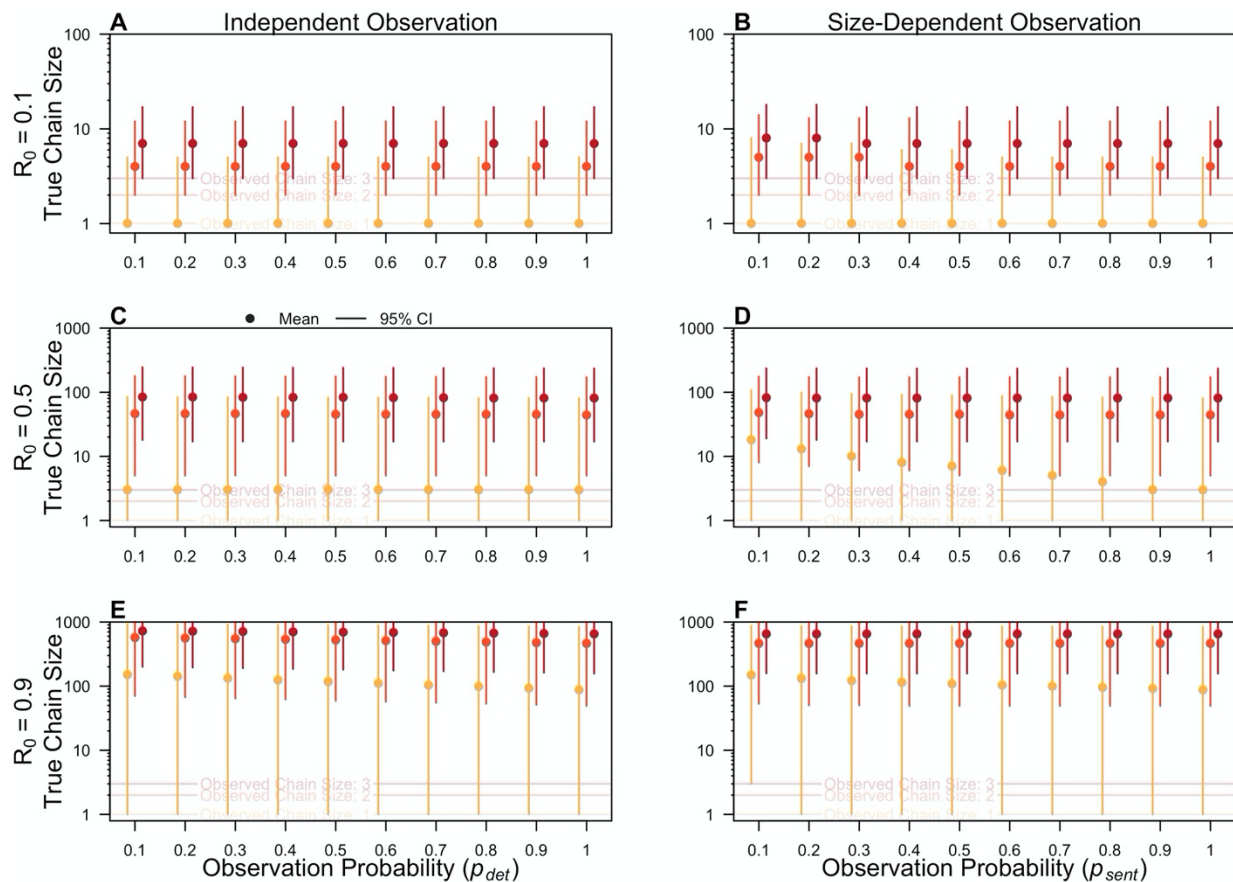
172 **RESULTS**

173 Assuming complete observation of infections (i.e., p_{det} and p_{sent} equal to 1), misdiagnosis of *P.*
174 *knowlesi* infections would underestimate the true *P. knowlesi* burden, with the magnitude of this
175 effect depending upon R_0 (Fig. 1). For a scenario in which R_0 was 0.1, the expected true size of a
176 transmission chain is one infection (95% CI: 1 – 5) if the observed size is one case, four
177 infections (2 – 12) if the observed size is two cases, and seven infections (3 – 17) if the observed
178 size is three cases. Under an alternative scenario in which R_0 was 0.9, the expected true size of
179 the transmission chains increased to 87 (1 – 838), 461 (49 – 965), and 650 (157 – 983)
180 infections, respectively.

181 The effect of incomplete observation (i.e., p_{det} or $p_{sent} < 1$) on the expected burden was
182 most apparent at an intermediate R_0 of 0.5 and with the model of size-dependent observation

183 (Fig. 1D). Under this scenario, given a transmission chain of size one, the expected true
184 transmission chain was 18 infections (1 – 107) if p_{sent} was equal to 0.1, compared to 3 infections
185 (1 – 80) if p_{sent} was equal to 1. In all other scenarios, the expected burden did not change
186 significantly with p_{det} or p_{sent} . This occurred because, even with complete observation (i.e.,
187 p_{det} and p_{sent} equal to 1), 99.881% of *P. knowlesi* cases were expected to be misdiagnosed,
188 given a sensitivity of 0.119%. Therefore, irrespective of the observation probability, only a
189 subset of true transmission chain sizes is consistent with the sizes of the observed transmission
190 chains, given this high percentage of false negatives. If I instead assumed perfect sensitivity of
191 the method, I observed a greater effect of the observation probability on the expected burden
192 (Fig. S1).

193



194

195 **Figure 1. Effect of misdiagnosis and imperfect observation on the expected burden.** The mean
196 true transmission chain size (dots) and 95% CI (segments) are shown conditional upon on an
197 observed transmission chain size of one (yellow), two (orange), or three (red) cases and an R_0 of
198 0.1 (A,B), 0.5 (C,D), and 0.9 (E,F). The horizontal axis is the observation probability,
199 representing p_{det} for the Model of Independent Observation (A, C, E) and p_{sent} for the Model of
200 Size-Dependent Observation (B, D, F).

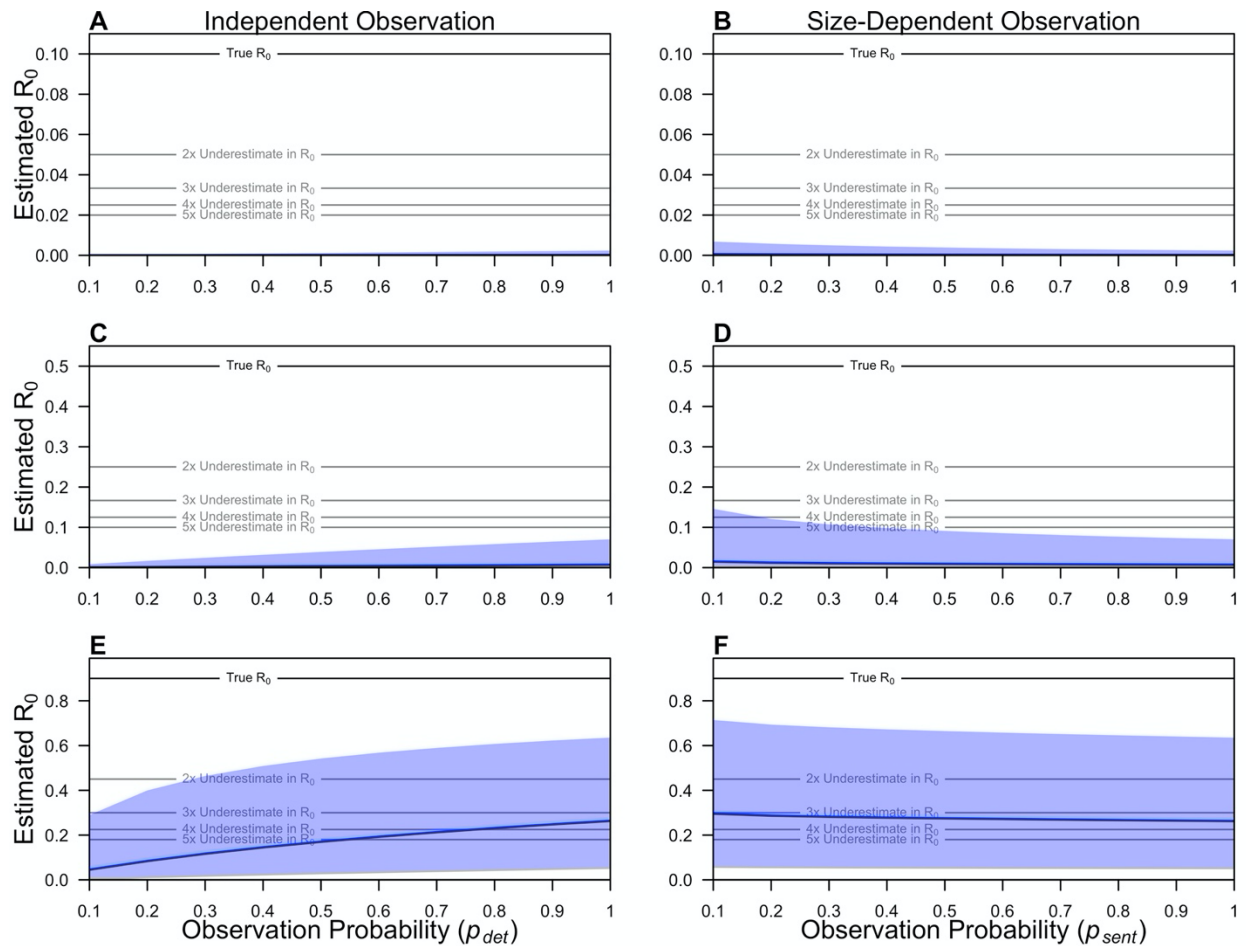
201

202 Given that misdiagnosis underestimated the burden of *P. knowlesi* in this analysis, I assessed its
203 effect on my estimates of transmission, \hat{R}_0 . Because misdiagnosis caused the average observed
204 size of transmission chains to be less than the average true size of transmission chains, I
205 consistently underestimated \hat{R}_0 , with the effect being more severe at lower R_0 (Fig. 2). For
206 example, with an R_0 of 0.1 and assuming perfect observation of infections (i.e., p_{det} and p_{sent}
207 equal to 1), my median \hat{R}_0 estimate was 1.9×10^{-4} (95% PPI: $2.0 \times 10^{-5} - 2.9 \times 10^{-3}$),
208 corresponding to a 520-fold (34 – 4900) underestimate of transmission. Under an alternative
209 scenario in which R_0 was 0.9, my median \hat{R}_0 was 0.26 (0.046 – 0.66), corresponding to a 3.5-
210 fold (1.4 – 19.4) underestimate in transmission.

211 My estimates of transmission were sensitive to the simulated observation probability,
212 though the direction of the effect depended upon the assumed model of observation. For the
213 model of independent observation, \hat{R}_0 estimates increased with increasing p_{det} , because the
214 average observed size of transmission chains increased as more infections were observed (Fig. 2,
215 left column). By contrast, \hat{R}_0 estimates decreased with increasing p_{sent} for the model of size-
216 dependent observation (Fig. 2, right column). This counterintuitive effect can be explained by the
217 observation that, if p_{sent} is low, larger transmission chains have a greater probability that at least

218 one infection is a sentinel infection. This causes a bias in the size of the transmission chains that
 219 are observed at low values of p_{sent} , increasing the mean observed transmission chain size
 220 relative to that at higher values of p_{sent} and thus inflating the \hat{R}_0 estimate.

221



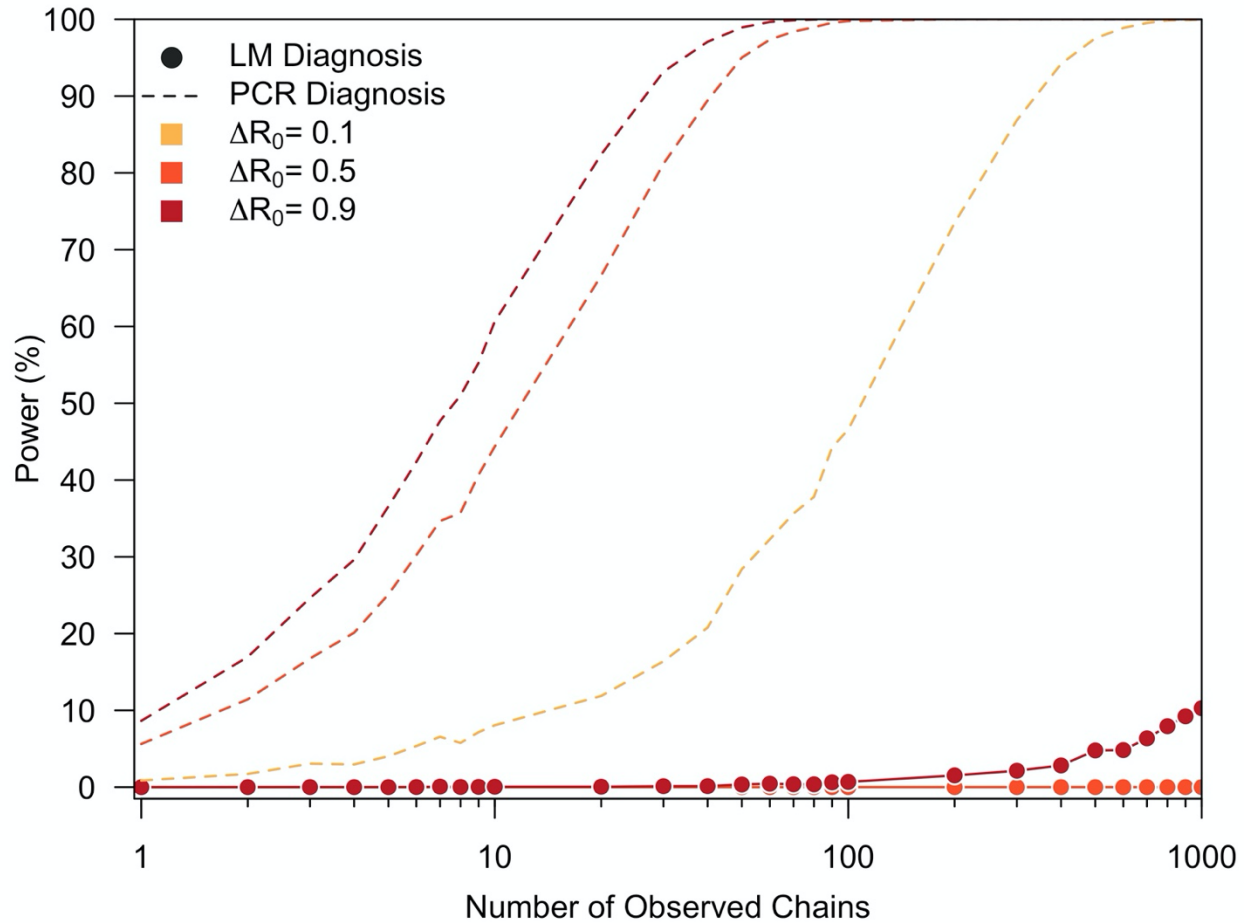
222

223 **Figure 2. Effect of misdiagnosis and imperfect observation on estimates of transmission.** The
 224 posterior median (blue line) and 95% posterior prediction interval (blue shaded region) of
 225 maximum-likelihood estimates of R_0 are shown as a function of the observation probability. The
 226 observation probability represents p_{det} for the Model of Independent Observation (A, C, E) and
 227 p_{sent} for the Model of Size-Dependent Observation (B, D, F). The solid black denotes the true R_0
 228 in each panel, and the grey lines denote two-to-five-fold underestimates of R_0 in each panel.

229

230 The underestimates of burden (Fig. 1) and transmission (Fig. 2) indicated that
231 misdiagnosis of *P. knowlesi* may affect the statistical power to detect changes in transmission
232 based on the size of observed transmission chains. To test this, I simulated changes in
233 transmission and measured the statistical power to detect that change. I observed that, under
234 scenarios in which R_0 increased by 0.9, data on 1,000 observed transmission chains provided
235 only 10.3% power using an imperfect diagnostic method, compared to 100% if using a perfect
236 diagnostic method (Fig. 3). At smaller increases in R_0 , data on observed transmission chain sizes
237 obtained using an imperfect diagnostic method had effectively no power to detect a change in
238 transmission.

239



240

241 **Figure 3. Effect of misdiagnosis on the statistical power to detect changes in transmission.** The
242 statistical power (%) to detect an increase in transmission is shown as a function of the number
243 of observed chains for a transmission increase (ΔR_0) of 0.1 (yellow), 0.5 (orange), and 0.9 (red).
244 Solid lines and points represent an imperfect diagnostic method (i.e., LM) and the dotted lines
245 represent perfect diagnosis (i.e., PCR).

246

247 DISCUSSION

248 Obtaining accurate estimates of transmission and burden is important for monitoring the
249 emergence of infectious diseases. Previous studies have explored the extent to which incomplete
250 observation of infections affects the estimates of transmission for such pathogens (1,2,7,11). In

251 this work, I built upon these studies by considering the effect of misdiagnosis on estimates of
252 pathogen transmission and burden. Using the zoonotic *P. knowlesi* malaria as a case study, I
253 found that misdiagnosis— independent of incomplete observation of infections— may cause us to
254 underestimate the transmission and burden of pathogens with sub-critical dynamics and hinders
255 effective, prospective monitoring of changes in transmission.

256 My results demonstrate that, for pathogens with sub-critical transmission, misdiagnosis
257 leads to an underestimate of overall pathogen burden. Depending upon the R_0 simulated, I found
258 that there could be as many as 86 misdiagnosed cases, on average, for each correctly diagnosed
259 case of *P. knowlesi*, even when assuming complete observation of infections. This effect
260 increased under select settings when the assumption of complete observation of infections was
261 relaxed. The underestimation of burden due to misdiagnosis has the potential to shape our
262 epidemiological understanding of an emerging pathogen. For instance, singleton cases of a
263 zoonotic pathogen, such as *P. knowlesi*, are commonly assessed as dead-end spillover events
264 from the zoonotic reservoir (17). However, my simulations suggest that such singleton cases
265 could instead represent a broad range of epidemiological outcomes, spanning dead-end spillover
266 events to larger transmission chains.

267 Due to its effect on observed pathogen burden, misdiagnosis contributed a downward
268 bias in estimates of transmission. Except for scenarios in which the true simulated R_0 was close
269 to one, my maximum-likelihood estimates of R_0 approached zero, representing situations in
270 which we would incorrectly conclude that human-to-human transmission of the pathogen was
271 unlikely to be occurring. For every scenario considered, the estimate of R_0 was less than the true
272 value, indicating that bias due to misdiagnosis exceeds the competing positive bias from
273 incomplete observation when assuming size-dependent observation (1). For pathogens such as *P.*

274 *knowlesi*, these simulation results suggest that, in settings where misdiagnosis is common, the
275 extent of human-to-human transmission could be greater than previously thought. To date, it has
276 been believed that nearly all cases of *P. knowlesi* in humans are caused by spillover from long-
277 tailed and pig-tailed macaques, the zoonotic reservoir (17). The lack of observed human-to-
278 human transmission may be explained by multiple factors, including low parasite densities in
279 humans (16) and restricted vector habitat preference (15), and is supported by a lack of genetic
280 diversity across human *P. knowlesi* infections (22). Nevertheless, human-to-human transmission
281 of *P. knowlesi* has been demonstrated experimentally (23), and these results suggest that, if or
282 when human-to-human transmission occurs, misdiagnosis could cause us to underestimate its
283 magnitude.

284 Finally, I demonstrated that data on the sizes of transmission chains diagnosed using a
285 diagnostic with realistic sensitivity would be insufficient to monitor changes in transmission.
286 Even with 1,000 observed transmission chains, I calculated a power of only 10% to detect an
287 increase in R_0 from 0.1 to 1. This empirical power calculation assumed complete observation of
288 infections, so it represents an upper bound on the statistical power that we might expect if a
289 diagnostic with realistic sensitivity was used. Therefore, more sensitive diagnostics, such as
290 polymerase chain reaction, may be needed to detect changes in transmission that could result
291 from pathogen evolution (24), among other factors (25–27).

292 This analysis is subject to a number of limitations. First, the conclusions that I reached
293 were based upon simulated data only. I used simulations representative of *P. knowlesi* to
294 illustrate possible outcomes that may occur due to misdiagnosis (19,20), yet I lacked empirical
295 data on the distribution of transmission chain sizes for *P. knowlesi*. As such, this analysis is not
296 estimating the true extent of human-to-human transmission of *P. knowlesi*. Second, methods

297 exist to account for incomplete observation of infections in estimates of R_0 (11). However, as
298 noted by Blumberg *et al.* (2), it is challenging to estimate the proportion of infections that are
299 captured by the surveillance system. Consequently, these calculations were conditioned upon the
300 assumption of complete observation and perfect diagnoses, so violations therein should be
301 interpreted as the upper bound on the bias that would likely be observed. Finally, I considered a
302 single pathogen in isolation, though misdiagnosis is commonly due to co-circulation of related
303 pathogens (12,13,18). Accounting for the upward bias due to false-positive diagnoses from other
304 pathogens and exploring the magnitude of this effect across epidemiological settings could be
305 important directions for future work.

306

307 **ACKNOWLEDGMENTS**

308 I acknowledge funding from a Graduate Research Fellowship from the National Science
309 Foundation and a Richard and Peggy Notebaert Premier Fellowship from the University of Notre
310 Dame. The funders had no role in the study design, in the collection, analysis, and interpretation
311 of data, in the writing of the report, or in the design to submit the article for publication. I thank
312 Alex Perkins for helpful comments on this manuscript.

313

314 **SUPPLEMENT**

315 **Methods**

316 *Mean Transmission Chain Size*

317 *Complete Observation and Correct Diagnosis*

318 For a pathogen with sub-critical transmission dynamics, I modeled the number of offspring
319 caused by a single infection through one generation of transmission as a negative binomial

320 distribution with mean R_0 and dispersion parameter κ (1,11). Therefore, it follows that
321 transmission chains of size j occur with probability,

322

$$323 \quad r_j = \frac{\Gamma(\kappa j + j - 1)}{\Gamma(\kappa j)\Gamma(j + 1)} \frac{\left(R_0/\kappa\right)^{j-1}}{\left(1 + \left(R_0/\kappa\right)\right)^{\kappa j + j - 1}}. \quad (S1)$$

324

325 In eq. (S1), $\Gamma(\cdot)$ is the gamma function. Because $R_0 < 1$, the mean transmission chain size μ can
326 be calculated as the mean of a geometric series with common ratio R_0 and is equal to $\frac{1}{1-R_0}$.

327

328 *Incomplete Observation and Imperfect Diagnosis*

329 In the case of *P. knowlesi* and many other pathogens, the size of transmission chains that are
330 identified by a surveillance system will be affected by two factors. First, infections in a
331 transmission chain may not present within the health system, due to a lack of symptoms or
332 access to treatment. Second, infections in the transmission chain that do present within the health
333 system may be misdiagnosed and thus inaccurately recorded within the surveillance system.
334 Both factors will make the observed transmission chain appear smaller in size than the true
335 transmission chain. Previous work by Blumberg and Lloyd-Smith (1) has quantified the effect of
336 two models of incomplete observation on estimates of transmission and burden. Here, I build
337 upon this work by integrating the effect of (mis)diagnosis of infections that occurs secondary to
338 the observation of infections within the health system.

339

340 Model of Independent Observation

341 The first model of incomplete observation and diagnosis assumes that each individual is subject
342 to an independent probability p_{ind} equal to the product of observation probability, p_{det} , and the
343 sensitivity of the diagnostic method, se . Therefore, the probability that we observe and correctly
344 diagnose j cases from a transmission chain is equal to

345

346
$$s'_j = \sum_{k=j}^{\infty} r_k \cdot \binom{k}{j} \cdot p_{ind}^j \cdot (1 - p_{ind})^{k-j}, \quad (S2)$$

347

348 where r_k is the probability that a transmission chain is of true size k , calculated using eq. (S1).

349 The probability that a transmission chain is of observed size j is equal to the normalized

350 probability of s'_j , computed as

351

352
$$r'_j = \frac{s'_j}{1 - s'_0}. \quad (S3)$$

353

354 In eq. (S3), s'_0 is the probability that a transmission chain is neither observed nor correctly

355 diagnosed. Due to incomplete observation and misdiagnosis, the probability that a transmission

356 chain is of observed size j , r'_j , is not equal to the probability that a transmission chain is of true

357 size j , r_j . Finally, because each infection within the transmission chain is subject to an

358 independent probability of observation and correct diagnosis, the probability p_{obs} that a

359 randomly sampled infection is observed and correctly diagnosed is equal to p_{ind} .

360

361 Model of Size-Dependent Observation

362 The alternative model assumes that transmission chains are observed through a sentinel
363 infection, such that, if the sentinel infection is observed, then all other infections in the
364 transmission chain will be observed. Incorporating the effect of imperfect diagnosis, the
365 probability that we do not observe a transmission chain of size j is equal to

366

$$367 \quad v_j = (1 - p_{sent})^j + (1 - (1 - p_{sent})^j)(1 - se)^j. \quad (S4)$$

368

369 In eq. (S4), the first term in the summation is the probability that none of the j infections of the
370 transmission chain are a sentinel infection, and the second term in the summation is the product
371 of the probability that at least one infection is a sentinel infection (i.e., the probability that we
372 observe the transmission chain) and the probability that all j infections are misdiagnosed. Using
373 eq. (S4), I calculated the probability that a transmission chain is neither observed nor correctly
374 diagnosed as

375

$$376 \quad s'_0 = \sum_{k=1}^{\infty} r_k \cdot v_k. \quad (S5)$$

377

378 The probability that a transmission chain is of observed size j is then equal to

379

$$380 \quad r'_j = \frac{\sum_{k=j}^{\infty} r_k \cdot (1 - (1 - p_{sent})^k) \cdot \binom{k}{j} \cdot se^j \cdot (1 - se)^{k-j}}{1 - s'_0}, \quad (S6)$$

381

382 and the probability that a randomly chosen infection is observed and correctly diagnosed is equal
383 to

384

385
$$p_{obs} = \frac{\sum_{j=1}^{\infty} j \cdot r_j \cdot (1 - (1 - p_{sent})^j) \cdot se}{\mu}. \quad (S7)$$

386

387 Eq. (S7) accounts for the probability that non-sentinel infections are detected, a quantity that
388 increases as a function of the transmission chain size.

389

390 Mean Transmission Chain Size

391 For both the model of independent observation and the model of size-dependent observation, the
392 mean observed size of transmission chains is calculated as

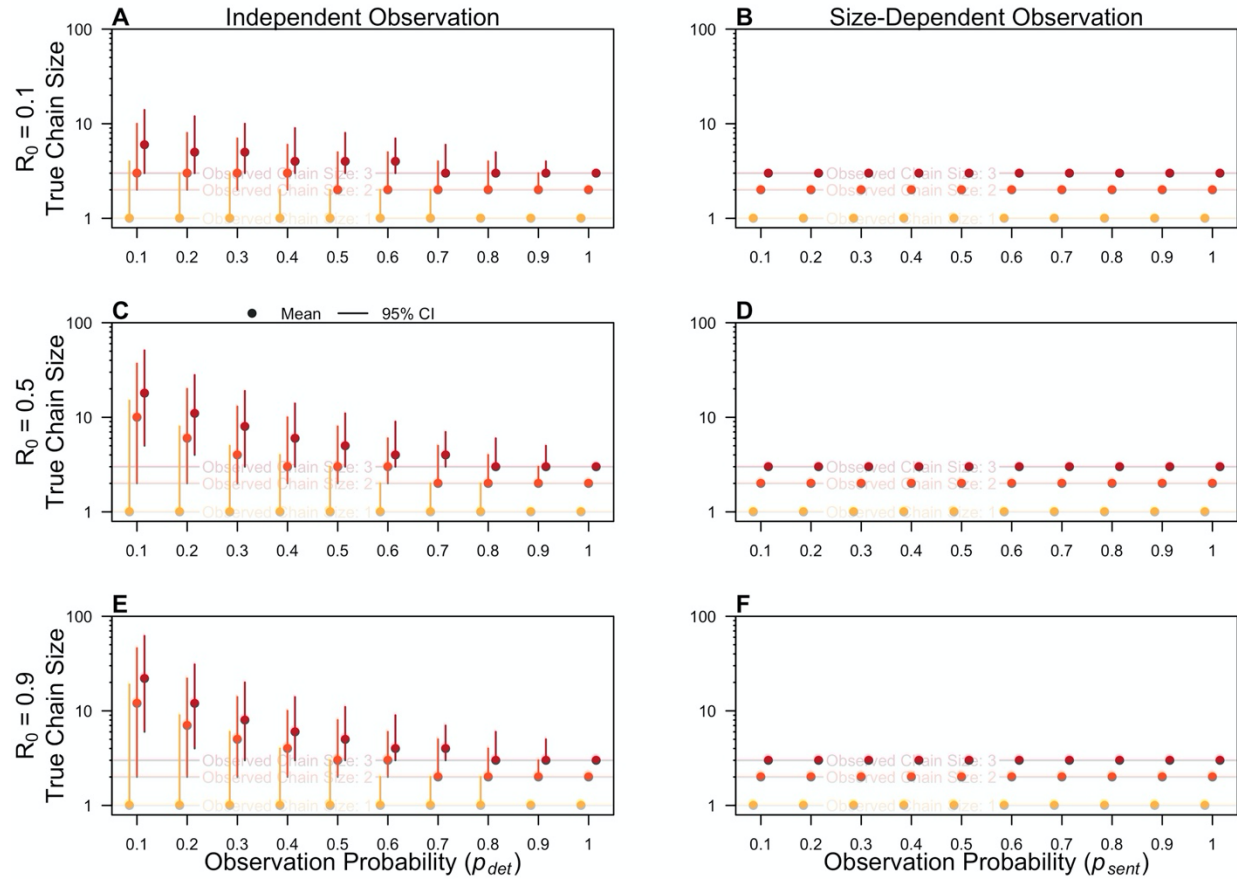
393

394
$$\mu^* = \sum_{j=1}^{\infty} j \cdot r'_j = \frac{p_{obs} \cdot \mu}{1 - s'_0}. \quad (S8)$$

395

396 **Results**

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).



397

398 **Figure S1. Effect of misdiagnosis and imperfect observation on the expected true pathogen**

399 **burden assuming perfect diagnosis. The mean true transmission chain size (dots) and 95% CI**

400 **(segments) are shown conditional upon on an observed transmission chain size of one (yellow),**

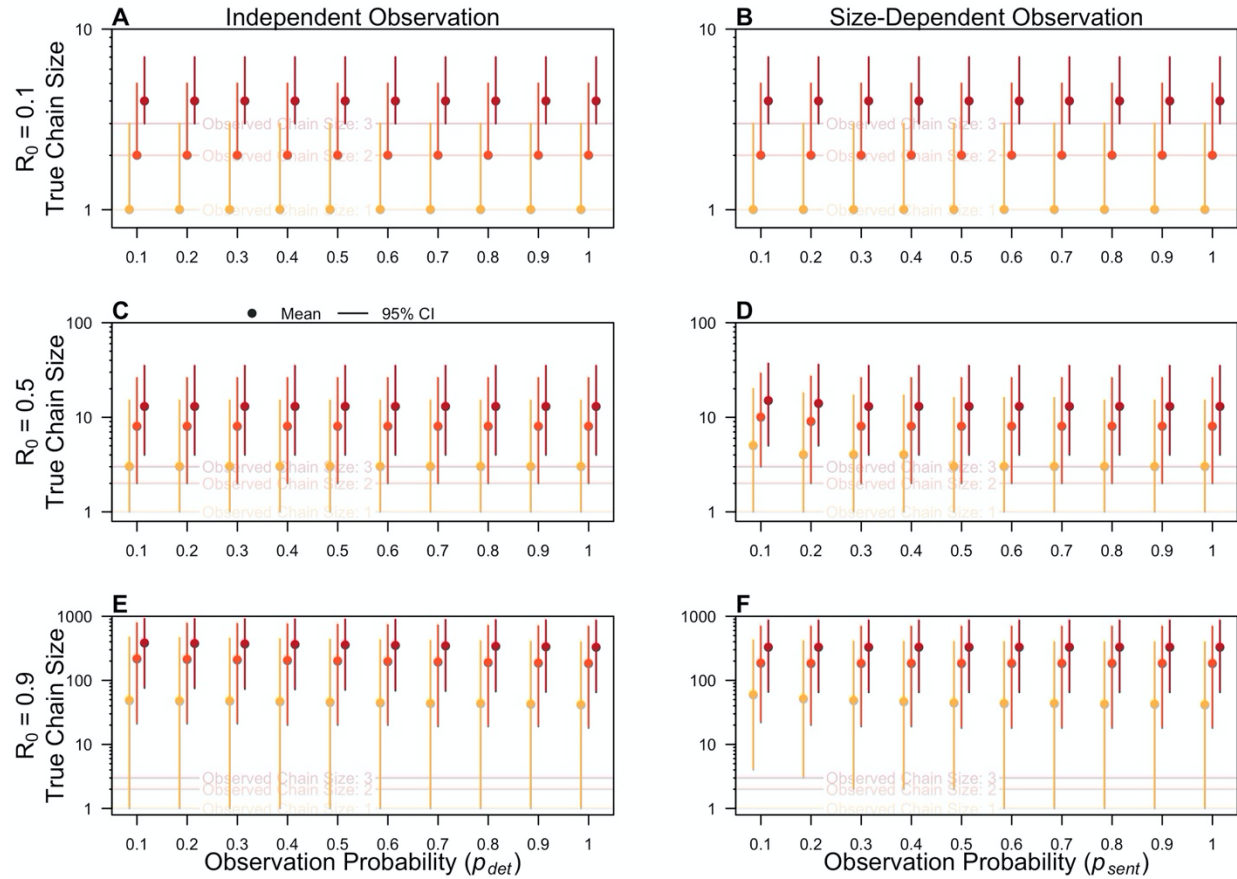
401 **two (orange), or three (red) cases and an R_0 of 0.1 (A,B), 0.5 (C,D), and 0.9 (E,F). The**

402 **horizontal axis is the observation probability, representing p_{det} for the Model of Independent**

403 **Observation (A, C, E) and p_{sent} for the Model of Size-Dependent Observation (B, D, F).**

404

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

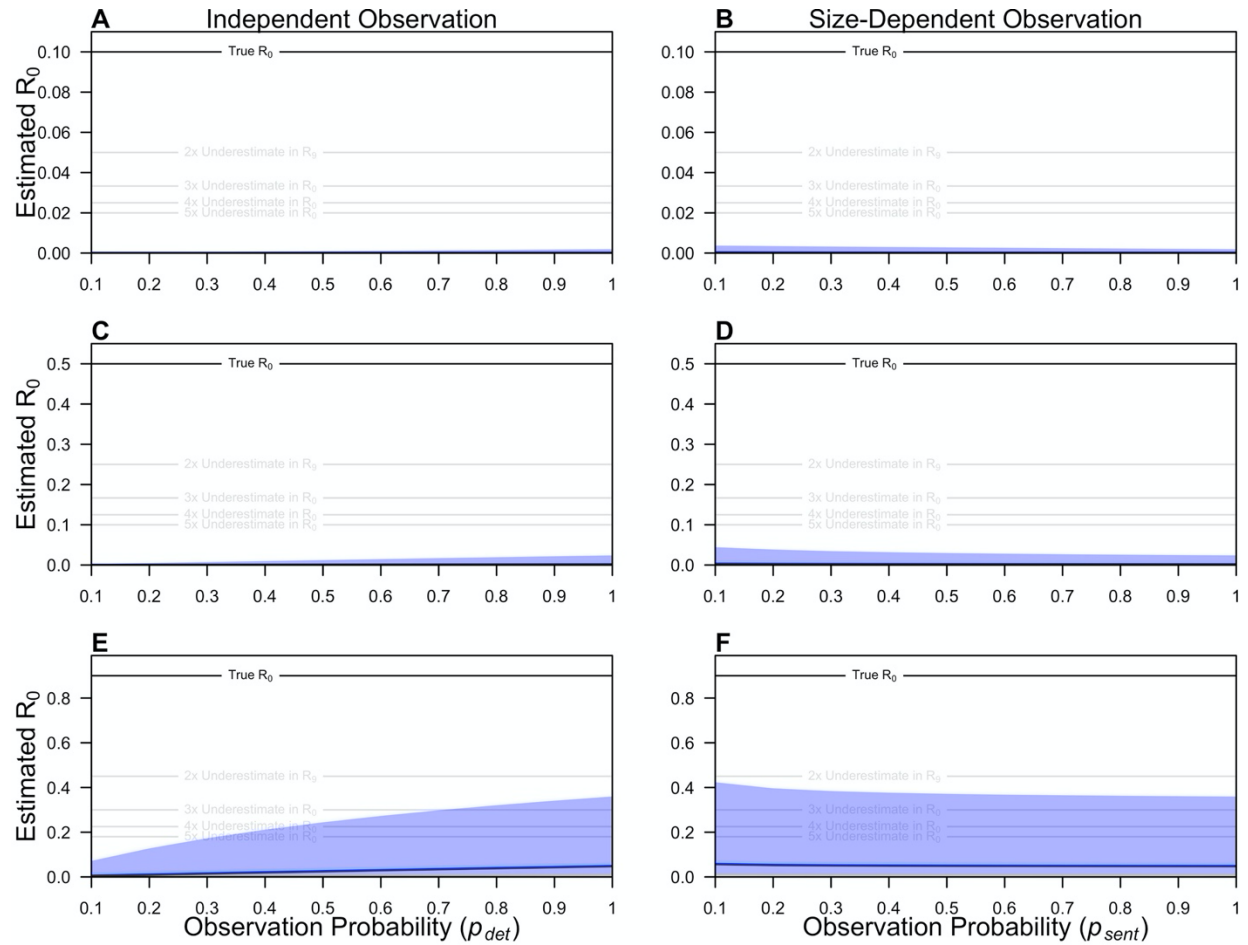


405

406 **Figure S2. Effect of misdiagnosis and imperfect observation on the expected true pathogen**
 407 **burden when $\kappa \rightarrow \infty$. The mean true transmission chain size (dots) and 95% CI (segments) are**
 408 **shown conditional upon on an observed transmission chain size of one (yellow), two (orange), or**
 409 **three (red) cases and an R_0 of 0.1 (A,B), 0.5 (C,D), and 0.9 (E,F). The horizontal axis is the**
 410 **observation probability, representing p_{det} for the Model of Independent Observation (A, C, E)**
 411 **and p_{sent} for the Model of Size-Dependent Observation (B, D, F).**

412

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).



413

414

415

416

417

418

419

420

421

Figure S3. Effect of misdiagnosis and imperfect observation on estimates of transmission

when $\kappa \rightarrow \infty$. The posterior median (blue line) and 95% posterior prediction interval (blue

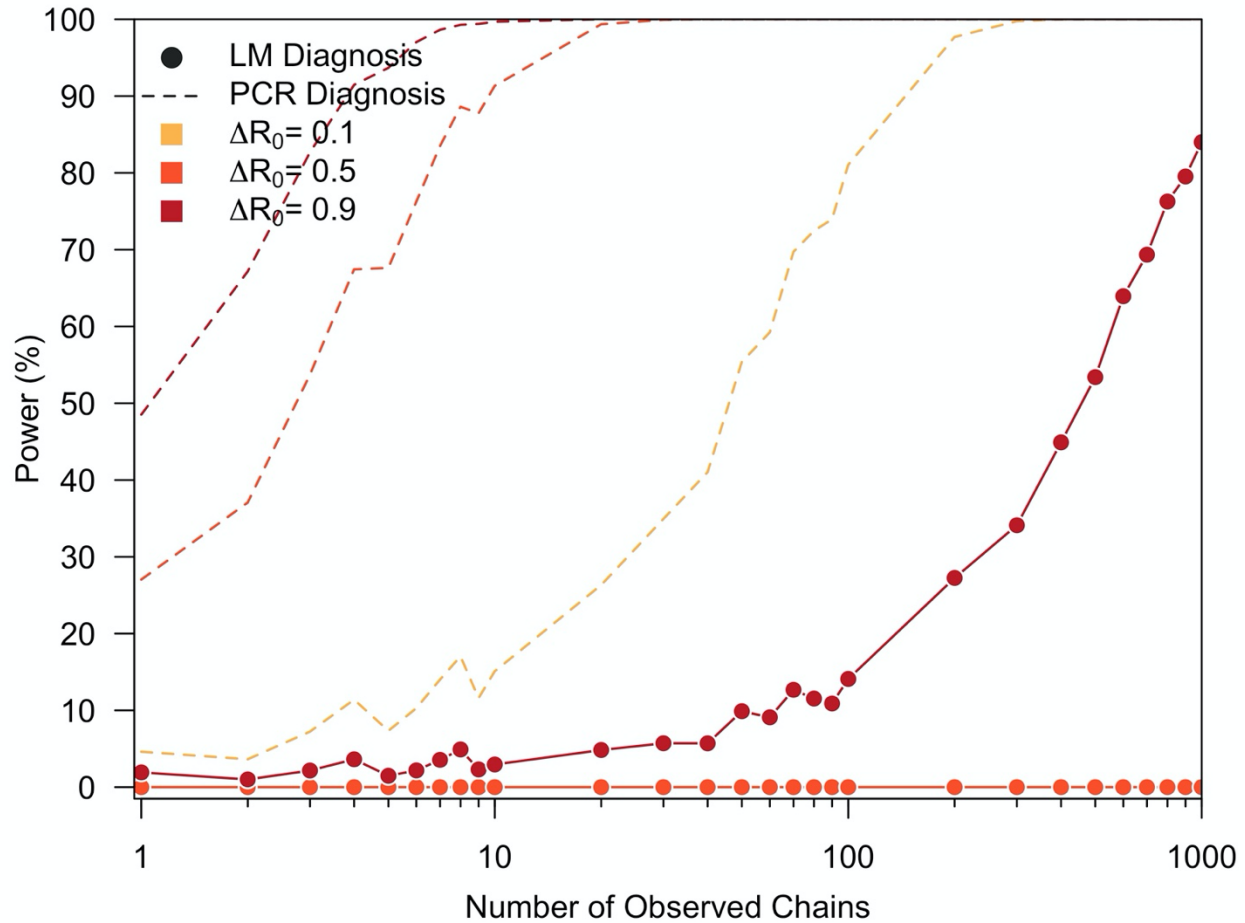
shaded region) of maximum-level estimates of R_0 are shown as a function of the observation

probability. The observation probability represents p_{det} for the Model of Independent

Observation (A, C, E) and p_{sent} for the Model of Size-Dependent Observation (B, D, F). The solid

black denotes the true R_0 in each panel, and the grey lines denote two-to-five-fold underestimates

of R_0 in each panel.



422

423 **Figure S4. Effect of misdiagnosis on the statistical power to detect changes in transmission**

424 **when $\kappa \rightarrow \infty$.** The statistical power (%) to detect an increase in transmission is shown as a

425 **function of the number of observed chains for a transmission increase (ΔR_0) of 0.1 (yellow), 0.5**

426 **(orange), and 0.9 (red). Solid lines and points represent an imperfect diagnostic method (i.e.,**

427 **LM) and the dotted lines represent perfect diagnosis (i.e., PCR).**

428

429 REFERENCES

- 430 1. Blumberg S, Lloyd-Smith JO. Inference of $R(0)$ and transmission heterogeneity from the size
431 distribution of stuttering chains. PLoS Comput Biol. 2013;9(5):e1002993.

- 432 2. Blumberg S, Funk S, Pulliam JRC. Detecting Differential Transmissibilities That Affect the
433 Size of Self-Limited Outbreaks. Wilke CO, editor. PLoS Pathog. 2014 Oct
434 30;10(10):e1004452.
- 435 3. Blumberg S, Worden L, Enanoria W, Ackley S, Deiner M, Liu F, et al. Assessing Measles
436 Transmission in the United States Following a Large Outbreak in California. PLoS Curr
437 [Internet]. 2015 [cited 2021 Jul 27]; Available from:
438 [https://currents.plos.org/outbreaks/article/assessing-measles-transmission-in-the-united-](https://currents.plos.org/outbreaks/article/assessing-measles-transmission-in-the-united-states-following-a-large-outbreak-in-california/)
439 [states-following-a-large-outbreak-in-california/](https://currents.plos.org/outbreaks/article/assessing-measles-transmission-in-the-united-states-following-a-large-outbreak-in-california/)
- 440 4. Farrington CP. Branching process models for surveillance of infectious diseases controlled
441 by mass vaccination. Biostatistics. 2003 Apr 1;4(2):279–95.
- 442 5. Churcher TS, Cohen JM, Novotny J, Ntshalintshali N, Kunene S, Cauchemez S. Measuring
443 the path toward malaria elimination. Science. 2014 Jun 13;344(6189):1230–2.
- 444 6. Plowright RK, Parrish CR, McCallum H, Hudson PJ, Ko AI, Graham AL, et al. Pathways to
445 zoonotic spillover. Nat Rev Microbiol. 2017 Aug;15(8):502–10.
- 446 7. Ferguson NM. PUBLIC HEALTH: Enhanced: Public Health Risk from the Avian H5N1
447 Influenza Epidemic. Science. 2004 May 14;304(5673):968–9.
- 448 8. Ambrose MR, Kucharski AJ, Formenty P, Muyembe-Tamfum J-J, Rimoin AW, Lloyd-Smith
449 JO. Quantifying transmission of emerging zoonoses: Using mathematical models to
450 maximize the value of surveillance data [Internet]. Epidemiology; 2019 Jun [cited 2021 Jul
451 27]. Available from: <http://biorxiv.org/lookup/doi/10.1101/677021>
- 452 9. Cauchemez S, Nouvellet P, Cori A, Jombart T, Garske T, Clapham H, et al. Unraveling the
453 drivers of MERS-CoV transmission. Proc Natl Acad Sci USA. 2016 Aug 9;113(32):9081–6.
- 454 10. De Serres G, Gay NJ, Farrington CP. Epidemiology of Transmissible Diseases after
455 Elimination. American Journal of Epidemiology. 2000 Jun 1;151(11):1039–48.
- 456 11. Blumberg S, Lloyd-Smith JO. Comparing methods for estimating R_0 from the size
457 distribution of subcritical transmission chains. Epidemics. 2013 Sep;5(3):131–45.
- 458 12. Oidtman RJ, España G, Perkins TA. Co-circulation and misdiagnosis led to underestimation
459 of the 2015–2017 Zika epidemic in the Americas. Lacerda MVG, editor. PLoS Negl Trop
460 Dis. 2021 Mar 1;15(3):e0009208.
- 461 13. Glennon EE, Jephcott FL, Oti A, Carlson CJ, Bustos Carillo FA, Hranac CR, et al.
462 Syndromic detectability of haemorrhagic fever outbreaks [Internet]. Epidemiology; 2020
463 Mar [cited 2021 Jul 27]. Available from:
464 <http://medrxiv.org/lookup/doi/10.1101/2020.03.28.20019463>
- 465 14. Shearer FM, Huang Z, Weiss DJ, Wiebe A, Gibson HS, Battle KE, et al. Estimating
466 Geographical Variation in the Risk of Zoonotic Plasmodium knowlesi Infection in Countries

- 467 Eliminating Malaria. Churcher TS, editor. PLoS Negl Trop Dis. 2016 Aug
468 5;10(8):e0004915.
- 469 15. Collins WE. *Plasmodium knowlesi* : A Malaria Parasite of Monkeys and Humans. Annu Rev
470 Entomol. 2012 Jan 7;57(1):107–21.
- 471 16. Kantele A, Jokiranta TS. Review of cases with the emerging fifth human malaria parasite,
472 *Plasmodium knowlesi*. Clin Infect Dis. 2011 Jun;52(11):1356–62.
- 473 17. Feachem RGA, Chen I, Akbari O, Bertozzi-Villa A, Bhatt S, Binka F, et al. Malaria
474 eradication within a generation: ambitious, achievable, and necessary. The Lancet. 2019
475 Sep;394(10203):1056–112.
- 476 18. Barber BE, William T, Grigg MJ, Yeo TW, Anstey NM. Limitations of microscopy to
477 differentiate *Plasmodium* species in a region co-endemic for *Plasmodium falciparum*,
478 *Plasmodium vivax* and *Plasmodium knowlesi*. Malar J. 2013 Jan 8;12:8.
- 479 19. Huber JH, Elliott M, Koepfli C, Perkins A. The impact of emerging *Plasmodium knowlesi*
480 on accurate diagnosis by light microscopy: a systematic review and modelling analysis
481 [Internet]. Epidemiology; 2021 Sep [cited 2021 Sep 13]. Available from:
482 <http://medrxiv.org/lookup/doi/10.1101/2021.09.08.21263294>
- 483 20. Imai N, White MT, Ghani AC, Drakeley CJ. Transmission and Control of *Plasmodium*
484 *knowlesi*: A Mathematical Modelling Study. Churcher TS, editor. PLoS Negl Trop Dis. 2014
485 Jul 24;8(7):e2978.
- 486 21. Akaike H. A new look at the statistical model identification. IEEE Trans Automat Contr.
487 1974 Dec;19(6):716–23.
- 488 22. Lee K-S, Divis PCS, Zakaria SK, Matusop A, Julin RA, Conway DJ, et al. *Plasmodium*
489 *knowlesi*: Reservoir Hosts and Tracking the Emergence in Humans and Macaques. Kazura
490 JW, editor. PLoS Pathog. 2011 Apr 7;7(4):e1002015.
- 491 23. Alpert E, Collins WE, Jeter MH, Chin W, Contacos PG. Experimental Mosquito-
492 Transmission of *Plasmodium Knowlesi* to Man and Monkey. The American Journal of
493 Tropical Medicine and Hygiene. 1968 May 1;17(3):355–8.
- 494 24. Geoghegan JL, Senior AM, Di Giallonardo F, Holmes EC. Virological factors that increase
495 the transmissibility of emerging human viruses. Proc Natl Acad Sci USA. 2016 Apr
496 12;113(15):4170–5.
- 497 25. Faust CL, McCallum HI, Bloomfield LSP, Gottdenker NL, Gillespie TR, Torney CJ, et al.
498 Pathogen spillover during land conversion. Ostfeld R, editor. Ecol Lett. 2018 Apr;21(4):471–
499 83.
- 500 26. Baeza A, Santos-Vega M, Dobson AP, Pascual M. The rise and fall of malaria under land-
501 use change in frontier regions. Nat Ecol Evol. 2017 May;1(5):0108.

502 27. Morse SS. Factors in the Emergence of Infectious Diseases. *Emerg Infect Dis.* 1995
503 Mar;1(1):7–15.

504