

1 **Assessing physical and environmental predictors of**  
2 **bovine *Schistosoma japonicum* infection in rural China**

3

4 **Short title:** Predictors of bovine schistosomiasis in rural China

5

6 **Elise Grover<sup>1</sup>, Sara Paull<sup>1</sup>, Katerina Kechris<sup>2</sup>, Andrea Buchwald<sup>1,3</sup>, Katherine James<sup>1</sup>, Yang Liu<sup>4</sup>, Elizabeth**  
7 **J. Carlton<sup>1\*</sup>**

8

9 <sup>1</sup> Department of Environmental and Occupational Health, Colorado School of Public Health, University of  
10 Colorado Anschutz Medical Campus, Aurora, USA

11 <sup>2</sup> Department of Biostatistics and Informatics, Colorado School of Public Health, University of Colorado Anschutz  
12 Medical Campus, Aurora, USA

13 <sup>3</sup> Center for Vaccine Development and Global Health, University of Maryland School of Medicine

14 <sup>4</sup> Institute of Parasitic Diseases, Sichuan Center for Disease Control and Prevention, Chengdu, China

15

16 \* Corresponding author

17 Email: [elizabeth.carlton@cuanschutz.edu](mailto:elizabeth.carlton@cuanschutz.edu)

18

19

20

## 21 **Abstract**

### 22 **Background**

23 Bovines have been repeatedly highlighted as a major reservoir for human *Schistosoma japonicum* infection in  
24 rural farming villages in China. However, little is known about the individual and environmental risk factors for  
25 bovine schistosomiasis infection. The current body of literature on individual-level risk factors features  
26 inconsistent, and sometimes contradictory results, and to date, few studies have assessed the broader  
27 environmental conditions that predict bovine schistosomiasis.

### 28 **Methodology/Principal Findings**

29 Using data collected as a part of a longitudinal study in 39 rural villages in Sichuan, China from 2007 to 2016, we  
30 aimed to identifying the strongest individual, household and village-level predictors of bovine *S. japonicum*  
31 infection. Candidate predictors for this assessment included: 1) physical/biological characteristics of bovines, 2)  
32 potential human sources of environmental schistosomes, 3) socio-economic indicators, 4) potential animal  
33 reservoirs, and 5) agricultural risk factors. A Random Forests machine learning approach was used to determine  
34 which of our candidate predictors serve as the best predictors of bovine schistosomiasis infection in each survey  
35 year. Of the five categories of predictors, high-risk agricultural practices and animal reservoirs, specifically, bovine  
36 density at the village-level, were repeatedly found to be among the top predictors of bovine *S. japonicum*  
37 infection.

### 38 **Conclusion/Significance**

39 Our findings highlight the potential utility of presumptively treating bovines residing in villages and households that  
40 engage in high-risk agricultural practices, or bovines belonging to villages with particularly high levels of bovine  
41 ownership. Additionally, village-level predictors were stronger predictors of bovine infection than household-level  
42 predictors, suggesting future investigations and interventions may need to apply a broad ecological lens in order  
43 to successfully extricate and address environmental sources of ongoing transmission.

44

## 45 **Author Summary**

46 Schistosomiasis is a burdensome global disease that is frequently transmitted between humans and animals. The  
47 parasite that causes schistosomiasis is released into water by snails that become infected via contact with eggs  
48 from human or animal feces, allowing other human and animal hosts to become infected when they come in  
49 contact with contaminated water. In China, bovines are believed to be the most common animal source of human  
50 infections, though little is known about what factors promote bovine infections. Because schistosomiasis is a  
51 sanitation-related, water-borne disease transmitted by many animals, we hypothesized that several environmental  
52 factors – such as the lack of improved sanitation systems, or participation in agricultural production that is water  
53 or fertilizer-intensive – could promote schistosomiasis infection in bovines. Our study investigated this using data  
54 collected in 39 villages in a region of China where bovine and human schistosomiasis both occur. We found that  
55 several agriculture-related factors and bovine density in the village were predictive of bovine infection status.  
56 These findings highlight the importance of assessing environmental sources of disease transmission across large  
57 geographic scales, and suggest that preventative treatment of bovines residing in high risk villages may help to  
58 control local transmission.

59

60

61

62

63

64

65

66

67

68

69

70

71

## 72 Introduction

73 Schistosomiasis is among the most burdensome helminth infections, with transmission being reported in  
74 a total of 78 countries in 2018 and approximately 230 million people in need of preventative treatment (1).  
75 Although great strides have been made in the last several decades in the control of schistosomiasis in several  
76 countries worldwide (1), pockets of reemergent or persistent transmission within such areas highlight the need for  
77 careful consideration of possible local drivers of transmission (2, 3). A poignant example of this is found when  
78 looking at the transmission of *Schistosoma japonicum* in China, where despite well-established control programs  
79 and great progress towards elimination since the mid-1950s (4), a 2018 national report highlighted that there  
80 remains 450 endemic counties where transmission interruption has yet to be achieved (5). *S. japonicum* has been  
81 found to be transmitted by at least 40 species of wild and domestic mammals (6), and animal activities near likely  
82 transmission sites may be important sources of reemergence and persistence. In China, several domesticated  
83 and wild animals have been identified as being capable of carrying and transmitting *S. japonicum*, including  
84 bovines, pigs, goats, dogs, cats and rodents (7, 8). Estimates from Jiangxi Province of Eastern China suggest that  
85 bovines may be responsible for as much as 75% of human transmission (9). This substantial contribution is  
86 thought to be related to the high degree of environmental overlap between humans and bovines during  
87 agricultural production, as well as the large amount of fecal output of bovines, which has been estimated to be as  
88 high as 100 times that of human fecal production each day (6, 10, 11). Additionally, the high frequency of  
89 livestock movement via the livestock trade within mountainous regions of China further highlights the important  
90 role that bovines may be playing in the *S. japonicum* transmission cycle in endemic areas (12, 13).

91 Despite increasing awareness that bovines may be an important driver of human schistosomiasis  
92 infection, little is known about what factors are likely to be influencing infection within bovine populations. Studies  
93 have highlighted several risk factors associated with other bovine infectious diseases, including individual  
94 characteristics like old-age, male sex, a range of breeds and uses (e.g. dairy, beef or agricultural work), group  
95 characteristics like herd size and herd density, and environmental characteristics like contact with other animals  
96 and the presence/absence of irrigation systems (14-19). While recent assessments of bovine risk of *S. bovis*  
97 infection in Eastern Africa have also studied individual-level risk factors such as bovine sex, age, breed and body  
98 condition, the results are contradictory (20-26). Reasons for such discrepancies have not been fully elucidated,

99 though Defersha & Belete (2018) hypothesize that it may be related to variations in management practices for  
100 different bovine groups (e.g. separation of sexes or of age groups) and different grazing ranges or grazing  
101 patterns allowed on different farms (e.g. smaller grazing area of very young and very old bovines) (23).

102 Outside of eastern Africa, few studies have set out to characterize predictors of bovine schistosomiasis  
103 infection. One study from Malaysia found that low weight, male sex, and older age were all risk factors for *S.*  
104 *spindale* infection in a range of different cattle species, though notably, no water buffalo species were included in  
105 this study (27). By comparison, in Southern China, a study conducted primarily among water buffaloes (96.2%  
106 water buffaloes, 3.8% cattle) found that infection intensity was highest in bovines under the age of two (28). These  
107 seemingly contradictory results may potentially be explained by isolation and limited grazing for calves in  
108 Malaysia (27), as well as potential genus-related differences in acquired immunity and self-cure rates (29). Of the  
109 two main types of bovines found in *S. japonicum* endemic areas, yellow cattle are believed to be more susceptible  
110 to infection than water buffalo based on studies assessing worm establishment success in the two genera (30,  
111 31). Nevertheless, He et al. (2001) also point out that water buffaloes may still act as important hosts in  
112 marshland areas of China, as they are more likely to spend time in water, and therefore more likely to be involved  
113 in the *S. japonicum* transmission cycle (10).

114 As *S. japonicum* is the primary species responsible for schistosomiasis infection in both humans and  
115 bovines in China, and given the considerable role that bovines are posited to have in contributing to human  
116 infection risk, studies aimed at assessing potential predictors of *S. japonicum* infection in bovines are of  
117 paramount importance. Not only is there a great deal of disagreement in the current body of literature over the key  
118 risk factors for bovine schistosomiasis infection, the limited studies to date have almost exclusively focused on  
119 physical/individual-level characteristics rather than broader environmental conditions. As such, this study set out  
120 to assess potential individual and environmental predictors of bovine *S. japonicum* infection in 2007, 2010 and  
121 2016 at the individual, household and village-levels in a region where schistosomiasis persistence has been  
122 demonstrated to exist in both humans and bovine populations.

123

## 124 **Methods**

### 125 **Village selection**

126 A longitudinal assessment of human and bovine infection was conducted in villages of Sichuan province  
127 in 2007, 2010 and 2016. Villages were located in the hilly regions of rural Sichuan and ranged from ~20-150  
128 households and a population of ~50-200 people. Village selection has been described previously (32). Briefly, to  
129 identify villages with evidence of *S. japonicum* reemergence, county surveillance records were reviewed from the  
130 year that transmission control was achieved in the county through March 2007. Out of eight Sichuan counties  
131 where schistosomiasis had been identified despite control efforts (33), three were selected for inclusion based on  
132 surveillance record availability and the local control stations' willingness and capacity to collaborate (32).  
133 However, due to a 7.9 magnitude earthquake in May 2008 that severely impacted one of the study counties (34),  
134 follow-up surveys were conducted in 36 villages in the two remaining counties in 2010. Based on infection rates, 7  
135 of the original 36 villages were surveyed again in 2016, in addition to 3 newly reemerging villages, giving a total of  
136 36 villages included in the analysis in 2007 and 2010, and 10 villages in 2016.

### 137 **Demographic, household and GPS surveys**

138 A village census was conducted in each collection year and all residents over the age of five were invited  
139 to participate in surveys and stool sample screenings for *S. japonicum* infection. In addition, attempts were made  
140 to survey all bovines in the village for *S. japonicum* infection. In the summers of 2007, 2010 and 2016, the head of  
141 each household was asked to complete a household survey that contained closed-ended questions related to  
142 socioeconomic status, domestic and farm animal ownership, sanitation and water access and agricultural  
143 practices. Bovine age, type and sex were collected at the time of the bovine infection surveys in 2007 and 2010  
144 (these data were not collected in 2016). Trained staff from the Sichuan Center for Disease Control and Prevention  
145 and the county Schistosomiasis Control Stations piloted and conducted all surveys in the local Sichuan dialect.

### 146 **Ethics statement**

147 This study was approved by the Sichuan Institutional Review Board, the University of California, Berkeley,  
148 Committee for the Protection of Human Subjects, and the Colorado Multiple Institutional Review Board. All  
149 participants provided written, informed consent. The collection of bovine samples we determined to be exempt

150 from review by the Animal Care and Use Committee at the University of California, Berkeley and the Institutional  
151 Animal Care and Use Committee at the University of Colorado Anschutz.

## 152 **Infection surveys**

153 Infection surveys were conducted by attempting to test three stool samples on three consecutive days  
154 from eligible humans and all bovines in the village. Infection surveys were conducted in November and December  
155 of 2007 and 2010, and July 2016. Individual bovines were isolated in a pen or tied up until a stool was produced  
156 on three separate days (consecutive, when possible). All stool samples were transported to the central laboratory  
157 soon after collection to be examined using the miracidium hatching test, following standard protocols (35). To  
158 account for the short survival and rapid hatching of bovine miracidia, the bovine samples were examine for  
159 miracidia at one, three and five hours after preparation for at least two minutes each time, whereas human  
160 samples were assessed at two, five and eight hours after preparation. One sample from each human was also  
161 examined using the Kato Katz thick smear procedure in 2007 and 2010 (36). A bovine was considered positive for  
162 *S. japonicum* if any miracidium hatching test was positive. A human was considered positive for *S. japonicum* if  
163 any miracidium hatching test or the Kato-Katz test was positive.

164 For each data collection period, the proportion of bovines in the village that were captured by infection  
165 surveys was assessed by comparing the total number of bovines reported in household surveys to the total  
166 number of bovines that participated in the infection survey in each village. Between 2007 and 2016, bovine  
167 infection status was assessed in 35/36 villages where residents reported owning bovines in 2007, 31/35 villages  
168 in 2010, and 8/8 villages in 2016. Details about participation and infection survey completeness are provided in  
169 S1 Table.

## 170 **Predictor selection and definitions**

171 The primary outcome of interest in this analysis was bovine *S. japonicum* infection in 2007, 2010 and  
172 2016. All candidate predictors were defined using either the household surveys or the human and bovine infection  
173 surveys and were divided into five categories: 1) biological/physical characteristics; 2) potential human sources of  
174 environmental schistosomes; 3) socio-economic indicators; 4) potential animal reservoirs/sources of infection; 5)  
175 agricultural risk factors. We included agriculture as its own category because bovines are frequently employed in  
176 agricultural work in China (13), and because different crop types and agricultural practices have their own

177 inherent exposure risks (e.g. planting wet crops like rice may increase the likelihood of contact with snail habitat  
178 and exposure to cercariae (37)). Variables identified as predictors with hypothesized similar mechanisms of  
179 transmission risk were aggregated where possible (e.g. wet vs. dry crops). Three crop type categories were  
180 created: winter crops, summer dry crops and summer wet crops (i.e. rice). Night soil use – that is, the collection of  
181 either treated or untreated human and/or animal waste for use as fertilizer – was also included as an agricultural  
182 risk factor and divided into three categories: night soil use on winter crops, dry summer crops and wet summer  
183 crops.

184           There were minor variations in the household survey content and question formulation across the study  
185 period. Namely, some variables were not available in all the study years (e.g. pig ownership was not assessed in  
186 2007). Where possible, continuous/discrete predictors were included over binary measures of a predictor for the  
187 household-level predictors. For binary variables, we excluded variables from the analysis of a given collection  
188 year if they represented very rare (<10%) or very common conditions (>90%). For continuous variables, variables  
189 were excluded when >90% of the observations took a single value. For example, household dog and pig  
190 ownership were both excluded in 2016 because >90% of the households owned one or more dogs (a binary  
191 variable), while >90% did not own any pigs. A composite household asset score (0-9) was developed for use in  
192 this assessment, which included eight household assets assessed in all three collection years (washing machine,  
193 television, air conditioner, refrigerator, computer, car, motorcycle), as well as a binary measure indicating that the  
194 home was made from either concrete, wood or bricks (vs. adobe).

195           Because prior work has demonstrated that group-level measures can serve as important predictors of  
196 schistosomiasis infection in humans (34), we also generated village-level candidate predictors from the household  
197 survey data. Village-level variables represent all households that participated in the household survey from a  
198 given village, even if they didn't own bovines. Village-level variables were either the village-average value of  
199 continuous household measures, or for binary variables, the proportion of the village population reporting the  
200 condition. Notably, the village-level variables excluded all observations from the bovine's own household, and  
201 instead used only the data from the other households in the village that participated in the household survey. This  
202 allowed for an assessment of how the surrounding village environment impacts individual bovine infection risk,  
203 independent of the home environment, whereas the household-level variables aim to unpack the influence of the



204 unique household environment on bovine infection status. The aforementioned predictor definitions and exclusion  
 205 criteria led to a total of 31 predictors of bovine infection, which are summarized in Table 1.

206

207 **Table 1. Summary of predictor variables included in the analysis.**

Predictor list	Scale of analysis	Years available	Variable type	Predictor category <sup>a</sup>
Bovine type	Individual	2007, 2010	Binary	Physical/biological
Bovine sex	Individual	2007, 2010	Binary	Physical/biological
Bovine age	Individual	2007, 2010	Continuous	Physical/biological
Number of hatch tests <sup>b</sup>	Individual	All	Discrete	Physical/biological
County of residence	Household	All	Binary	Physical/biological
Number of infected human household members	Household	All	Discrete	Potential human sources
Household has improved sanitation (y/n) <sup>c</sup>	Household	All	Binary	Socio-economic indicators
Household asset score (0-9)	Household	All	Discrete	Socio-economic indicators
Household cat ownership	Household	All	Binary	Animal reservoirs/sources
Household dog ownership	Household	2007, 2010	Binary	Animal reservoirs/sources
Household pig ownership	Household	2010	Discrete	Animal reservoirs/sources
Household bovine ownership	Household	All	Discrete	Animal reservoirs/sources
Household rice area	Household	All	Continuous	Agricultural risk factors
Household dry summer crop area	Household	All	Continuous	Agricultural risk factors
Household winter crop area	Household	All	Continuous	Agricultural risk factors
Household night soil rice: # buckets	Household	All	Discrete	Agricultural risk factors
Household night soil summer dry crop: # buckets	Household	All	Discrete	Agricultural risk factors
Household night soil winter crop: # buckets	Household	All	Discrete	Agricultural risk factors
Village prevalence of human infection	Village	All	Continuous	Potential human sources
Village prevalence of improved sanitation	Village	All	Continuous	Socio-economic indicators
Village mean asset score (0-9)	Village	All	Continuous	Socio-economic indicators
Village prevalence of cat ownership	Village	All	Continuous	Animal reservoirs/sources
Village prevalence of dog ownership	Village	All	Continuous	Animal reservoirs/sources
Village mean number of pigs owned	Village	2010, 2016	Continuous	Animal reservoirs/sources
Village mean number of bovines owned	Village	All	Continuous	Animal reservoirs/sources
Village mean rice area	Village	All	Continuous	Agricultural risk factors
Village mean dry summer crop area	Village	All	Continuous	Agricultural risk factors
Village mean winter crop area	Village	All	Continuous	Agricultural risk factors
Village night soil rice: mean # buckets	Village	All	Continuous	Agricultural risk factors
Village night soil summer dry crop: mean # buckets	Village	All	Continuous	Agricultural risk factors
Village night soil winter crop: mean # buckets	Village	All	Continuous	Agricultural risk factors

208

209 <sup>a</sup> Predictors were grouped into five categories relevant to bovine *S. japonicum* infection risk probability: 1) physical/biological  
 210 characteristics (e.g. old-age); 2) potential human sources of environmental schistosomes (e.g. human *S. japonicum* infection  
 211 prevalence in the bovine's household); 3) socio-economic indicators (e.g. prevalence of improved sanitation systems in the  
 212 surround village); 4) potential animal reservoirs/sources of infections (e.g. prevalence of dog ownership in the surrounding  
 213 village); 5) agricultural risk factors (e.g. a household's total rice crop area).

214 <sup>b</sup> Because not all bovines produced three stool samples, and examination of a greater number of stool samples can increase  
 215 the probability of detecting infection, the number of hatch tests used on a given bovine was also included as a predictor in our  
 216 analyses.

217 <sup>c</sup> Improved sanitation was defined as access to an improved toilet in the household, including a biogas digester or a three-  
 218 compartment toilet.

219

## 220 Analysis

221 Across the collection years, 67 bovines with infection data were excluded from this analysis due to lack of  
222 household survey data (30/503 bovines in 2007; 36/233 bovines in 2010; and 1/72 in 2016). Infection prevalence  
223 was similar among the excluded bovines (11/67, 16.4%) as compared to those included in this analysis (111/741,  
224 15.0%). Among the remaining bovines included, household and village-level variables generally had low levels of  
225 missing data (all <20% missing). By contrast, the individual-level bovine data was recorded with less consistency:  
226 the variable with the most missing data was bovine sex in 2007 (21.6% missing). Missing values were imputed  
227 separately for each collection year for all variables with <25% missing using the `rflmpute` function from the  
228 “`randomForest`” package in R (38, 39).

229 Spatial patterns of bovine infection prevalence were inspected using ESRI’s ArcGIS ArcMap software  
230 release 10.5.1 (40). Categorical versions of each of the individual, household and village-level candidate  
231 predictors were generated and compared between *S. japonicum* infected and uninfected bovines to investigate  
232 potential changes in predictor distribution patterns by infection status across the study period.

233 To determine which of our candidate predictors serve as the best predictors of schistosomiasis  
234 transmission in 2007, 2010 and 2016, a random forests (RF) machine learning approach was used. For each  
235 collection year, 25% of the data was reserved for validation, while the remaining 75% was used for model  
236 construction. To address class imbalance in our outcome of interest (bovine *S. japonicum* prevalence of 13.3%,  
237 17.3%, and 19.7% in 2007, 2010 and 2016, respectively), over sampling of the minority class was conducted. For  
238 model tuning, 10-fold cross validation was performed using the `Caret` package in R to help select the optimal  
239 maximum node size and the number of variables to try at each branch. Once the optimal value of each of these  
240 parameters was determined, a final model was run using 5000 trees per forest (41).

241 For each collection year, we conducted a total of ten rebalancing and model tuning iterations to assess  
242 the degree of stability in our variable importance rankings. The mean decrease in accuracy (MDA) value was  
243 used to rank the top ten predictors from each model on a scale of ten to one from most important (highest MDA)  
244 to least important (lowest MDA). These variable rankings were then summed across the 10 rebalancing iterations  
245 to give a 10-model summary score of variable importance, ranging from 100-1 and the ten highest scoring  
246 variables from the ten-model summary score were then reassigned a final ranking of 1-10. Next, using only those  
247 predictors ranked first through tenth within each collection year, we performed an additional ten iterations of the

248 aforementioned balancing and tuning process to create “lean” prediction model summary score of variable  
249 importance, thereby reducing excess noise in the variable ranking assessment caused by including a large  
250 number of candidate predictors. Because we hypothesized that the inclusion of human infection as a predictor of  
251 bovine infection would strongly influence the predictive capacity of our RF models due to a presumed association  
252 between bovine and human infection, we also conducted ten iterations of a sensitivity analysis for each collection  
253 year that excluded the human infection variables from the assessment. The ability of the full, lean and sensitivity  
254 RF models to predict infection status was then assessed using ROC area under the curve and accuracy. In the  
255 case of disagreement or a tie when comparing our chosen performance metrics, the sensitivity, kappa and  
256 specificity were subsequently compared to select the top performing model for each year.

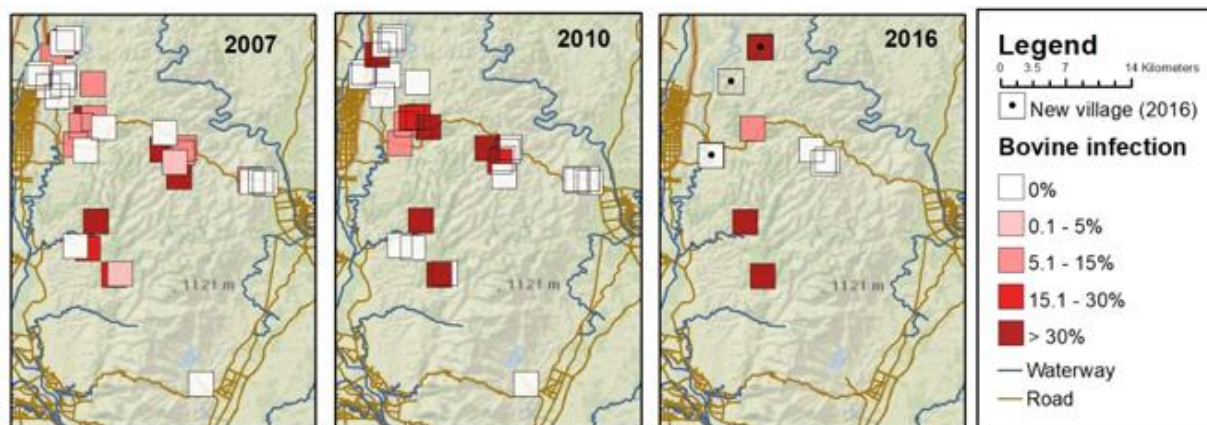
257 Each of the full model, lean model and sensitivity model summary scores were used to generate heat  
258 maps highlighting variable importance scores within each collection year, their change over time, and the  
259 frequency with which the different levels of analysis (individual, household or village) were each found to be  
260 among the top ten most important predictors. Simple logistic regression analyses were performed to assess the  
261 direction of association between the top predictors and bovine infection, dividing continuous variables into tertiles  
262 to assess for potential non-linear relationships. The direction of association was recorded for the top predictors  
263 within each collection year, using a p-value of <0.2 to indicate weak evidence of a between-group difference. In  
264 the case that no difference was indicated between tertile groups at the p<0.2 level, the predictors were further  
265 divided into quartiles and re-assessed. If still no evidence for a between group difference was identified using  
266 quartiles, this point was noted in the results. Density plots by infection status were also examined for a subset of  
267 predictors that were found to have a change in the direction of association across the collection years. Stata 15  
268 and R Studio 4.0 were used for all analyses (42, 43).

269

## 270 **Results**

271 This analysis included bovines from 37 villages across the study period, with a total of 473 bovines from  
272 35 villages in 2007, 197 bovines from 31 villages in 2010, and 71 bovines from 8 villages in 2016. The overall  
273 bovine infection prevalence was 13.3%, 17.3% and 19.2% in 2007, 2010 and 2016, respectively. Figure 1 shows  
274 a map of bovine infection distribution by village across two counties.

275 **Figure 1. Village-level prevalence of schistosomiasis in bovines in 2007, 2010 and 2016.** Unshaded squares  
276 indicate study villages where no bovines were tested. Service Layer Credits: National Geographic, ESRI,  
277 DeLorme, HERE, UNEP-WCMC, USGS, NASA, METI, NOAA, increment P Corp, and OpenStreetMap  
278 Contributors, Geofabrik GmbH, Copyright 2018.



## 281 **Bovine infection prevalence and individual-level characteristics**

282 Of the individual characteristics assessed in this analysis, none were consistently associated with  
283 infection status (Table 2). For example, in 2007, water buffalo were less likely to be infected than cattle (8.1% in  
284 water buffalo, 14.2% in cattle), while in 2010, the prevalence of infection was ~17% in both groups. Similarly, in  
285 2007 bovines over 5 years of age were three times more likely to be infected than younger bovines, but in 2010  
286 there wasn't a clear pattern of infection by age. Notably, there were fewer water buffaloes than cattle in both  
287 assessment years (18.5% of all bovines were water buffalo in 2007; 15.5% in 2010), bovines were predominantly  
288 female (86.3% in 2007; 87.2% in 2010), and ranged widely in age from less than a year to 26 years old.

289

290 **Table 2. Tabulation of individual-level predictors by bovine infection status.**  
291

		2007			2010			2016		
		N Positive	Total tested	% Positive	N Positive	Total tested	% Positive	N Positive	Total tested	% Positive
<b>County</b>	<b>Total</b>	<b>63</b>	<b>473</b>	<b>13.32%</b>	<b>34</b>	<b>197</b>	<b>17.26%</b>	<b>14</b>	<b>71</b>	<b>19.72%</b>
	1	15	191	7.85%	17	71	23.94%	3	21	14.29%
	2	48	282	17.02%	17	126	13.49%	11	50	22.00%
<b>Bovine sex</b>	Female	50	320	15.63%	29	163	17.79%	ND <sup>a</sup>	ND <sup>a</sup>	
	Male	9	51	17.65%	4	24	16.67%	ND <sup>a</sup>	ND <sup>a</sup>	
	Missing	4	102	3.92%	1	10	10.00%			
<b>Bovine age<sup>b</sup></b>	<2	3	43	6.98%	6	35	17.14%	ND <sup>a</sup>	ND <sup>a</sup>	
	2 to 4	26	220	11.82%	9	59	15.25%	ND <sup>a</sup>	ND <sup>a</sup>	
	5+	28	137	20.44%	18	93	19.35%	ND <sup>a</sup>	ND <sup>a</sup>	
	Missing	6	73	8.22%	1	10	10.00%			
<b>Bovine type</b>	Water buffalo	7	86	8.14%	5	29	17.24%	ND <sup>a</sup>	ND <sup>a</sup>	
	Cattle	54	380	14.21%	28	158	17.72%	ND <sup>a</sup>	ND <sup>a</sup>	
	Missing	2	7	28.57%	1	10	10.00%			
<b>Hatch tests (N)</b>	1	3	59	5.08%	0	29	0.00%	1	10	10.00%
	2	8	80	10.00%	6	41	14.63%	2	9	22.22%
	3	52	334	15.57%	28	127	22.05%	11	52	21.15%

292 <sup>a</sup> No data (ND). Data was not collected on bovine sex, age or type in 2016.

293 <sup>b</sup> Although age is broken into categories to facilitate a comparison of the distributions between 2007 and 2010, it was  
294 assessed as a continuous variable in our RF models.  
295

## 296 **Bovine infection prevalence and household-level characteristics**

297 Bovine infection prevalence was highest among bovines in households where one or more humans were  
298 infected, and in households that did not own pigs (Table 3). Relationships between bovine infection and the rest  
299 of the household predictors were inconsistent across years. For example, access to improved sanitation and  
300 infection status by sanitation group shifted across our study period, rising from 22.8% of households reporting  
301 improved sanitation in 2007 and roughly equal infection prevalence in the two sanitation groups, to 52.1% with  
302 improved sanitation by 2016 and a higher probability of infection in the households with unimproved sanitation  
303 (23.5%) compared to the households with improved sanitation (16.2%). Across the study period, there was a  
304 steady increase in the prevalence of households reporting planting rice (69.1% in 2007; 71.5% in 2010; 81.7% in  
305 2016), other summer crops (77.2% in 2007; 98% in 2010; 100% in 2016) and winter crops (97.7% in 2007; 99% in  
306 2010; 100% in 2016). For rice crops in 2007 and 2010, the prevalence of bovine infection increased as the area of  
307 rice crop planted increased, whereas in 2016, this pattern did not hold. Other noteworthy changes in agricultural  
308 production across the study period includes a decrease in night soil use on rice and winter crops: the proportion of

309 households reporting any night soil use on rice crops dropped from 35.6% to 11.7% between 2007 and 2016, and  
 310 for winter crops it dropped from 58.3% to 12.6%. By contrast, the proportion of night soil users for summer crops  
 311 remained relatively constant across the study period (52.4% in 2007; 53.5% in 2010; 50.7% in 2016).

312

313 **Table 3. Household predictors by bovine infection status.**

		2007			2010			2016		
		# Positive	# tested	% Positive	# Positive	# tested	% Positive	# Positive	# tested	% Positive
<b># of infected humans in the household</b>	0	38	366	10.38%	19	154	12.34%	12	61	19.67%
	1+	19	84	22.62%	13	29	44.83%	2	10	20.00%
	Missing	6	23	26.09%	2	14	14.29%	0	0	--
<b>Household asset score (0-9)</b>	0 - 2	22	184	11.96%	4	24	16.67%	2	10	20.00%
	2 - 3	31	211	14.69%	15	89	16.85%	3	14	21.43%
	4+	10	78	12.82%	15	84	17.86%	9	47	19.15%
<b>Improved sanitation</b>	No	47	365	12.88%	27	129	20.93%	8	34	23.53%
	Yes	16	108	14.81%	7	68	10.29%	6	37	16.22%
<b>Household owns dogs<sup>a</sup></b>	No	6	67	8.96%	3	29	10.34%	2	7	28.57%
	Yes	57	404	14.11%	31	168	18.45%	12	64 <sup>a</sup>	18.75%
	Missing	0	2	0.00%	0	0	--	0	0	--
<b>Household owns cats</b>	No	17	155	10.97%	12	69	17.39%	5	23	21.74%
	Yes	46	312	14.74%	22	128	17.19%	9	48	18.75%
	Missing	0	6	0.00%	0	0	--	0	0	--
<b># Pigs owned by household<sup>a, b</sup></b>	0	ND <sup>b</sup>	ND <sup>b</sup>		27	116	23.28%	14	64 <sup>a</sup>	21.88%
	1+	ND <sup>b</sup>	ND <sup>b</sup>		7	81	8.64%	0	7	0.00%
<b># Other bovines owned by household</b>	0	49	334	14.67%	20	126	15.87%	8	38	21.05%
	1+	14	139	10.07%	14	71	19.72%	6	33	18.18%
<b>Total rice crop area (mu)</b>	0	16	146	10.96%	5	56	8.93%	5	13	38.46%
	<2	24	195	12.31%	16	84	19.05%	4	31	12.90%
	2+	23	132	17.42%	13	57	22.81%	5	27	18.52%
<b>Total dry summer crop area (mu)</b>	0	8	108	7.41%	0	4	0.00%	0	0	--
	<3	39	242	16.12%	18	99	18.18%	2	26	7.69%
	3+	16	123	13.01%	16	94	17.02%	12	45	26.67%
<b>Total winter crop area (mu)</b>	<2	14	142	9.86%	9	52	17.31%	7	29	24.14%
	2-3	33	223	14.80%	18	97	18.56%	7	29	24.14%
	4+	16	108	14.81%	7	48	14.58%	0	13	0.00%
<b># buckets night soil on rice</b>	0	40	304	13.16%	29	168	17.26%	8	53	15.09%
	1+	22	168	13.10%	5	29	17.24%	1	7	14.29%
	Missing	1	1	100.0%	0	0	--	5	11	45.45%



# buckets night soil on dry summer crops	0	31	248	12.50%	24	105	22.86%	6	36	16.67%
	1-25	6	53	11.32%	5	32	15.63%	5	22	22.73%
	>25	26	172	15.12%	5	60	8.33%	3	13	23.08%
# buckets night soil household winter crops	0	27	197	13.71%	31	137	22.63%	9	62	14.52%
	1-26	5	79	6.33%	0	16	0.00%	1	3	33.33%
	>26	31	197	15.74%	3	44	6.82%	4	6	66.67%

314

315 <sup>a</sup> Variables with >90% of observations taking on a single value were excluded from the RF assessment. This exclusion criteria  
 316 applied twice in 2016: >90% of the included households reported owning at least one dog, and >90% of households owned  
 317 zero pigs.

318 <sup>b</sup> ND. Data was not collected on pig ownership in 2007.

319

## 320 Bovine infection prevalence and village-level characteristics

321 Bovine infection prevalence was highest in villages with high levels of bovine ownership (Table 4). For the  
 322 remaining village-level predictors however, the infection patterns were inconsistent. For example, bovine infection  
 323 prevalence was highest in 2007 and 2010 for bovines residing in villages where a high percentage of the human  
 324 population was infected, whereas in 2016, that pattern did not hold. In 2007, infection prevalence incrementally  
 325 decreases as the percent of households in the village that own dogs increases, but infection prevalence was  
 326 higher among bovines residing in villages with higher dog ownership in 2010 and 2016. Similarly, in 2007, the  
 327 prevalence of bovine infection is highest in villages where more night soil is used on rice crops, dry summer crops  
 328 and winter crops, but in 2010, bovine infection prevalence decreases as the surrounding village's night soil use  
 329 increases. Notably, the average amount of night soil being applied to crops dropped across the study period for all  
 330 crop types.

331

332 **Table 4. Village-level predictors by bovine infection status.**

		2007			2010			2016		
		# Positive	# tested	% Positive	# Positive	# tested	% Positive	# Positive	# tested	% Positive
Human infection prevalence (%)	0-2.53%	14	147	9.52%	0	76	0.00%	4	27	14.81%
	2.54- 11.11%	13	158	8.23%	11	57	19.30%	8	29	27.59%
	≥11.12%	36	168	21.43%	23	64	35.94%	2	15	13.33%
Mean household asset score (0-9)	<2	39	225	17.33%	0	0	0.00%	0	0	0.00%
	2-3	24	248	9.68	26	143	18.18%	2	23	8.70%
	≥4	0	0	0.00%	8	54	14.81%	12	48	25.00%

<b>% households with an improved toilet</b>	<10%	14	180	7.78%	3	24	12.50%	0	0	0.00%
	10-50%	38	215	17.67%	29	133	21.80%	6	24	25.00%
	≥50%	11	78	14.10%	2	40	5.00%	8	47	17.02%
<b>% of households that own at least one dog</b>	<70%	17	98	17.35%	2	51	3.92%	7	50	14.00%
	70-85%	34	217	15.67%	18	94	19.15%	7	21	33.33%
	≥85%	12	158	7.59%	14	52	26.92%	0	0	0.00%
<b>% of households that own at least one cat</b>	<45%	12	109	11.01%	12	71	16.90%	4	36	11.11%
	45-60%	21	124	16.94%	14	73	19.18%	5	13	38.46%
	≥60%	30	240	12.50%	8	53	15.09%	5	22	22.73%
<b>Mean number of pigs owned</b>	0	ND <sup>a</sup>	ND <sup>a</sup>		1	21	4.76%	11	20	55.00%
	0.01 - 1	ND <sup>a</sup>	ND <sup>a</sup>		26	119	21.85%	1	45	2.22%
	≥ 1	ND <sup>a</sup>	ND <sup>a</sup>		7	57	12.28%	2	6	33.33%
<b>Mean number of bovines owned</b>	< 0.5	5	116	4.31%	4	59	6.78%	10	57	17.54%
	0.5 – 1	40	300	13.33%	19	99	19.19%	4	14	28.57%
	≥ 1	18	57	31.58%	11	39	28.21%	0	0	0.00%
<b>Mean area of rice planted (mu)</b>	<0.75	14	143	9.79%	9	90	10.00%	7	18	38.89%
	0.75-1.5	19	183	10.38%	12	48	25.00%	5	41	12.20%
	≥1.5	30	147	20.41%	13	59	22.03%	2	12	16.67%
<b>Mean area of dry summer crops planted (mu)</b>	<1	11	184	5.98%	1	14	7.14%	0	3	0.00%
	1 -2.5	44	195	22.56%	23	110	20.91%	2	26	7.69%
	≥2.5	8	94	8.51%	10	73	13.70%	12	42	28.57%
<b>Mean area of winter crops plants (mu)</b>	<2	21	149	14.09%	17	77	22.08%	12	53	22.64%
	2-2.75	18	178	10.11%	11	66	16.67%	0	2	0.00%
	≥2.75	24	146	16.44%	6	54	11.11%	2	16	12.50%
<b>Mean # buckets of night soil used on rice crops</b>	<1	3	58	5.17%	29	146	19.86%	8	41	19.51%
	1-9.9	21	153	13.73%	4	38	10.53%	6	27	22.22%
	≥10	39	262	14.89%	1	13	7.69%	0	3	0.00%
<b>Mean # buckets of night soil used on dry summer crops</b>	<10	13	190	6.84%	24	99	24.24%	1	15	6.67%
	10-30	18	109	16.51%	7	47	14.89%	13	47	27.66%
	≥30	32	174	18.39%	3	51	5.88%	0	9	0.00%
<b>Mean # buckets of night soil used on winter crops</b>	<10	1	26	3.85%	24	121	19.83%	14	71	100.0%
	10-30	19	218	8.72%	7	41	17.07%	0	0	--
	≥30	43	229	18.78%	3	35	8.57%	0	0	--

333

334 <sup>a</sup> ND. Data was not collected on pig ownership in 2007.

335

336



## 337 Predictors of bovine infection

338 The full models, lean models and sensitivity models within a given collection year all resulted in relatively  
339 stable rankings, while more variability in predictor rankings is seen when comparing across collection years  
340 (Figure 2). Within each model type for a given year, the ten iterations of re-balancing and tuning led to some  
341 variation in the MDA scores across the top ten predictors, with the top five more consistent in their high rankings.  
342 This is particularly prominent in the 2007 and 2010 models, whereas 2016 showed more variation overall. With  
343 few exceptions (4/60 model iterations), variables that scored in the top five within any of the ten iterations of either  
344 the full or sensitivity models were among the top ten predictors using the 10-model summary scores.

345 Agricultural variables were most frequently ranked in the top ten across all years. Specifically, the  
346 household area of winter crops planted, the mean area of rice planted in the surrounding village, and the mean  
347 amount of night soil applied to dry summer crops in the surrounding village were all ranked in the top ten for all  
348 collection years and the full, lean and sensitivity analyses. Additionally, the total household area of summer crops  
349 planted, the village mean area of winter crops and the mean number of bovines owned by the surrounding village  
350 were also all among the top ten predictors in at least one of the three model types used for 2007, 2010 and 2016.  
351 Of those predictors that ranked in the top ten in at least one collection year, four were scaled to the village-level,  
352 and two were assessed at the household-level. Because the full list of predictors changed slightly across the  
353 collection years, a supplemental analysis was conducted in which only predictors that were available in all three  
354 collection years were included in the RF models. This analysis demonstrated that 1) intra-year rankings and extra-  
355 year patterns did not change substantially, and 2) agricultural variables remained the most prominent predictor  
356 category when comparing across the entire study period. See S1 Figure for details of the supplemental analysis.

357

358 **Figure 2. Variable importance rankings and direction of association for candidate predictors of bovine *S.***  
359 ***japonicum* infection in 2007, 2010 and 2016.** Variable importance rankings are based on a composite of mean  
360 decrease in accuracy scores for 10 random forest (RF) models for each model type (full, lean and sensitivity) and  
361 collection year. The direction of association was determined through logistic regression, using tertile categories  
362 for continuous variables to assess evidence for non-linearity. A p-value of <0.2 was used to indicate evidence of a  
363 between-group difference, and, when a between group difference was found, the direction of association is  
364 indicated. See S2 Table for detailed logistic regression results.

Predictor categories	2007			2010			2016																				
	Full	Lean	Sens.	Full	Lean	Sens.	Full	Lean	Sens.																		
<b>Physical/biological characteristics</b>																											
Bovine type							ND	ND	ND																		
Bovine sex							ND	ND	ND																		
Bovine age	3 ↑	1 ↑	2 ↑				ND	ND	ND																		
Number of hatch tests				2 ↑	2 ↑	1 ↑																					
County of residence																											
<b>Potential human sources of schistosomes</b>																											
Human infection prevalence in the village	5 ↑	3 ↑	Excl.	1 ↑	1 ↑	Excl.			Excl.																		
Number of infected human household members			Excl.	6 ↑	7 ↑	Excl.			Excl.																		
<b>Socio-economic indicators</b>																											
Village prevalence of improved sanitation						9 ∩																					
Household has improved sanitation (y/n)																											
Village mean asset score (0-9)			10 ↓				8 ↑	5 ↑	8 ↑																		
Household asset score (0-9)						10 ∅																					
<b>Potential animal reservoirs/sources of infection</b>																											
Village prevalence of cat ownership							9* ↑	10 ↑	10 ↑																		
Household cat ownership																											
Village prevalence of dog ownership							9* ↑	9 ↑																			
Household dog ownership							Excl.	Excl.	Excl.																		
Village mean number of pigs owned	ND	ND	ND				3 ↓	7 ↓	4 ↓																		
Household pig ownership	ND	ND	ND	7 ↓	8 ↓	6 ↓	Excl.	Excl.	Excl.																		
Village mean number of bovines owned	8 ↑	8 ↑	6 ↑				2 ↑	1 ↑	1 ↑																		
Household bovine ownership																											
<b>Agricultural risk factors</b>																											
Village mean rice area	6 ↑	5 ↑	3 ↑	5 ↑	5 ↑	3* ↑	5 ↓	8 ↓	5 ↓																		
Household rice area				8 ↑	10 ↑	7 ↑																					
Village mean dry summer crop area	2 ↑	7 ↑	4 ↑				6 ↑	6 ↑	7 ↑																		
Household dry summer crop area	4 ↑	6 ↑	5 ↑	4 ↓	4 ↓	3* ↓			9 ↑																		
Village mean winter crop area	1 ↑	2 ↑	1 ↑	10 ↓	6 ↓		1 ↓	2 ↓	2 ↓																		
Household winter crop area	7 ↑	4 ↑	7 ↑	9 ∅	9 ∅	5 ∅	7 ↓	4 ↓	6 ↓																		
Village night soil rice: mean # buckets																											
Household night soil rice: # buckets																											
Village night soil summer dry crop: mean # buckets	9 ↑	10 ↑	8 ↑	3 ↓	3 ↓	2 ↓	4 ↑	3 ↑	3 ↑																		
HH night soil summer dry crop: # buckets																											
Village night soil winter crop: mean # buckets	10* U	9 U	9 U																								
Household night soil winter crop: # buckets	10* U																										
<p><b>Symbol key:</b></p> <ul style="list-style-type: none"> <li>↑ Positive association</li> <li>↓ Negative association</li> <li>∩ Non-linear association (Rise-Fall)</li> <li>U Non-linear association (Fall-Rise)</li> <li>ND Data not collected</li> <li>Excl. Variable excluded from model</li> <li>* Tied for importance rank</li> <li>∅ No evidence for a significant between - group difference found</li> </ul> <p><b>Color key: Variable importance rankings (1<sup>st</sup>-10<sup>th</sup>) by scale</b></p> <table border="1" style="margin-left: 20px;"> <tr> <td style="background-color: #fff9c4;"><b>Individual</b></td> <td style="background-color: #fff9c4;">1<sup>st</sup> – 2<sup>nd</sup></td> <td style="background-color: #fff9c4;">3<sup>rd</sup> – 4<sup>th</sup></td> <td style="background-color: #fff9c4;">5<sup>th</sup> – 6<sup>th</sup></td> <td style="background-color: #fff9c4;">7<sup>th</sup> – 8<sup>th</sup></td> <td style="background-color: #fff9c4;">9<sup>th</sup> – 10<sup>th</sup></td> </tr> <tr> <td style="background-color: #e1f5fe;"><b>Household</b></td> <td style="background-color: #e1f5fe;">1<sup>st</sup> – 2<sup>nd</sup></td> <td style="background-color: #e1f5fe;">3<sup>rd</sup> – 4<sup>th</sup></td> <td style="background-color: #e1f5fe;">5<sup>th</sup> – 6<sup>th</sup></td> <td style="background-color: #e1f5fe;">7<sup>th</sup> – 8<sup>th</sup></td> <td style="background-color: #e1f5fe;">9<sup>th</sup> – 10<sup>th</sup></td> </tr> <tr> <td style="background-color: #e1f5fe;"><b>Village</b></td> <td style="background-color: #e1f5fe;">1<sup>st</sup> – 2<sup>nd</sup></td> <td style="background-color: #e1f5fe;">3<sup>rd</sup> – 4<sup>th</sup></td> <td style="background-color: #e1f5fe;">5<sup>th</sup> – 6<sup>th</sup></td> <td style="background-color: #e1f5fe;">7<sup>th</sup> – 8<sup>th</sup></td> <td style="background-color: #e1f5fe;">9<sup>th</sup> – 10<sup>th</sup></td> </tr> </table> <p style="margin-left: 20px;">= Not ranked 1-10 for the given model/year</p>										<b>Individual</b>	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>	<b>Household</b>	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>	<b>Village</b>	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>
<b>Individual</b>	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>																						
<b>Household</b>	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>																						
<b>Village</b>	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>																						

365

366

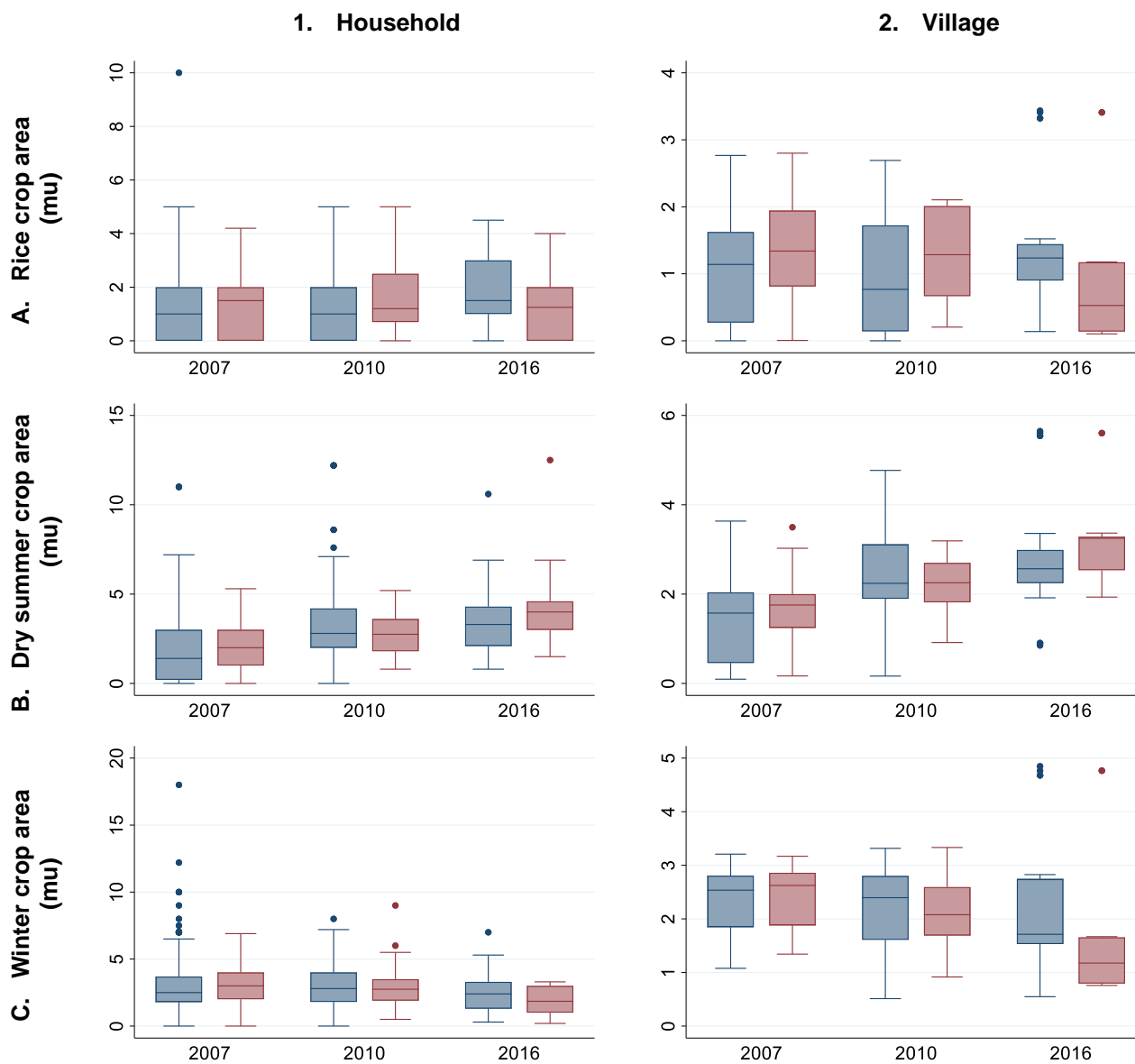
367 Despite the inter-year agreement for several of the agricultural variables' high importance rankings, the  
368 direction of association between the top agricultural predictors and bovine infection was not consistent across the  
369 three collection years. For example, the logistic regression assessments suggest that the direction of association  
370 with bovine infection flips from positive to negative for village rice crop area (2007 & 2010 = ↑; 2016 = ↓), and  
371 winter crop area (2007 = ↑; 2010 & 2016 = ↓), while for household summer crop area and village night soil use on  
372 summer crops, the direction of association flips from positive to negative to positive (2007 = ↑; 2010 = ↓; 2016 =  
373 ↑). Notably, in 2007 increases in all the key agricultural predictors were associated with an increase in bovine  
374 infection risk, apart from night soil use on winter crops. By contrast, in 2010 and 2016 our models indicate a  
375 mixture of positive and negative associations across the key agricultural predictors, and in one instance  
376 (household winter crop area in 2010), no evidence of a relationship was found.

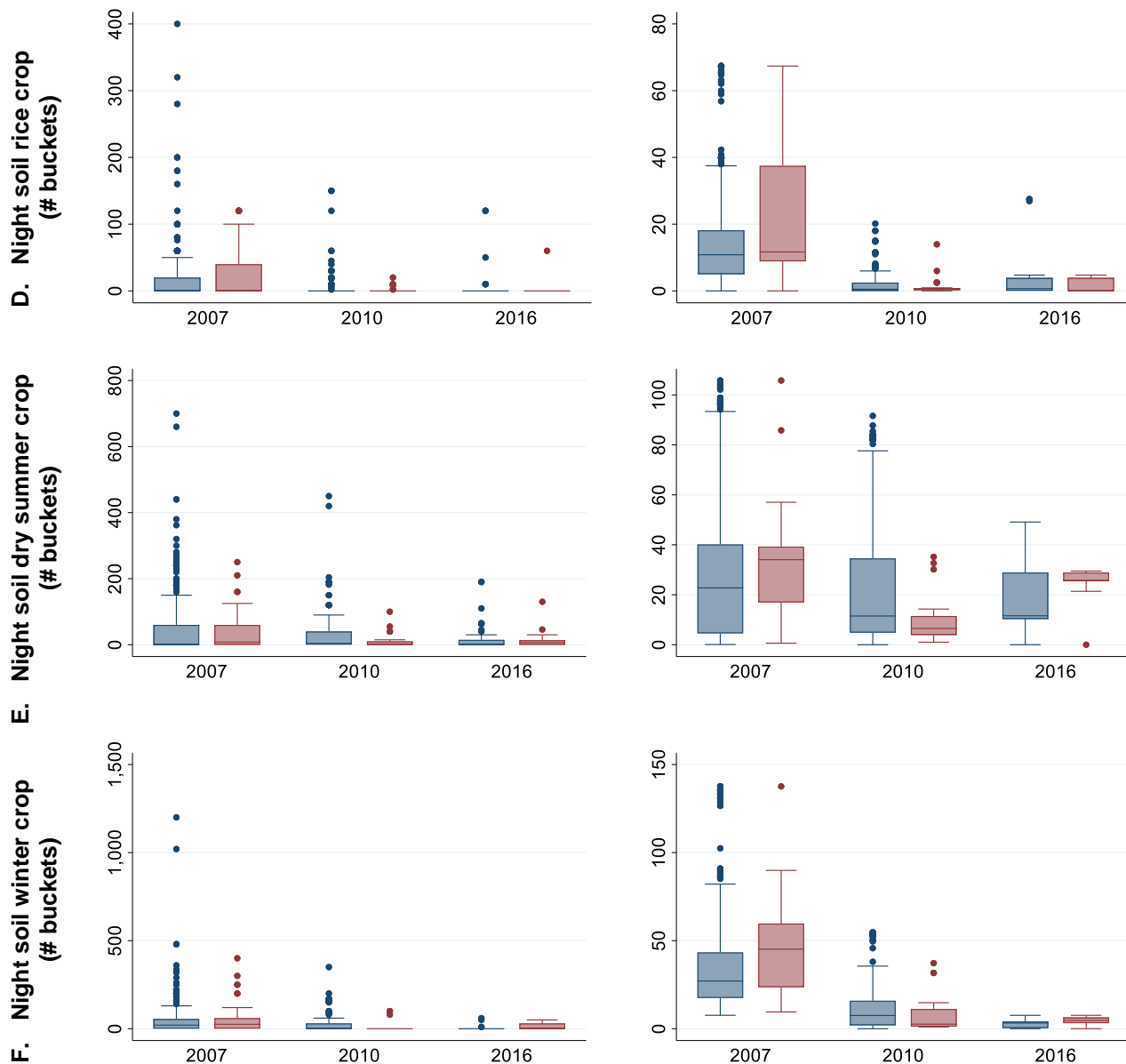
377 As mentioned above, the proportion of households planting rice, dry summer crops and winter crops, and  
378 the proportion of households reporting night soil use on rice and winter crops (but not dry summer crops) all  
379 shifted over the study period. Figure 3 depicts these shifting patterns over time, illustrating changes in the  
380 distribution of different agricultural practices by bovine infection status between 2007 and 2016. Despite the  
381 previously noted rise in the prevalence of households farming rice, dry summer crops and winter crops over the  
382 study period, panels A-C of Figure 3 show that only dry summer crop farming saw a notable increase in the total  
383 and mean area of crop being planted by households and villages between 2007 and 2016. On the other hand, a  
384 general decrease in the overall range and mean number of buckets of night soil being applied to rice, winter crops  
385 and, to a lesser extent, dry summer crops, can be observed when comparing between 2007 and 2016 (Figure 3,  
386 panels D-F).

387

388

389 **Figure 3. Changes in agricultural practices and the relationship between bovine infection and agricultural**  
390 **predictors over time.** For each of the agricultural predictors included in this analysis, boxplots are used to  
391 represent the distribution of uninfected (blue), infected (red), for household-level (left) and village-level variables  
392 (right) in 2007, 2010 and 2016.





*S. japonicum* negative bovines

*S. japonicum* positive bovines

393

394

395

396

397

398

399

In addition to the agricultural variables, there are also some other notable predictors that stand out in one or more collection year. Village bovine ownership is among the top ten predictors in at least one RF analysis from each year, all of which indicate that an increase in bovine ownership in the surrounding village corresponds with an increase in bovine infection risk. Human infection prevalence in the surrounding village was among the top five predictors of bovine infection in 2007 and 2010, and the number of infected humans within the household was among the top ten predictors in 2010. For both the household and village human infection predictors, an increase

400 in human infections was associated with an increase in bovine infections. When the human infection predictors  
401 were removed for the sensitivity analysis, the rankings of the remaining predictors did not shift substantially in any  
402 collection year. Of the physical/biological characteristics assessed, bovine age was among the top predictors in  
403 2007, with the logistic regression results suggesting that a bovine's infection risk increased with age. In all of the  
404 2010 analyses, the number of hatch tests was an important predictor of bovine infection, a feature not shared by  
405 the 2007 and 2016 analyses. This may be related to the relatively high proportion of bovines that had less than  
406 three hatch test results in 2010 (35.5%), as compared to 2007 (29.4%) and 2016 (26.8%).

407         Of the three different analyses performed (full, lean and sensitivity) for each collection year, the full  
408 models (i.e. those that included the full list of predictors available in a given year) tended to perform the best, as is  
409 highlighted in Table 5. Overall our models had high accuracy values, with the top performing models producing a  
410 maximum accuracy of 0.864 (95% CI: 0.79 – 0.92) in 2007, 0.816 (95% CI: 0.68 – 0.91) in 2010, and 1.0 (0.81 –  
411 1.0), in 2016. However, due to class imbalance in our reserved test datasets (see the no information rate (NIR) in  
412 Table 5), the Kappa value is a useful performance metric for our models, as this takes class imbalance into  
413 account. According to the benchmarks laid out by Landis and Koch (1977), the Kappa statistics from our 2007  
414 analyses suggest a “Fair” level of agreement (0.21 – 0.40) between our best RF models and the true known  
415 values in 2007. For 2010, the highest Kappa statistic came from the full predictor analysis, with a Kappa of 0.463,  
416 indicating a “Moderate” level of agreement (0.41 – 0.60) between the prediction model and the reserved test  
417 dataset (44). In 2016, both the full and sensitivity models achieved perfect prediction (Kappa = 1) for the test  
418 dataset in at least one of the ten model iterations, whereas the Kappa statistic for the top performing lean model  
419 was 0.853, or “Almost Perfect”, according to Landis & Koch (44).

420

421

422 **Table 5. Comparison of model performance metrics for the top performing model from the full, lean and**  
 423 **sensitivity analyses in 2007, 2010 and 2016.** The top performing model was defined as the one with the highest  
 424 accuracy for each analysis type (full, lean & sensitivity) and collection year (2007, 2010, 2016). In the case of a tie  
 425 for the highest accuracy value, the sensitivity, kappa and specificity were subsequently compared to select the top  
 426 performing model for each analysis type and year.  
 427

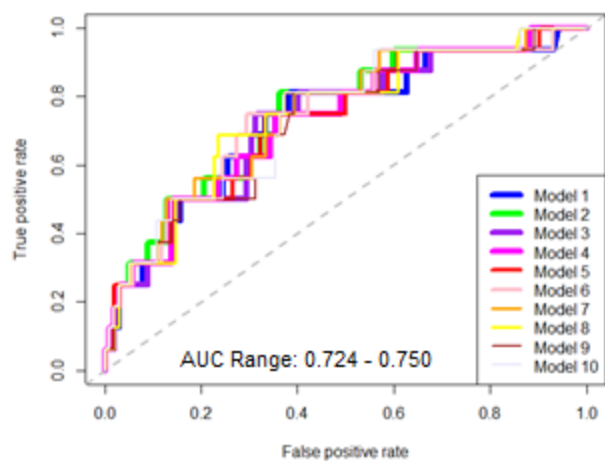
Performance metrics	2007 Models			2010 Models			2016 Models		
	Full	Lean	Sens.	Full	Lean	Sens.	Full	Lean	Sens.
Accuracy	0.864	0.864	0.856	0.816	0.816	0.816	1	0.944	1
95% CI	0.79 - 0.92	0.79 - 0.92	0.78 - 0.91	0.68 - 0.91	0.68 - 0.91	0.68 - 0.91	0.81 - 1	0.73 - 1	0.81 - 1
No info rate (NIR)	0.864	0.864	0.864	0.837	0.837	0.837	0.778	0.778	0.778
p-value (Acc > NIR)	0.566	0.566	0.667	0.729	0.779	0.729	0.011	0.067	0.011
Kappa <sup>a</sup>	0.313	0.313	0.246	0.463	0.360	0.416	1	0.853	1
Sensitivity	0.313	0.313	0.25	0.75	0.5	0.505	1	1	1
Specificity	0.951	0.951	0.951	0.829	0.878	0.625	1	0.929	1
Pos pred Value	0.5	0.5	0.444	0.462	0.444	0.854	1	0.8	1
Neg Pred Value	0.898	0.898	0.890	0.944	0.9	0.455	1	1	1

428 <sup>a</sup> Due to the high degree of imbalance between the outcome classes across the study period, the Kappa values is a useful  
 429 metric for our models, as it helps to correct bias that results when rewarding the prediction of the majority class. The  
 430 benchmark values outlined by Landis & Koch (1977) are useful here for determining the relative strength of the predictive  
 431 models: <0.00 = Poor; 0.00 – 0.20= Slight; 0.21 – 0.40 = Fair; 0.41 – 0.60 = Moderate; 0.61 – 0.81 = Substantial; 0.81 – 1.0 =  
 432 Almost Perfect.  
 433

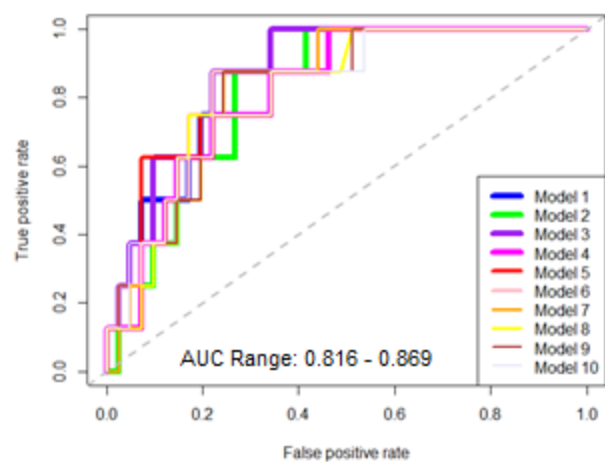
434 While there was some variation in model performance across the ten iterations of RF models for each  
 435 analysis year, overall the models were relatively stable. For the ten iterations of full analyses conducted for each  
 436 collection year, the AUC ranged from 0.724 – 0.75 in 2007, 0.816 – 0.819 in 2010, and 0.982 – 1.0 in 2016.  
 437 Figure 5 illustrates the ROC curve and corresponding best and worst AUC for each of the ten RF models of the  
 438 full predictor list analyses.

439  
 440 **Figure 5. Receiver operator curves for each of the ten full RF model iterations conducted for 2007, 2010**  
 441 (A) The ten RF models ROC curves for 2007 are shown in the top panel. The AUC in the 2007 full models ranged  
 442 from 0.724 – 0.75. (B) The ten RF models ROC curves for 2010 are shown in the middle panel. The AUC in the  
 443 2010 full models ranged from 0.816 – 0.869. (C) The ten RF models ROC curves for 2016 are shown in the  
 444 bottom panel. The AUC in the 2016 full models ranged from 0.982 – 1.0.

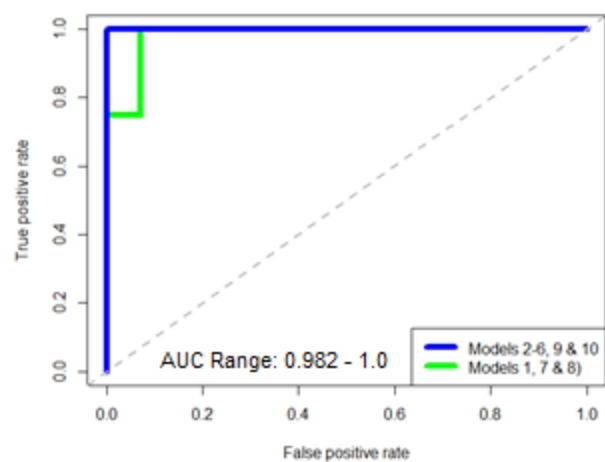
### A) 2007 Full Models



### B) 2010 Full Models



### C) 2016 Full Models



445  
446



## 447 Discussion

448 Of the five categories that were assessed as potential predictors of bovine infection in this study  
449 (physical/biological characteristics, human infection-related, socio-economic, potential animal reservoirs and  
450 agricultural factors), agricultural factors were important predictors of bovine *S. japonicum* infection in all collection  
451 years. Night soil use on summer crops, the village-level area of rice crops, and both the household and village-  
452 level areas of summer and winter crops were each ranked among the top five predictors for one or more  
453 collection years in our RF models. Interestingly, for 2007, all of the ranked agricultural variables except one were  
454 associated with an increase in bovine infection risk in our logistic regression assessments, whereas in 2010 and  
455 2016, these agricultural factors were found to be variably positively and negatively associated with infection risk.  
456 This finding may be related to changing norms and interventions that have taken hold in recent years as a result  
457 of increasing awareness of the potential risks posed by both bovines as a reservoir of schistosomiasis, and  
458 specific agricultural practices. For example, across our study period, we saw a steady increase in the prevalence  
459 of households planting rice (69.1% in 2007; 71.5% in 2010; 81.7% in 2016), and a simultaneous decrease in the  
460 prevalence of households applying any night soil to their rice crops (35.6% in 2007, 14.7% in 2010; 11.7% in  
461 2016). These shifting norms in rice production and night soil use likely resulted in a decrease in the overall  
462 concentration of night soil on rice crops within our study villages, which in turn, may help to explain why the  
463 village-level rice crop area shifts from having a positive association with bovine infection in 2007 and 2010, to a  
464 negative association by 2016.

465 Assessments conducted in China early in the new millennium repeatedly highlighted bovines as a key  
466 source of environmental contamination and as the main animal reservoir of *S. japonicum* in the country (9, 28,  
467 45). Beginning in 2004, a new government-led approach to eliminating schistosomiasis transmission in China was  
468 adopted, which – in conjunction with infrastructure improvements in rural areas and several new schistosomiasis  
469 elimination interventions – featured replacing bovines with machinery in agricultural production (46). Thus, the  
470 negative associations that were found intermittently between bovine infection and some of our agricultural  
471 variables in 2010 and 2016 may be linked to the added precautions that were being adopted when bovines were  
472 being used for agriculture, or because bovines were being reallocated for other purposes (e.g. beef production) as  
473 machinery became the norm for large crop areas or those deemed high risk (e.g. wet rice crops). Increasing

474 recognition of the potential risks posed by night soil use during our study period (32) may have also contributed to  
475 some decreases in environmental contamination as a result of decreases in night soil applications and/or the  
476 more careful treatment of night soil prior to field applications. Indeed, a downward trend in the range of reported  
477 night soil use (total and mean number of buckets) on crops can be observed in Figure 4 (parts D-F), though  
478 notably, we do not see any substantial shift in the overall proportion of households that reported any night soil use  
479 on summer crops over the years (52.4% in 2007; 53.5% in 2010; 50.7% in 2016) (Table 3). The continued  
480 prominence of applying some amount of night soil to summer crops, paired with the steady increase in the total  
481 area of summer crops being planted by villagers over the study period (see Figure 4, part B) may help to explain  
482 why night soil use on summer crops returns to being positively association with bovine infection status in 2016.

483 Bovine ownership in the surrounding village was in the top ten predictors of RF models and bovine  
484 density in a village was positively associated with bovine infection in our regression models in all collection years.  
485 These findings align well with the existing literature that points to bovines as the most important reservoir of *S.*  
486 *japonicum* infection in China (9, 45), and suggests that being in close proximity to higher densities of bovine hosts  
487 may correspond with increasing infection risk, as has been found for other bovine pathogens (47, 48). However, it  
488 is worth noting that household-level bovine ownership was not among the top predictors in any of our RF models,  
489 highlighting that the larger-scale lens (i.e. village-level analysis scale) may be particularly important to future  
490 investigations and control strategies. Likewise, recent informal interviews with locals from our study sites have  
491 revealed that bovines are infrequently kept near the home, as allowing bovines to graze (and defecate) freely is  
492 an economical and efficient way of raising bovines, further illustrating that the household scale may not always be  
493 broad enough to capture larger scale trends. Instead, villagers opt to bring their bovines to the mountains to graze  
494 during the day, which subsequently presents more opportunities for contact between bovines from different  
495 households, and may ultimately result in more widespread environmental contamination (e.g. bovine feces  
496 washed into nearby irrigation ditches after precipitation).

497 In the developmental stages of this analysis, we hypothesized that human infection prevalence and the  
498 number of infected people in the household would be among the top predictors of bovine infection status, given  
499 the known link between human schistosomiasis and bovine reservoirs (e.g. 45). It was therefore somewhat  
500 surprising to find that household-level human infection was only ranked as important from RF models in 2007, and  
501 village human infection prevalence was only ranked as important in 2007 and 2010. One potential explanation for

502 the apparent drop in the importance of human infection status as a predictor of bovine infection could be related  
503 to the aforementioned bovine-removal phenomenon, in which bovines are increasingly being removed from the  
504 village area and brought to alternative mountain locations for grazing, resulting in less frequent contact between  
505 bovines and humans, but more opportunities for contact with other bovines. In fact, the drop in the important  
506 rankings of human infection status in 2016 coincides with a jump in the variable importance rankings for village-  
507 level bovine ownership (6<sup>th</sup> – 8<sup>th</sup> in 2007 and 2010; 1<sup>st</sup> -2<sup>nd</sup> in 2016), providing further support of the theory that  
508 bovines may be becoming increasing important reservoirs of continued schistosomiasis infection. On the other  
509 hand, an altogether different explanation for the differences in the 2016 rankings compared to 2007 and 2010 is  
510 that the 2016 data collection simple didn't have a large enough sample size to allow for the detection of a true  
511 relationship between relatively rare events.

512 As such, one limitation of this assessment was the relatively small sample sizes, particularly in 2016  
513 (N=71), though to a lesser extent, 2010 (N=197) and 2007 (N=473), given the correspondingly large number of  
514 predictors that were included in the full predictor models (N=29, N=31, N=26, in 2007, 2010 and 2016  
515 respectively). While RF models are generally acknowledged as being able to handle assessments of high  
516 dimensional data even with relatively small sample sizes (49), it remains that small samples sizes can still give  
517 rise to the aforementioned issue of non-detection of rare events. Another limitation to this assessment is that RF  
518 models tends to favor continuous predictors over categorical measures, as they allow for a wider range of  
519 potential split points for classifying observations. For this reason, it is not particularly surprising that age was the  
520 only predictor from the individual/physical characteristics predictor group that was ranked among the top ten  
521 predictors, as the remaining individual characteristics were binary measures. Another notable limitation of the  
522 variable importance rankings used in RF models is that they become less reliable when predictors are highly  
523 correlated with one another (50). This may be particularly important to the rankings ascribed to the agricultural  
524 variables, as correlation between the area of the different crop types planted and the amount of night soil used on  
525 each crop tended to be high across all collection years, with the highest predictor correlations found in the 2016  
526 collection year (See S2 – S4 Figures for correlation matrices). This is notable, as a higher degree of instability in  
527 the variable importance rankings was also found for 2016 as compared to 2010 or 2007, suggesting predictor  
528 correlation may be responsible. We therefore recommend that the variable rankings presented from this analyses  
529 be interpreted more holistically (e.g. agricultural variables are strong predictors of bovine infection), and advise

530 caution when comparing unique variable ranking values against one another (e.g. rice crop area is less important  
531 than winter crop area).

532 Our main interests in this assessment were to 1) identify the best predictors of bovine *S. japonicum*  
533 infection within rural farming communities in Sichuan China, and 2) to ascertain whether there are broader trends  
534 in bovine infection distribution across individual, household or village-levels scales or over time. Our RF  
535 assessments have highlighted several key patterns that were repeated across multiple collection years and  
536 multiple iterations of three different models. Agricultural factors and high levels of bovine ownership at the village-  
537 level were repeatedly found to be among the top predictors of bovine *S. japonicum* infection, highlighting the  
538 potential utility of presumptively treating bovines belonging to villages with particularly high levels of bovine  
539 ownership, or those who engage in high-risk agricultural practices such as planting rice. Additionally, village-level  
540 predictors tended to be better predictors of bovine infection than household-level predictors, suggesting that  
541 interventions may need to take a multipronged approach to address broader ecological sources of ongoing  
542 transmission.

543

544

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

546 **Data Availability:** Due to the inclusion of potentially identifying information (e.g. socio-economic  
547 indicators and human infection status), access to study data must be requested through the Carlton Lab group to  
548 ensure the protection of research subjects and compliance with all ethical guidelines. Please contact Elizabeth  
549 Carlton at [elizabeth.carlton@cuanschutz.edu](mailto:elizabeth.carlton@cuanschutz.edu) for more information.

550 **Funding:** This research was supported by grants from the National Institute of Allergy and Infectious  
551 Diseases: R01AI134673 (EJC, PI), R21AI115288 (EJC, PI) and R01AI068854 (Robert Spear, PI). The content is  
552 solely the responsibility of the authors and does not necessarily represent the official views of the National  
553 Institutes of Health. The funders had no role in study design, data collection and analysis, decision to publish, or  
554 preparation of the manuscript.

555 **Author Contributions:** Conceptualization: EC, YL, SP, AB; Data Curation: DL, EG; Formal  
556 Analysis: EG, SP, KK, AB, EC; Funding Acquisition: EC, YL, SP; Investigation: YL, DL, EC; Methodology: EG,  
557 SP, EC, AB, KK; Project Administration: YL, DL; Resources: EC, YL, DL; Validation: EG, SP, EC, KK; Writing –  
558 Original Draft Preparation: EG, SP, EC, KJ; Writing – Review & Editing: EG, SP, KK, YL, DL, EC, AB, KJ.

559 **Acknowledgments:** We are grateful for the support and efforts of the field research team members  
560 from the Institute of Parasitic Diseases and the county anti-schistosomiasis control stations for their efforts in  
561 collecting the data presented here.

## References

562

563

- 564 1. World Health Organization. Schistosomiasis: Key Facts 2020 [Available from: <https://www.who.int/news-room/fact-sheets/detail/schistosomiasis>.  
565
- 566 2. Kittur N, King CH, Campbell CH, Kinung'hi S, Mwinzi PNM, Karanja DMS, et al. Persistent Hotspots in  
567 Schistosomiasis Consortium for Operational Research and Evaluation Studies for Gaining and Sustaining  
568 Control of Schistosomiasis after Four Years of Mass Drug Administration of Praziquantel. *The American*  
569 *journal of tropical medicine and hygiene*. 2019;101(3):617-27.
- 570 3. Song LG, Wu XY, Sacko M, Wu ZD. History of schistosomiasis epidemiology, current status, and  
571 challenges in China: on the road to schistosomiasis elimination. *Parasitology research*.  
572 2016;115(11):4071-81.
- 573 4. Xu J, Steinman P, Maybe D, Zhou XN, Lv S, Li SZ, et al. Evolution of the National Schistosomiasis  
574 Control Programmes in The People's Republic of China. *Adv Parasitol*. 2016;92:1-38.
- 575 5. Zhang LJ, Xu ZM, Guo JY, Dai SM, Dang H, Lü S, et al. [Endemic status of schistosomiasis in People's  
576 Republic of China in 2018]. *Zhongguo xue xi chong bing fang zhi za zhi = Chinese journal of*  
577 *schistosomiasis control*. 2019;31(6):576-82.
- 578 6. Gray DJ, Williams GM, Li Y, McManus DP. Transmission Dynamics of *Schistosoma japonicum* in the  
579 Lakes and Marshlands of China. *PloS one*. 2009;3(12):e4058.
- 580 7. Li H, Dong GD, Liu JM, Gao JX, Shi YJ, Zhang YG, et al. Elimination of schistosomiasis japonica from  
581 formerly endemic areas in mountainous regions of southern China using a praziquantel regimen.  
582 *Veterinary parasitology*. 2015;208(3-4):254-8.
- 583 8. Van Dorssen CF, Gordon CA, Li Y, Williams GM, Wang Y, Luo Z, et al. Rodents, goats and dogs – their  
584 potential roles in the transmission of schistosomiasis in China. *Parasitology*. 2017;144(12):1633-42.
- 585 9. GUO J, LI Y, GRAY D, NING A, HU G, CHEN H, et al. A DRUG-BASED INTERVENTION STUDY ON  
586 THE IMPORTANCE OF BUFFALOES FOR HUMAN SCHISTOSOMA JAPONICUM INFECTION  
587 AROUND POYANG LAKE, PEOPLE'S REPUBLIC OF CHINA. *The American Journal of Tropical*  
588 *Medicine and Hygiene*. 2006;74(2):335-41.
- 589 10. He YX, Salafsky B, Ramaswamy K. Host--parasite relationships of *Schistosoma japonicum* in mammalian  
590 hosts. *Trends Parasitol*. 2001;17(7):320-4.
- 591 11. Ross AG, Sleight AC, Li Y, Davis GM, Williams GM, Jiang Z, et al. Schistosomiasis in the People's  
592 Republic of China: prospects and challenges for the 21st century. *Clinical microbiology reviews*.  
593 2001;14(2):270-95.
- 594 12. Zhou YB, Liang S, Jiang QW. Factors impacting on progress towards elimination of transmission of  
595 schistosomiasis japonica in China. *Parasites & vectors*. 2012;5:275.
- 596 13. Zheng J, Guo JG, Wang XF, Zhu HQ. Relationship of the livestock trade to schistosomiasis transmission  
597 in mountainous area. *Zhongguo ji sheng chong xue yu ji sheng chong bing za zhi = Chinese journal of*  
598 *parasitology & parasitic diseases*. 2000;18(3):146-8.
- 599 14. Soomro A. A Study on Prevalence and Risk Factors of Brucellosis in Cattle and Buffaloes in District  
600 Hyderabad, Pakistan. *Journal of Animal Health and Production*. 2014;2:33-7.
- 601 15. Mugizi DR, Boqvist S, Nasinyama GW, Waiswa C, Ikwap K, Rock K, et al. Prevalence of and factors  
602 associated with *Brucella* sero-positivity in cattle in urban and peri-urban Gulu and Soroti towns of  
603 Uganda. *J Vet Med Sci*. 2015;77(5):557-64.
- 604 16. Skuce RA, Allen AR, McDowell SWJ. Herd-Level Risk Factors for Bovine Tuberculosis: A Literature  
605 Review. *Veterinary Medicine International*. 2012;2012:621210.
- 606 17. Nzalawahe J, Kassuku AA, Stothard JR, Coles GC, Eisler MC. Trematode infections in cattle in Arumeru  
607 District, Tanzania are associated with irrigation. *Parasites & vectors*. 2014;7:107.

- 608 18. Deka RP, Magnusson U, Grace D, Lindahl J. Bovine brucellosis: prevalence, risk factors, economic cost  
609 and control options with particular reference to India-a review. *Infection Ecology & Epidemiology*.  
610 2018;8(1):1556548.
- 611 19. Tempia S, Salman M, Keefe T, Morley P, Freier J, DeMartini J, et al. A sero-survey of rinderpest in  
612 nomadic pastoral systems in central and southern Somalia from 2002 to 2003, using a spatially integrated  
613 random sampling approach. *Revue scientifique et technique*. 2010;29(3):497.
- 614 20. Chanie M, Dejen B, Fentahun T. Prevalence of cattle schistosomiasis and associated risk factors in  
615 Fogera cattle, south Gondar zone, Amhara national regional state, Ethiopia. *Journal of Advanced  
616 Veterinary Research*. 2012;2:153-6.
- 617 21. Kebede A, Dugassa J, Haile G, Wakjira BM. Prevalence of bovine of schistosomosis in and around  
618 Nekemte, East Wollega zone, Western Ethiopia. *Journal of Veterinary Medicine and Animal Health*.  
619 2018;10:123-7.
- 620 22. Gebremeskel AK, Simeneh ST, Mekuria SA. Prevalence and Associated Risk Factors of Bovine  
621 Schistosomiasis in Northwestern Ethiopia. *World*. 2017;7(1):01-4.
- 622 23. Defersha T, Belete B. The Neglected Infectious Disease, Bovine Schistosomiasis: Prevalence and  
623 Associated Risk Factors for its Occurrence among Cattle in the North Gulf of Lake Tana, Northwest  
624 Ethiopia. *J Vet Med Health*. 2018;2:112.
- 625 24. Yihunie A, Urga B, Alebie G. Prevalence and risk factors of bovine schistosomiasis in Northwestern  
626 Ethiopia. *BMC veterinary research*. 2019;15(1):12.
- 627 25. Tsega M, Derso S. Prevalence of bovine schistosomiasis and its associated risk factor in and around  
628 Debre Tabor town, north west of Ethiopia. *Europ J Biol Sci*. 2015;7:108-13.
- 629 26. Lulie B, Guadu T. Bovine schistosomiasis: A threat in public health perspective in Bahir Dar town,  
630 northwest Ethiopia. *Acta Parasitologica Globalis*. 2014;5(1):1-6.
- 631 27. Tan TK, Low VL, Lee SC, Panchadcharam C, Kho KL, Koh FX, et al. Detection of *Schistosoma spindale*  
632 ova and associated risk factors among Malaysian cattle through coprological survey. *Japanese Journal of  
633 Veterinary Research*. 2015;63(2):63-71.
- 634 28. Gray DJ, Williams GM, Li Y, Chen H, Li RS, Forsyth SJ, et al. A cluster-randomized bovine intervention  
635 trial against *Schistosoma japonicum* in the People's Republic of China: design and baseline results. *The  
636 American journal of tropical medicine and hygiene*. 2007;77(5):866-74.
- 637 29. Li YS, McManus DP, Lin DD, Williams GM, Harn DA, Ross AG, et al. The *Schistosoma japonicum* self-  
638 cure phenomenon in water buffaloes: potential impact on the control and elimination of schistosomiasis in  
639 China. *International journal for parasitology*. 2014;44(3-4):167-71.
- 640 30. Xu S, Shi F, Shen W, Lin J, Wang Y, Lin B, et al. Vaccination of bovines against schistosomiasis japonica  
641 with cryopreserved-irradiated and freeze-thaw schistosomula. *Veterinary parasitology*. 1993;47(1-2):37-  
642 50.
- 643 31. He Y, Xu S, Shi F, Shen W, Hsü S, Hsü H. Comparative studies on the infection and maturation of  
644 schistosoma japonicum in cattle and buffaloes. *Current Zoology*. 1992;38(3):266-71.
- 645 32. Carlton EJ, Bates MN, Zhong B, Seto EYW, Spear RC. Evaluation of Mammalian and Intermediate Host  
646 Surveillance Methods for Detecting Schistosomiasis Reemergence in Southwest China. *PLoS Negl Trop  
647 Dis*. 2011;5(3).
- 648 33. Liang S, Yang C, Zhong B, Qiu D. Re-emerging schistosomiasis in hilly and mountainous areas of  
649 Sichuan, China2006. 139-44 p.
- 650 34. Carlton EJ, Liu Y, Zhong B, Hubbard A, Spear RC. Associations between Schistosomiasis and the Use of  
651 Human Waste as an Agricultural Fertilizer in China. *PLoS Negl Trop Dis*. 2015;9(1).
- 652 35. Control DoD. Textbook for Schistosomiasis Control. Shanghai Shanghai Publishing House for Science  
653 and Technology. 2000:72-6.



- 654 36. Katz N, Chaves A, Pellegrino J. A simple device for quantitative stool thick-smear technique in  
655 Schistosomiasis mansoni. *Rev Inst Med Trop Sao Paulo*. 1972;14(6):397-400.
- 656 37. Gordon CA, Kurscheid J, Williams GM, Clements ACA, Li Y, Zhou XN, et al. Asian Schistosomiasis:  
657 Current Status and Prospects for Control Leading to Elimination. *Trop Med Infect Dis*. 2019;4(1).
- 658 38. Breiman LC, A.; Liaw, A.; Wiener, M. . Breiman and Cutler's Random Forests for Classification and  
659 Regression. <https://cran.r-project.org/web/packages/randomForest/randomForest.pdf>; 2018.
- 660 39. RColorBrewer S, Liaw MA. Package 'randomForest'.
- 661 40. Environmental Systems Research Institute (ESRI). ArcGIS Desktop Release. 10.5.1 ed. Redlands, CA.  
662 2017
- 663 41. Breiman L, Cutler A. Manual—setting up, using, and understanding random forests V4. 0. 2003. URL  
664 [https://www.stat.berkeley.edu/~breiman/Using\\_random\\_forests\\_v4\\_0.pdf](https://www.stat.berkeley.edu/~breiman/Using_random_forests_v4_0.pdf). 2011.
- 665 42. StataCorp. Stata Statistical Software: Release 15. College Station, TX: StataCorp LP; 2015.
- 666 43. RStudio Team. RStudio: Integrated Development for R. RStudio, PBC, Boston, MA, 2020. 2020.
- 667 44. Landis JR, Koch GG. The Measurement of Observer Agreement for Categorical Data. *Biometrics*.  
668 1977;33(1):159-74.
- 669 45. Gray DJ, Williams GM, Li Y, Chen H, Forsyth SJ, Li RS, et al. A cluster-randomised intervention trial  
670 against *Schistosoma japonicum* in the Peoples' Republic of China: bovine and human transmission. *PloS*  
671 *one*. 2009;4(6):e5900.
- 672 46. Liu Y, Zhong B, Wu Z-S, Liang S, Qiu D-C, Ma X. Interruption of schistosomiasis transmission in  
673 mountainous and hilly regions with an integrated strategy: a longitudinal case study in Sichuan, China.  
674 *Infectious Diseases of Poverty*. 2017;6(1):79.
- 675 47. Spencer SE, Besser TE, Cobbold RN, French NP. 'Super' or just 'above average'? Supershedders and  
676 the transmission of *Escherichia coli* O157:H7 among feedlot cattle. *J R Soc Interface*. 2015;12(110):0446.
- 677 48. Meadows AJ, Mundt CC, Keeling MJ, Tildesley MJ. Disentangling the influence of livestock vs. farm  
678 density on livestock disease epidemics. *Ecosphere*. 2018;9(7):e02294.
- 679 49. Biau G, Scornet E. A random forest guided tour. *TEST*. 2016;25(2):197-227.
- 680 50. Strobl C, Boulesteix A-L, Kneib T, Augustin T, Zeileis A. Conditional variable importance for random  
681 forests. *BMC Bioinformatics*. 2008;9(1):307.
- 682



676 **Supporting Information**

677  
 678 **S1 Table. Completeness of bovine infection and household surveys.** Some differences between the number  
 679 of bovines reported by households and the number of bovines tested may have arisen due to the lag time  
 680 between the household surveys, which were completed in June/July in 2007 and 2010 and the infection surveys,  
 681 which were conducted in November and December in 2007 and 2010. Both the household surveys and infection  
 682 surveys were conducted during June/July of 2016.

	2007			2010			2016		
	N	Total possible	%	N	Total possible	%	N	Total possible	%
Number of villages where 1+ household reported owning bovines	36	36	100%	35	36	97.2%	8	10	80%
Number of villages where 1+ bovine tested	35	36	97.2%	31	35	88.6%	8	8	100%
Number of villages where:									
≥80% of bovines tested	18	35	51.4%	13	31	41.9%	4	8	50%
40 – 80 % tested	14	35	40.0%	14	31	45.2%	4	8	50%
<40% tested	3	35	8.6%	4	31	12.9%	0	8	0%
Total number of bovines owned, as reported on household surveys	675	--	--	371	--	--	95	--	--
Total number of bovines tested	503	675	74.5%	233	371	62.8%	72	95	58.6%
Hatch tests									
1	68	503	13.5%	34	233	14.6%	10	72	8.6%
2	91	503	18.1%	49	233	21.0%	9	72	20.7%
3	344	503	68.4%	150	233	64.4%	53	72	73.6%
Total number of bovines with household survey data	473	503	94.0%	197	233	84.5%	71	72	98.6%

683

684

685

686

687

688

689

690 **S2 Table. Simple logistic regression analyses to determine the direction of association between bovine**  
 691 **infection status and each predictor, by collection year.** Tertiles (and sometimes quartiles) by year were used  
 692 in simple logistic regression analyses to help investigate potential non-linearity. Results highlighted in gray  
 693 indicate that the predictor was one of the top ten predictors in one or more RF analyses for a given collection  
 694 year.

	2007				2010				2016			
	Tertile	Point estimate	SE	P-value	Tertile	Point estimate	SE	P-value	Tertile	Point estimate	SE	P-value
<b>Individual characteristics</b>												
Bovine age	≤ 3	ref			≤ 3	Ref						
	3.1 - 5	0.58	0.34	0.094*	3.1 - 6	0.10	0.46	0.824				
	≥ 5.1	0.92	0.36	0.011*	≥ 6.1	0.01	0.47	0.988				
# of hatch tests	1	Ref			1	p.f.p.	p.f.p.	p.f.p.*	1	Ref		
	2	0.73	0.70	0.297	2	Ref			2	0.94	1.32	0.476
	3	1.24	0.61	0.043	3	0.50	0.49	0.308	3	0.88	1.11	0.426
<b>Human sources</b>												
Human infection: household count	0	Ref			0	Ref			0	Ref		
	1+	0.93	0.31	0.003	1+	1.75	0.45	<0.001*	1+	0.02	0.85	0.981
Human infection: village prevalence	≤ 2.6%	Ref			0%	p.f.p.	p.f.p.	p.f.p.*	≤ 2.53%	Ref		
	2.7 - 11.9%	-0.24	0.40	0.558	0.1 - 10.8%	Ref			2.54 - 4.49%	0.99	0.71	0.164
	≥ 12%	1.06	0.33	0.001*	≥ 10.9%	1.11	0.43	0.010*	≥ 4.5%	-0.10	0.82	0.907
<b>Socio-economic</b>												
Vil: Imp. sanitation	≤ 7.2%	Ref			≤ 25%	Ref			≤ 48%	Ref		
	7.3 - 23.9%	1.29	0.35	<0.001	26 - 42%	0.87	0.43	0.046*	49 - 53%	p.f.p.	p.f.p.	p.f.p.
	≥ 24%	0.36	0.41	0.379	≥ 43%	-0.99	0.62	0.112*	≥ 54%	0.61	0.65	0.346
Vil. asset score (0-9)	≤ 1.654	Ref			≤ 2.77	Ref			≤ 4.03	Ref		
	1.655 - 2.5	-1.09	0.35	0.002*	2.78 - 3.79	0.38	0.44	0.384	4.04 - 4.59	1.40	0.86	0.104*
	≥ 2.51	-1.15	0.34	<0.001*	≥ 3.8	-0.30	0.50	0.549	≥ 4.6	1.23	0.90	0.170*
HH asset score (0-9)	≤ 1	Ref			≤ 2	Ref			≤ 3	Ref		
	2 - 3	0.24	0.30	0.427	3	0.32	0.51	0.531	4	-1.06	0.89	0.234
	≥ 4	0.08	0.41	0.845	4	0.04	0.59	0.951	≥ 5	0.51	0.68	0.453
					≥ 5	0.41	0.53	0.440				
<b>Animal reservoir</b>												
Vil.: cat ownership	≤ 47.9%	Ref			≤ 44.9%	Ref			≤ 43.9%	Ref		
	48 - 64.9%	-0.61	0.37	0.095	45 - 55.9%	0.94	0.47	0.044	44 - 58.9%	2.81	1.11	0.012*
	≥ 65%	0.31	0.31	0.324	≥ 56%	0.13	0.53	0.813	≥ 59%	2.07	1.14	0.069*
Vil.: dog ownership	≤ 75%	Ref			≤ 71.9%	Ref			≤ 64%	Ref		
	76 - 85.2%	-0.43	0.31	0.166	72 - 80%	2.40	0.77	0.002	64.1 - 74.9%	0.24	0.82	0.769
	≥ 85.2%	-0.95	0.36	0.009	≥ 81%	2.38	0.77	0.002	≥ 75%	1.25	0.77	0.104*
HH: pigs owned					0	Ref						
					1	-0.42	0.67	0.534				
					>1	-1.50	0.56	0.008*				
Vil. mean pigs owned					≤ 0.49	Ref			≤ 0.14	Ref		
					.5 - 0.84	0.34	0.43	0.436	0.21 - 0.43	-2.05	1.09	0.060*
					≥ 0.85	-0.38	0.52	0.470	≥ 0.44	-1.30	0.83	0.119*
Vil.: mean bovines owned	≤ 0.65	Ref			≤ 0.55	Ref			≤ 0.30	Ref		
	0.66 - 0.89	1.19	0.38	0.002*	0.56 - 0.74	-0.10	0.51	0.839	0.30 - 0.44	0.63	0.96	0.513
	≥ 0.89	1.16	0.40	0.003*	≥ 0.75	0.77	0.45	0.086*	≥ 0.441	2.24	0.86	0.009*

<b>Agriculture</b>												
Vil. mean rice area	≤ 0.84	Ref			≤ 0.61	Ref			≤ 0.95	Ref		
	0.84 – 1.47	-0.33	0.37	0.380	0.62 – 1.29	2.33	0.77	0.003*	0.96 – 1.34	-0.97	0.70	0.166*
	≥ 1.47	0.55	0.32	0.081*	≥ 1.30	2.35	0.77	0.002*	≥ 1.35	-1.61	0.86	0.061*
HH rice area	≤ 0.5	Ref			≤ 0.5	Ref			≤ 1.2	Ref		
	0.6 – 1.5	0.36	0.35	0.301	0.6 – 1.5	0.57	0.49	0.253	1.2 – 2.3	-0.49	0.67	0.463
	≥ 1.6	0.54	0.34	0.114	≥ 1.6	0.76	0.48	0.116*	≥ 2.4	-1.25	0.87	0.149
Vil sum. crop area	≤ 0.73	Ref			≤ 2.03	Ref			≤ 2.412	Ref		
	0.74 – 1.96	1.51	0.39	<0.001*	2.03– 2.76	0.02	0.44	0.966	2.413 – 2.98	0.88	0.92	0.336
	≥ 1.97	0.90	0.42	0.031*	≥ 2.77	-0.52	0.49	0.285	≥ 2.99	1.81	0.86	0.034*
HH sum. crop area	≤ 1	Ref			≤ 2.2	Ref			≤ 2.8	Ref		
	1.1 – 2.2	0.69	0.34	0.041*	2.3 – 3.8	-0.09	0.43	0.830	2.81 – 4.1	1.70	0.86	0.046*
	≥ 2.3	0.40	0.33	0.230	≥ 3.9	-0.69	0.50	0.167*	≥ 4.2	0.84	0.92	0.362
Vil. win. crop area*	≤ 1.869	Ref			≤ 1.869	Ref			≤ 1.58	Ref		
	1.87 – 2.55	0.01	0.40	0.981	1.87 – 2.58	-0.50	0.45	0.269	1.59 – 1.89	-0.86	0.70	0.219
	≥ 2.82	0.48	0.37	0.193*	≥ 2.59	-0.60	0.46	0.193*	≥ 1.9	-1.55	0.86	0.071*
HH win. crop area*	1.79	Ref			≤ 1.9	Ref			≤ 1.2	Ref		
	1.8 – 2.5	-0.03	0.41	0.944	2 – 2.8	-0.07	0.53	0.895	1.3 – 2.1	-0.62	0.81	0.442
	2.6 – 3.7	0.54	0.38	0.152*	2.9 – 3.7	0.36	0.51	0.617	2.2 – 3.2	-0.76	0.80	0.340
Vil. night soil sum.	≥ 3.8	0.32	0.39	0.413	≥ 3.8	-0.23	0.54	0.678	≥ 3.3	-1.17	0.90	0.193*
	≤ 7.3	Ref			≤ 6.39	Ref			≤ 11.2	Ref		
	7.4 – 37	1.18	0.38	0.002*	6.4 – 14.4	-0.12	0.41	0.766	11.3 – 27.9	2.48	1.11	0.025*
Vil. night soil winter	≥ 38	1.06	0.38	0.005*	≥ 14.5	-1.80	0.66	0.006*	≥ 28	1.95	1.14	0.087*
	≤ 23.6	Ref			≤ 2.589	Ref			≤ 2.9	Ref		
	23.7 – 42.6	-0.69	0.44	0.120*	2.59 – 11.2	-0.82	0.44	0.066	3 – 4.2	0.41	0.87	0.642
HH night soil winter*	≥ 42.7	1.04	0.32	0.001*	≥ 11.3	-1.39	0.54	0.010	≥ 4.3	1.79	0.76	0.019
	0	0.90	0.56	0.105*	0	Ref			0	Ref		
	1-20	Ref			>0	-1.71	0.63	0.006	>0	2.00	0.76	0.009
	21 – 59	1.11	0.59	0.060*								
	≥ 60	0.96	0.58	0.097*								

695

696 \*In the case that the tertiles did not show evidence of any moderate difference ( $p < 0.2$ ) between one or more groups, quartiles

697 were tried. When still no difference was found between groups, this was noted in results Table 4.

698 p.f.p = perfect failure predicted.

699

700

701

702

703

704

705 **S1 Figure. Supplemental analysis assessing changes over time.** Two additional RF model iterations were run  
 706 for each collection year that only included those predictors that were available in all three of the collection years.  
 707 The top ten predictors for these two iterations were given a score of 1-10, and the summed scores were used to  
 708 determine the variable ranking 1<sup>st</sup> – 10<sup>th</sup> for each collection year, as well as a final variable ranking “all year score”  
 709 that summed the rankings across all six iterations (two per collection year) conducted.

710

	2007	2010	2016	All year score																		
<b>Physical characteristics</b>																						
Number of hatch tests		2		9																		
County of residence																						
<b>Infection</b>																						
Human infection: household count		9																				
Human infection: village prevalence	5	1		4*																		
<b>Socio-economic indicators</b>																						
Village prevalence of improved sanitation			10																			
Household has improved sanitation (y/n)																						
Village mean asset score (0-9)	10		4	10																		
Household asset score (0-9)																						
<b>Animal ownership</b>																						
Village prevalence of cat ownership			9																			
Household cat ownership																						
Village prevalence of dog ownership			6																			
Household bovine ownership																						
Village prevalence of bovine ownership	8		1*	8																		
<b>Agriculture</b>																						
Village mean rice area	6	5	7*	6																		
Household rice area		6																				
Village mean dry summer crop area	2	10	7*	7																		
Household dry summer crop area	4	4		4*																		
Village mean winter crop area	1	8	3	2																		
Household winter crop area	3	7	5	3																		
Village night soil rice: mean # buckets																						
Household night soil rice: # buckets																						
Village night soil summer dry crop: mean # buckets	7	3	1*	1																		
Household night soil summer dry crop: # buckets																						
Village night soil winter crop: mean # buckets																						
Household night soil winter crop: # buckets	9																					
<p>* = Tied for importance rank                      = Not ranked 1-10 for the given model/year</p> <p><b>Color key: Variable importance rankings (1<sup>st</sup>-10<sup>th</sup>) by scale</b></p> <table border="1"> <thead> <tr> <th>Individual</th> <th>1<sup>st</sup> – 2<sup>nd</sup></th> <th>3<sup>rd</sup> – 4<sup>th</sup></th> <th>5<sup>th</sup> – 6<sup>th</sup></th> <th>7<sup>th</sup> – 8<sup>th</sup></th> <th>9<sup>th</sup> – 10<sup>th</sup></th> </tr> </thead> <tbody> <tr> <td>Household</td> <td>1<sup>st</sup> – 2<sup>nd</sup></td> <td>3<sup>rd</sup> – 4<sup>th</sup></td> <td>5<sup>th</sup> – 6<sup>th</sup></td> <td>7<sup>th</sup> – 8<sup>th</sup></td> <td>9<sup>th</sup> – 10<sup>th</sup></td> </tr> <tr> <td>Village</td> <td>1<sup>st</sup> – 2<sup>nd</sup></td> <td>3<sup>rd</sup> – 4<sup>th</sup></td> <td>5<sup>th</sup> – 6<sup>th</sup></td> <td>7<sup>th</sup> – 8<sup>th</sup></td> <td>9<sup>th</sup> – 10<sup>th</sup></td> </tr> </tbody> </table>					Individual	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>	Household	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>	Village	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>
Individual	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>																	
Household	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>																	
Village	1 <sup>st</sup> – 2 <sup>nd</sup>	3 <sup>rd</sup> – 4 <sup>th</sup>	5 <sup>th</sup> – 6 <sup>th</sup>	7 <sup>th</sup> – 8 <sup>th</sup>	9 <sup>th</sup> – 10 <sup>th</sup>																	

711

712 **S2 Figure. Correlation matrix for 2007 predictors.** A correlation matrix for predictors included in the 2007 RF models is provided to highlight those  
 713 predictors whose relative variable ranking positions may be less reliable due to correlation with other influential predictors. Only predictors with a  
 714 correlation coefficient of < -0.499 or > 0.499 are included. The 2007 correlation matrix demonstrates that there are some strongly correlated predictors,  
 715 particularly in the agricultural predictor category, that may be impacting their relative importance rankings.

	Count y	Hatch tests	HH asset	HH rice crop area	HH sum. crop area	HH win. crop area	HH NS sum. crop	Vil. asset	Vil cat own	Vil. bov. own	Vil hum. Inf. prev.	Vil rice area	Vil. sum. crop area	Vil. win. crop area	Vil. NS sum. crop	Vil. NS win. crop	Vil NS Rice crop
County	1.0																
Hatch tests	0.553	1.0															
HH asset			1.0														
HH rice area				1.0													
HH sum. crop area	0.578				1.0												
HH win. crop area					0.559	1.0											
HH NS sum. Crop					0.594		1.0										
Vil Asset	-0.728		0.576		-0.501			1.0									
Vil Cat Own								-0.565	1.0								
Vil Bov. Own								-0.520		1.0							
Vil Hum. Inf. Prev.											1.0						
Vil rice area	-0.664			0.691	-0.575			0.544				1.0					
Vil sum. crop area	0.799				-0.699			-0.632				-0.696	1.0				
Vil win. crop area													0.648	1.0			
Vil NS sum. crop	0.674				0.644			-0.626				-0.599	0.814	0.712	1.0		
Vil NS win. Crop												0.547				1.0	
Vil NS rice crop											0.598						1.0

**Color Key:**

Moderate positive correlation	0.50 – 0.599	Moderate negative correlation	-0.599 – -0.50
Strong positive correlation	0.60 – 0.799	Strong negative correlation	-0.799 – -0.60
Very strong positive correlation	≥ 0.80	Very strong negative correlation	≤ -0.80

717 **S3 Figure. Correlation matrix for 2010 predictors.** A correlation matrix for predictors included in the 2010 RF models is provided to highlight those  
 718 predictors whose relative variable ranking positions may be less reliable due to correlation with other influential predictors. Only predictors with a  
 719 correlation coefficient of < -0.499 or > 0.499 are included. The 2010 correlation matrix demonstrates that there are just a few strongly correlated predictors  
 720 in the agricultural predictor category. As well as the socio-economic indicator category that may be impacting relative importance rankings.  
 721

	County	HH rice crop area	HH sum. crop area	HH win. crop area	HH NS sum. crop	HH NS win. crop	Vil. asset	Vil cat own	Vil. bov. own	Vil pigs	Vil rice area	Vil. sum. crop area	Vil. win. crop area	Vil. NS sum. crop	Vil. NS win. crop
County	1.0														
HH rice area		1.0													
HH sum. crop area			1.0												
HH win. crop area			0.575	1.0											
HH NS sum. Crop					1.0										
HH NS win. Crop					0.565	1.0									
Vil Asset	-0.771						1.0								
Vil Cat Own	0.589						-0.533	1.0							
Vil Bov. Own	0.556						-0.525		1.0						
Vil pigs										1.0					
Vil rice area		0.553									1.0				
Vil sum. crop area	0.576						-0.569	0.535				1.0			
Vil win. crop area								0.789		0.515		0.675	1.0		
Vil NS sum. crop											0.578			1.0	
Vil NS win. crop						0.509								0.878	1.0

**Color Key:**

Moderate positive correlation	0.50 – 0.599	Moderate negative correlation	-0.599 – -0.50	742
Strong positive correlation	0.60 – 0.799	Strong negative correlation	-0.799 – -0.60	743
Very strong positive correlation	≥ 0.80	Very strong negative correlation	≤ -0.80	744

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

746 **S4 Figure. Correlation matrix for 2016 predictors.** A correlation matrix for predictors included in the 2016 RF models is provided to highlight those  
 747 predictors whose relative variable ranking positions may be less reliable due to correlation with other influential predictors. Only predictors with a  
 748 correlation coefficient of  $< -0.499$  or  $> 0.499$  are included. The 2016 correlation matrix demonstrates that there are several strongly correlated  
 749 predictors across the different predictor categories that may be impacting relative importance rankings for the 2016 RF models.

	Count y	Hatch tests	HH hum. inf. prev	HH rice crop area	HH win. crop area	HH NS sum. crop	HH NS rice crop	Vil hum inf prev	Vil. asset	Vil toilet imp	Vil. dog own	Vil cat own	Vil pigs own	Vil rice crop area	Vil. sum. crop area	Vil. win. crop area	Vil. NS sum. crop	Vil. NS win. crop	Vil. NS rice crop
County	1.0																		
Hatch tests	0.615	1.0																	
HH hum. Inf.	-0.625		1.0																
HH rice area				1.0															
HH win. crop area				0.806	1.0														
HH NS sum. crop						1.0													
HH NS rice crop						0.834	1.0												
Vil. hum. inf prev.	-0.839	-0.641	0.661					1.0											
Vil asset									1.0										
Vil toilet Imp.										1.0									
Vil. dog own											1.0								
Vil cat Own											0.800	1.0							
Vil pigs	-0.502		0.611					0.838	-0.563				1.0						
Vil rice area								0.669	-0.635			0.522	0.858	1.0					
Vil sum. crop area								0.501	-0.612			0.721	0.801	0.678	1.0				
Vil win. crop area								0.590	-0.653			0.640	0.820	0.955	0.774	1.0			
Vil NS sum. crop										-0.645							1.0		
Vil NS win. crop	0.649							-0.546		-0.814							0.578	1.0	
Vil. NS rice crop																	0.581		1.0

**Color Key:**

Moderate positive correlation	0.50 – 0.599	Moderate negative correlation	-0.599 – -0.50
Strong positive correlation	0.60 – 0.799	Strong negative correlation	-0.799 – -0.60
Very strong positive correlation	$\geq 0.80$	Very strong negative correlation	$\leq -0.80$

750