

1 Seasonal Dynamics of *Anopheles stephensi* and its 2 Implications for Mosquito Detection and Emergent Malaria 3 Control in the Horn of Africa

4 Charles Whittaker^{1,†}, Arran Hamlet¹, Ellie Sherrard-Smith¹, Peter Winskill¹, Gina Cuomo-Dannenburg¹,
5 Patrick G.T. Walker¹, Marianne Sinka², Samuel Pironon^{3,4}, Ashwani Kumar⁵, Azra Ghani¹, Samir
6 Bhatt^{1,6†}, Thomas S. Churcher^{1†}

7 ¹MRC Centre for Global Infectious Disease Analysis & Abdul Latif Jameel Institute for Disease and Emergency Analytics,
8 School of Public Health, Imperial College London, London, UK

9 ²Department of Biology, University of Oxford, Oxford, UK

10 ³Royal Botanic Gardens Kew, Richmond, Surrey, UK

11 ⁴United Nations Environment Program World Conservation Monitoring Centre, Cambridge, UK

12 ⁵Vector Control Research Centre, Indira Nagar, Puducherry, India

13 ⁶Section of Epidemiology, Department of Public Health, University of Copenhagen, Copenhagen, Denmark

14 †These authors jointly supervised this work.

15 *Corresponding Author: Charles Whittaker, Department of Infectious Disease Epidemiology, School of Public Health, Imperial
16 College, London W2 1PG, United Kingdom. Email: charles.whittaker16@imperial.ac.uk

17 **Keywords:** *Anopheles stephensi*; malaria ecology; urban malaria; population dynamics;
18 seasonality; epidemiology.

19 Contents

20
21 **Supplementary Information 1:** Description of Systematic Review, Data
22 Extraction and Initial Processing

23
24 **Supplementary Information 2:** Description of Statistical Methodologies
25 Utilised

26
27 **Supplementary Information 3:** Additional Figures and Results

28
29 **References**

30

31 **Outline of Document**

32 In this supplementary document we outline the methods and data used to explore and analyse
33 the patterns and drivers of *Anopheles stephensi* population dynamics across South Asia and
34 the Middle East. In **Supplementary Information 1**, we present an overview of the systematic
35 search strategy used to collate the references containing the extracted and analysed data, as
36 well as details about the initial pre-processing steps applied to said data. In **Supplementary**
37 **Information 2**, we describe the statistical methodologies used to process and analyse this
38 extracted data. The output of these analyses forms the basis for the results presented in the
39 main text. Finally, in **Supplementary Information 3**, we present a number of additional figures
40 and tables to support the work detailed in the main text.

41

42 **Supplementary Information 1: Description of Systematic Review:** 43 **Data Extraction and Initial Pre-Processing**

44 **Systematic Review: Search Procedure and Record Screening**

45 We collated references from two previously published systematic reviews of literature relating
46 to *Anopheles stephensi* (focusing on its presence/absence across a wide geographical range¹
47 and its seasonal dynamics in India² respectively), and updated these previous searches (both
48 conducted in 2017) by searching *Web of Science* and *PubMed* databases from January 2017
49 for further relevant references containing temporally disaggregated *Anopheles stephensi*
50 catch data. Key words for this search were:

51 ((anophel* AND ((India) OR (BURMA) OR (MYANMAR) OR (BANGLADESH) OR
52 (THAILAND) OR (ISLAMIC REPUBLIC OF IRAN) OR (ETHIOPIA) OR (DJIBOUTI)
53 OR (SUDAN))) AND (("2017"[Date - Publication] : "3000"[Date - Publication])) OR
54 ((anophel* AND ((Pakistan) OR (Iran) OR (Afghanistan)) AND (("1990"[Date -
55 Publication] : "3000"[Date - Publication]))

56 with references for Pakistan, Iran and Afghanistan searched for over an extended time-period
57 (i.e. date range of 1990-2020 rather than 2017-2020) to ensure completeness of the collated
58 references, and fill in countries not included during previous reviews. Our searches identified
59 a total of 926 records, which were screened according to the following Inclusion/Exclusion
60 criteria:

61 **Inclusion Criteria:**

- 62 • Reference contains temporally disaggregated adult mosquito catch data for
63 *An. stephensi*, at a temporal resolution of monthly or finer.
- 64 • The time-period spanned by the survey must be at least 10 consecutive months in
65 duration and have caught at least a total of 25 *An. stephensi* over the period for which
66 catches were being carried out.

67 **Exclusion Criteria:**

- 68 • Mosquito catch data is not temporally disaggregated to a sufficient extent (e.g. catches
69 were done yearly or seasonally rather than monthly).
- 70 • Mosquito catch data was collected as part of a trial assessing a vector control
71 intervention (which would perturb the natural dynamics of the vector, rendering the
72 data unrepresentative of the population dynamics in the absence of control).
- 73 • The reference only contained information on immature/larval mosquito life cycle stages
74 rather than mature adults.
- 75 • The reference contained insufficient information to geolocate the area in which the
76 study was conducted to at least the administrative unit 2 level.

77 Overall, a total of 34 references were collated containing 62 time-series from catch surveys
78 carried out in distinct locations from across Afghanistan (n=2), Djibouti (n=1), India (n=30),
79 Iran (n=17), Myanmar (n=5) and Pakistan (n=7) from the systematic review. These were
80 further supplemented with 2 references (from Pakistan and India respectively, yielding a total
81 of 3 time-series) collated as part of a review of the bionomics of *An. stephensi* previously
82 carried out³, yielding a total of 65 time-series from these 36 references. The next section
83 describes in further detail about extraction and collation of the data associated with each study.

84 **Systematic Review: Data Extraction, Collation and Initial Processing**

85 Entomological Data Extraction

86 For each reference, we extracted all relevant entomological catch data provided that pertained
 87 specifically to *An. stephensi*. Where data were presented in a table, data was copied directly
 88 from the table. Where the data were in a graph, data were extracted using the DataThiefTM
 89 software. This yielded a total of 65 time series of monthly mosquito catch data (no reference
 90 presented data at a finer temporal resolution), ranging in length from 10 – 60 months, with a
 91 mean time-period of 15.6 months and a median time-period of 12 months, a mean catch size
 92 of 758 and a median catch size of 289.

93 **Supplementary Table 1: Number of time series collated according to method of** 94 **mosquito collection.**

	Landing Catch	Resting Collections	Pit Collections	Light Traps	Pyrethrum Spray Catch
# Time-Series	3	33	2	4	14

95

96 Of the collated studies, the majority sampled mosquitoes via resting collections (n=33), though
 97 there was significant variation between surveys as to where mosquitoes had been sampled
 98 (e.g. human dwellings or cattlesheds), when sampling had been carried out (daytime, night-
 99 time or overnight) and for the small number of landing catch studies collated (n=3), which bait
 100 had been used (cattle or humans). Of the 65 collated time-series, 56 presented results arising
 101 from a survey carried out using 1 catch methodology (described in **Supplementary Table 1**
 102 above). 9 time-series represented results which presented the total number of *Anopheles*
 103 *stephensi* mosquitoes caught across all methods of collection and could not be disaggregated
 104 by catch-type. They have not been counted in **Supplementary Table 1** above.

105 The primary focus of these analyses was to characterise annual and seasonal patterns of
 106 variation in *An. stephensi* abundance. Given this, and also that variations in time-series length
 107 are a factor known to affect their statistical properties⁴ (and therefore limit the comparability of
 108 the time series gathered and analysed here), all time-series were standardised to be 12
 109 months in length. For time series containing more than 12 time points (i.e. time series that
 110 spanned longer than a single year), we averaged the recorded catches for a given month. For
 111 time-series containing less than 12 months of data, this was not carried out. Where the study
 112 has been initiated in a month other than January, and concluded in a month other than
 113 December, the recorded counts were reordered to yield a complete time series running from
 114 January to December (and then subsequently adjusted so that the month of peak vector
 115 density is arbitrarily set to month 7 when plotting the time-series, to enable graphical
 116 comparability, see **Fig.2A**).

117 The results presented in the collated references were frequently presented in the form
 118 standardised by sampling effort, such as Man-Hour Density (MHD). They do not therefore
 119 represent the total number of mosquitoes caught each month (required for the statistical
 120 framework utilised to characterise temporal properties) and therefore, where information on
 121 sampling effort was present (e.g. number of hours spent sampling/catching *An. stephensi*,
 122 number of households or cattlesheds searched, number of trap nights etc), we used this
 123 information to convert MHD back to the raw counts. In the small number of instances where
 124 there was variable sampling effort across the time series (which would bias the conversion
 125 away from the underlying population abundance), we conservatively used the lowest sampling
 126 effort recorded across the time series in the conversion. Together, this allowed us to produce

127 an estimate of the number of mosquitoes sampled (a raw count, based on equal sampling
128 effort across the time series). See **Supplementary Data: *Extracted Entomological Data*** for
129 more information about how each time-series was processed).

130 **Study Geolocation and Environmental Covariate Extraction**

131 For each study where geolocation was possible, we recorded the location at both the
132 administrative unit 1 and 2 level, based on information provided in the reference. A number of
133 the references identified in our review had previously been utilised as part of previous
134 reviews^{1,2} – where this data was available, these descriptions of study location were used. For
135 each location, we then extracted a suite of satellite-derived environmental covariates. These
136 environmental covariates consist of raster layers spanning all countries in which studies had
137 been conducted in (i.e. Afghanistan, Djibouti, India, Iran, Myanmar and Pakistan) at a 2.5 arc-
138 minute (~5km by 5km, depending on the exact location and distance from the equator) spatial
139 resolution. The covariates utilised here were initially selected from a set of 19 derived from the
140 *BioClimatic* variables (a suite of biological relevant covariates defined from monthly rainfall
141 and temperature satellite data⁵, making the strong assumption that these variables, which
142 represent location averages over the period 1970-2000, adequately describe the climactic
143 factors present in the periods spanned by our studies, which were predominantly conducted
144 after 2000) as well as measures of landcover and urbanicity⁶, population density^{7,8} and
145 enhanced vegetation index^{9,10}. This provided a total of 43 covariates, many of which were
146 highly correlated with one another. To reduce the degree of this multicollinearity, we generated
147 a reduced subset of covariates using tools available in the *tidymodels* collection of R
148 packages¹¹ that aim to minimise the Spearman rank correlation coefficients between retained
149 covariates, and also exclude covariates where there is minimal variation for that covariate
150 across the full dataset, leaving **19** covariates in total. In addition to the environmental
151 covariates described above, for each of the administrative units a survey had been carried out
152 in, we also collated daily rainfall estimates for the time-period the survey had been conducted
153 in, using the “*The Climate Hazards Group Infrared Precipitation With Stations*” (CHIRPS)
154 dataset¹². These data were aggregated up to the same temporal resolution as the *An.*
155 *stephensi* catch data (i.e. monthly). These rainfall data were used to calculate the cross-
156 correlation coefficient between mosquito catches and rainfall.

157

158 **Supplementary Table 2: The Complete Suite of Covariates Collated and Subsequently**
 159 **Reduced for Modelling and Prediction of Seasonal Population Dynamics**

#	Variable	Temporal Resolution	Source
1	BioClimatic - Annual Mean Temperature	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
2	BioClimatic - Mean Diurnal Range	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
3	BioClimatic - Isothermality	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
4	BioClimatic - Temperature Seasonality	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
5	BioClimatic - Max Temperature of Warmest Month	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
6	BioClimatic - Min Temperature of Coldest Month	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
7	BioClimatic - Temperature Annual Range	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
8	BioClimatic - Mean Temperature of Wettest Quarter	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
9	BioClimatic - Mean Temperature of Driest Quarter	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
10	BioClimatic - Mean Temperature of Warmest Quarter	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
11	BioClimatic - Mean Temperature of Coldest Quarter	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
12	BioClimatic - Annual Precipitation	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
13	BioClimatic - Precipitation of Wettest Month	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
14	BioClimatic - Precipitation of Driest Month	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
15	BioClimatic - Precipitation Annual Coefficient of Variation	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
16	BioClimatic - Precipitation of Wettest Quarter	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
17	BioClimatic - Precipitation of Driest Quarter	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
18	BioClimatic - Precipitation of Warmest Quarter	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
19	BioClimatic - Precipitation of Coldest Quarter	Annual Average, 1970 - 2000	https://www.worldclim.org/bioclim
20	Population Density	1992-2020, using the closest year value to the year of the study	http://www.worldpop.org.uk
21	Enhanced Vegetation Index	1992-2020, using the closest year value to the year of the study	https://modis.gsfc.nasa.gov/data/dataproduct/mod13.php
22	Landcover	1992-2020, using the closest year value to the year of the study	https://maps.elie.ucl.ac.be/CCI/viewer/index.php
23	Average Monthly Catch	Calculated empirically for each time-series	NA
24	Maximum proportion of total annual rainfall in any consecutive 4 month period	Calculated empirically for each location and time-series	NA
25	Country survey had been carried out in	Calculated empirically for each time-series (grouped into "India", "Iran" and "other").	NA

160 **Note:** There are **43** covariates total here, as Landcover contains **19** distinct covariates (each
 161 describing the proportion of cover attributable to a particular landcover class in a given area).

162 **Note:** All WorldClim data is from Version 2 of the datasets.

163

164 **Supplementary Information 2: Description of Statistical** 165 **Methodologies**

166 **Negative Binomial Gaussian Process – Fitting and Inference:**

167 In-line with previously work modelling the seasonal dynamics of different *Anophele* mosquito
168 species from across India², we utilise a flexible Gaussian Process modelling framework to
169 temporally interpolate between the monthly-catch datapoints and smooth the raw, noisy and
170 overdispersed catch data. Gaussian processes specify a distribution over functions such that
171 any finite set of function values $f(x_1), f(x_2), \dots, f(x_N)$ have a joint Gaussian distribution¹³. The
172 Gaussian process is entirely specified by its mean function:

$$173 \quad E[f(x)] = \mu(x)$$

174 and by its covariance function:

$$175 \quad Cov[f(x), f(x')] = k(x, x')$$

176 The covariance function is also known as the kernel and defines, based on the Euclidean
177 distance between any two points, their covariance (and thus the covariance matrix of the
178 Gaussian Process when all pairwise combinations of points are considered). Many different
179 forms of the kernel are possible that each encode different prior information about how we
180 expect two datapoints (x and x' in this instance) to be similar, and the distance over which we
181 expect this similarity to persist. Given that mosquito population dynamics are typically
182 characterised by seasonally repeating patterns occurring either, a periodic kernel function was
183 used to define the covariance between pairs of points:

$$184 \quad k(x, x') = \alpha^2 \exp\left(-\frac{2}{l^2} \sin^2\left(\frac{\pi|x - x'|}{p}\right)\right)$$

185 where p represents the period over which we would expect points to show similar dynamics
186 (i.e. a period of twelve would imply we expect points separated by 12 months to be most
187 similar), α specifies the magnitude of the covariance, and l represents a lengthscale
188 parameter further constraining the extent to which two values separated by a given time can
189 co-vary.

190 Bayesian inference and fitting of normal Gaussian Processes typically follow this hierarchical
191 formulation:

$$192 \quad \theta \sim \pi(\theta)$$

$$193 \quad f \sim GP(0, K_\theta(x))$$

$$194 \quad y_i \sim MVN(f(x_i), \sigma^2) \forall i \in \{1, \dots, N\}$$

195 where θ represents a vector of hyperparameters involved in defining the kernel's properties,
196 f is a distribution of functions from a zero-mean Gaussian Process with covariance function
197 K_θ , $f(x)$ are function evaluations at times x , and y the observed data. We modify this structure
198 to account for specific characteristics of the mosquito data being utilised – specifically that the
199 data are integer counts, that mosquito catch data is rarely normally distributed and frequently
200 displays high levels of overdispersion (a common property of biological systems generally).
201 We therefore adapted the above framework to accommodate a Negative Binomial likelihood,
202 leading to the following inferential framework:

$$203 \quad p, \alpha, l \sim \pi(p, \alpha, l)$$

$$204 \quad \mathbf{f} \sim GP(0, K_{\theta}(x))$$

$$205 \quad \text{where:} \quad k(x, x') = \alpha^2 \exp\left(-\frac{2}{l^2} \sin^2\left(\frac{\pi|x - x'|}{p}\right)\right)$$

$$206 \quad y_i \sim \text{Negative Binomial}(e^{f(x_i)}, \sigma) \forall i \in \{1, \dots, N\}$$

207 where $e^{f(x)}$ is used to reflect the fact that we use a log link between the observed counts and
 208 the underlying latent process reflecting the population dynamics, and σ represents the
 209 overdispersion parameter of the Negative Binomial distribution.

210 **Prior Specification**

211 Per previous work², prior distributions for the estimated parameters were defined as follows:

$$212 \quad l \sim \text{Normal}(2, 1^2)$$

$$213 \quad \alpha \sim \text{Half - Normal}(0, \sqrt{SD(y)})$$

$$214 \quad p \sim \text{Normal}(12, 4^2)$$

$$215 \quad \sigma \sim \text{Half - Normal}(0, 8^2)$$

216 Weakly informative priors were set on the scaling factor α , the period, p , and the
 217 overdispersion parameter, σ . The prior for the kernel period (p) was centred on 12 (a value of
 218 the period that would represent annual variation being the dominant temporal modality) to
 219 reflect our prior belief that observed variation in mosquito abundance is likely to cycle annually.
 220 However, recognising that other temporal patterns of fluctuating abundance are possible, we
 221 placed a large standard deviation on p to allow the model to accommodate instances of
 222 bimodality or periods operating across timescales longer than a year. We placed lower and
 223 upper bounds on p at 4 and 18 months respectively, to avoid identifiability issues arising from
 224 the lack of data at temporal resolutions below and above these bounds.

225 **Model Fitting and Parameter Inference**

226 This Negative Binomial Gaussian Process were fitted using a Bayesian framework
 227 implemented in STAN, a probabilistic programming language for statistical inference written
 228 in C++ that employs the No-U-Turn sampler, a variant of the gradient-based Hamilton Monte
 229 Carlo algorithm for inference¹⁴. For each time-series, 2 chains of 20,000 iterations were run
 230 for purposes of model fitting and parameter inference. Half of each chain's iterations were
 231 discarded as burn-in/the adaptive phase of the sampling, leaving a total of 20,000 iterations
 232 available for inference. Measures of MCMC convergence such as the Gelman-Rubin statistic
 233 were monitored in all cases and were all consistently < 1.02 .

234 **Fitted Time Series Normalisation and Von Mises Distribution Fitting**

235 After having fitted and smoothed the mosquito catch time-series, we normalised each in the
 236 following way:

$$237 \quad p_i = \frac{y_i}{\sum y_i}$$

238 where p_i is the proportion of the annual catch recorded at timepoint i . This was done in order
 239 to establish comparability across the time series (which varied substantially in the absolute

240 numbers of *Anopheles stephensi* caught). We then further characterised the periodic
 241 properties of these time series by fitting Von Mises distribution to the time-series. The Von-
 242 Mises distribution is a continuous probability distribution that exists on the circle, with range 0
 243 to 2π . It is the circular analogue of the normal distribution (which exists on the line), with the
 244 probability density function for the angle x given by:

$$245 \quad f(x|\mu, \kappa) = \frac{e^{\kappa \cos(x-\mu)}}{2\pi I_0(\kappa)}$$

246 where $I_0(\kappa)$ is the modified Bessel function of order 0, the parameter μ is a measure of location
 247 (analogous to the mean of the normal distribution, describing where on the circle the
 248 distribution is clustered around) and κ describes the concentration of density around μ (and
 249 thus its inverse is a measure of dispersion, analogous to σ^2 for the normal distribution. We
 250 fitted two sets of Von Mises densities to the normalised time series, the first containing a single
 251 component:

$$252 \quad f(x|\mu_1, \kappa_1) = f_1(x|\mu_1, \kappa_1)$$

253 and another with two-components, formulated as:

$$254 \quad f(x|\mu_1, \kappa_1, \mu_2, \kappa_2, \omega) = \omega f_1(x|\mu_1, \kappa_1) + (1 - \omega) f_2(x|\mu_2, \kappa_2)$$

255 where x represents the normalised mosquito count formulated as a random variable on the
 256 circle (i.e. $x = \frac{2\pi p_i}{12}$). Fitting was undertaken using the *optim* function in R, with the root mean
 257 squared error as the loss function. The outputs from this fitting were then included in the
 258 process generating aggregate summaries of the temporal properties of the time-series, a
 259 process described in further detail below.

260 Time Series Characterisation and Analysis

261 To characterise the temporal properties of each time-series, we calculated a series of
 262 summary statistics for each, drawing on previous work carried out exploring the empirical
 263 structure of time series¹². In doing this, we can make explicit comparisons between time-series
 264 about key aspects of their temporal properties (e.g., the degree or timing of seasonality), and
 265 in doing so, identify time-series with similar statistical and temporal properties. These
 266 summary statistics were the following:

- 267 1. **Periodic Kernel Median:** Fitting the Negative Binomial Gaussian Process with a
 268 periodic kernel allowed inference of the period, p , providing us with an estimate of the
 269 frequency of repeating patterns in the monthly abundance of mosquitoes. An estimate
 270 of p was calculated for each fitted time series, with the median value of p across the X
 271 HMC iterations for each time-series used here
- 272 2. **Kullback-Leibler Divergence:** Also known as the relative entropy, the Kullback-
 273 Liebler divergence represents a measure of how different one probability distribution
 274 is from a second probability distribution (where a value of 0 indicates that the two
 275 distributions are identical). It is specified in the following manner:

$$276 \quad E_i = p_i \log_2 \left(\frac{p_i}{q_i} \right)$$

277

$$E = \sum_{i=1}^{12} p_i \log_2 \left(\frac{p_i}{q_i} \right)$$

278

where p_i is the average value of the normalised time series for month i , and $q_i = 1/12$ for $i = 1, \dots, 12$. This operation therefore measures the deviation of a normalised time series from a uniform distribution, in doing so, informing about the extent to which a seasonal peak (or peaks) is present in the time series.

279

280

281

282

3. **Time Difference Between Vector Density Peak and Rainfall Peak Timings:** The time difference between the highest recorded vector density and the highest recorded rainfall for that year.

283

284

285

4. **Proportion of Points Greater Than 1.65x the Mean:** For each fitted, normalised time series, the proportion of points greater than 1.65x the time-series' mean was calculated, informing the degree and width of any seasonal peaks.

286

287

288

5. **Number of Peaks:** Estimates of the parameters governing the fitted two component Von Mises distribution were used to infer the number of peaks in each time series. Specifically, and in-keeping with previous work², a time series was deemed to possess one peak if the value of the Von Mises component weighting was either < 0.3 or > 0.7 and the difference in means was $< \frac{2\pi}{3}$ or $> \frac{4\pi}{3}$, indicating that the majority of the density could be attributed to one of the two components, and that the two means identified during the fitting were temporally close to one another. Otherwise, a time series was judged to possess two peaks.

289

290

291

292

293

294

295

296

6. **Von Mises 1 Component Mean:** If a 1 component Von Mises distribution was preferred, then the Von Mises mean corresponding to the maximum likelihood predicted value was used. If the 2 component Von Mises distribution was preferred, the value for this operation for that particular time series is set to -5.

297

298

299

300

7. **Von Mises Two Component Weight:** Estimates of the weight parameter governing the two component Von Mises distribution were also used to infer the bimodality of the time series. The weight specifies the proportion of each component that is used to fit the time series and thus a very high (or very low weight) indicates the dominance of a single component and the comparatively small contribution of the other.

301

302

303

304

305

8. **Maximum Percentage of Total Annual Catch In Any 3 Month Period:** In-keeping with previous, operationally aligned estimates of malaria seasonality¹⁵, we calculated using a sliding 3-month window the maximum percentage of the total annual catch that was caught in any 3 month period.

306

307

308

309

Principal Components Analysis and Clustering

309

310

Principal Components Analysis (PCA) is a statistical procedure that utilises an orthogonal transformation to convert a set of correlated variables (in this case the outputs of the 7 mathematical operations described above for each of the time series) into a set of linearly uncorrelated variables (known as the "principal components"). In doing so, this allows us to summarise this set of variables with a smaller number of representative variables that together explain the majority of the variability in the variables. Reducing the dimensionality of the dataset in this way facilitates visualisation of time series properties (as defined by the mathematical operations) as well as clustering of the time series into groups which share similar properties (clustering algorithms typically perform poorly in high dimensional settings, necessitating the use of PCA as described here). Clustering was then undertaken using the k-means clustering algorithm, using the first four PCA components that together described 85% of the total variation present in the data.

311

312

313

314

315

316

317

318

319

320

321

322 **Random Forest Modelling and Prediction of Seasonality**

323 Random Forests are a machine learning, ensemble-based method that work by constructing
324 a collection of decision trees that together explain the results (where results are either a
325 continuous outcome variable in the regression context, or a binary indicator in the classification
326 context)¹⁶. The outputs of these decision trees are subsequently aggregated in a statistically
327 principled and coherent way to produce a “forest” (or ensemble) of trees that together produce
328 predictions for comparison with data. They have previously been shown to provide significant
329 improvements in accuracy over traditional linear regression based approaches, particularly in
330 contexts where non-linear relationships or interactions between covariates are likely present
331 and to be relevant to prediction of an outcome¹⁷.

332 We used a Random Forest based approach to either 1) classify time-series cluster
333 membership (i.e. predict whether a time-series belonged to either Cluster 1 or Cluster 2, as
334 defined via the PCA and k-means clustering analysis described above); or 2) predict *An.*
335 *stephensi* time-series seasonality (defined as the percentage of total annual vector density in
336 any continuous 3-month period). These models were fitted using the software package
337 *Ranger*¹⁸, implemented in the *tidymodels* framework for R¹¹, with 6-fold cross-validation
338 utilised to optimise hyperparameter combinations; presented results are based on averaging
339 the results of 25 separate iterations of cross-validation and model fitting (to account for
340 stochasticity in model fitting), and any predictions made using out-of-bag model estimates in
341 all instances. Due to significant imbalances in class size across the time-series clusters (49
342 time-series in Cluster 1 compared to only 16 time-series in Cluster 2, we carried out
343 upsampling using the SMOTE (synthetic minority over-sampling technique¹⁹) algorithm. We
344 also carried out model fitting without this upsampling, the results of which are presented in
345 **Supp Fig. 6.**

346 In all instances, out-of-sample predictive accuracy was assessed using 6-fold cross-validation
347 (CV) and used to optimise the hyperparameters associated with the Random Forest method
348 algorithm. Random Forest models were fitted to the training dataset (i.e. the full dataset minus
349 one of the CV folds) and then model accuracy assessed on the remaining fold of data not
350 included in model training. In the case of the cluster classification example, the metric used to
351 evaluate model performance was the area under the curve (AUC). In the case of the
352 regression prediction of seasonality, the metric used to evaluate model performance was the
353 root mean squared error (RMSE). The Random Forest hyperparameters providing the best
354 out-of-sample AUC/RMSE were then selected, and a final Random Forest model then fitted
355 on the full set of data available. Predictive accuracy (assessed via AUC/RMSE) was then
356 calculated for the entire dataset by using out-of-bag predictions for each sample i.e.
357 predictions on each training sample using only the trees that did not have that training sample
358 in their bootstrap sample. We also calculated both permutation variable importance and
359 generated partial dependency plots²⁰ for each model to assess the contribution of specific,
360 individual environmental covariates to whether a time-series had a single seasonal peak or
361 not. Together these methods allow evaluation of the importance of each included covariate to
362 model predictive accuracy, and in turn, allows us to “rank” covariates according to their
363 contribution to the predictive performance of the model. This entire process was repeated 25
364 times in order to average over the stochasticity and variation inherent in the Random Forest
365 fitting process.

366 We also carried out an additional sensitivity analysis where a set of the available data (n=12
367 time-series) was held-out at the onset, and the random forest model trained (using 6-fold

368 cross-validation) on the remaining available data (n=53 time-series total, with 43 time-series
 369 used in model fitting and 10 time-series used for performance evaluation in each of the cross-
 370 validation folds). Optimal hyperparameters were selected in the same way as described
 371 above, and then a final model fitted to the full, non-held out data (n=53 time-series), and model
 372 predictive accuracy assessed by evaluating performance on the held-out data (n=12 time-
 373 series).

374 **Probability of Detecting *Anopheles stephensi* With Different Surveillance and** 375 **Monitoring Strategies**

376 We explore the implications of seasonal variation in *An. stephensi* abundance on the
 377 probability of detecting the vector in entomological surveillance and monitoring using human
 378 landing catches. Note that what follows below assumes there is no seasonal variation in
 379 factors other than mosquito abundance (such as resting preferences) that might influence the
 380 probability of *An. stephensi* being caught in a human landing catch. In the absence of
 381 estimates of overall mosquito population size, we first start by considering an arbitrary
 382 Entomological Inoculation Rate (*EIR*, the number of infectious bites an individual receives
 383 each year) and Sporozoite Rate (*SR*, the prevalence of sporozoites in the mosquito
 384 population), which together define an overall annual biting rate (*ABR*).

$$385 \quad ABR = \frac{EIR}{SR}$$

386 For the purposes of the results in the main text, we select an *EIR* of 1 and an *SR* of 0.05 to
 387 give an *ABR* of 20, though we stress these choices are arbitrary and meant to be illustrative
 388 only, and that the methods below could be used to calculate the results for any combination
 389 of *EIR* and *SR*. For a given *ABR* and for each *An. stephensi* time-series *i*, this *ABR* is
 390 proportionally divided up over the course of 365 days according to the normalised vector
 391 density at each timepoint, such that the biting rate *b* for time-series *i* on day *d* is given by:

$$392 \quad b_{i,d} = ABR \left(\frac{D_{i,d}}{\sum_{d=1}^{365} D_{i,d}} \right)$$

393 where $D_{i,d}$ is the normalised vector-density on day *d* for time-series *i*. Because the sampling
 394 resolution of the studies we collated was never finer than monthly, we then use $b_{i,d}$ to calculate
 395 an average daily biting rate for each month, $b_{i,m}$. $b_{i,m}$ therefore describes the expected daily
 396 number of bites an individual receives in month *m*. The number of bites an individual would
 397 receive during a specific day of human landing catch sampling during month *m* can then be
 398 considered a draw from a Poisson distribution with rate as follows:

$$399 \quad C_{i,m} \sim \text{Poisson}(\lambda = b_{i,m})$$

400 The expected number of mosquitoes caught over multiple days and months of mosquito
 401 sampling can then be calculated by exploiting the following property of the Poisson distribution:

$$402 \quad X_1 \sim \text{Poisson}(\lambda_1) \ \& \ X_2 \sim \text{Poisson}(\lambda_2)$$

$$403 \quad \text{then } X_1 + X_2 \sim \text{Poisson}(\lambda_1 + \lambda_2)$$

404 Given this, for time-series *i*, carrying out mosquito sampling for n_m consecutive months
 405 starting at month *j*, and within each month carrying out n_d days-worth of sampling, the total
 406 number of *An. stephensi* expected to be caught is given as follows:

$$C_{i,j,n_m,n_d} \sim \text{Poisson} \left(\sum_{m=j}^{m=(j+n_m-1)} n_d b_m \right)$$

407
408 from which the probability of not sampling *An. stephensi* (i.e. the total number of *An. stephensi*
409 caught is equal to 0) during that sampling period can be calculated.

410 For each time-series, we then identified the month in which monthly rainfall peaked, and the
411 month in which vector density was highest (noting that these months were very rarely the
412 same month). We then calculated the cumulative probability of *An. stephensi* detection under
413 a range of different surveillance strategies. Specifically three strategies were simulated:

- 414 • **Vector-Peak Timed:** Starting the survey at the month with peak vector density (noting
415 that in the absence of pre-existing detailed entomological information this is largely a
416 hypothetical quantity, meant to illustrate the maximum detection probability that could
417 be achieved).
- 418 • **Rainfall-Peak Timed:** Starting the survey at the month with peak rainfall.
- 419 • **Random Month Timed:** The expected cumulative probability of detection achieved if
420 the survey was started during a random month (calculated in practice by simulating
421 survey starting in each of the year's 12 months and then calculating the average
422 cumulative probability of these surveys).

423 In addition to varying the timing of the survey (which varies according to the surveillance
424 strategy considered, as described directly above), we also varied the amount of sampling effort
425 (number of days sampled within each month) and the overall duration of the (i.e. how many
426 consecutive months were sampled). Note that the aim here is not to describe the exact
427 probability of missing *An. stephensi* in any given entomological survey, as this will depend on
428 a wide array of other, poorly defined and heterogeneous factors (such as type of catch
429 methodology used etc). Instead, the aim is to highlight how variation in seasonal dynamics
430 can influence the nature of surveillance required to successfully *An. stephensi*.

431 **Modelling of Malaria Transmission and the Impact of *Anopheles stephensi***

432 We integrated the temporal profiles of *An. stephensi* abundance into a well-established
433 deterministic compartmental model of *Plasmodium falciparum* malaria transmission and
434 disease^{21–23} to explore the implications of the vector's establishment and seasonality on the
435 dynamics of malaria transmission, with a particular focus on areas where malaria transmission
436 is currently low or absent. What follows is a description of the mathematical modelling
437 framework in general terms, followed by specific details about how exactly this framework was
438 used to model malaria transmission underpinned by *An. stephensi* in settings where malaria
439 is currently absent or only minimally present.

440 The deterministic malaria model used here considers both human and mosquito populations.
441 Humans begin as Susceptible (*S*), and upon infection (at a rate which is dependent on the
442 force of infection they experience), progress to either Asymptomatic (*A*) or clinical disease,
443 with the comparative probability of these two outcomes depending on the degree of acquired
444 natural immunity due to previous exposure to the parasite. If an individual progresses to
445 clinical disease, they enter either a Treated (*T*) or Clinical Disease (*D*) state that depends on
446 the probability of receiving treatment. For those treated, individuals progress through a period
447 of prophylactic protection following treatment (*P*), and then return to the susceptible
448 compartment. For those developing clinical disease, they remain symptomatic for the duration

449 of the disease, before moving to an asymptomatic state (*A*, detectable by light microscopy),
450 before subsequently moving to a submicroscopically infected state (*U*, not detectable by light
451 microscopy). Individuals who are currently asymptotically infected (including individuals in
452 both the *A* and *U* states) can be reinfected and develop clinical disease once again – if this
453 does not occur, they subsequently clear the infection and return to the susceptible state.

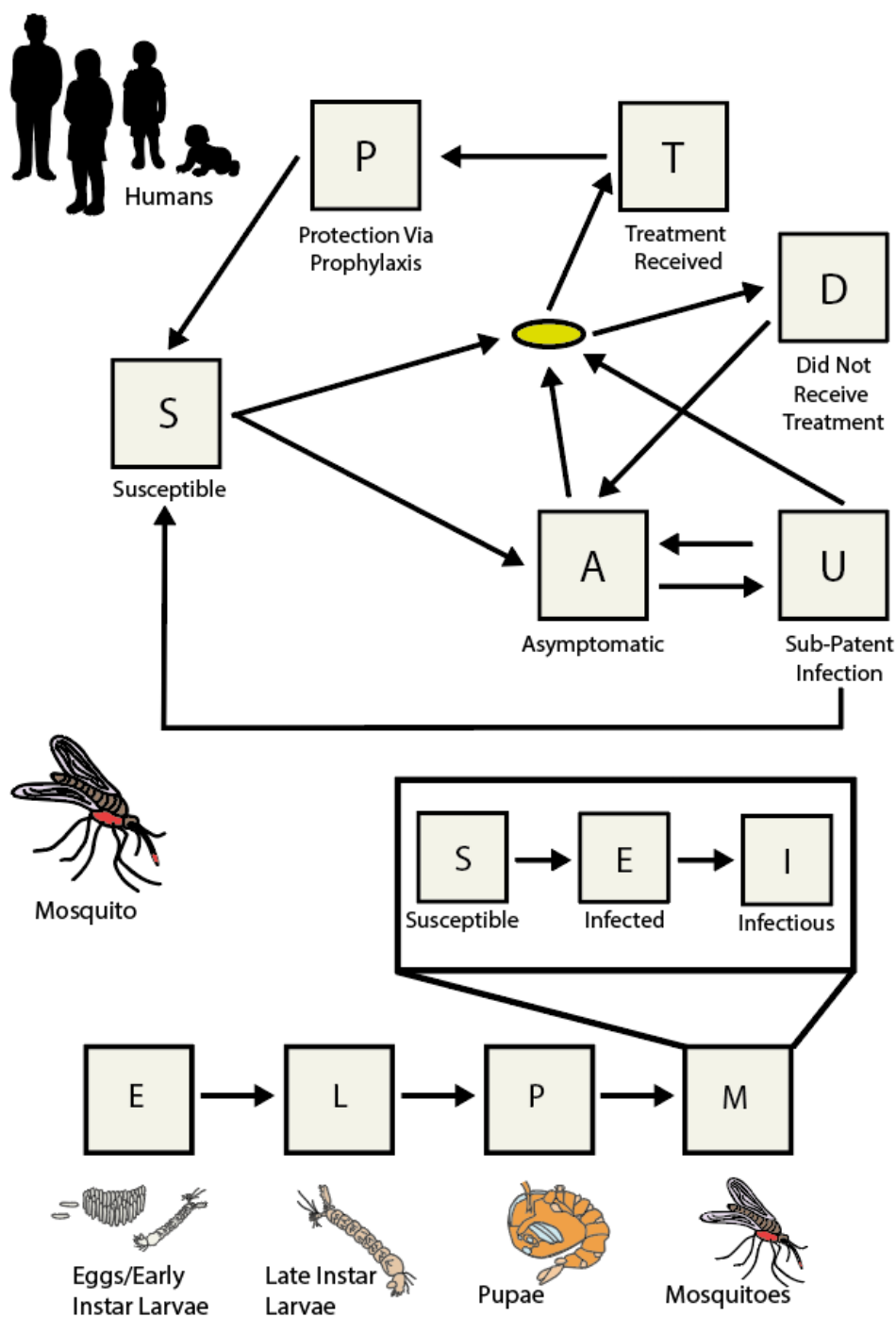
454 Adult mosquito populations and their preceding juvenile stages are also explicitly modelled.
455 Immature mosquitoes start off as larvae, divided into early and late stage (*Es* and *Ls*
456 respectively) which then mature into pupae (*P*) before eventually maturing into adult
457 mosquitoes. Adult mosquitoes are further stratified according to infection with *P.falciparum*
458 status – they begin as susceptible (*Sm*) and upon infection, progress to an exposed (but un-
459 infectious, *Em*) state, and then onto the infectious state (*Im*) following the extrinsic incubation
460 period (EIP, assumed to be constant over time). Mosquitoes are infected through exposure to
461 humans currently possessing transmissible infections i.e. the treated (*T*), clinical disease (*D*),
462 asymptomatic (*A*) and submicroscopic (*U*) infection states.

463 Seasonality in mosquito abundance is incorporated through a flexible, time-varying carrying
464 capacity that in broad terms describes temporal variation in the ability of the local environment
465 to support mosquito breeding. The value of this carrying capacity relative to the size of the
466 mosquito population influences the mortality of early and late-stage larvae, which as previous
467 modelling work has shown, enables the model to accurately and adequately capture temporal
468 fluctuations in mosquito abundance²⁴. We integrate each of the seasonal profiles of *An.*
469 *stephensi* density implied by the corresponding time-series of catch data into the model,
470 matching the carrying capacity to the empirically observed temporal variation in *An. stephensi*
471 abundance. Estimates of the bionomic properties of *An. stephensi* (specifically the mosquito's
472 daily mortality, degree of anthropophagy, degree of endophily and the proportion of bites taken
473 on individuals indoors and/or in bed) were taken from previous work that reviewed the
474 properties³, and the vector to human ratio arbitrarily set to 20, which corresponds to
475 approximately 9% malaria prevalence in a setting where the risk of malaria is constant year
476 round (i.e. a perennial setting). Indoor residual spraying (IRS) is assumed to reduce malaria
477 burden primarily by both killing adult mosquitoes and deterring them from biting and feeding –
478 specifically, IRS can either repel before biting and feeding, or kill following biting (when the
479 vector rests on a sprayed wall). The efficacy of IRS decays over time due to a loss of
480 insecticide. The efficacy of the different IRS compounds considered (bendiocarb, clothianidin
481 and pirimiphos methyl), as well as the different rates of efficacy decay were parameterised
482 using Sherrard-Smith et al 2018²⁵. We modelled the impact of a single round of IRS, timed
483 according to a range of different strategies that largely mirror the strategies described in the
484 section on surveillance and entomological monitoring above. Specifically, these were:

- 485 • **Optimal-Timing:** Starting the survey at the timepoint where the reduction in incidence
486 is maximised (noting that in the absence of pre-existing detailed entomological
487 information on the timing of peak vector abundance, this is a hypothetical quantity,
488 meant to illustrate the maximum impact that could be achieved with perfect
489 information).
- 490 • **Rainfall-Peak Based Timing:** Starting the survey at the midpoint of the month with
491 peak rainfall.
- 492 • **Random Month:** The expected reduction in malaria incidence achieved if the IRS
493 campaign was started during a random month (calculated in practice by simulating

494 survey starting in each of the year's 12 months and then calculating the average
495 cumulative probability of these surveys).

496 In all cases, the impact was calculated by comparing the reduction in malaria burden (as
497 measured by total annual incidence in the 12-month period following spraying) compared to a
498 counterfactual of no IRS.
499



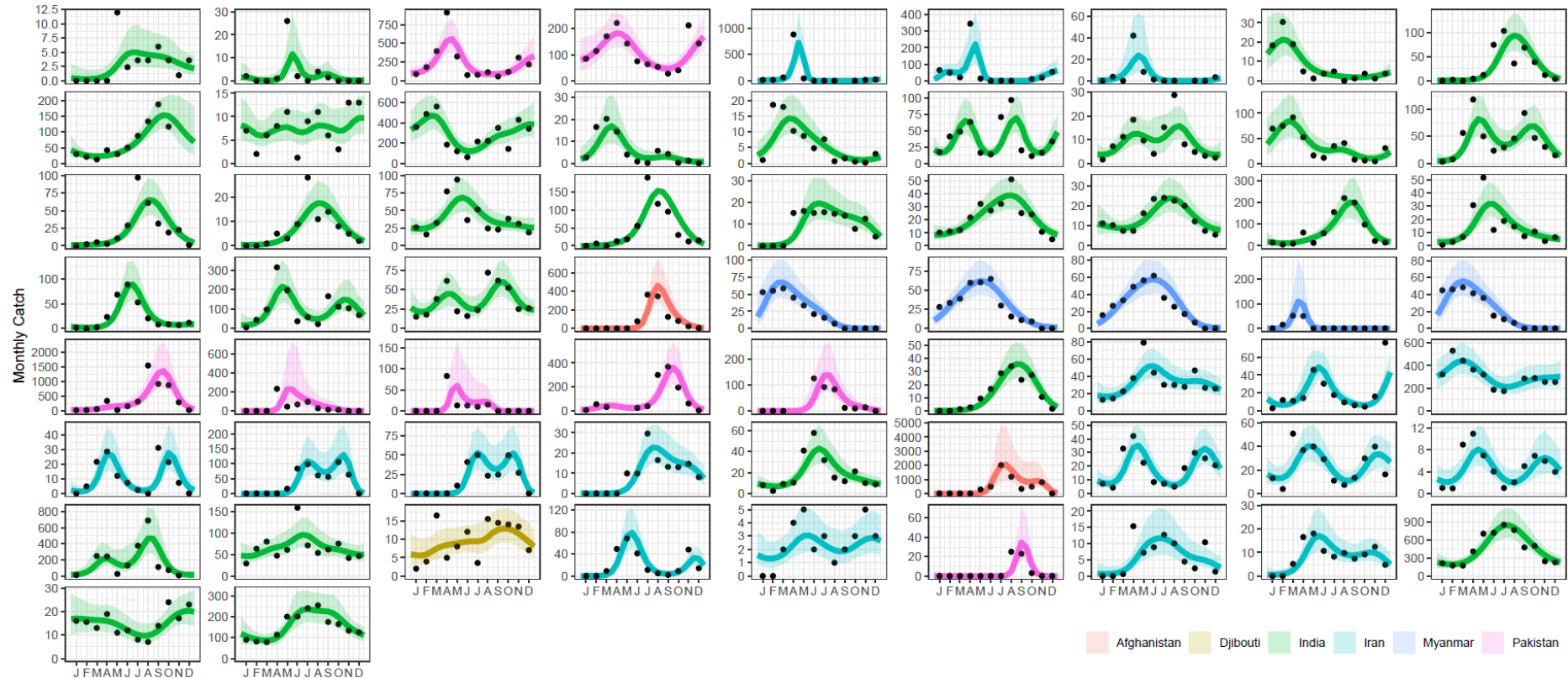
500

501 **Supplementary Figure Model Schematic:** Humans exist in either S (Susceptible), A
 502 (asymptomatic infection), T (infected and treated), D (infected and have clinical disease), U
 503 (submicroscopically infected) or P (prophylactically protected from infection by treatment
 504 received). The full-life cycle of the mosquito is modelled, with states including E (eggs/early
 505 larvae), L (late instar larvae), P (pupae), and finally M (mature adult mosquitoes). Mosquitoes
 506 in the M state begin in state S (susceptible) and upon infection more to a latently infected (but
 507 not yet infectious state) denoted by E. Upon becoming infectious they transition to the I state.
 508 Arrows show transitions between states, with the yellow oval indicating a decision point in the
 509 human part of the model based on whether the individual receives treatment or not.

510

511 **Supplementary Information 3: Additional Figures and Results**

512

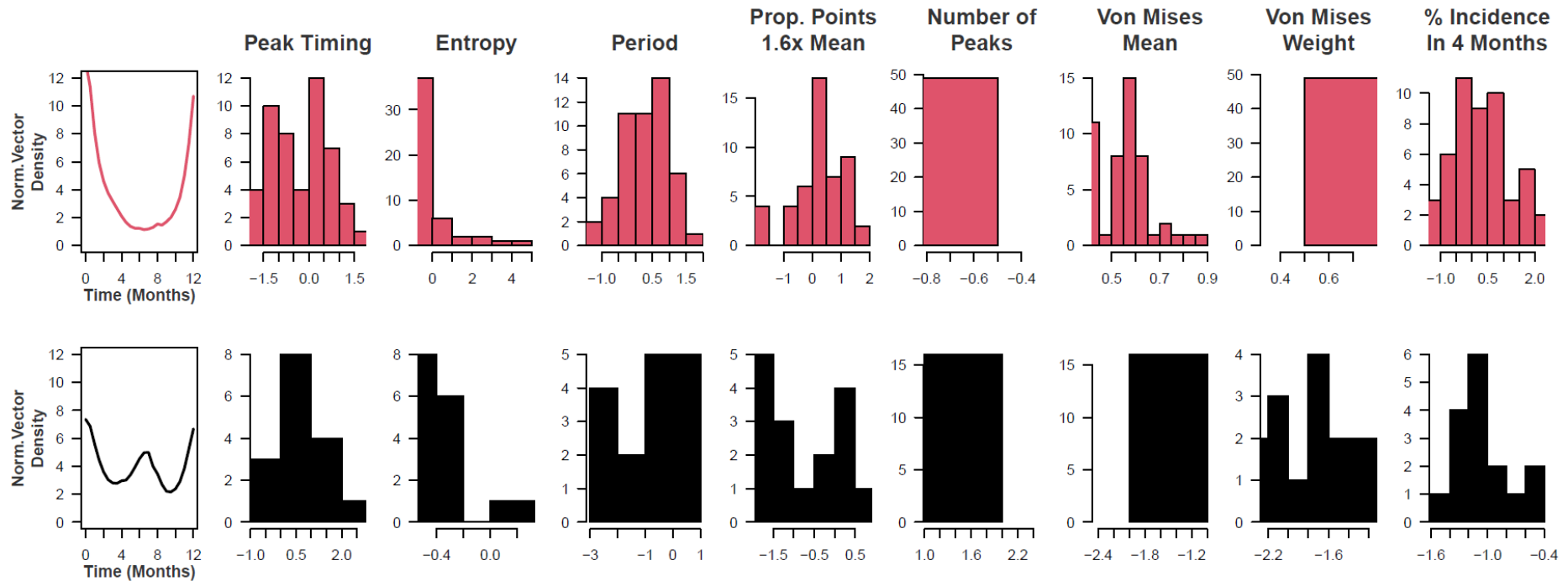


513

514 **Supplementary Figure 1: Results of model fitting to the longitudinal entomological data collated in this study.** Reviews of the literature
 515 in tandem with previously published databases of entomological data identified 65 *Anopheles stephensi* time-series matching the inclusion criteria
 516 (>10 months of catch data at monthly temporal resolution or finer), and a negative binomial gaussian process with period kernel fitted to each
 517 time-series. For the results presented above, black points are the data, and the lines represent the model output, coloured according to the
 518 country in which the study was conducted. Line indicates the mean model output, with the shaded ribbon delineating the 95% credible interval
 519 (CI).

520

521

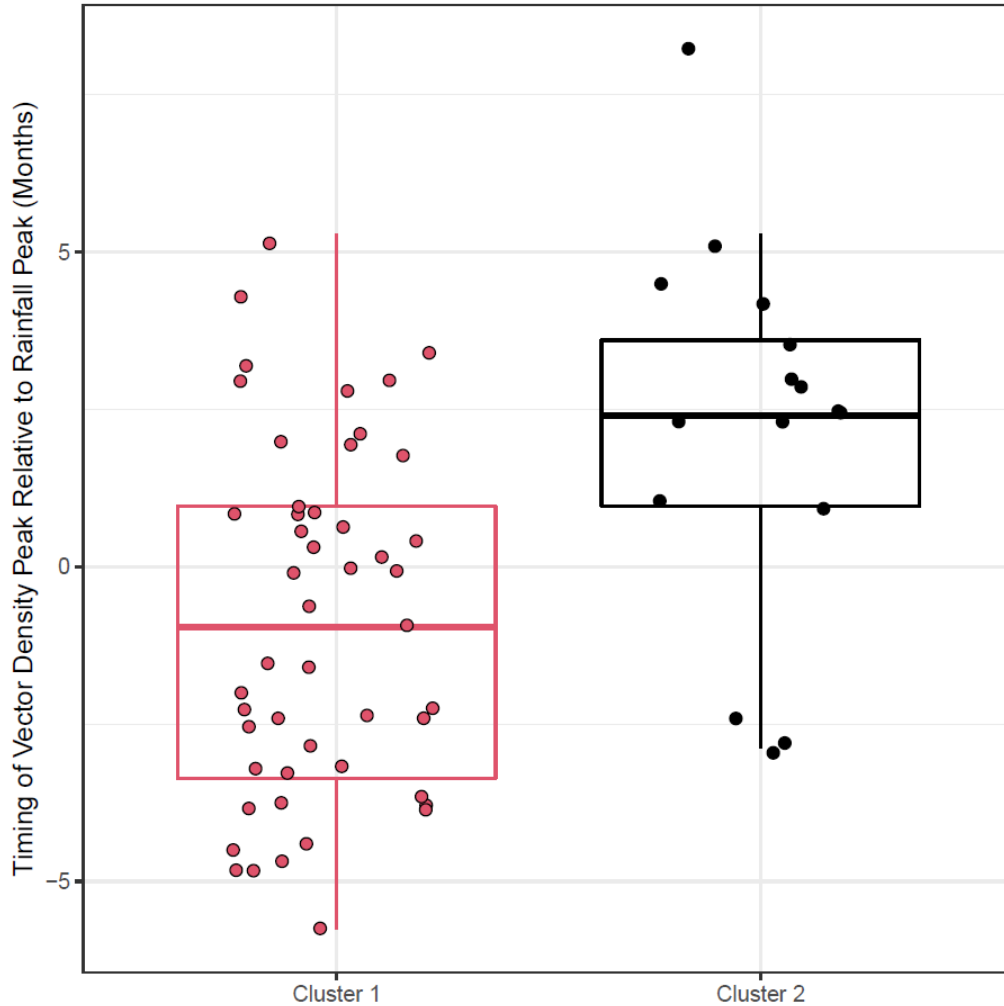


522

523 **Supplementary Figure 2: Archetype/Cluster Temporal Properties.** A series of mathematical operations were applied to the fitted time-series
 524 to characterise and explore their temporal properties. The results of this characterisation were then clustered using the k-means algorithm. For
 525 each cluster, the mean temporal profile is displayed, as well as the underlying distribution of values for each temporal property for each cluster
 526 (where the values for a given temporal properties for all time-series have first been normalised and standardised to have mean 0 and unit
 527 variance). For further information on each of these operations, see **Supplementary Information: Time Series Characterisation and Analysis.**

528

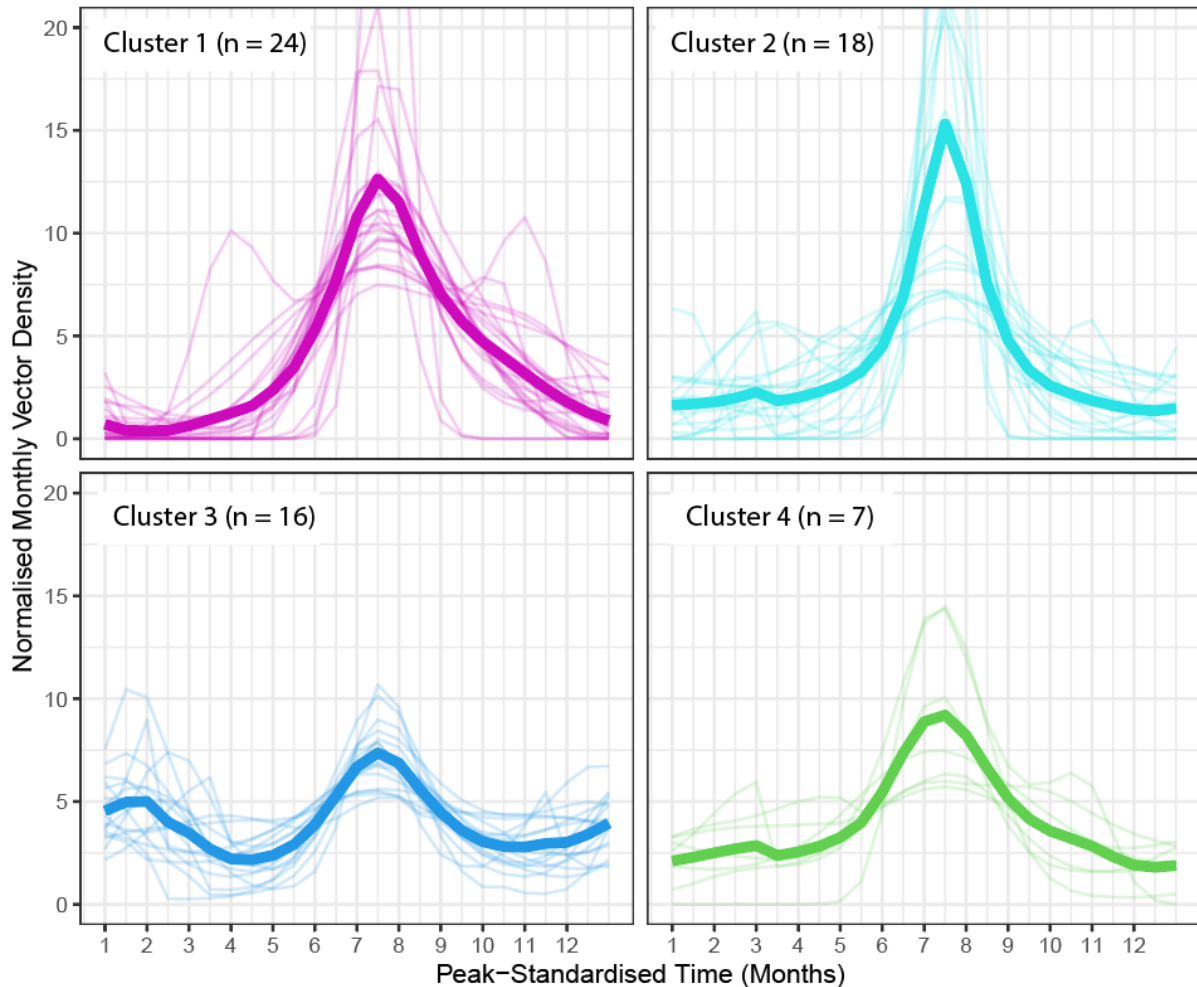
529
530
531
532



533
534
535
536
537
538
539

Supplementary Figure 3: Timing of Vector Density Peak Relative to Rainfall Peak, Stratified by Cluster. For each study location and time-series, we calculated the time-difference (in months) between the peak vector density and peak monthly rainfall. Our results highlighted systematic differences between clusters, but also significant variation within clusters.

540
541
542
543



544

545 **Supplementary Figure 4: Results of Clustering For 4 Clusters Instead of 2.** In order to
546 further investigate the different patterns of temporal dynamics present in the collated dataset,
547 we re-ran the k-means clustering algorithm this time specifying 4 clusters. The less seasonal
548 cluster from the 2 cluster analysis in the main text (Cluster 2 in the main text results) was
549 retained (here Cluster 3), and Cluster 1 from the main text was further disaggregated into 3
550 different clusters (here, Clusters 1, 2 and 4), each defined by different peak timings (mean
551 timing of vector density peak 7, 8.25 and 5.86 months after January for Clusters 1, 2 and 4
552 respectively) and the timing of the vector peak relative to peaks in rainfall (rainfall peak on
553 average 1.03 and 2.32 months before vector density peak for Clusters 1 and 2, 1.09 months
554 after vector density peak on average for Cluster 4).

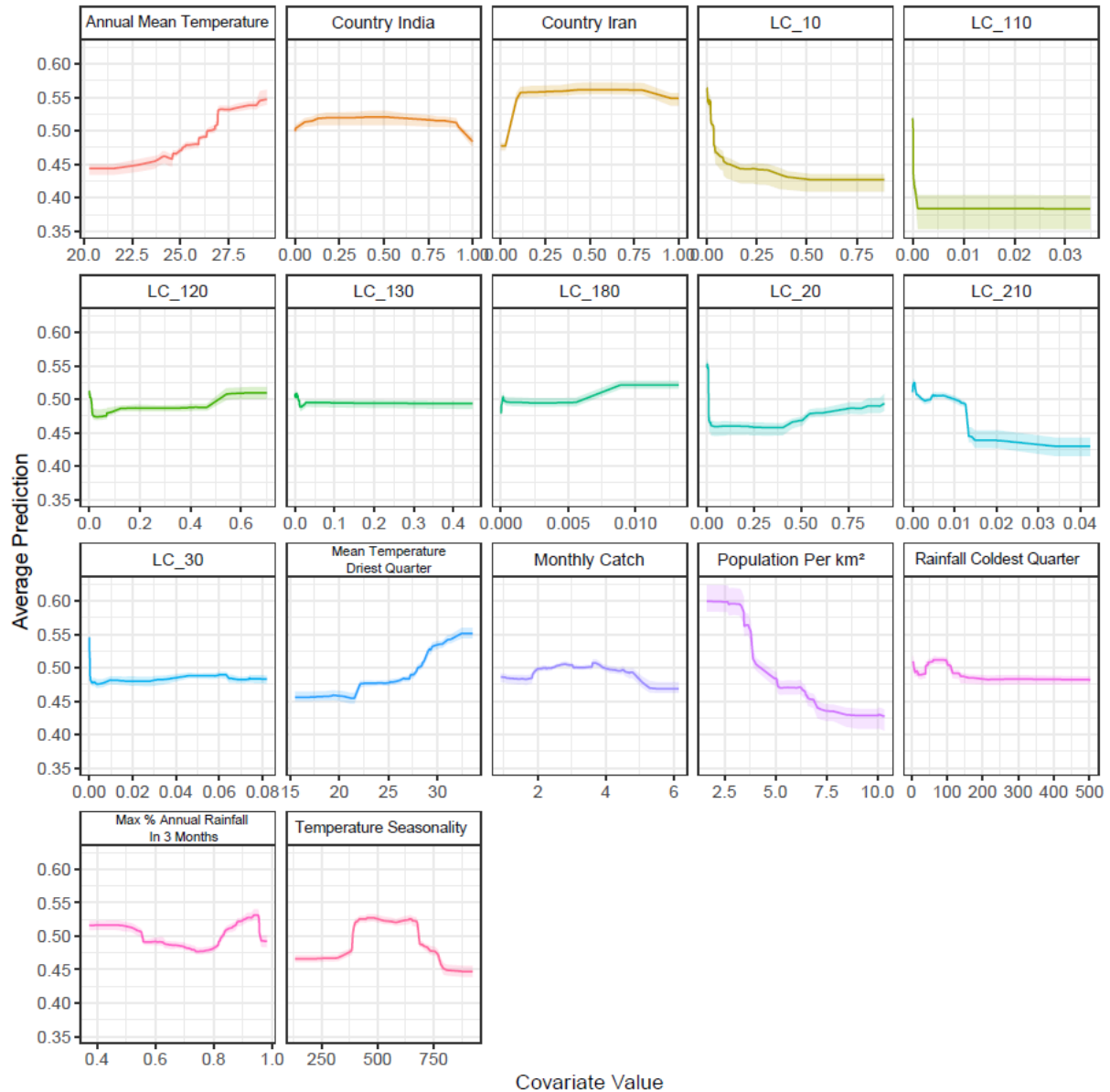
555

556

557

558

559



560

561 **Supplementary Figure 5: Partial Dependence Plots for Covariates Used in the Random**
 562 **Forest Classification Modelling.** The y-axis on the left shows the probability of the time-series
 563 belonging to Cluster 2 (i.e. a high probability indicates the time-series is predicted to
 564 likely belong to Cluster 2, a low probability indicates the time-series likely belongs to Cluster
 565 1). The x-axis describes the value of the (scaled, normalised) covariate.

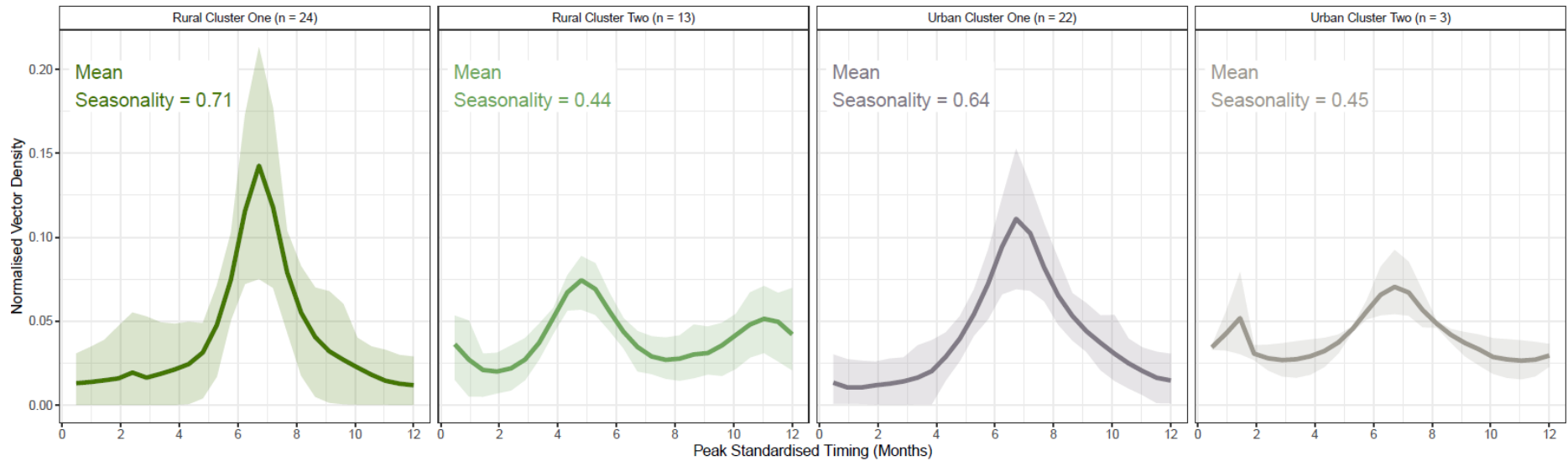
566

567

568

569

570

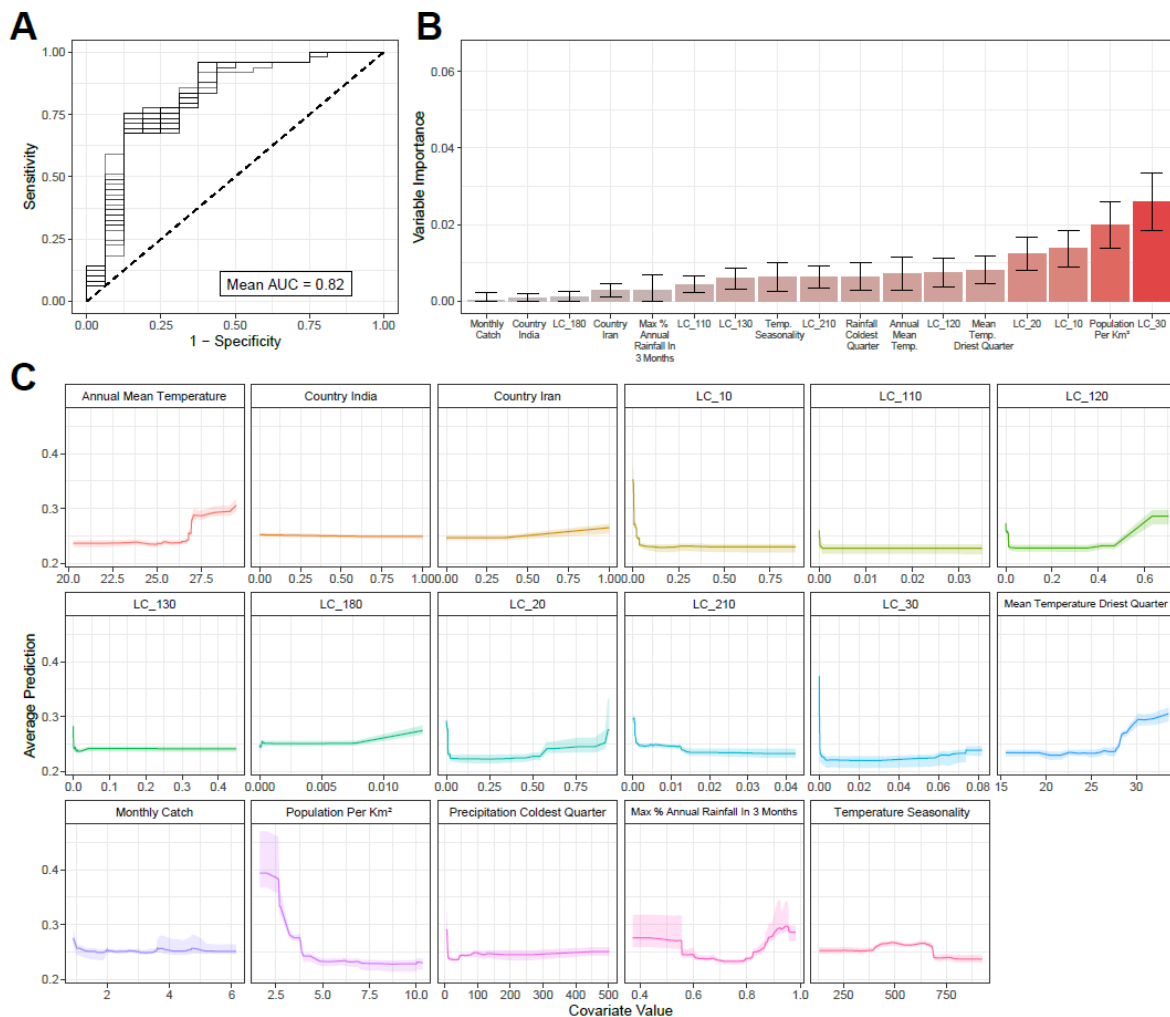


571

572 **Supplementary Figure 6:** Collated *An. stephensi* time-series, disaggregated according to urbanicity and cluster membership. Cluster 1 and
573 Cluster 2 time-series from rural locations and urban locations are plotted separately. Coloured line indicates the mean and ribbon indicates the
574 90% range spanned by the group of time series belonging to each displayed grouping. The average seasonality (defined as the maximum
575 percentage of total annual incidence in any continuous 4-month period) is also displayed for each group.

576

577
578
579



580

581 **Supplementary Figure 7: Random Forest Classification Results Without Upsampling**

582 **Cluster 2.** Due to the extreme class-imbalance of Clusters 1 and 2 (49 vs 16 time-series

583 respectively), the results presented in the main text are following upsampling of the Cluster 2

584 time-series to create a dataset with equal numbers of time-series belonging to each cluster.

585 As a sensitivity analysis, we also carried out the random forest fitting without upsampling and

586 assessed both model fit (as measured by AUC, **(A)**) and variable importance **(B)**. Model

587 performance was somewhat reduced compared to the upsampled data (mean AUC of 0.81 vs

588 mean AUC >0.9 for the upsampled dataset), whilst variable importance results were broadly

589 consistent across both analyses, with population per square kilometre and various land-cover

590 measures all emerging as important predictive variables. We also present partial dependence

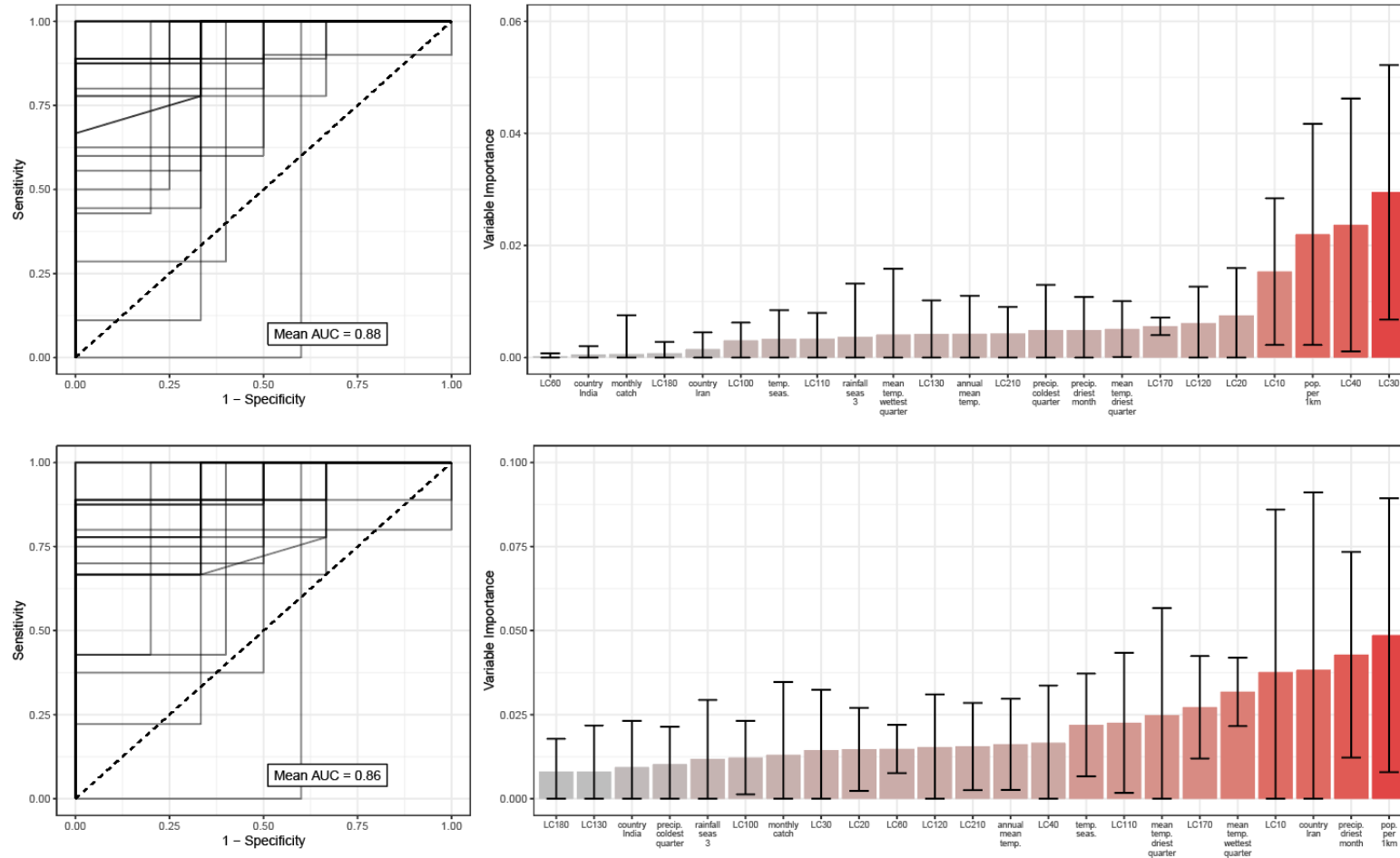
591 plots for all of the included covariates **(C)**. The y-axis on the left shows the probability of the

592 time-series belonging to Cluster 2 (i.e. a high probability indicates the time-series is predicted

593 to likely belong to Cluster 2, a low probability indicates the time-series likely belongs to Cluster

594 1). The x-axis describes the value of the (scaled, normalised) covariate.

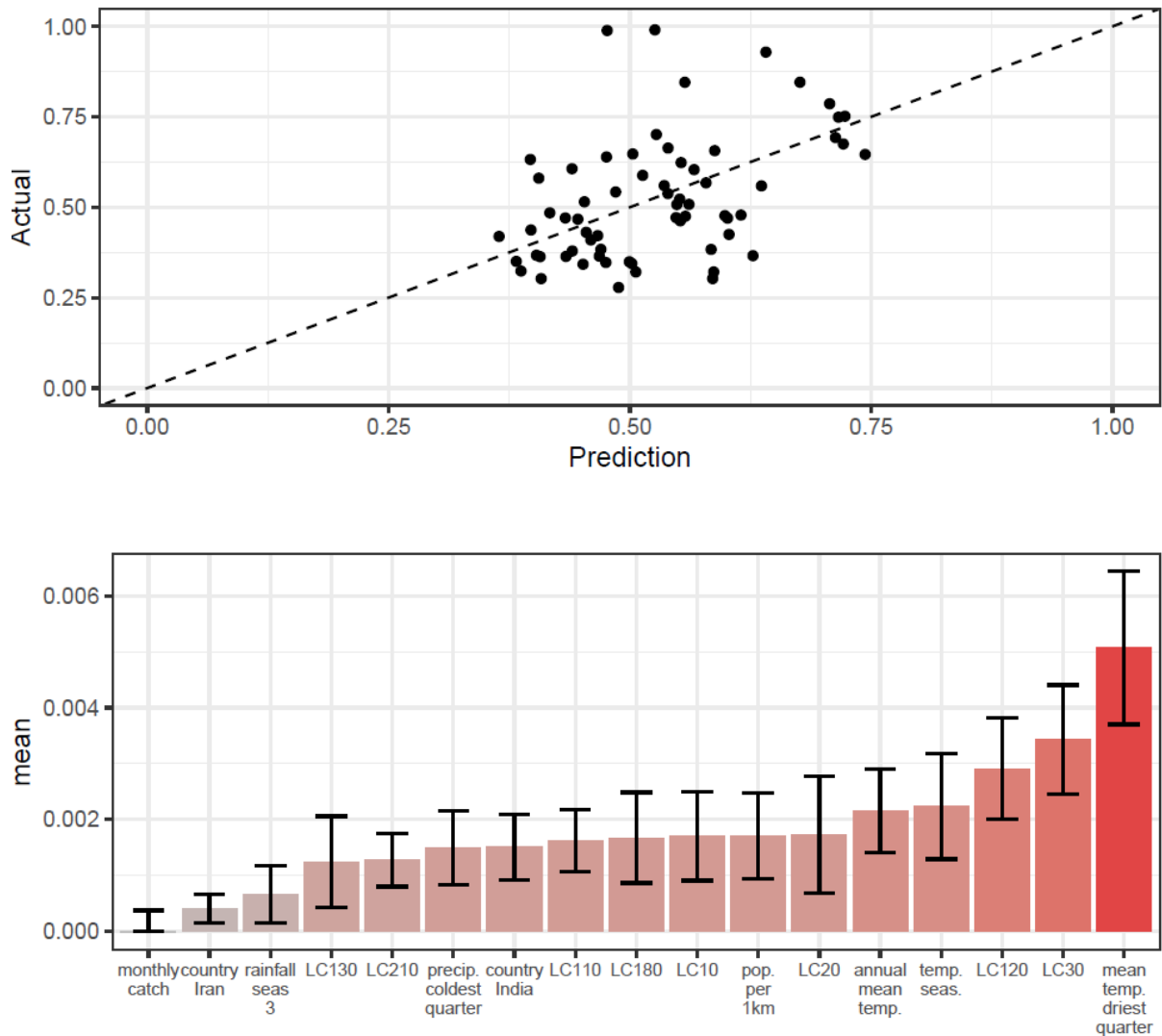
595



596

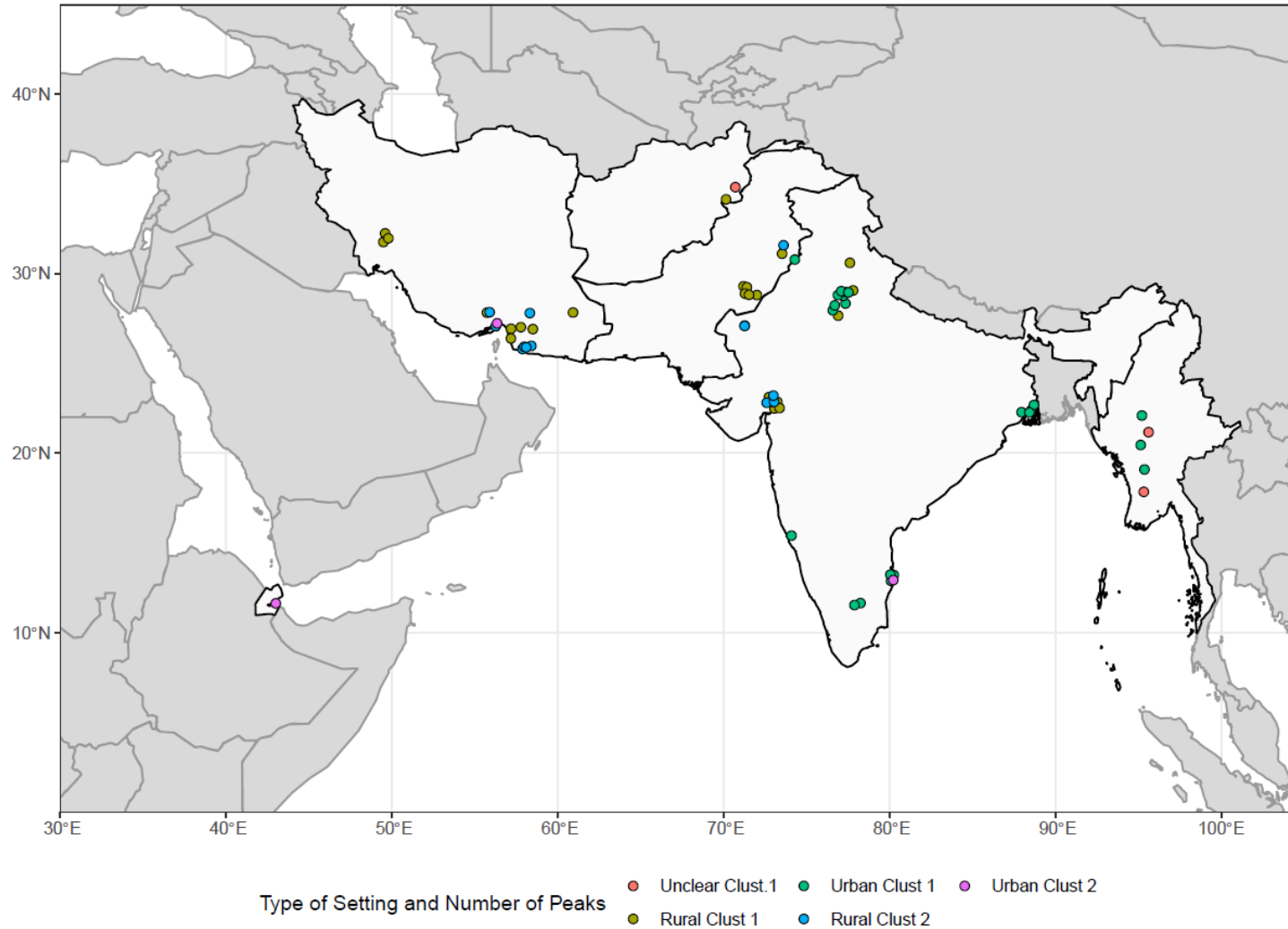
597 **Supplementary Figure 8: Random Forest Classification Results With Hold-Out Data.** Due to the overall sample size ($n = 65$ time-series),
 598 the results presented in the main text were generated using a random forest-based workflow where final model fitting (using hyperparameters
 599 tuned using 6-fold cross-validation) utilised the entirety of the dataset. As a sensitivity analysis, we also carried out the random forest fitting
 600 holding out a small portion of the dataset ($n = 9$) during model fitting, with model performance subsequently evaluated on this held-out data.
 601 Results presented above are in the case where data was upsampled to address class imbalance (top) and where no upsampling was carried out
 602 (bottom).

603
604
605
606
607



608
609
610
611
612
613
614
615

Supplementary Figure 9: Random Forest Prediction of Percentage of Vector Density In Any 3 Month Period. As a further sensitivity analysis, we used a random forest modelling approach to predict the percentage of vector density occurring in a single continuous 3-month period. Results presented above are the average of 25 independent random forest model fittings, with no upsampling of the data carried out, and the final model fitted (using hyperparameters tuned using 6-fold cross-validation) to the full dataset. Model predictive power was moderate, with correlation between predicted and actual values = 0.43.



616

617 **Supplementary Figure 10: Sources and Locations of *Anopheles stephensi* Time-Series Data According to Urban/Rural Assignment.**

618 Collated time-series are displayed above coloured according to 1) whether or not the study was carried out in an urban or rural location; and 2)

619 which cluster they were assigned to.

620 **References**

- 621 1. Sinka, M. E. *et al.* A new malaria vector in Africa: Predicting the expansion range of *Anopheles*
622 *stephensi* and identifying the urban populations at risk. *Proc. Natl. Acad. Sci. U. S. A.* **117**,
623 24900–24908 (2020).
- 624 2. Whittaker, C. *et al.* The ecological structure of mosquito population seasonal dynamics. *bioRxiv*
625 (2021) doi:10.1101/2021.01.09.21249456.
- 626 3. Hamlet, A. *et al.* The potential impact of *Anopheles stephensi* establishment on the transmission
627 of *Plasmodium falciparum* in Ethiopia and prospective control measures. *BMC Med.* **20**, 135
628 (2022).
- 629 4. Fulcher, B. D., Little, M. A. & Jones, N. S. Highly comparative time-series analysis: the empirical
630 structure of time series and their methods. *J. R. Soc. Interface* **10**, 20130048 (2013).
- 631 5. Fick, S. E. & Hijmans, R. J. WorldClim 2: new 1-km spatial resolution climate surfaces for global
632 land areas. *Int. J. Climatol.* **37**, 4302–4315 (2017).
- 633 6. European Space Agency (ESA). ESA CCI Land Cover time-series v2.0.7 (1992 - 2015).
- 634 7. Gaughan, A. E., Stevens, F. R., Linard, C., Jia, P. & Tatem, A. J. High resolution population
635 distribution maps for Southeast Asia in 2010 and 2015. *PLoS One* **8**, e55882 (2013).
- 636 8. Linard, C., Gilbert, M., Snow, R. W., Noor, A. M. & Tatem, A. J. Population distribution,
637 settlement patterns and accessibility across Africa in 2010. *PLoS One* **7**, e31743 (2012).
- 638 9. Justice, C. O. *et al.* An overview of MODIS Land data processing and product status. *Remote*
639 *Sens. Environ.* **83**, 3–15 (2002).
- 640 10. Didan, K. MOD13Q1 MODIS/Terra Vegetation Indices 16-Day L3 Global 250m SIN Grid V006
641 NASA EOSDIS Land Processes DAAC. doi:10.5067/MODIS/MOD13Q1.006.
- 642 11. Kuhn, M. & Wickham, H. Tidymodels: a collection of packages for modeling and machine
643 learning using tidyverse principles. *Boston, MA, USA*. [(accessed on 10 December 2020)]
644 (2020).
- 645 12. Funk, C. *et al.* The climate hazards infrared precipitation with stations—a new environmental
646 record for monitoring extremes. *Scientific Data* **2**, 1–21 (2015).
- 647 13. Rasmussen, C. E. Gaussian processes in machine learning. in *Advanced Lectures on Machine*
648 *Learning* 63–71 (Springer Berlin Heidelberg, 2004).
- 649 14. Carpenter, B. *et al.* Stan: A probabilistic programming language. *J. Stat. Softw.* **76**, 1–32 (2017).
- 650 15. Cairns, M. *et al.* Estimating the potential public health impact of seasonal malaria
651 chemoprevention in African children. *Nat. Commun.* **3**, 881 (2012).
- 652 16. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
- 653 17. Biau, G. Analysis of a random forests model. *J. Mach. Learn. Res.* **13**, 1063–1095 (2012).
- 654 18. Wright, M. N. & Ziegler, A. ranger: A Fast Implementation of Random Forests for High
655 Dimensional Data in C++ and R. *arXiv [stat.ML]* (2015).
- 656 19. Chawla, N. V., Bowyer, K. W., Hall, L. O. & Kegelmeyer, W. P. SMOTE: Synthetic Minority Over-
657 sampling Technique. *J. Artif. Intell. Res.* **16**, 321–357 (2002).
- 658 20. Molnar, C. *Interpretable Machine Learning*. (Lulu.com, 2020).
- 659 21. Griffin, J. T., Ferguson, N. M. & Ghani, A. C. Estimates of the changing age-burden of
660 *Plasmodium falciparum* malaria disease in sub-Saharan Africa. *Nat. Commun.* **5**, 3136 (2014).
- 661 22. Challenger, J. D. *et al.* Predicting the public health impact of a malaria transmission-blocking
662 vaccine. *Nat. Commun.* **12**, 1494 (2021).
- 663 23. Griffin, J. T. *et al.* Reducing *Plasmodium falciparum* malaria transmission in Africa: a model-
664 based evaluation of intervention strategies. *PLoS Med.* **7**, e1000324 (2010).
- 665 24. White, M. T. *et al.* Modelling the impact of vector control interventions on *Anopheles gambiae*
666 population dynamics. *Parasit. Vectors* **4**, 153 (2011).
- 667 25. Sherrard-Smith, E. *et al.* Systematic review of indoor residual spray efficacy and effectiveness
668 against *Plasmodium falciparum* in Africa. *Nat. Commun.* **9**, 4982 (2018).