

1
2
3
4 **Flying In-formation: A computational method for the classification of host**
5 **seeking mosquito flight patterns using path segmentation and**
6 **unsupervised machine learning**
7
8

9 **Mark T Fowler^{1*}, Anthony J Abbott¹, Gregory PD Murray^{1,2}, Philip J McCall^{1*}**

10

11

12 ¹ Vector Biology Department, Liverpool School of Tropical Medicine, Liverpool, L3 5QA, UK.

13 ² Centre de Recherches Interdisciplinaires, U1284 INSERM, Université de Paris, Paris, France.

14

15

16 *Address for correspondence:

17 Mathematical modelling: Mark.Fowler@lstmed.ac.uk

18 Vector biology: Philip.McCall@lstmed.ac.uk

19

20

21

22 **Keywords:** *Anopheles* mosquito, bed nets, vector behaviour, Behavioural Change Point Analysis, machine
23 learning, unsupervised clustering.

24 Abstract

25 The rational design of effective vector control tools requires detailed knowledge of vector behaviour. Yet,
26 behavioural observations, interpretations, evaluations and definitions by even the most experienced
27 researcher are constrained by subjectivity and perceptual limits. Seeking an objective alternative to
28 'expertise', we developed and tested an unsupervised method for the automatic identification of video-
29 tracked mosquito flight behaviour. This method unites path-segmentation and unsupervised machine
30 learning in an innovative workflow and is implemented using a combination of R and python. The workflow
31 (1) records movement trajectories; (2) applies path-segmentation; (3) clusters path segments using
32 unsupervised learning; and (4) interprets results. Analysis of the flight patterns of *An. gambiae* s.s.,
33 responding to human-baited insecticide-treated bednets (ITNs), by the new method identified four distinct
34 behaviour modes: with 'swooping' and 'approaching' modes predominant at ITNs; increased 'walking'
35 behaviours at untreated nets; similar rates of 'reacting' at both nets; and higher overall activity at treated
36 nets. The method's validity was tested by comparing these findings with those from a similar setting using
37 an expertise-based method. The level of correspondence found between the studies validated the accuracy
38 of the new method. While researcher-defined behaviours are inherently subjective, and prone to corollary
39 shortcomings, the new approach's mathematical method is objective, automatic, repeatable and a validated
40 alternative for analysing complex vector behaviour. This method provides a novel and adaptable analytical
41 tool and is freely available to vector biologists, ethologists and behavioural ecologists.

42 Author summary

43 Vector control targets the insects and arachnids that transmit 1 in every 6 communicable diseases worldwide.
44 Since the effectiveness of many vector control tools depends on exploiting or changing vector behaviour, a
45 firm understanding of this behaviour is required to maximise the impact of existing tools and design new
46 interventions. However, current methods for identifying such behaviours are based primarily on expert
47 knowledge, which can be inefficient, difficult to scale and limited by perceptual abilities. To overcome this,
48 we present, detail and validate a new method for categorising vector behaviour. This method combines
49 existing path segmentation and unsupervised machine learning algorithms to identify changes in vector

50 movement trajectories and classify behaviours. The accuracy of the new method is demonstrated by
51 replicating existing, expert-derived, findings covering the behaviour of host-seeking mosquitos around
52 insecticide treated bednets, compared to nets without insecticide. As the method found the same changes
53 in mosquito activity as previous research, it is said to be validated. The new method is significant, as it
54 improves the analytical capabilities of biologists working to reduce the burden of vector-borne diseases, such
55 as malaria, through an understanding of behaviour.

56 Introduction

57 Vector-borne diseases (VBDs) are illnesses caused by Protozoa, viruses and nematodes and transmitted by
58 infected arthropods, such as mosquitoes and ticks. VBDs threaten 80% of the planet's population, and are
59 responsible for an estimated 17% of all human communicable diseases and over 700 000 deaths annually [1–
60 3]. Many effective strategies to reduce the burden of VBDs target the arthropod vector. Such an approach
61 involves the development and use of interventions that control or exploit vector behaviour and prevent
62 human contact with pathogens. For example, tools that exploit a vector's host-seeking behaviour include
63 decoys or targets for *Glossina sp.* (tsetse fly, vectors of human animal trypanosomiasis) [4,5] and insecticide-
64 treated bednets (ITNs) for *Anopheles sp.* (the mosquitoes that transmit malaria) and *Aedes sp.* (the principal
65 vector for dengue fever) [6]. Significantly, although these devices are now essential tools for their respective
66 disease control or elimination programmes in sub-Saharan Africa [7], both continue to undergo further
67 research to improve their performance and applicability [8,9]. For example, efforts to improve ITNs have
68 entailed analysis of mosquito net responses through the segregation of flight paths around a human-baited
69 bednet into distinct movements and behaviours. These behaviours were based on flight characteristics
70 detected and defined by the researchers and interpreted as responses to the human, the net itself and/or
71 the presence of any insecticide treatment on the net [9–11]. However, investigations into distinguishing,
72 defining and classifying vector behaviour in such contexts are still principally based on researchers' expertise
73 and experience with the target species' biology and ethology [9,12–14]. Nevertheless, reliance on such a
74 solely subjective method is problematic. Expert knowledge is intrinsically inefficient to apply at scale, it is

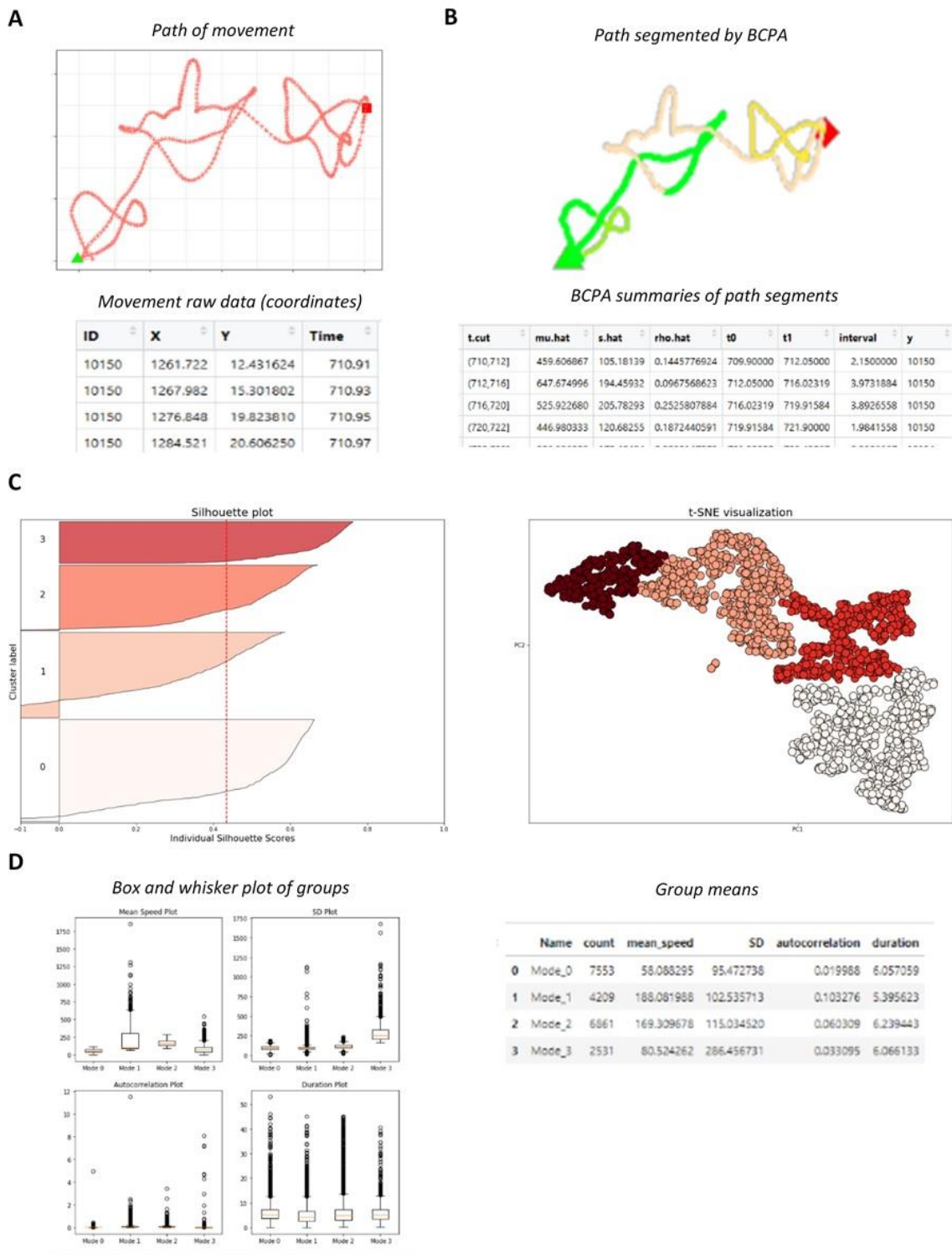
75 domain-specific, subject to cognitive biases and constrained by the physical limits of human perception [15–
76 17].

77 Much of this subjectivity can be eliminated by the application of objective computational processes. These
78 processes have the potential to analyse the movement paths of arthropod vectors to isolate and define
79 behaviours, but do so in an automated, repeatable and objective way. Computational algorithms that can
80 investigate animal movement and behaviour in this manner are already available. More specifically,
81 Behavioural Change Point Analysis (BCPA) is a form of path segmentation that splits movements into distinct
82 behavioural ‘bursts’ at significant changes in activity [18–20], thereby isolating movements. For example,
83 BCPA has been used to identify the timing of animal movements [21,22] and to quantify animal behavioural
84 shifts from seasonal environmental changes [23,24]. Secondly, unsupervised machine learning is a statistical
85 approach that can identify hidden patterns present within datasets and has the potential to cluster, and
86 therefore define, movements. These clustering algorithms group datapoints into distinct collections based
87 on any latent structure present within data [25,26] and have been used to identify patterns in neuronal
88 ensembles in the brain [27] and to identify subgroups within patient populations [28]. However, BCPA and
89 unsupervised machine learning are yet to successfully classify insect behaviour from movement trajectories
90 alone. Their application has been restricted in this context, as path segmentation can only identify changes
91 in behaviour rather than behaviours themselves [18], while clustering requires a sufficiently high signal-to-
92 noise ratio to be successful (something raw movement trajectories do not possess) [26].

93 This study proposes and tests that a solution to the problem of a total subjective base for the classification
94 of vector behaviour is possible by combining the above two identified computational processes. That is, path
95 segmentation and unsupervised machine learning can be brought together to discriminate and categorise
96 vector movements into distinct behavioural modes. However, the combination and application of these
97 algorithms requires a workflow to collect, prepare and analyse trajectories. In this report, we present,
98 describe and test such a workflow. This is a novel method that was devised to support complex behavioural
99 analyses, specifically concerning resource location by mosquitoes, in which: (1) detailed spatial and time-
100 series data covering the movement trajectory of a vector in a domain-specific setting is collected [10,29,30]

101 (Fig 1A); (2) movement trajectories are segmented into behavioural ‘bursts’ through BCPA [19,20] (Fig 1B);
102 (3) these behavioural ‘bursts’ are grouped through an optimised clustering algorithm [31,32] (Fig 1C); and
103 (4) results are interpreted through the analysis of descriptive statistics and examination of representative
104 samples [33–35] (Fig 1D).

105 Combining path segmentation and unsupervised machine learning into a single unique workflow provides a
106 novel method to overcome the inherent subjectivity and perceptual limitations of any investigator-led
107 alternative. In a first application of the method, we analysed the flight paths of the primary African malaria
108 vector mosquito, *An. gambiae* s.s., during host location around an ITN with a single human occupant,
109 recorded under experimental conditions in the laboratory. We report that the new workflow distinguished
110 four behavioural types that varied in frequency depending on net treatment. These findings corresponded
111 well with those in a previous investigator-led interpretation [9], but were achieved in a more objective,
112 repeatable manner.



113
114
115
116
117
118
119
120
121
122
123
124

Fig 1. Example of workflow process. (A) Time series data detailing a single vector movement trajectory. For each observation that comprise the movement an identifier, an x-coordinate, a y-coordinate and a time are required. The triangle is the start of the movement, the square the end of the movement. (B) BCPA is used to segment the movement into distinct behaviours, here based on significant changes in persistence velocity. Three significant changes in persistence velocity are identified in this example, giving four tokens of behaviour. BCPA segmentation produces a data frame summarising each phase. (C) The movement segments are grouped using the optimum clustering algorithm and initial parameters, as defined by internal validation. Clustering is internally validated through silhouette score, silhouette plot and manual inspection of a t-SNE visualisation. (D) A label is attached to behavioural groups by interpreting the clustering results. Interpretation is systemised through analysis of group statistics and examination of representative examples from each cluster (i.e., those found at the centre of each groups' t-SNE plot).

125 Results

126 To assess the accuracy of the new workflow, we applied the method to the activity of *An. gambiae*, a principal
 127 vector of malaria in sub-Saharan Africa, around either an insecticide-treated net (as approved by the World
 128 Health Organisation, hereafter ‘treated’) or an untreated polyester net (‘untreated’). A strain of mosquito
 129 susceptible to all insecticides, Kisumu, was used in both the untreated and treated arms of the experiment.
 130 The results of this application were then compared with those from a previous, expert-derived, study to
 131 validate the accuracy of the workflow.

132 Data acquisition, cleaning and assessment

133 Activity rates, based on observations from the raw data, were found to be much higher around an untreated
 134 net, with the number and length of movements significantly lower when an ITN was used (Table 1). When
 135 the autocorrelation of the datasets was assessed, it was found that movement velocity was positively
 136 autocorrelated through 50 time-lags in both the untreated (Fig 2A) and treated (Fig 2B) data. Accordingly,
 137 the data were taken to be suitable for analysis.

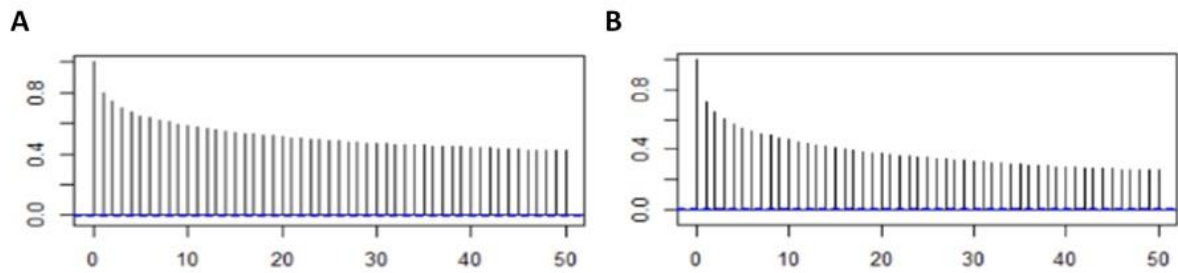
138

		Untreated	Treated
Tracking - Raw	Replicates	5	5
	Total Length	10 hrs	10 hrs
	Strain	Kisumu	Kisumu
	Observations	3 514 999	491 708
	Trajectories	9 076	1 472
	Max trajectory length	215.3 secs	124.0 secs
	Mean trajectory length	7.4 secs	6.7 secs
	Mode trajectory length	0.5 secs	0.5 secs
Tracking - Cleaned	Observations	3 320 055	453 201
	Trajectories	5 475	798
	Max trajectory length	215.3 secs	124.0 secs
	Mean trajectory length	12.1 secs	11.4 secs
	Mode trajectory length	2.2 secs	2.1 secs
Segmentation (BCPA)	Total phases	33 350	3 979
	Max trajectory phases	61	40
	Mean trajectory phases	6	5
	Mode trajectory phases	2	2

139

140 **Table 1. Untreated and Treated tracking and segmentation figures for evaluation.** All BCPA set to detect a significant
 141 change in persistence velocity (‘Velocity*cos(Turning Angle)’), using a window size of 30, a window step of 1, a sensitivity
 142 value of 2 and a cluster width of 1.

143



144
 145 **Fig 2. Correlograms.** (A) Untreated velocity correlogram. Movement speed is autocorrelated with its recent past (through
 146 a maximum of 50 lags). This association becomes weaker as the lag increases (from 0.8 at lag 1 to 0.5 at lag 50). (B)
 147 Treated velocity correlogram. Movement speed is autocorrelated with its recent past (through a maximum of 50 lags).
 148 This association becomes weaker as the lag increases (from 0.7 at lag 1 to 0.3 at lag 50).

149 **Path segmentation**

150 Results of the path segmentation are found in Table 1. Activity of *An. gambiae* s.s. was again found to be
 151 significantly higher in the untreated trial.

152 **Clustering**

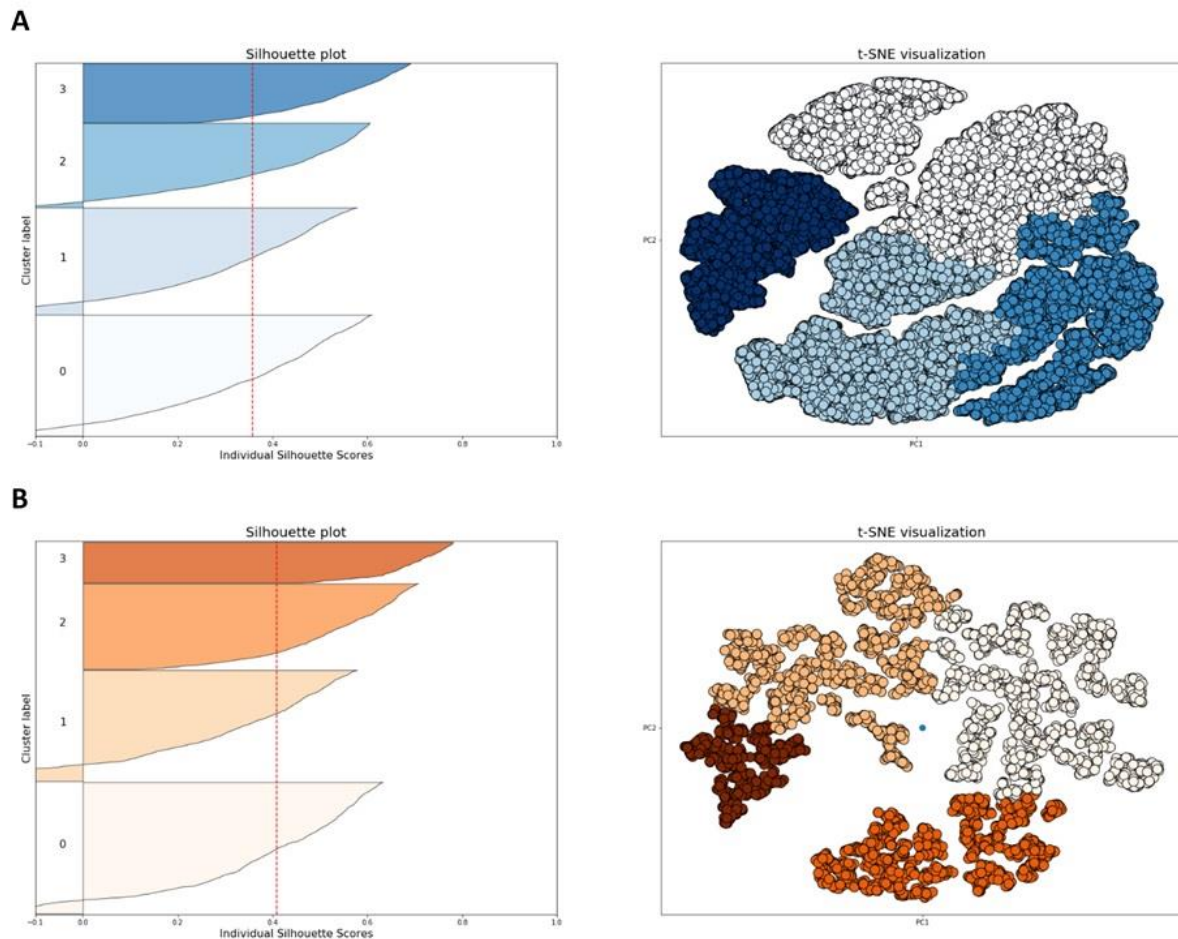
153 Internal validation indicated that the optimal algorithm and parameters to cluster both the untreated and
 154 treated data was an agglomerative clustering algorithm using Ward's method for linkage and four clusters.
 155 The untreated clustering produced a silhouette score of 0.36 (Fig 3A), while the treated grouping's silhouette
 156 score was 0.41 (Fig 3B). Results of this analysis are shown in Table 2.

157

Label	Untreated				Treated			
	Swoop	Approach	React	Walk	Swoop	Approach	React	Walk
Count	5 332	7 572	10 850	9 596	924	1 198	1 419	438
Count PCT (%)	15.99	22.70	32.53	28.77	23.22	30.11	35.66	11.01
Duration (s)	24 625	46 918	63 280	56 534	33 326	8 259	8 287	3 241
Duration PCT (%)	12.87	24.52	33.07	29.54	14.39	35.73	35.85	14.02
Mean Speed (mm/s)	301.06	92.80	122.49	51.42	369.17	90.23	182.32	33.59
Mean SD (±)	131.40	93.46	181.82	93.83	139.05	102.54	107.16	63.93
Mean Autocorrelation	0.17	0.03	0.04	0.02	0.27	0.03	0.08	0.02
Mean Duration (s)	4.62	6.20	5.83	5.89	3.60	6.89	5.84	7.40

158

159 **Table 2. Cluster summaries and interpretation labels for evaluation.** 'Count' is the total number of discrete phases in
 160 each cluster; 'Count PCT' is percentage of phases; 'Total Duration' is the total time, in seconds, for each group. 'Duration
 161 PCT' is percentage of duration. 'Mean Speed' and 'Mean SD' are group means given in mm/s. 'Mean Autocorrelation' and
 162 'Mean Duration' are group means. An autocorrelation of 1.0 represents a perfect correlation and 0.0 represents no
 163 correlation.



164
165
166 **Fig 3. Internal validation and clustering.** A) Untreated silhouette score and t-SNE plot. 33 350 datapoints on a t-SNE plot
167 with a perplexity of 75 in four clusters using agglomerative clustering with Ward's linkage, giving a silhouette score of
168 0.36. (B) Treated silhouette score and t-SNE plot. 3 979 datapoints on a t-SNE plot with a perplexity of 25 in four clusters
169 using agglomerative clustering with Ward's linkage, giving a silhouette score of 0.41.

170

171 Interpreting results

172 Mean statistics of each group were investigated to interpret the results (Table 2). Similarly broad behavioural
173 types were found in both arms of the study. After interpretation, these groups were labelled 'swooping',
174 'approaching', 'reacting' and 'walking.' 'Swooping' captures fast, short and highly autocorrelated
175 movements; 'approaching' slower, less variable behaviour with low autocorrelation; 'reacting' faster, more
176 variable actions with some autocorrelation; and 'walking' encompasses long, slow movements that are not
177 autocorrelated. As only information about vector movements is used in classification, environmental
178 interactions (e.g., net contact) cannot be included in the definition of behaviours. The labels, and the broad

179 nature of each grouping, were then confirmed through investigation of representative samples from each
180 group.

181 **Conclusions**

182 Seven principal conclusions can be drawn from this analysis into *An. gambiae* s.s. activity around untreated
183 and treated nets: (1) in any fixed time period, mosquito flight activity is significantly greater when the human
184 host is protected within an untreated net compared to a treated net (two sample Z-test, $P < 0.01$); (2) four
185 behavioural modes are exhibited around both treated and untreated nets; (3) The proportion of ‘swooping’
186 behaviour increases significantly around a treated net (two sample Z-test, $P < 0.01$); (4) The proportion of
187 ‘approaching’ increases significantly in the presence of a treated net (two sample Z-test, $P < 0.01$); (5) The
188 proportion of ‘reacting’ increases significantly around a treated net (two sample Z-test, $P < 0.01$); (6) The
189 proportion of ‘walking’ decreases significantly around a treated net (two sample Z-test, $P < 0.01$); and (7) *An.*
190 *gambiae* s.s. ‘swoop’ faster in experiments with a treated net (GLM, $P < 0.001$).

191 **External Validation**

192 A similar study of the effect of bednet treatment on vector behaviour had previously been conducted in
193 Tanzania using wild *An. arabiensis*, a sibling species closely related to *Anopheles gambiae* s.s. and that
194 exhibits many of the same host seeking behaviour characteristics [9]. This previous study used expert
195 knowledge to identify behaviour types, determining that mosquitos exhibited four behaviours around both
196 an untreated and treated net (‘swooping’, ‘visiting’, ‘bouncing’ and ‘resting’) and that total activity levels
197 dropped significantly at ITNs compared to untreated bednet (from a geometric mean time of 73.5 mins to
198 23.8 mins). Where particular behaviours are concerned, and comparing total mean times, the study found
199 that ‘swooping’ (where “tracks do not contact the bednet”), ‘visiting’ (“long periods of flight are interspersed
200 with infrequent net contacts’) and ‘resting’ (“mosquito movement is under 1.33 mm /s”) all increased in the
201 presence of a treated net (however, this increase in swooping was not found to be statistically significant).
202 The study also found that ‘bouncing’ (“rapid contacts with the bednet surface... include[ing] walking”)
203 reduced significantly around the treated net (when evaluating geometric mean times).

204 Comparing findings from this study and those of [9], several replications are clear: (1) both studies recognised
205 four types of mosquito behaviour; (2) total vector activity fell at a treated bednet; (3) ‘swooping’ behaviour
206 increased with a treated net; and (4) ‘walking’ / ‘bouncing’ is decreased when using a treated net. However,
207 although ‘swooping’ and ‘bouncing’ from [9] are acceptable analogues to the behaviours ‘swooping’ and
208 ‘walking’ from this study, it is not possible to align the prior study’s ‘resting’ and ‘visiting’ with the
209 ‘approaching’ or ‘reacting’ categories of this study, due to divergent definitions.

210 As such, comparison with previous results can only be said to validate four of this study’s conclusions (i.e.,
211 (1), (2), (3) and (6) from the [Conclusions](#) section). Although there are slight differences between [9] and the
212 current study (i.e., in vectors observed and the definition of behaviour modes), these differences are minor,
213 potentially explicable by the use of a wild population, which is inherently more genetically diverse. With this
214 knowledge, the similarities are such that [9] can be said to support several major findings from this study in
215 the given setting. Consequently, the external validity of the new method was deemed to be proven.

216 Discussion

217 In this study, we present an automated, generalised method for the identification and classification of the
218 behaviour of vectors based on their movement trajectories. This new workflow combines BCPA [19,20] and
219 unsupervised machine learning [31,32] and offers a new solution to current challenges faced by vector
220 biologists and for vector control [1–3]. Although a similar methodology has been proposed for the
221 investigation of marine animal behaviour [36], to our knowledge this is the first use of such an approach
222 within entomology. The method has particular relevance in vector biology, where an automated, repeatable
223 and generalisable means of identifying and defining behaviour that has been validated against vector activity
224 is most pertinent.

225 Here we supply a preliminary application of the new method, analysing the behaviour of *An. gambiae* s.s., in
226 the presence of both baited untreated bednets and baited ITNs. As the study replicated previous findings,
227 the method is deemed to be an innovative, validated and productive approach that improves and expands
228 the existing toolkit available to vector biologists. Furthermore, the method is repeatable, as any individual

229 with the same dataset will produce the same behavioural tokens and behavioural types; it is generalisable,
230 as it is not limited to a single domain, but can be applied to any vector in any setting; and its foundation in
231 mathematical processes ensures it is immune to observer bias. The accuracy of the new method is confirmed
232 using both internal and external validation [37,38]. The former ensures the correct algorithm and initial
233 parameters are applied, while the latter tests the accuracy of the approach itself. Internal validity is measured
234 in two ways: (1) formal metrics of similarity between datapoints (i.e., silhouette scores) are studied; and (2)
235 t-SNE visualisations of cluster assignment are manually inspected [39–41]. External validation is achieved by
236 replicating known results [42]. The new method is used to compare the behaviour of an insecticide
237 susceptible strain of *An. gambiae* s.s. (Kisumu) around both Long-lasting Insecticidal Nets (ITNs) and
238 untreated nets, producing findings that are corroborated by previous research [9].

239 The new method offers advantages over alternative, objective approaches that are theoretically automated,
240 repeatable and generalisable. One such method, Hidden Markov Models (HMMs) are probabilistic models
241 that determine the underlying hidden states (e.g., behavioural modes) that cause an observed process (e.g.,
242 movement trajectories) [43]. However, for HMMs to apply in this instance, a vector's behavioural states must
243 be a first-order Markov process. That is, a vector's behaviour at time t must be determined solely by their
244 behavioural state at time $t-1$ [43–45]. Nevertheless, it is reasonable to assert that vector behaviour is
245 influenced by internal and external drivers acting over greater periods of time than this and that vector
246 ecology is determined by a wider range of datapoints that cannot be described by a first-order Markov
247 process and an HMM [46]. To capture this more nuanced conception of vector behaviour, a sliding window,
248 such as is applied in BCPA, is needed. Similarly, although several path segmentation methods exist for
249 detecting changes in animal movements other than BCPA [47–50], a form of segmentation that can account
250 for the particular difficulties encountered when tracking vectors must be used in this instance. As the key
251 difficulty here is the frequency of lost frames (caused by the recording system momentarily losing track of
252 the small vector), a method that can handle an irregular dataset is required. As BCPA is a likelihood-based
253 form of path segmentation, which sweeps an analysis window over an entire movement path to identify
254 significant shifts in a parameter value, it provides a robust method for dissecting vector activity into

255 behavioural tokens that can account for irregular temporal measurement intervals and does so without any
256 *a priori* assumptions [19,20].

257 Although offering several advancements, the method presented here is subject to its own limitations. One
258 such constraint concerns the clarity of the silhouette created by any grouping of behaviours. That is, as
259 behavioural units are nebulous concepts, any silhouette of their classification will be equally unclear and
260 datapoints from different behaviours will not necessarily have a high separateness [51–53]. For example, the
261 distinction between fast walking and slow running is not clear. Consequently, the identification of strong
262 patterns when assessing the clustering of behaviour is unlikely. This is shown in the contiguous silhouettes
263 and the silhouette values produced by movement data (Fig 3). Additionally, it is important to make explicit
264 the assumptions on which this study is based. These assumptions are that vectors are always in some
265 behavioural state, that vectors have more than one potential behavioural mode and that these modes are
266 discrete and expressed over a period of time. Finally, it needs to be clarified that the method presented here
267 is not totally objective. Since the workflow’s Interpretation stage requires experts to attach a label to clusters,
268 a level of subjectivity is still required to implement this analysis. Although this labelling is not theoretically
269 necessary to produce and compare results (as clusters can be described by their characteristics alone), a level
270 of subjectivity is still needed to interpret these results and apply them to everyday discourse concerning
271 behaviour [51–53].

272 In conclusion, we present and test a new workflow that represents an innovative use of path segmentation
273 and unsupervised machine learning to classify vector behaviour and expands the analytic toolkit available to
274 researchers. This represents a promising development that can improve the evidence base available to
275 vector biologists and open new avenues for the exploitation of vector behaviour to improve intervention
276 performance. Given that global vector control is currently facing a raft of challenges – including
277 environmental and species distribution changes [2], limited resources [3] and an increase in insecticide
278 resistance [54] – novel methodological approaches are more important than ever. Furthermore, it is likely
279 that developments can be made to improve performance and applicability. For example, an analysis of
280 transitions between behaviours could be undertaken, potentially providing additional insights into vector

281 activity and ensuring ecological limits to behavioural transitions have been captured. Finally, the output from
282 this workflow could be used as input to a supervised machine learning algorithm, increasing the efficiency of
283 future analyses.

284 **Materials and Methods**

285 We present a four-stage workflow in which vector movement trajectories are first collected and pre-
286 processed via BCPA. The most appropriate unsupervised clustering algorithm, and initial parameters, are
287 then identified and applied before the workflow concludes with the interpretation of results, decoding and
288 attaching a behavioural label to each group. The whole workflow is then validated by measuring the accuracy
289 of its results.

290 **Resources**

291 The workflow presented here is implemented using a combination of R and Python. R is used for pre-
292 processing, utilising the BCPA package built for that language. Python, through a Jupyter notebook, is used
293 at the clustering stage to exploit the `scikit-learn` library. We recommend that the Anaconda platform
294 be used to access RStudio and JupyterLab, as up-to-date installations for Windows, Linux and Mac can all be
295 found in that single distribution. Code, and further details, needed to run the workflow can be accessed
296 through a public GitLab repository: https://gitlab.com/MTFowler/lstm_flightcluster. All analysis found here
297 was performed on a ThinkPad X1 Carbon, using an Intel i7-7500u CPU.

298 All procedures associated with the collection of mosquito flight data are as described in [9,11,55,56]. Briefly,
299 the ‘Kisumu’ laboratory strain of *An gambiae*, a primary malaria vector across sub-Saharan Africa and
300 susceptible to all insecticides, was used in both the untreated and treated arms of the experiment. All
301 mosquito flight assays were completed in a purpose-built climate-controlled insectary in Liverpool.

302 **Data acquisition, cleaning and assessment**

303 Vector movement paths were represented by spatial identifiers ordered sequentially via a time variable [18–
304 20,30]. Each event was captured by a unique identifier, an x (longitude or easting) coordinate, a y (latitude
305 or northing) coordinate and a time variable (Fig 1A). This data was collected using an optical imaging and

306 flight-tracking system detailed in [55,56]. This system allowed for multiple vectors to move unconstrained
307 within an enclosed area, a subset of this space being within the field of view of the recording system, creating
308 the recording volume. After collection, movement trajectory data was cleaned and assessed (Table 1).
309 Movements considered noise were removed. This ‘noise’ included short tracks deemed to be isolated
310 fragments from a larger track, or disturbance that has been missed during video cleaning [9–11].
311 Furthermore, although BCPA accounts for semi-regular sampling [19,20], allowing for some irregularity in
312 the dataset, movement tracks were removed from the analysis if they contained two datapoints at the same
313 time or if they had especially large time gaps (i.e., greater than 10 seconds). Finally, as path segmentation
314 assumes that all time series data displays serial dependence, it was confirmed that the dataset was
315 autocorrelated (i.e., that the velocity of each datapoint is statistically correlated with its recent past) [20].
316 This was accomplished in R using the ‘Autocorrelation and Cross-Correlation Function Estimation’, $ACF()$.

317 **Path segmentation**

318 With a correctly formatted dataset, that had been cleaned and assessed, BCPA was applied. BCPA is a form
319 of path segmentation that identifies changes in animal behaviour, at the path-level, based on significant
320 shifts in a parameter value of an organism’s movement trajectory. As BCPA accepts movement paths as
321 sequentially ordered step lengths, turning angles and velocities, rather than the spatial identifiers collected
322 by tracking technology, spatial values were converted into the required variables using the $GetVT()$
323 function from R’s *BCPA* package [57]. Within BCPA there are four user defined parameters: (1) the
324 ‘Parameter Value’ (the response time-series variable in which significant changes will identify a behavioural
325 change point); (2) the ‘Window Size’ (the number of datapoints the window will capture when sweeping);
326 (3) a sensitivity parameter ‘K’; and (4) the ‘Cluster Width.’ For arthropod activity, it was determined that
327 optimal segmentation occurs at a significant change in persistence velocity ($Velocity * \cos(\text{Turning Angle})$),
328 using a window size of 30, a window step of 1, a sensitivity value of 2 and a cluster width of 1. These initial
329 parameters were determined following BCPA documentation recommendations [19,20,57] and to maximise
330 sensitivity to behavioural shifts. (Note, however, that this increase in sensitivity amplifies the chances of
331 spurious shifts being detected which will ultimately result in transitions to the same behaviour in the final

332 output. However, as the alternative is to lower sensitivity and potentially miss legitimate changes in
333 behaviour, a high sensitivity, with corollary spurious shifts, is preferred.)

334 **Clustering**

335 To determine the optimal unsupervised learning algorithm and initial parameters for clustering, internal
336 validation was undertaken. Following [37,38], the form of internal validation used was silhouette scores
337 [58,59] and visual inspection of t-SNE plots [39–41]. A silhouette score measures how well data had been
338 grouped, comparing each object’s similarity to others within its own cluster (group tightness) and those from
339 other clusters (group separation) and was calculated using Python’s `silhouette_score()` function from
340 the `metrics` module of the `scikit-learn` library. This measure gives a score between -1.00 and +1.00,
341 with a silhouette value below 0.20 showing no structure is present in the data and the grouping is invalid; a
342 figure over 0.70 representing a strong structure and a valid grouping; and a silhouette score around 0.50
343 illustrating that a reasonable structure has been found within the data and that the clustering is acceptable
344 [59]. Detailed silhouette coefficients for each sample was then visualised using a silhouette plot in Python
345 with the `silhouette_samples()` function from the `sklearn.metrics` module (Fig 1C). As all
346 clusterings require manual review to validate appropriateness [37], the high-dimensional data was reduced
347 and positioned in a two-dimension map using t-distributed Stochastic Neighbour Embedding (t-SNE) [39].
348 Once mapped, the appropriateness of the clustering was verified through manual visual inspection. Review
349 ensured there were acceptable levels of cohesion between members of the same group and separateness
350 between members of different groups. t-SNE was undertaken using the `TSNE()` function from the
351 `sklearn.manifold` module found within Python. When performing t-SNE, the user needs to define the
352 perplexity (an estimate number of nearest neighbours for each datapoint), with larger datasets requiring a
353 larger perplexity [41]. Consequently, perplexity was fine-tuned to show global geometry.

354 Once the optimum algorithm and parameters were determined by silhouette score and inspection of t-SNE
355 plot, findings were applied to the BCPA output. Python’s `scikit-learn` library was used as it is an efficient
356 means of building standard machine learning models [60]. (Other packages, such as R’s `class`, are available

357 when implementing unsupervised machine learning, however different packages may produce different
358 results.)

359 **Interpreting results**

360 The final stage of analysis is to interpret and label clusters [35]. Although the naming of individual clusters
361 can be systematised, its final interpretation is ultimately a somewhat subjective process. Interpretation here
362 entailed scrutiny of the attributes of each cluster. Taking each group's mean velocity, standard deviation,
363 autocorrelation and duration, an expert analysed and named the behaviour associated with the movement
364 classes. This initial understanding was then confirmed and refined through visual inspection of
365 representative samples. 'Representative' was defined as those samples found at the centre point of a group
366 and such datapoints were deemed to be the most typical of that behaviour class [33,34]. Consequently,
367 interpretation of results was bolstered through centroid analysis, with those datapoints at the heart of each
368 group's t-SNE mapping, or as close to the centre as possible, isolated and visually inspected by an expert.
369 Multiple such examples close to the centre were isolated and inspected, thereby confirming, or refining, the
370 initial analysis.

371 Using both the analysis of descriptive statistics and centroid analysis, an expert was able to interpret the
372 broad behavioural type of each cluster and attach a sensible label. If no label was able to be attached, either
373 because no behaviour is being demonstrated or behaviours are spread between clusters, the full workflow
374 was undertaken again. By manipulating the user defined settings during BCPA or dimensionality reduction,
375 significant changes in clustering can result. Consequently, these parameters were fine-tuned to optimise
376 performance and final behaviour identification.

377 **External Validation**

378 To ascertain the accuracy of the new method, its performance was externally validated by comparing results
379 concerning the difference in flight patterns of *An. gambiae* s.s. around human-baited insecticide-treated
380 bednets (ITNs) and untreated bednets. Findings generated by the new, computation-derived, method were
381 contrasted with those from a previous study that employed the existing, expert-defined, method for

382 behavioural classification. Although a standard method of external validation is through comparison against
383 an *a priori* dataset, testing whether a known, true, clustering can be recreated [38], this was not possible in
384 this instance. There is no known, incontrovertibly true, grouping for this *An. gambiae* s.s. behaviour and
385 therefore no such comparison can be made. Consequently, external validation was established via
386 confirmation with previous results. That is, the workflow's accuracy as a method was corroborated by
387 comparing its conclusions to those already present in the literature [9] to verify whether prior findings could
388 be replicated.

389

390 Acknowledgements

391 We are grateful to Elizabeth Bandason, Amy Guy and Josie Parker for advice on interpreting
392 mosquito behaviour and James Maas for advice on statistical methods.

393 Author Contributions

394 Conceptualisation, formal analysis, methodology and writing (draft preparation) by M.F. Data
395 curation and writing (review & editing) by A.A. and G.M. Software by V.V., C.T. and D.T. Supervision
396 and writing (review & editing) by P.M. All authors have read and agreed to the published version of
397 the manuscript.

398 Data availability statement

399 All data and code used for running experiments, model fitting, and plotting is available on a GitLab
400 repository at https://gitlab.com/MTFowler/lstm_flightcluster.

401 Funding

402 This research was funded in part by Medical Research Council of the UK (grant number
403 MR/P027873/1) through the Global Challenges Research Fund and the Bill & Melinda Gates
404 Foundation under Grant Agreement No OPP1200155. The findings and conclusions contained within
405 are those of the authors and do not necessarily reflect positions or policies of the Bill & Melinda
406 Gates Foundation.

407 References

- 408 1. Wilson AL, Courtenay O, Kelly-Hope LA, Takken W, Lindsay SW. The importance of vector control for
409 the control and elimination of vector-borne diseases. *PLoS Neglected Tropical Diseases*. 2020;14.
410 doi:<https://doi.org/10.1371/journal.pntd.0007831>
- 411 2. World Health Organization. *Global Vector Control Response 2017–2030*. Geneva, Switzerland: World
412 Health Organization; 2017.
- 413 3. Golding N, Wilson AL, Moyes CL, Cano J, Pigott DM, Velayudhan R, et al. Integrating vector control
414 across diseases. *BMC Medicine*. 2015;13. doi:<https://doi.org/10.1186/s12916-015-0491-4>
- 415 4. Shaw APW, Tirados I, Mangwiro CTN, Esterhuizen J, Lehane MJ, Torr SJ, et al. Costs of using “tiny
416 targets” to control *Glossina fuscipes fuscipes*, a vector of gambiense sleeping sickness in Arua District
417 of Uganda. *PLoS Neglected Tropical Diseases*. 2015;9.
418 doi:<https://doi.org/10.1371/journal.pntd.0003624>
- 419 5. Tirados I, Hope A, Selby R, Mpenbele F, Miaka EM, Boelaert M, et al. Impact of tiny targets on *Glossina*
420 *fuscipes quanzensis*, the primary vector of human African trypanosomiasis in the Democratic Republic
421 of the Congo. *PLoS Neglected Tropical Diseases*. 2020;14.
422 doi:<https://doi.org/10.1371/journal.pntd.0008270>
- 423 6. Pryce J, Richardson M, Lengeler C. Insecticide-treated nets for preventing malaria. *Cochrane Database*
424 *of Systematic Reviews*. 2018;11. doi:<https://doi.org/10.1002/14651858.CD000363.pub3>
- 425 7. Bhatt S, Weiss DJ, Cameron E, Bisanzio D, Mappin B, Dalrymple U, et al. The effect of malaria control on
426 *Plasmodium falciparum* in Africa between 2000 and 2015. *Nature*. 2015;526: 207–211.
427 doi:<https://doi.org/10.1038/nature15535>
- 428 8. Torr SJ, Vale GA. Know your foe: Lessons from the analysis of tsetse fly behaviour. *Trends in*
429 *Parasitology*. 2015;31. doi:<https://doi.org/10.1016/j.pt.2014.12.010>
- 430 9. Parker JEA, Angarita-Jaimes NC, Gleave K, Mashauri F, Abe M, Martine J, et al. Host-seeking activity of a
431 Tanzanian population of *Anopheles arabiensis* at an insecticide treated bed net. *Malaria Journal*.
432 2017;16. doi:<https://doi.org/10.1186/s12936-017-1909-6>
- 433 10. Parker JEA, Angarita-Jaimes NC, Abe M, Towers CE, Towers DP, McCall PJ. Infrared video tracking of
434 *Anopheles gambiae* at insecticide-treated bed nets reveals rapid decisive impact after brief localised
435 net contact. *Scientific Reports*. 2015;5. doi:<https://doi.org/10.1038/srep13392>
- 436 11. Murray GPD, Lissenden N, Jones J, Voloshin V, Toé KH, Sherrard-Smith E, et al. Barrier bednets target
437 malaria vectors and expand the range of usable insecticides. *Nature Microbiology*. 2020;5: 40–47.
438 doi:<https://doi.org/10.1038/s41564-019-0607-2>
- 439 12. Gibson G, Torr SJ. Visual and olfactory responses of haematophagous Diptera to host stimuli. *Medical*
440 *and Veterinary Entomology*. 1999;13: 2–23. doi:<https://doi.org/10.1046/j.1365-2915.1999.00163.x>
- 441 13. Mathis A, Mamidanna P, Cury KM, Abe T, Murthy VN, Mathis MW, et al. DeepLabCut: Markerless pose
442 estimation of user-defined body parts with deep learning. *Nature Neuroscience*. 2018;21: 1281–1289.
443 doi:<https://doi.org/10.1038/s41593-018-0209-y>
- 444 14. Nath T, Mathis A, Chen AC, Patel A, Bethge M, Mathis MW. Using DeepLabCut for 3D markerless pose
445 estimation across species and behaviors. *Nature Protocols*. 2019;14: 2152–2176.
446 doi:<https://doi.org/10.1038/s41596-019-0176-0>

- 447 15. Kennedy JS. *The New Anthropomorphism*. Cambridge: Cambridge University Press; 1992.
- 448 16. Kuhnert PM, Martin TG, Griffiths SP. A guide to eliciting and using expert knowledge in Bayesian
449 ecological models. *Ecology Letters*. 2010;13. doi:<https://doi.org/10.1111/j.1461-0248.2010.01477.x>
- 450 17. Krueger T, Page T, Hubacek K, Smith L, Hiscock K. The role of expert opinion in environmental
451 modelling. *Environmental Modeling and Software*. 2012;36: 4–18.
452 doi:<https://doi.org/10.1016/j.envsoft.2012.01.011>
- 453 18. Edelhoff H, Signer J, Balkenhol N. Path segmentation for beginners: An overview of current methods for
454 detecting changes in animal movement patterns. *Movement Ecology*. 2016;4.
455 doi:<https://doi.org/10.1186/s40462-016-0086-5>
- 456 19. Gurarie E, Bracis c., Delgado M, Meckley TD, Kojola I, Wagner CM. What is the animal doing? Tools for
457 exploring behavioural structure in animal movements. *Journal of Animal Ecology*. 2016;85: 69–84.
458 doi:<https://doi.org/10.1111/1365-2656.12379>
- 459 20. Gurarie E, Andrews RD, Laidre KL. A novel method for identifying behavioural changes in animal
460 movement data. *Ecology Letters*. 2009;12: 395–408. doi:<https://doi.org/10.1111/j.1461-0248.2009.01293.x>
- 462 21. Chamberlain MJ, Cohen BS, Bakner NW, Collier BA. Behavior and movement of wild turkey broods. *The
463 Journal of Wildlife Management*. *The Journal of Wildlife Management*. 11582;84: 1139.
464 doi:<https://doi.org/10.1002/jwmg.21883>
- 465 22. Requier F, Henry M, Decourtye A, Brun F, Aupinel P, Bretagnolle V. Measuring ontogenetic shifts in
466 central-place foraging insects: a case study with honey bees. *Journal of Animal Ecology*. 2020;89: 1860–
467 1871. doi:<https://doi.org/10.1101/2020.03.31.017582>
- 468 23. Garstang M, Davis RE, Leggett K, Frauenfeld OW, Greco S, Zipser E, et al. Response of African elephants
469 (*loxodonta africana*) to seasonal changes in rainfall. *PLoS One*. 2014;9.
470 doi:<https://doi.org/10.1371/journal.pone.0108736>
- 471 24. Kidd-Weaver A, Hepinstall-Cymerman J, Welch CN, Murray MH, Adams HC, Ellison TJ, et al. The
472 movements of a recently urbanized wading bird reveal changes in season timing and length related to
473 resource use. *PLoS One*. 2020. doi:<https://doi.org/10.1371/journal.pone.0230158>
- 474 25. Tarca AL, Carey VJ, Chen XW, Romero R, Draghici S. Machine learning and its applications to biology.
475 *PLOS Computational Biology*. 3. doi:<https://doi.org/10.1371/journal.pcbi.0030116>
- 476 26. Valletta JJ, Torney C, Kings M, Thornton A, Madden J. Applications of machine learning in animal
477 behaviour studies. *Animal Behaviour*. 2017;124: 203–220.
478 doi:<https://doi.org/10.1016/j.anbehav.2016.12.005>
- 479 27. Grossberger L, Battaglia FP, Vinck M. Unsupervised clustering of temporal patterns in high-dimensional
480 neuronal ensembles using a novel dissimilarity measure. *PLoS Computational Biology*. 2018;14.
481 doi:<https://doi.org/10.1371/journal.pcbi.1006283>
- 482 28. Lopez C, Tucker S, Salameh S, Tucker C. An unsupervised machine learning method for discovering
483 patient clusters based on genetic signatures. *Journal of Biomedical Informatics*. 2018;85: 30–39.
484 doi:<https://doi.org/10.1016/j.jbi.2018.07.004>

- 485 29. Lynd A, McCall PJ. Clustering of host-seeking activity of *Anopheles gambiae* mosquitoes at the top
486 surface of a human-baited bed net. *Malaria Journal*. 2013;12. doi:[https://doi.org/10.1186/1475-2875-](https://doi.org/10.1186/1475-2875-12-267)
487 12-267
- 488 30. Getz WM, Saltz D. A framework for generating and analyzing movement paths on ecological
489 landscapes. *Proceedings of the National Academy of Science*. 2008;105: 19066–19071.
490 doi:<https://doi.org/10.1073/pnas.0801732105>
- 491 31. Murtagh F, Contreras P. Algorithms for hierarchical clustering: an overview. *Data Mining and*
492 *Knowledge Discovery*. 2012;2: 86–97. doi:<https://doi.org/10.1002/widm.1219>
- 493 32. Chidananda Gowda K, Krishna G. Agglomerative clustering using the concept of mutual nearest
494 neighbourhood. *Pattern Recognition*. 1978;10: 105–112. doi:[https://doi.org/10.1016/0031-](https://doi.org/10.1016/0031-3203(78)90018-3)
495 3203(78)90018-3
- 496 33. Alexiou A, Riddlesden D, Singleton A. The geography of online retail behaviour.) *Consumer data*
497 *research*. UCL Press; 2018. pp. 97–109.
- 498 34. Kang J, Ryu KR, Kwon HC. Using cluster-based sampling to select initial training set for active learning in
499 text classification. *Advances in Knowledge Discovery and Data Mining*. Berlin, Heidelberg: Springer;
500 2004. pp. 384–388.
- 501 35. Labrinidis A, Jagadish HV. Challenges and opportunities with big data. *Proceedings of the VLDB*
502 *Endowment*. 2012;5: 2032–2033. doi:<https://doi.org/10.14778/2367502.2367572>
- 503 36. Zhang J, O'Reilly KM, Perry GLW, Taylor GA, Dennis TE. Extending the functionality of behavioural
504 change-point analysis with k-means clustering: A case study with the Little Penguin (*Eudyptula minor*).
505 *PLoS ONE*. 2015. doi:<https://doi.org/10.1371/journal.pone.0122811>
- 506 37. Halkidi M, Batistakis Y, Vazirgiannis M. On clustering validation techniques. *Journal of Intelligent*
507 *Information Systems*. 2001;17: 107–145. doi:<https://doi.org/10.1023/A:1012801612483>
- 508 38. Arbelaitz O, Gurrutxaga I, Muguerza J, Pérez JM, Perona I. An extensive comparative study of cluster
509 validity indices. *Pattern Recognition*. 2013;46: 243–256.
510 doi:<https://doi.org/10.1016/j.patcog.2012.07.021>
- 511 39. Van der Maaten L, Hinton G. Visualizing data using t-SNE. *Journal of Machine Learning Research*.
512 2008;9: 2579–2605.
- 513 40. Hinton GE, Roweis ST. Stochastic neighbor embedding. *Advances in neural information processing*
514 *systems*. Cambridge, MA, USA: The MIT Press; 2002. pp. 833–840.
- 515 41. Wattenberg M, Viégas F, Johnson I. How to use t-SNE effectively. *Distill*. 2016;1.
516 doi:<https://doi.org/10.23915/DISTILL.00002>
- 517 42. Muppalaneni NB, Gunjan VK. *Computational intelligence techniques for comparative genomics*.
518 Dordrecht: Springer Singapore; 2015.
- 519 43. Rabiner LR. An introduction to Hidden Markov Models. *IEEE ASSP Magazine*. 1986;3: 4–16.
520 doi:<https://doi.org/10.1109/MASSP.1986.1165342>
- 521 44. Wang G. Machine learning for inferring animal behaviour from location and movement data. *Ecological*
522 *Informatics*. 2019;49: 69–76. doi:<https://doi.org/10.1016/j.ecoinf.2018.12.002>

- 523 45. Roever C, Beyer H, Chase M, Aarde R. The pitfalls of ignoring behaviour when quantifying habitat
524 selection. *Diversity and Distributions*. 2014;20: 322–333. doi:<https://doi.org/10.1111/ddi.12164>
- 525 46. Vinauger C, Lahondère C, Wolff GH, Locke LT, Liaw JE, Parrish JZ, et al. Modulation of host learning in
526 *Aedes aegypti* mosquitoes. *Current Biology*. 2018;28: 333–344.
527 doi:<https://doi.org/10.1016/j.cub.2017.12.015>
- 528 47. Lavielle M. Detection of multiple changes in a sequence of dependent variables. *Stochastic Processes
529 and their Applications*. 1999;83: 79–102. doi:[https://doi.org/10.1016/S0304-4149\(99\)00023-X](https://doi.org/10.1016/S0304-4149(99)00023-X)
- 530 48. Barraquand F, Benhamou S. Animal movements in heterogeneous landscapes: identifying profitable
531 places and homogeneous movement bouts. *Ecology Letters*. 2008;89: 3336–3348.
532 doi:<https://doi.org/10.1890/08-0162.1>
- 533 49. Byrne RW, Noser R, Bates LA, Jupp PE. How did they get here from there? Detecting changes of
534 direction in terrestrial ranging. *Animal Behaviour*. 2009;77: 619–31.
535 doi:<https://doi.org/10.1016/j.anbehav.2008.11.014>
- 536 50. Buchin M, Kruckenberg H, Kölzsch A. Segmenting trajectories by movement states. *Advances in spatial
537 data handling, geospatial dynamics, geosimulation and exploratory visualization*. Berlin, Heidelberg:
538 Springer; 2013. pp. 15–25.
- 539 51. Drummond H. The nature and description of behaviour patterns. *Perspectives in ethology*. Boston, MA:
540 Springer; 1981. pp. 1–33.
- 541 52. Drickamer LC, Snowdon CT. The emerging science: Defining the goals, approaches, and methods.
542 *Animal behaviour with commentaries*. Chicago: University of Chicago press; 1996. pp. 71–86.
- 543 53. Lehner PN. *Handbook of ethological methods*. Cambridge: Cambridge University Press; 1999.
- 544 54. Ranson H, Lissenden N. Insecticide resistance in African Anopheles mosquitoes: a worsening situation
545 that needs urgent action to maintain malaria control. *Trends in Parasitology*. 2016;32: 187–96.
546 doi:<https://doi.org/10.1016/j.pt.2015.11.010>
- 547 55. Angarita-Jaimes NC, Parker JEA, Abe M, Mashauri F, Martine J, Towers CE, et al. A novel videotracking
548 system to quantify the behaviour of nocturnal mosquitoes attacking human hosts in the field. *Journal
549 of the Royal Society Interface*. 2016;13. doi:<http://dx.doi.org/10.1098/rsif.2015.0974>
- 550 56. Voloshin V, Kröner C, Seniya C, Murray GPD, Guy A, Towers CE, et al. Diffuse retro-reflective imaging
551 for improved video tracking of mosquitoes at human baited bednets. *Royal Society Open Science*.
552 2020;7. doi:<http://dx.doi.org/10.1098/rsos.191951>
- 553 57. Gurarie E. Behavioral change point analysis in R: The BCPA package. CRAN; 2013. Available:
554 <https://cran.r-project.org/web/packages/bcpa/vignettes/bcpa.pdf>
- 555 58. Rousseeuw PJ. A graphical aid to the interpretation and validation of cluster analysis. *Journal of
556 computational and applied mathematics*. 1987;20: 53–65. doi:[https://doi.org/10.1016/0377-
557 0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7)
- 558 59. Kaufman L, Rousseeuw PJ. *Finding groups in data: An introduction to cluster analysis*. John Wiley &
559 Sons; 2009.
- 560 60. Patel P. *Hands-on unsupervised learning using Python: How to build applied machine learning solutions
561 from unlabelled data*. Farnham: O'Reilly; 2019.

