

# A PÓLYA TREE MODELLING FRAMEWORK FOR BATCH-MARK DATA

BY IOANNIS ROTOUS<sup>1,a</sup>, ALEX DIANA<sup>2,c</sup>, AND ELENI MATECHOU<sup>1,b</sup>

<sup>1</sup>Statistical Ecology @ Kent, School of Mathematics, Statistics and Actuarial Science, University of Kent, <sup>a</sup>ir237@kent.ac.uk;  
<sup>b</sup>e.matechou@kent.ac.uk

<sup>2</sup>Department of Mathematics, Statistics and Actuarial Science, University of Essex, <sup>c</sup>ad23269@essex.ac.uk

Wildlife population surveys typically consist of multiple sampling occasions, where individuals are followed over time, enabling estimation of population size and, in open populations, of entry and exit patterns. Batch-mark (BM) surveys, where newly sampled individuals are given the same marking, unique for each sampling occasion but not for each individual, provide the only viable monitoring tool for many species of amphibians, birds and fish. Modelling BM data for open populations has proven more challenging than modelling data where individuals are uniquely marked, and approaches proposed in the literature thus far rely on approximate inference or do not scale well with the number of individuals, and do not readily extend to the joint modelling of different observation processes often employed in practice. In this paper, we propose a novel approach for modelling BM data, by defining a bivariate grid for modelling the latent entry and exit patterns, as well as population size. We employ the Bayesian nonparametric Pólya Tree (PT) prior for defining a model on the grid cells, which enables exact and highly efficient Bayesian inference on the number of individuals in each cell, and hence of the entry/exit pattern. Our approach scales with the number of sampling occasions, instead of the number of individuals, and allows us to easily write the likelihood function for BM data under different observation processes. We demonstrate our new Pólya Tree prior Batch Mark (PTBM) approach using extensive simulations and two case studies, comparing its performance with two recently proposed approaches.

**1. Introduction.** Ecological monitoring is of paramount importance in safeguarding our natural ecosystems. By studying wildlife populations and accurately estimating essential demographic parameters, such as population size, birth/death rates or phenological patterns, we can gain invaluable insights into their dynamics. One approach that can be employed for population monitoring is batch-marking (BM), which involves repeated sampling occasions when individuals are physically caught, with newly caught individuals marked using a sampling occasion-specific mark, that is a mark that is typically different among sampling occasions but not among individuals. BM surveys are particularly useful in cases where individual marking cannot be employed, such as with species that are highly abundant or small in size, for example fish or insects, and have been extensively utilized providing valuable data for population monitoring (Cowen et al., 2017; Vavassori, Saddler and Müller, 2019; Davidson et al., 2019; Doll et al., 2021; Rosser, Willden and Loeb, 2022).

BM data on each sampling occasion consist of the number of individuals caught that are unmarked and the number of individuals caught that were first marked on previous sampling occasions. Therefore, as opposed to standard capture-recapture studies (McCrea and Morgan, 2014; Seber and Schofield, 2019; King and McCrea, 2019), where individuals are uniquely marked, it is not known how many times and on which sampling occasions each individual is re-caught. This aggregated nature of BM data means that likelihood-based inference is

---

*Keywords and phrases:* Bayesian nonparametrics, Stopover model, Population size, Removal models, Survival probability.

more challenging, and the literature on appropriate models for BM data is limited. To discuss existing approaches we need to introduce the concept of individual presence histories and capture histories. Both are vectors of length equal to the number of sampling occasions. The former indicates when the corresponding individual is present at the surveyed site, while the latter indicates when the individual has been caught. Individual presence histories are always latent in ecological data, while individual capture histories are observed in capture-recapture data but latent in BM data. Current approaches for modelling BM data include the work by Huggins, Wang and Kearns (2010), who derived estimating equations and closed-form solutions for survival and capture probabilities, along with abundance estimates using a Horvitz-Thompson-type estimator. Another important contribution came from Cowen et al. (2014), who developed a tractable likelihood function for marked individuals only, but with an associated computational cost that increases substantially with more sampling occasions, since the calculation of the likelihood involves nested summations of all possible individual latent presence histories. The Cowen et al. (2014) work was further extended by Cowen et al. (2017) who employed Hidden Markov Models (HMM) to model BM data for both marked and unmarked individuals, but this HMM approach does not scale well with the number of individuals, since it involves high dimensional state-transition and state-dependent probability matrices. More recently, Zhang, Bonner and McCrea (2023) introduced an innovative approach to approximating the likelihood function for BM data under the robust design (RD, Kendall and Pollock, 1992). When the RD is employed, it is assumed that the population is open between primary periods, e.g. months, but closed within secondary periods, e.g. days within a month, so that individuals can enter/exit between primary periods but not between secondary periods. Zhang, Bonner and McCrea (2023) proposed an approach that relies on the saddlepoint approximation (SPA Butler, 2007) to the likelihood function, which requires reconstructing the latent capture history as well as the latent presence history for each individual, so the computational burden increases considerably with the number of sampling occasions.

In this paper, we present a novel Bayesian nonparametric approach utilizing the Pólya Tree (PT) prior for BM data analysis. This approach, which builds on the work by Diana et al. (2023) for count and ring-recovery data, offers numerous benefits. The fundamental idea is that we build a lower-right triangular grid with  $\frac{(K+1)(K+2)}{2}$  latent cells, as shown in Figure 1, where  $K$  is the number of sampling occasions. The grid cells represent the latent number of individuals with a specific combination of entry and exit intervals, where the intervals are defined by pairs of consecutive sampling occasions, so that  $n_{i,j}$  is the number of individuals with entry between sampling occasions  $i$  and  $i + 1$  and exit between sampling occasions  $j$  and  $j + 1$ . We note that entry and exit can correspond to birth/death, arrival/departure, or any other process that introduces/removes individuals from the population. By using our proposed grid approach and corresponding PT prior, as we discuss below, we naturally account for the aggregated nature of BM data because we do not need to infer, impute or marginalise over latent individual presence histories or capture histories, and instead only need to infer the latent cells of the grid. Consequently, within this framework, we overcome previous challenges related to BM model inference and establish a tractable likelihood function and a corresponding efficient model-fitting algorithm for standard BM data, as well as BM data collected under different sampling designs, as we discuss in the two case studies of the paper.

Inference of the cells  $n_{i,j}$  is efficient and flexible within a Bayesian framework based on the PT prior, which allows us to define and infer the grid probabilities. By relying on the PT prior, we can build a model directly on the distributions of entry and exit patterns, with minimal parametric assumptions on the shape of these distributions. Additionally, the replicate PT framework, first introduced in Diana et al. (2023), allows us to impose constraints on the entry and exit processes, leading to more parsimonious models, as we discuss in Section 3.

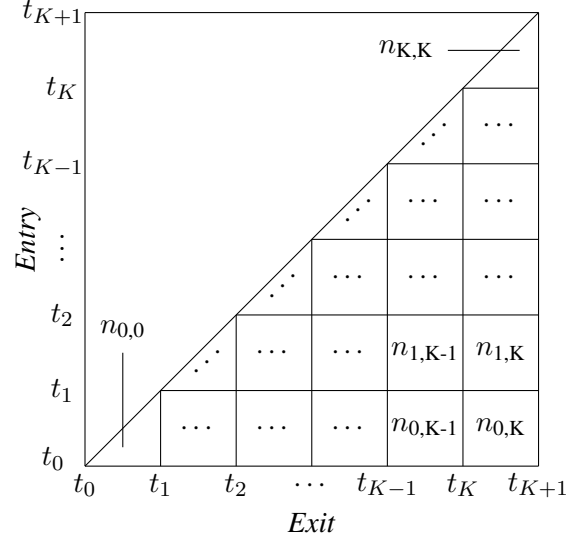


Fig 1: Entry and exit sample space, for  $K$  sampling occasions, taking place at times  $t_1, t_2, \dots, t_K$  with convention that  $t_0 = -\infty, t_{K+1} = \infty$ . The latent number of individuals in cell  $(i, j)$ , that is with entry between the  $i$ th and  $(i + 1)$ th and exit between the  $j$ th and  $(j + 1)$ th sampling occasions are denoted by  $n_{i,j}$ , with  $i = 0, 1, \dots, K$  and  $j = i, \dots, K$ .

We present results for two case studies previously analysed in the literature. Both consider open populations and fit models of Jolly Seber (JS Jolly, 1965; Seber, 1965) type, which are models that allow direct estimation of population size while accounting for entry and exit from the population. The data from the two case studies result from different observation processes. The first case study considers data on weather loaches (*Misgurnus anguillicaudatus*), which are freshwater fish, monitored in South-Eastern Australia. In this case, during the BM survey, individuals caught as unmarked are either marked and returned to the population in the standard BM approach, or are removed from the population, akin to removal sampling (Matechou et al., 2016). We show how our framework can easily be adapted to model the two processes jointly, which was not possible using the approach proposed by Cowen et al. (2017). A corresponding simulation study demonstrates the resulting bias when the removal process is ignored. The second case study considers data on golden mantella frogs (*Mantella aurantiaca*), collected in Central Madagascar under the RD across six primary periods. In this case, we demonstrate how our framework can easily account for BM data collected under a RD sampling, using the exact likelihood, and compare our results to those obtained by Zhang, Bonner and McCrea (2023). A corresponding simulation study explores issues around allocation of effort within a RD framework for BM surveys.

The article is organized as follows: In Section 2 we introduce our modelling approach of the latent entry/exit pattern and population size from BM data, and in Section 3 we describe the PT prior, which provides the foundation of our modelling framework. We present simulation and real data results for each of the two case studies in Sections 4 and 5. The article concludes with a discussion of the findings and outlines potential avenues for future research in Section 6.

**2. Model for batch-mark data.** First, we introduce the notation and modelling of standard BM data, with  $K$  sampling occasions, taking place at times  $t_1, \dots, t_K$ , with  $t_0 = -\infty, t_{K+1} = \infty$ , for convenience. Additional notation for the two different observation processes employed in the case studies is introduced in the corresponding sections. In what follows we refer to individuals in (grid) cell  $(i, j)$  as those individuals that entered between sampling occasions  $i$  and  $(i + 1)$  and exited between sampling occasions  $j$  and  $(j + 1)$ .

### Data

$u_t$  : observed number of unmarked individuals caught on sampling occasion  $t$ .

$m_{k,t}$  : observed number of individuals that were marked on sampling occasion  $k$  and were recaptured on sampling occasion  $t$ .

### Parameters

$N$  : population size, with  $\sum_{i=0}^K \sum_{j=i}^K n_{i,j} = N$  (see Figure 1). We model  $N$  as a Poisson random variable with mean  $\omega$ ,  $N \sim \text{Poisson}(\omega)$ .

$w_{i,j}$  : probability of an individual belonging to cell  $(i, j)$ .

$p_t$  : probability of capturing an individual that is present on sampling occasion  $t$ .

### Latent

$n$  :  $(K + 1) \times (K + 1)$  matrix where  $n_{i,j}$  corresponds to the latent number of individuals in cell  $(i, j)$ .

The latent cells  $n_{i,j}$ , conditional on  $N$  are modelled as

$$n|N \sim \text{Multinomial}(N, \{w_{i,j}\}_{i=0,\dots,K,j=i,\dots,K})$$

so that, since  $N \sim \text{Poisson}(\omega)$ , the unconditional distribution of cells in matrix  $n$  is given by

$$(1) \quad n \sim \text{Poisson}(\omega \times \{w_{i,j}\}_{i=0,\dots,K,j=i,\dots,K})$$

$U^k$  :  $(K + 1) \times (K + 1)$  matrix where  $U_{i,j}^k$  corresponds to the latent number of individuals in cell  $(i, j)$  that are present and unmarked on sampling occasion  $k$ .

$M^k$  :  $(K + 1) \times (K + 1)$  matrix where  $M_{i,j}^k$  corresponds to the number of individuals in cell  $(i, j)$  that were caught as unmarked, and subsequently marked, on sampling occasion  $k$ . Each of the latent  $U_{i,j}^k$  individuals can be caught with probability  $p_k$ , independent of other individuals, so that

$$M_{i,j}^k \sim \text{Binomial}(U_{i,j}^k, p_k)$$

We note that the latent number of individuals in cell  $(i, j)$  that are unmarked on sampling occasion  $t$ ,  $U_{i,j}^t$ , is equal to the difference between the total latent number of individuals in that cell and the latent number of individuals in that cell that were marked before sampling occasion  $t$ , so that  $U_{i,j}^t = n_{i,j} - \sum_{k=1}^{t-1} M_{i,j}^k$ , with  $i < t \leq j$ . Note that  $M_{i,j}^k = 0$  when  $i \geq k$ , since individuals must enter before they are caught, and when  $t = 1$  there are no marked individuals yet and hence  $U_{i,j}^1 = n_{i,j}$ .

Therefore, conditional on the latent matrices  $M^k$ , the observed number of unmarked individuals caught on sampling occasion  $t$ ,  $u_t$ , are

$$u_t = \sum_{i=0}^{t-1} \sum_{j=t}^K M_{i,j}^t.$$

The number of individuals marked on sampling occasion  $k$  that are still present on sampling occasion  $t$  is the sum of  $M_{i,j}^k$  terms for individuals that entered before  $k$  and have not yet exited at  $t$ , i.e. with  $i < k$  and  $j \geq t$ , and each of these individuals is recaptured with probability  $p_t$ , independent of other individuals, so that the observed number of individuals marked on sampling occasion  $k$  and recaptured on sampling occasion  $t$  are distributed as

$$m_{k,t} \sim \text{Binomial} \left( \sum_{i=0}^{k-1} \sum_{j=t}^K M_{i,j}^k, p_t \right), \quad t > k$$

The complete model for the standard BM data is given in Equation 2.

(2)

latent number of individuals in cell  $(i, j)$   $n_{i,j} \sim \text{Poisson}(\omega \times w_{i,j})$

latent number of individuals in cell  $(i, j)$  that are unmarked on sampling occasion  $t$   $U_{i,j}^t = \begin{cases} n_{i,j}, & t = 1 \\ n_{i,j} - \sum_{k=1}^{t-1} M_{i,j}^k, & t > 1 \end{cases}$

latent number of individuals in cell  $(i, j)$  caught as unmarked and subsequently marked on sampling occasion  $t$   $M_{i,j}^t \sim \text{Binomial}(U_{i,j}^t, p_t)$

observed number of unmarked individuals caught on occasion  $t$   $u_t = \sum_{i=0}^{t-1} \sum_{j=t}^K M_{i,j}^t$

observed number of individuals that were marked on sampling occasion  $k$  and were recaptured on sampling occasion  $t$   $m_{k,t} \sim \text{Binomial} \left( \sum_{i=0}^{k-1} \sum_{j=t}^K M_{i,j}^k, p_t \right), \quad t > k$

118 We note that in our model  $N$  is not the equivalent of the super-population of individuals,  
 119  $N^S$ , that became available for capture at least once (Schwarz and Arnason, 1996) but instead  
 120 also accounts for individuals that entered and exited without ever becoming available for  
 121 capture (Matechou and Caron, 2017) that is individuals in cells with  $i = j$ . Inference on these  
 122 individuals is possible within this framework thanks to the PT prior, as also discussed in  
 123 Diana, Griffin and Matechou (2019). However, inference on the super-population size is also  
 124 readily available from the  $n$  matrix, since  $N^S = N - \sum_{i=0}^K n_{i,i}$ , i.e. by excluding individuals  
 125 that never became available for capture, and in the case studies we report both quantities.

126 A simulation study for this standard BM model is presented in Web Appendix B.

**3. Pólya Tree Prior.** Inferring the grid matrix  $n$  relies on inferring the grid probabilities  $w_{i,j}$ , through the application of the PT prior. The PT prior requires partitioning the sample space and specifying a sequence of positive numbers to assign probabilities to **each set in** these partitions.

In our case, following the methodology outlined in Section 2, we need to partition the entry and exit space, the latter conditional upon entry. This partitioning is achieved by first splitting the sample space into the sets  $B_i = (t_i, t_{i+1}] \times (t_i, t_{K+1})$  for  $i = 0, 1, 2, \dots, K$ , corresponding to the individuals entering in the interval between the  $i$ th and  $(i + 1)$ th sampling occasions and exiting after entry. The probability of an individual belonging to  $B_i$  is  $V_i$ , where  $(V_0, V_1, \dots, V_K) \sim \text{Dirichlet}(\alpha_0, \alpha_1, \dots, \alpha_K)$ . The Dirichlet distribution on the  $V$ s can be centred on any distribution  $G_0$  that expresses our prior expectation for the entry pattern, by assuming  $\mathbb{E}[V_i] = G_0(B_i)$ , which can be achieved by defining the sequence of positive numbers such that  $a_i \propto G_0(B_i)$ . If we do not have prior knowledge about the entry pattern we can choose a uniform distribution for  $G_0$ , which corresponds to setting  $a_0 = a_1 = \dots = a_K = 1$ . Further examples about partitions of the entry/exit space and PT prior centring can be found in Diana et al. (2023).

Next, we split each entry set,  $B_i = (t_i, t_{i+1}) \times (t_i, t_{K+1})$  into the sets  $B_{i,j} = (t_i, t_{i+1}) \times (t_j, t_{j+1})$ ,  $j = i, \dots, K$ , where  $(t_i, t_{i+1}) \times (t_j, t_{j+1})$  corresponds to the individuals entering between the  $i$ th and  $(i + 1)$ th sampling occasion and exiting between the  $j$ th and  $(j + 1)$ th sampling occasion. The probability of an individuals belonging to  $B_{i,j}$  conditional on being in  $B_i$  is  $V_{i,j}$ , where  $V_{i,j}, V_{i,j+1}, \dots, V_{i,K} \sim \text{Dirichlet}(a_{i,j}, a_{i,j+1}, \dots, a_{i,K})$ ,  $j = i, \dots, K$ . We employ the replicate PT as described in Diana et al. (2023), which imposes a structure on the exit probabilities, and hence builds a more parsimonious PT prior with a smaller number of parameters. In our case, we use a time-dependent constraint for the exit probabilities and assume that  $V_{0,j} = V_{1,j} = \dots = V_{j,j}$ ,  $j = 1, \dots, K$ , i.e. we replicate the exit probabilities across the entry intervals, so that the probability of exit in a given interval is the same for all individuals, regardless of entry.

Finally, the PT prior distribution of the probabilities  $w_{i,j}$  is defined as

$$(3) \quad \{w_{i,j}\} = \{V_i V_{i,j}\} \sim PT(\{a_i\}, \{a_{i,j}\}), \quad i = 0, 1, 2, \dots, K \quad j = i, \dots, K$$

**3.1. Inference.** We describe our employed Markov Chain Monte Carlo (MCMC), with corresponding conditional distributions, where appropriate, in Web Appendix A. Briefly, the latent matrices  $\{M^k\}_{k=1}^K$  and  $n$  are updated using a standard Metropolis Hastings (MH) random walk (Robert et al., 2004) whereas the  $\{w_{i,j}\}_{i=0,j=i}^{K,K}$  parameters are updated using a Gibbs algorithm (Gelfand and Smith, 1990) since we exploit the conjugacy properties of the PT prior (Lavine, 1992). Parameters  $\{p_t\}_{t=1}^K$  and  $\omega$  are also updated with the use of a Gibbs sampler.

**4. Case Study 1.** BM data on weather-loch were collected across 11 sampling occasions, but in contrast to standard BM sampling, a proportion of unmarked individuals caught on each sampling occasion were removed from the population (Cowen et al., 2017). Here we demonstrate how our modelling framework introduced in section 2 can naturally account for removals on captures in BM surveys. From section 2, we make use of data  $\{u_t\}_{t=1}^K, \{m_{k,t}\}_{k=1,t=k+1}^{K-1,K}$ , parameters  $N, \{w_{i,j}\}_{i=0,j=i}^{K,K}, \{p_t\}_{t=1}^K$ , and latent  $n, \{M^k\}_{k=1}^K$  alongside we introduce additional data, parameters and latent corresponding to removals per sampling occasion.

169 4.1. *Data and Notation.* First we introduce some case-study-specific notation.

170 **Data**

171  $r_k$  : observed number of individuals that were removed from the population on sampling  
172 occasion  $k$ .

173 **Parameters**

174  $l_k$  : probability of removing an unmarked individual caught on sampling occasion  $k$ .

175 **Latent**

$Y^k$  :  $(K + 1) \times (K + 1)$  matrix where  $Y_{i,j}^k$  corresponds to the latent number of unmarked individuals in cell  $(i, j)$  caught on sampling occasion  $k$ . We note that this is not the same as  $M_{i,j}^k$  introduced in Section 2 since not all newly caught individuals are marked in this case, as we discuss below. Similarly to the standard BM model described in Section 2,

$$Y_{i,j}^k \sim \text{Binomial}(U_{i,j}^k, p_k)$$

$R^k$  :  $(K + 1) \times (K + 1)$  matrix where  $R_{i,j}^k$  corresponds to the latent number of individuals in cell  $(i, j)$  that were removed on sampling occasion  $k$ . Each of the  $Y_{i,j}^k$  individuals caught as unmarked has the same probability,  $l_k$ , independent of other individuals, of being removed instead of being marked and returned to the population, so that

$$R_{i,j}^k \sim \text{Binomial}(Y_{i,j}^k, l_k)$$

176 We note that now the number of individuals in cell  $(i, j)$  that are present and marked on  
177 sampling occasion  $t$  is  $M_{i,j}^t = Y_{i,j}^t - R_{i,j}^t$ , and the latent number of individuals in cell  
178  $(i, j)$  that are available as unmarked on sampling occasion  $t$  is  $U_{i,j}^t = n_{i,j} - \sum_{k=1}^{t-1} M_{i,j}^k -$   
179  $\sum_{k=1}^{t-1} R_{i,j}^k$ ,  $i < t \leq j$  since individuals in cell  $(i, j)$  are no longer available as unmarked  
180 on sampling occasion  $t$  if they were previously marked or removed from the population.

Therefore it follows that, conditional on  $R_{i,j}^k$ , the observed number of individuals removed from the population on sampling occasion  $k$  are

$$r_k = \sum_{i=0}^{k-1} \sum_{j=k}^K R_{i,j}^k$$

181 4.2. *Model.* The model of Equation (2) is extended to account for the removal process,  
182 in addition to the standard BM process, as shown in Equation (4).

(4)

latent number of individuals in cell  $(i, j)$   $n_{i,j} \sim \text{Poisson}(\omega \times w_{i,j})$ latent number of individuals in cell  $(i, j)$  that are unmarked on sampling occasion  $t$   $U_{i,j}^t = \begin{cases} n_{i,j}, & t = 1 \\ n_{i,j} - \sum_{k=1}^{t-1} M_{i,j}^k - \sum_{k=1}^{t-1} R_{i,j}^k, & t > 1 \end{cases}$ latent number of individuals in cell  $(i, j)$  caught as unmarked on sampling occasion  $t$   $Y_{i,j}^t \sim \text{Binomial}(U_{i,j}^t, p_t)$ observed number of unmarked individuals caught on sampling occasion  $t$   $u_t = \sum_{i=0}^{t-1} \sum_{j=t}^K Y_{i,j}^t$ latent number of individuals in cell  $(i, j)$  removed on sampling occasion  $k$   $R_{i,j}^k \sim \text{Binomial}(Y_{i,j}^k, l_k)$ observed number of individuals removed on sampling occasion  $k$   $r_k = \sum_{i=0}^{k-1} \sum_{j=k}^K R_{i,j}^k$ observed number of individuals that were marked on sampling occasion  $k$  and were recaptured on sampling occasion  $t$   $m_{k,t} \sim \text{Binomial}\left(\sum_{i=0}^{k-1} \sum_{j=t}^K M_{i,j}^k, p_t\right), t > k$ 

183 4.3. *Case study 1: Inference.* Inference for the parameters  $\{M^k\}_{k=1}^K$ ,  $n$ ,  $\{w_{i,j}\}_{i=0,j=i}^{K,K}$ ,  
 184  $\{p_t\}_{t=1}^K$  and  $\omega$  is described in Section 2. The update for the new parameter, the latent matrix  
 185  $\{R^k\}_{k=1}^K$ , is based on a MH random walk and is described in Web Appendix A

186 4.4. *Simulation study.* We conducted a simulation study with the aim of exploring how  
 187 ignoring removed individuals affects inference quality. For 100 replications, we simulated  
 188 BM data for 11 sampling occasions, involving 6000 individuals and **with** removal probabilities  
 189 of 0.05, 0.1, or 0.2. We used time varying entry/exit probabilities as well as capture  
 190 probabilities. The chosen values are given in Web Appendix C, together with the prior dis-  
 191 tribution choices for all parameters. In each case, we ran an MCMC algorithm for 500,000  
 192 iterations, with a burn-in of 20,000 iterations.

193 We denote the true parameter value in each case by  $\tilde{\theta}$  and the corresponding parameter  
 194 sampled on the  $b$ -th iteration of replication  $m$  after burn-in by  $\theta_b^m$ . We calculate and report  
 195 the posterior bias across simulations,  $\bar{B}_m$ , as  $\bar{B}_m = \frac{1}{B} \sum_{b=1}^B (\theta_b^m - \tilde{\theta})$ ,  $m = 1, 2, \dots, 100$ . In  
 196 Figure 2 we display the posterior bias obtained by our model when removed individuals are  
 197 accounted for (denoted by PTBM-R) and when they are not (denoted by PTBM). Inference  
 198 for  $N$  is summarised in Figure 2 and for all other parameters in Web Appendix C.

199 As expected, when the removed individuals are not accounted for, inference is biased,  
 200 especially for parameters concerning the number of individuals in the survey such as  $N$ ,  $U_t$   
 201 and  $N_t$ , which is defined as the number of individuals available on sampling occasion  $t$  (see  
 202 Web Appendix C). This bias becomes more severe as the proportion of unmarked individuals  
 203 removed increases. The rest of the parameters, such as  $V_i$ ,  $V_{.,j}$  and  $p_t$ , exhibit less severe bias,  
 204 especially when removal probability is low, since they are mainly informed by the marked  
 205 individuals that are followed over time, and hence are not affected as much by the removal



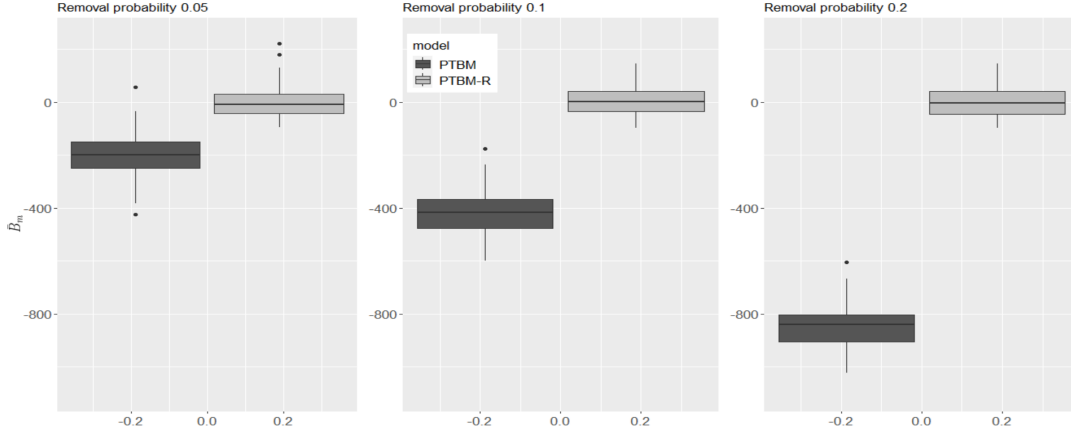


Fig 2: Posterior bias of population size,  $N$ . Each column corresponds to a different removal probability. We compare the posterior bias across the models PTBM-R and PTBM, i.e. when we account for removals and when we do not.

206 of unmarked individuals. Finally, as expected, when the correct model is fitted to the data,  
 207 estimation is unbiased.

208 **4.5. Weather-loch data.** We present results for this case study in Figures 3 and Table 1  
 209 and compare them to those obtained by Cowen et al. (2017). Following Cowen et al. (2017),  
 210 we assume that entry probabilities and exit probabilities are constant, that is  $V_i = \beta$ , for all  
 211  $i = 0, 1, \dots, K$  and  $V_{i,j} = 1 - \phi$  for all  $j = i, \dots, K$ . The prior distribution choices for our  
 212 parameters are described in Web Appendix D.

213 The overall patterns in our findings agree with those in the simulation study. Namely, in  
 214 Table 1, we observe that estimation of entry and exit probabilities for the PTBM-R, PTBM,  
 215 and HMM models is, on average, aligned. Additionally, in Figure 3, we observe that, on  
 216 average, our estimates for the number of unmarked individuals present on each sampling  
 217 occasion, and for the population size are larger when accounting for removals compared  
 218 to the other two models PTBM and HMM. Finally, the posterior median population and  
 219 super-population for the PTBM-R model are 3280 with a 95% posterior credible interval of  
 220 (2892, 3771) and 2805 with a 95% posterior credible interval of (2554, 3047), respectively,  
 221 whereas for the PTBM model are 2198 with a 95% posterior credible interval of (1958, 2419)  
 222 and 2182 with a 95% posterior credible interval of (2095, 2276) respectively. The super-  
 223 population for the Cowen et al. (2017) model is estimated as 2242 (no interval was reported)  
 224 which is in alignment with out PTBM model which is expected since both models do not  
 225 account for removals.

226 Figures 3 in this section and 8 in Web Appendix D demonstrate that the number of avail-  
 227 able individuals increases and peaks on the 5th sampling occasion before gradually decreas-  
 228 ing. This could potentially be interpreted as a result of seasonality. Notably, Huggins, Wang  
 229 and Kearns (2010) mentioned that they estimated the smallest number of available individuals  
 230 on the 7th sampling occasion, which they attributed to the winter season, with the population  
 231 then increasing during the spring. In contrast, Cowen et al. (2017) stated that they estimated  
 232 a decrease in the number of available individuals over time when using time-varying capture  
 233 probabilities, as opposed to Huggins, Wang and Kearns (2010). Interestingly, our PTBM-R  
 234 model manages to identify some seasonality, with population size peaking before the winter  
 235 time and subsequently estimating a decrease in the population while accounting for time-  
 236 varying probabilities.

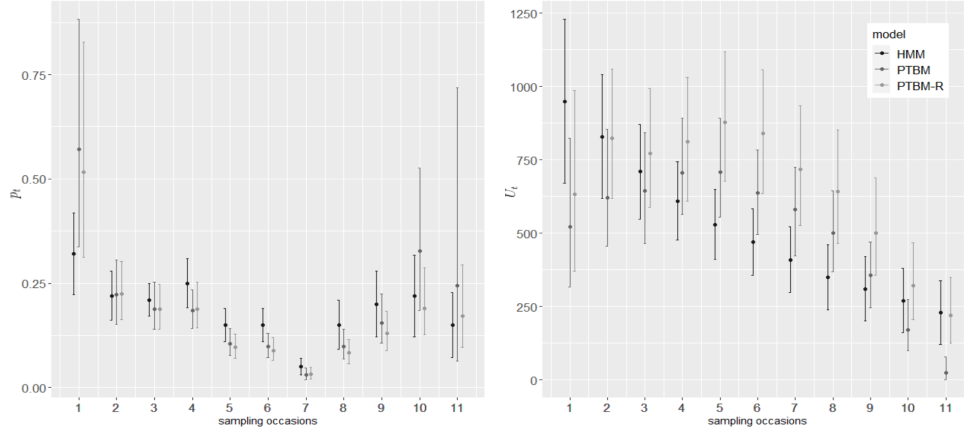


Fig 3: Summaries of capture probabilities  $p_t$ , (a), and of the inferred number of unmarked individuals present on each sampling occasion  $U_t$ ,  $t = 1, 2, \dots, K$ , (b). The HMM model of Cowen et al. (2017) was fitted classically, so the summaries correspond to the maximum likelihood estimate and corresponding 95% confidence interval in each case, whereas our PTBM and PTBM-R models are fitted in a Bayesian framework, so the summaries correspond to the posterior mean and 95% posterior credible interval.

TABLE 1

Summaries of entry probability  $\beta$ , (a) and exit probability  $1 - \phi$ , (b). In the HMM case, these correspond to the maximum likelihood estimate and corresponding 95% confidence interval, while in the PTBM and PTBM-R models, these correspond to the posterior median and corresponding 95% posterior credible interval.

Parameter		
Model	$\beta$	$\phi$
HMM	0.24 (0.18, 0.30)	0.63 (0.58, 0.69)
PTBM	0.22 (0.14, 0.35)	0.61 (0.27, 0.81)
PTBM-R	0.21 (0.14, 0.34)	0.63 (0.40, 0.87)

Finally, we conducted a goodness of fit assessment, by comparing the observed number of unmarked individuals  $\{u_t\}_{t=1}^K$  caught on each sampling occasion with summaries of the corresponding data simulated from each model (see Web Appendix D). It is evident that the PTBM-R model, which better describes the data-generating process, outperforms the other two models in generalizing observed patterns in the data.

**5. Case Study 2.** Next, we consider the case study first published in Zhang, Bonner and McCrea (2023) where sampling follows Pollock’s robust design (RD Pollock, 1982). The sampling of golden mantella took place over 6 primary periods with (3, 3, 3, 4, 4, 4) secondary sampling occasions for each primary period. During each primary period, captured individuals receive a distinct mark and are then released so can be recaptured on later secondary sampling occasions within the same or different primary periods.

**5.1. Data and Notation.** First we introduce some case-study-specific notation. In the case of the RD, there are  $K$  primary periods with  $T_k$  secondary sampling occasions within primary period  $k$ . The underlying latent process of entry and exit refers to intervals between primary periods, as described in Section 2, while the observation process has the same structure as in Equation (2), but in this case with multiple secondary sampling occasions per primary period, as we describe below.

**Data**

$u_{k,t}$  : observed number of unmarked individuals caught on secondary sampling occasion  $t$  of primary period  $k$ .

$m_{k,\nu,t}$  : observed number of individuals that were marked in primary period  $k$  and were recaptured on secondary sampling occasion  $t$  of primary period  $\nu$ .

**Parameters**

$p_{k,t}$  : probability of capturing an individual present on secondary sampling occasion  $t$  of primary period  $k$ .

**Latent**

$U^{k,t}$  :  $(K+1) \times (K+1)$  matrix where  $U_{i,j}^{k,t}$  corresponds to the latent number of unmarked individuals in cell  $(i,j)$  present on secondary sampling occasion  $t$  of primary period  $k$ .

$Y^{k,t}$  :  $(K+1) \times (K+1)$  matrix where  $Y_{i,j}^{k,t}$  corresponds to the latent number of individuals in cell  $(i,j)$  caught as unmarked on secondary sampling occasion  $t$  of primary period  $k$ . Similarly to Sections 2 and 4,

$$Y_{i,j}^{k,t} \sim \text{Binomial}(U_{i,j}^{k,t}, p_{k,t})$$

$M^{k,t}$ : is a  $(K+1) \times (K+1)$  matrix, where  $M_{i,j}^{k,t}$  corresponds to the latent number of individuals in cell  $(i,j)$  caught unmarked in primary period  $k$  and marked across any secondary sampling occasions from 1 up to  $t$  within primary period  $k$ .

$$M_{i,j}^{k,t} = \sum_{m=1}^t Y_{i,j}^{k,m}$$

Therefore, similarly to Sections 2 and 4, it follows that the observed number of unmarked individuals caught on secondary sampling occasion  $t$  of primary period  $k$  are

$$u_{k,t} = \sum_{i=0}^{k-1} \sum_{j=k}^K Y_{i,j}^{k,t}$$

while the number of individuals that were marked in primary period  $k$  and recaptured on secondary sampling occasion  $t$  of primary period  $\nu$  are modelled as

$$m_{k,\nu,t} \sim \begin{cases} \text{Binomial}\left(\sum_{i=0}^{k-1} \sum_{j=\nu}^K M_{i,j}^{k,t-1}, p_{k,t}\right), & \nu = k, t = 2, 3, \dots, T_\nu \\ \text{Binomial}\left(\sum_{i=0}^{k-1} \sum_{j=\nu}^K M_{i,j}^{k,T_k}, p_{k,t}\right), & \nu > k, t = 1, 2, \dots, T_\nu \end{cases}$$

In the above, we distinguish two cases. If recaptures take place in the same primary period ( $\nu = k$ ), individuals are caught from the available marked individuals  $M_{i,j}^{k,t-1}$ , which are the individuals caught for the first time on primary period  $k$  and marked on any secondary sampling occasion from 1 up to  $t-1$ , with  $i < k$  and  $j \geq \nu$ . If recaptures take place on subsequent primary periods ( $\nu > k$ ), individuals are caught from the available marked individuals  $M_{i,j}^{k,T_k}$ , which are the ones first caught on primary period  $k$  and marked on any secondary sampling occasion  $1, 2, \dots, T_k$ , with  $i < k$  and  $j \geq \nu$ , that is, that are available on primary period  $\nu$ .

5.2. *Model.* The complete model for BM data collected under a RD is given in Equation 5.

(5)

$$\begin{aligned}
 &\text{latent number individuals in cell } (i, j) & n_{i,j} &\sim \text{Poisson}(\omega \times w_{i,j}) \\
 &\text{latent number of individuals in cell } (i, j) \text{ that} & U_{i,j}^{k,t} &= \begin{cases} n_{i,j}, & k = 1, t = 1 \\ n_{i,j} - \sum_{\nu=1}^k \sum_{t'=1}^{t-1} Y_{i,j}^{\nu,t'}, & k \neq 1, t \neq 1 \end{cases} \\
 &\text{are unmarked on secondary sampling occa-} \\
 &\text{sion } t \text{ of primary period } k \\
 &\text{latent number of individuals in cell } (i, j) & Y_{i,j}^{k,t} &\sim \text{Binomial}(U_{i,j}^{k,t}, p_{k,t}) \\
 &\text{caught as unmarked on secondary sampling} \\
 &\text{occasion } t \text{ of primary period } k \\
 &\text{latent number of individuals in cell } (i, j) & M_{i,j}^k &= \sum_{t'=1}^t Y_{i,j}^{k,t'} \\
 &\text{marked in primary period } k \\
 &\text{observed number of unmarked individuals} & u_{k,t} &= \sum_{i=0}^{k-1} \sum_{j=k}^K Y_{i,j}^{k,t} \\
 &\text{caught on secondary sampling occasion } t \text{ of} \\
 &\text{primary period } k \\
 &\text{observed number of individuals that were} & m_{k,\nu,t} &\sim \begin{cases} \text{Binomial}\left(\sum_{i=0}^{k-1} \sum_{j=\nu}^K M_{i,j}^{k,t-1}, p_{k,t}\right), \\ \nu = k, t = 2, 3, \dots, T_\nu \\ \text{Binomial}\left(\sum_{i=0}^{k-1} \sum_{j=\nu}^K M_{i,j}^{k,T_k}, p_{k,t}\right), \\ \nu > k, t = 1, 2, \dots, T_\nu \end{cases} \\
 &\text{marked in primary period } k \text{ and were recap-} \\
 &\text{tured on secondary sampling occasion } t \text{ of} \\
 &\text{primary period } \nu
 \end{aligned}$$

278 5.3. *Case Study 2: Inference.* We employ an MCMC MH random walk for infer-  
 279 ring the latent matrices  $\{M^k\}_{k=1}^K$  and  $n$ , and a Gibbs sampler for the probabilities  
 280  $\{w_{i,j}\}_{i=0,j=i}^{K,K}$ ,  $\{p_{k,t}\}_{k=1,t=1}^{K,T_k}$  and for  $\omega$ . Details, about the conditional distributions are given  
 281 in Web Appendix A.

282 5.4. *Simulation study.* We simulated a population size of 5000 individuals across 5 pri-  
 283 mary periods, considering time varying entry/exit probabilities and capture probabilities.  
 284 Their values are given in Web Appendix E, together with the prior distribution choices for  
 285 the parameters. To assess the impact of the number of secondary sampling occasions on the  
 286 quality of inference, we simulated data using 1, 2, 4, 8, and 16 secondary sampling occasions.

287 We ran the MCMC for 500,000 iterations, with the 100,000 iterations as burn-in. The  
 288 results are presented in Table 2 and in Web Appendix E. Here, we summarize the effect,  
 289 in terms of percentage decrease in root mean squared error (RMSE), in the estimated entry  
 290 and exit probabilities and population size when the number of secondary occasions increases  
 291 (Table 2). Our findings indicate that, as expected, incorporating more sampling occasions  
 292 within each primary period leads to smaller RMSE for all parameters. However, the benefit in  
 293 terms of decrease in RMSE is not proportionate to the increase in effort for larger numbers of  
 294 secondary periods, and the effects of the additional effort practically diminish for more than  
 295 eight secondary periods, especially in the case of population size. These results highlight the  
 296 benefit of the RD within BM studies in comparison to the standard BM without any periods  
 297 of closure, and can support decisions around study design and allocation of effort in BM  
 298 surveys.

TABLE 2

The median decrease change in RMSE as we increase the number of secondary sampling occasions from 1 to 2, 2 to 4, 4 to 8 and lastly from 8 to 16. The displayed parameters are the population size,  $N$ , and the entry and exit probabilities  $V_i, V_{.,j}$  for  $i = 1, 2, \dots, K$  and  $j = 1, 2, \dots, K - 1$ , for  $K$  primary periods.

Secondary sampling occasions

Parameters			1→2	2→4	4→8	8→16
$N$			8%	5%	3%	0.3%
$V_i$			20%	19%	11%	4%
$V_{.,j}$			21%	20%	18%	14%

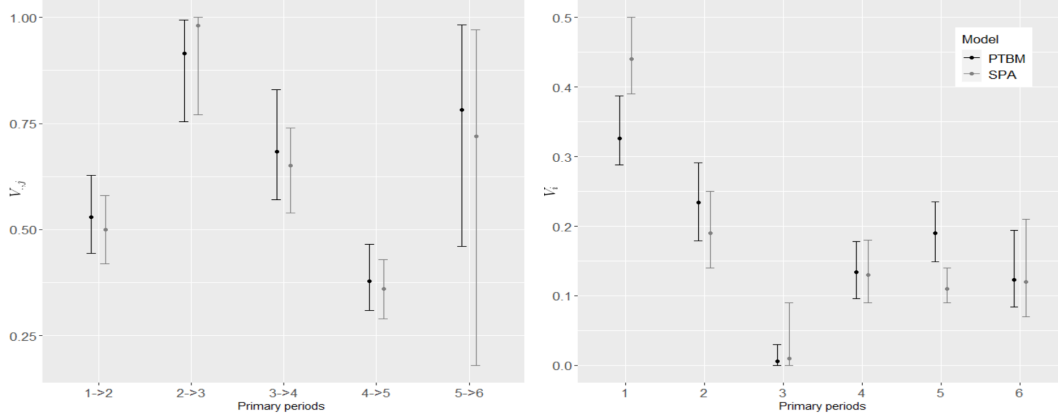


Fig 4: Summaries of exit probabilities  $V_{.,j}$  for  $j = 1, 2, \dots, K - 1$ , (a), and of entry probabilities  $V_i$ ,  $i = 1, 2, \dots, K$ , (b). The SPA model of Zhang, Bonner and McCrea (2023) was fitted classically, so the summaries correspond to the maximum likelihood estimate and corresponding 95% confidence interval in each case, whereas our PTBM and PTBM-R models are fitted in a Bayesian framework, so the summaries correspond to the posterior mean and 95% posterior credible interval.

5.5. *Golden mantella data.* We apply our model to the data presented in Zhang, Bonner and McCrea (2023) and compare the estimates of the entry/exit probabilities (Figure 4). Estimates of the capture probabilities and the number of available individuals can be found in Web Appendix F. We use the same prior distribution choices as the ones displayed in Web Appendix E.

The results demonstrate that the two approaches are in general agreement, but PTBM consistently gives narrower posterior credible intervals than the corresponding confidence intervals obtained by Zhang, Bonner and McCrea (2023) using SPA. We note here that we have chosen flat priors for all model parameters, to make our results as comparable as possible with the SPA results, and hence our interpretation of this difference in interval width is that it is due to the use of the exact, instead of approximate, likelihood in our case.

Figure 4 demonstrates that the exit probabilities are consistently larger for the midpoints of the breeding seasons, which is expected, as also reported by Zhang, Bonner and McCrea (2023). Additionally, the entry probabilities, on average, are higher in the early primary periods within each breeding season. The estimates of the number of individuals present on each primary period are in agreement for SPA and PTBM, with the PTBM having narrower intervals as shown in Figure 13 in Web Appendix F. Lastly, the super-population size is also estimated similarly between the two models, with PTBM leading to narrower posterior credible interval 5407, (5112, 5711) by PTBM, 5567 (5145, 6063) by SPA.

**6. Discussion.** BM surveys are an important monitoring approach for several species and our paper contributes to the, fairly limited thus far, statistical literature for modelling the corresponding data. Our introduced latent grid for modelling the entry and exit pattern, together with the associated PT prior for the grid probabilities, provides a new, general, flexible and computationally efficient modelling framework for BM data. Our model scales with the number of sampling occasions, which is typically much smaller than the number of individuals, as it infers the number of individuals with a specific entry/exit interval, rather than inferring latent individual presence or capture histories, and gives us access to exact inference under different observation processes. The PT model described in this article may appear similar to a Dirichlet-Multinomial model (Royle, Converse and Link, 2012), which is, in fact, a special case of the PT framework when considering the partition described in Section 3. However, it is not precisely the same and is less flexible since we utilize PT machinery, such as the replicate PT, to impose restrictions on the dynamic parameters of the model. The replicate PT structure allows us to place constraints on the model parameters to build parsimonious models that are also ecologically meaningfully.

In our case, we have assumed a replicate PT structure that leads to the assumption that the probability of exit depends on the sampling occasion, and not on the time of entry (time-varying exit). The PT framework gives us access to efficient approaches for model comparisons of different constraints, such as constant, time, or age-varying probabilities (Holmes et al., 2015), and hence future work could explore establishing the required methodology for building and comparing models with different constraints for the model parameters.

Additionally, in the case studies and corresponding simulations, we have employed flat prior distributions for all parameters, to make our results as comparable as possible to the classical models considered thus far in the literature. However, the PT can be centred on parametric distributions, such as the normal, that express our prior expectation of the entry/exit pattern, and in that case inference benefits from both the smoothness of the parametric curve and the flexibility of the nonparametric PT prior.

We considered two case studies with different observation processes to demonstrate the flexibility of our proposed PT prior approach to accommodate different data-generating processes. Our framework can easily be extended to jointly model BM data with capture-recapture, count or other types of ecological data that are often collected in practice. However, an open challenge that remains is the incorporation of covariates in the model parameters, and in particular when these are measured at the individual level, as in that case the grid-approach does not apply, at least not in its current form.

**Acknowledgements.** We are grateful to Professor Rachel S. McCrea, Dr Panagiotis (Takis) Besbeas, Professor Byron J. T. Morgan, Dr Wei Zhang, Dr Simon J. Bonner and Professor Richard Griffiths for providing data, insights and feedback on our analysis.

## REFERENCES

- BUTLER, R. W. (2007). *Saddlepoint Approximations with Applications*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- COWEN, L. L., BESBEAS, P., MORGAN, B. J. and SCHWARZ, C. J. (2014). A comparison of abundance estimates from extended batch-marking and Jolly–Seber-type experiments. *Ecology and Evolution* **4** 210–218.
- COWEN, L. L., BESBEAS, P., MORGAN, B. J. and SCHWARZ, C. J. (2017). Hidden Markov models for extended batch data. *Biometrics* **73** 1321–1331.
- DAVIDSON, J. R., SUDIRMAN, R., WAHID, I., BASKIN, R. N., HASAN, H., ARFAH, A. M., NUR, N., Hidayat, M. Y., SYAFRUDDIN, D. and LOBO, N. F. (2019). Mark-release-recapture studies reveal preferred spatial and temporal behaviors of *Anopheles barbirostris* in West Sulawesi, Indonesia. *Parasites & vectors* **12** 1–11.
- DIANA, A., GRIFFIN, J. and MATECHOU, E. (2019). A Polya tree based model for unmarked individuals in an open wildlife population. In *Bayesian Statistics and New Generations: BAYSM 2018, Warwick, UK, July 2-3 Selected Contributions* 3–11. Springer.

- DIANA, A., MATECHOU, E., GRIFFIN, J., ARNOLD, T., TENAN, S. and VOLPONI, S. (2023). A general modeling framework for open wildlife populations based on the Polya tree prior. *Biometrics* **79** 2171–2183.
- DOLL, J. C., WOOD, C. J., GOODFRED, D. W. and RASH, J. M. (2021). Incorporating batch mark–recapture data into an integrated population model of brown trout. *North American Journal of Fisheries Management* **41** 1390–1407.
- GELFAND, A. E. and SMITH, A. F. (1990). Sampling-based approaches to calculating marginal densities. *Journal of the American statistical association* **85** 398–409.
- HOLMES, C. C., CARON, F., GRIFFIN, J. E. and STEPHENS, D. A. (2015). Two-sample Bayesian nonparametric hypothesis testing.
- HUGGINS, R., WANG, Y. and KEARNS, J. (2010). Analysis of an extended batch marking experiment using estimating equations. *Journal of agricultural, biological, and environmental statistics* **15** 279–289.
- JOLLY, G. M. (1965). Explicit estimates from capture-recapture data with both death and immigration-stochastic model. *Biometrika* **52** 225–247.
- KENDALL, W. L. and POLLOCK, K. H. (1992). The robust design in capture-recapture studies: a review and evaluation by Monte Carlo simulation. *Wildlife 2001: populations* 31–43.
- KING, R. and MCCREA, R. (2019). Capture–Recapture methods and models: estimating population size. In *Handbook of statistics*, **40** 33–83. Elsevier.
- LAVINE, M. (1992). Some aspects of Polya tree distributions for statistical modelling. *The annals of statistics* 1222–1235.
- MATECHOU, E. and CARON, F. (2017). Modelling individual migration patterns using a Bayesian nonparametric approach for capture–recapture data.
- MATECHOU, E., MCCREA, R. S., MORGAN, B. J., NASH, D. J. and GRIFFITHS, R. A. (2016). Open models for removal data. *The Annals of Applied Statistics* **10** 1572–1589.
- MCCREA, R. S. and MORGAN, B. J. (2014). *Analysis of capture-recapture data*. CRC Press.
- POLLOCK, K. H. (1982). A capture-recapture design robust to unequal probability of capture. *The Journal of Wildlife Management* **46** 752–757.
- ROBERT, C. P., CASELLA, G., ROBERT, C. P. and CASELLA, G. (2004). The Metropolis—Hastings algorithm. *Monte Carlo statistical methods* 267–320.
- ROSSER, E., WILLDEN, S. A. and LOEB, G. M. (2022). Effects of SmartWater, a fluorescent mark, on the dispersal, behavior, and biocontrol efficacy of *Phytoseiulus persimilis*. *Experimental and Applied Acarology* **87** 163–174.
- ROYLE, J. A., CONVERSE, S. J. and LINK, W. A. (2012). Data augmentation for hierarchical capture-recapture models. *arXiv preprint arXiv:1211.5706*.
- SCHWARZ, C. J. and ARNASON, A. N. (1996). A general methodology for the analysis of capture-recapture experiments in open populations. *Biometrics* 860–873.
- SEBER, G. A. (1965). A note on the multiple-recapture census. *Biometrika* **52** 249–259.
- SEBER, G. A. F. and SCHOFIELD, M. R. (2019). *Capture-recapture: Parameter estimation for open animal populations*. Springer.
- VAVASSORI, L., SADDLER, A. and MÜLLER, P. (2019). Active dispersal of *Aedes albopictus*: a mark-release-recapture study using self-marking units. *Parasites & vectors* **12** 1–14.
- ZHANG, W., BONNER, S. J. and MCCREA, R. S. (2023). Latent multinomial models for extended batch-mark data. *Biometrics* **79** 2732–2742.

**Supporting Information.** Web Appendix A, referenced in Subsections 3.1, 4.3 and 5.3, Web Appendix B, referenced in Section 2, Web Appendix C, referenced in Subsection 4.4, Web Appendix D, referenced in Subsection 4.5, Web Appendix E, references in Subsections 5.4 and 5.5, Web Appendix F, in Subsection 5.5, are available with this paper at the Biometrics website on Wiley Online Library. The code for fitting all of the models presented in the paper is available on [https://github.com/IoannisRs/BatchMark\\_scripts.git](https://github.com/IoannisRs/BatchMark_scripts.git).