



**HAL**  
open science

# FlywayNet : A hidden semi-Markov model for inferring the structure of migratory bird networks from count data

Sam Nicol, Marie-josée Cros, Nathalie Peyrard, Régis Sabbadin, Ronan Trépos, Richard Fuller, Bradley Woodworth

## ► To cite this version:

Sam Nicol, Marie-josée Cros, Nathalie Peyrard, Régis Sabbadin, Ronan Trépos, et al.. FlywayNet : A hidden semi-Markov model for inferring the structure of migratory bird networks from count data. *Methods in Ecology and Evolution*, 2023, 14 (1), pp.265-279. 10.1111/2041-210X.14011 . hal-04092604

**HAL Id: hal-04092604**

**<https://hal.inrae.fr/hal-04092604>**

Submitted on 9 May 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.



L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

## RESEARCH ARTICLE

# FlywayNet: A hidden semi-Markov model for inferring the structure of migratory bird networks from count data

Sam Nicol<sup>1</sup>  | Marie-Josée Cros<sup>2</sup>  | Nathalie Peyrard<sup>2</sup>  | Régis Sabbadin<sup>2</sup>  |  
Ronan Trépos<sup>2</sup>  | Richard A. Fuller<sup>3</sup>  | Bradley K. Woodworth<sup>3</sup> 

<sup>1</sup>CSIRO Land and Water, Dutton Park, Queensland, Australia

<sup>2</sup>INRAE, UR MIAT, Castanet-Tolosan, France

<sup>3</sup>School of Biological Sciences, The University of Queensland, St. Lucia, Queensland, Australia

## Correspondence

Sam Nicol

Email: [sam.nicol@csiro.au](mailto:sam.nicol@csiro.au)

## Funding information

Agence Nationale de la Recherche, Grant/Award Number: ANR-21-CE40-005; Commonwealth Scientific and Industrial Research Organisation

**Handling Editor:** Marie Auger-Méthé

## Abstract

1. Every year, millions of birds migrate between breeding and nonbreeding habitat, but the relative numbers of animals moving between sites are difficult to observe directly.
2. Here we propose FlywayNet, a discrete network model based on observed count data, to determine the most likely migration links between regions using statistical modelling and efficient inference tools. Our approach advances on previous studies by accounting for noisy observations and flexible stopover durations by modelling using interacting hidden semi-Markov Models. In FlywayNet, individual birds sojourn in stopover nodes for a period of time before moving to other nodes with an unknown probability that we aim to estimate. Exact estimation using existing approaches is not possible, so we designed customised versions of the Monte Carlo expectation-maximisation and approximate Bayesian computation algorithms for our model. We compare the efficiency and quality of estimation of these approaches on synthetic data and an applied case study.
3. Our algorithms performed well on benchmark problems, with low absolute error and strong correlation between estimated and known parameters. On a case study using citizen science count data of the Far Eastern Curlew (*Numenius madagascariensis*), an endangered shorebird from the East Asian–Australasian Flyway, the ABC and MCEM algorithms generated contrasting recommendations due to a difference in optimisation criteria and noise in the data. For ABC, we recovered key features of population-level movements predicted by experts despite the challenges of noisy unstructured data.
4. Understanding connectivity places local conservation efforts and threat mitigation in the global context, yet it has proven difficult to rigorously quantify connectivity at the population level. Our approach provides a flexible framework to infer the structure of migratory networks in birds and other organisms.

## KEYWORDS

approximate Bayesian computation, connectivity, East Asian–Australasian flyway, eBird, hidden semi Markov model, Monte Carlo expectation-maximisation

Sam Nicol, Marie-Josée Cros, Nathalie Peyrard, Régis Sabbadin and Ronan Trépos contributed equally to this study.

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2022 Commonwealth of Australia and The Authors. *Methods in Ecology and Evolution* published by John Wiley & Sons Ltd on behalf of British Ecological Society.

## 1 | INTRODUCTION

The seasonal migration of animals around our planet is one of Earth's great natural spectacles. Apart from the inspiration humans draw from the endurance of migrants who undertake such arduous journeys across our world, migration is critical for connecting ecosystem processes and services across vast distances (Semmens et al., 2011; Wilcove & Wikelski, 2008). Sadly, the phenomenon of migration is threatened and many formerly abundant species are declining globally (Clemens et al., 2016; Rappole & McDonald, 1994; Robbins et al., 1989; Wilcove & Wikelski, 2008).

As well as understanding the drivers of decline, arresting declines in migratory species requires knowing the degree of geographical linkage between different stages of a species' annual range due to the movement trajectories of individuals as they complete their migration, referred to as migratory connectivity (Marra et al., 2019; Webster et al., 2002). Connectivity determines how changes in habitat at one part of a migratory network may influence others. For example, connectivity can explain how poor nonbreeding habitat quality will impact the breeding population (Silllett et al., 2000), the disproportionate impacts of the loss of migratory structure on population size (Iwamura et al., 2013; Runge et al., 2014) or how disease is likely to spread through a migration network (Webster et al., 2002). Understanding connectivity places local conservation efforts and threat mitigation in the global context. For example, if we know the main routes travelled by populations, we can prioritise management of threats in the parts of the flyway that are critical habitat for the largest number of migrants. Connectivity should also be the basis for informed reserve design for migratory species, yet recent analysis suggests that existing reserve networks rarely account for connectivity of migratory bird species across their annual cycle (Runge et al., 2015).

Measuring migratory connectivity is challenging due to the wide geographical areas and vast numbers of individuals involved (Webster et al., 2002). Great progress has been made in recent decades, with sophisticated advances in traditional mark-recapture/banding studies (Cohen et al., 2014), improved satellite tracking and geolocator technology, stable isotope analysis and genetic techniques all providing alternative ways to learn more about where individuals move (Webster et al., 2002). Accompanying these advances is an extensive literature on movement ecology, including models for analysing tracking data, which we do not attempt to review here. Of particular relevance to our study are the works of Joo et al. (2013), which used a hidden semi-Markov Model and tracking data to model foraging behaviour, and Kölzsch et al. (2018) which used tracking data to infer a migratory network structure. Despite these powerful methods and the increasingly clever ways that they are being combined, the expense and low scalability of tracking individuals (Webster et al., 2002) means that in most species our understanding of migratory routes is still drawn from a tiny subsample of the population, often just a few individuals of any given species.

A complementary approach to tracking individuals is to infer connectivity from count data at known aggregation sites. Count data, particularly for birds using the citizen science database eBird

(Sullivan et al., 2009), have been used to complement and boost inferences about connectivity from other methods such as geolocator data (Hallworth et al., 2015) and stable isotope analysis (Fournier et al., 2017). Unlike most tracking data and banding data, count data has the tremendous advantage of being widely and freely available, at least for birds, but increasingly for other organisms (e.g. Tonachella et al., 2012). There is an opportunity to develop more powerful estimation techniques that make greater use of count data in its own right.

Models of migratory networks have been posed to explore theoretical properties of migratory structure using assumed parameters (Taylor & Norris, 2010), and others have used tracking data to infer structure from a few individuals (Kölzsch et al., 2018). However, until recently there have been relatively few attempts to infer migratory network structure using count data, but a handful of studies have attempted to solve the problem for migratory bird networks. The most relevant to our study, (Jain & Dilkina, 2015; Sheldon et al., 2007), infer Markov transition probabilities for a migratory network from eBird data but make the unrealistic assumption of perfect observations. Our study advances previous attempts because it incorporates imperfect detection (i.e. it allows for error in the observed counts) and explicitly estimates the time spent at stopover locations.

As in previous models of migratory networks (e.g. Kölzsch et al. (2018), Jain and Dilkina (2015), Taylor and Norris (2010)), we model a migratory system as a network consisting of nodes (breeding nodes, nonbreeding nodes and stopover nodes) connected by edges. Individuals sojourn in stopover nodes for a period of time before moving to other nodes with an unknown probability that we aim to estimate. From the set of estimated transition probabilities we can reconstruct a weighted network which represents connections between stopovers and their relative strength. The model also enables us to estimate the mean duration of a bird's sojourn at each stopover (hereafter 'sojourn time').

Since animals are difficult to count precisely, to estimate the characteristics of the network, we introduce a hidden semi-Markov modelling (HSMM, Yu (2010), Joo et al. (2013)) approach to model imperfectly detected count data. The hidden part of the model is the position of each bird at each time step. The HSMM is an extension of the well-known hidden Markov model (HMM). The HMM assumes that the sojourn time in a given hidden state follows a geometric distribution; extending the HMM to the HSMM relaxes this assumption and allows explicit modelling of sojourn times. The geometric distribution assumes that the most probable sojourn time is always 1 time unit, which is a limiting assumption for birds that may spend a few weeks resting at sites before continuing their migration. To circumvent this limitation we use a HSMM.

Due to the dimension of the hidden variables, exact estimation of the model parameters using classical approaches is not feasible for even a small number of nodes. To overcome this, we present two dedicated estimation algorithms for our model: Monte Carlo expectation-maximisation (MCEM, Wei & Tanner, 1990) and approximate Bayesian computation (ABC, Csilléry et al., 2010). We present and compare the efficiency and quality of estimation of these

approaches on synthetic data before applying them to a case study of a migratory shorebird in the East Asian–Australasian flyway.

## 2 | MATERIALS AND METHODS

### 2.1 | The model of the migratory system

We assume that we are following the migration of a population of  $N$  birds over a set of  $l$  distinct sites (i.e. breeding, nonbreeding or stopover locations) over time. Sites are connected via migration links ('edges' in the following) for which we have some a priori knowledge, however, we do not know the strength of the connections, and our goal is to learn the most likely structure from count data. We name our model 'FlywayNet'.

#### 2.1.1 | A priori knowledge of the migratory network

We introduce some a priori knowledge on the presence or absence of an edge between two sites. First, since migration is a directed movement (from North to South or from South to North, depending on the season) we assume that birds do not fly backward. Although it is known that some birds do terminate migration and return to their place of origin (e.g. Driscoll & Ueta, 2002), the number of birds returning to their origin sites is very small compared to the number completing their migration. So we assume an ordering of the  $l$  sites such that if  $i < j$  then a bird cannot fly from site  $j$  to site  $i$ . The set of all potential connections is given by the set of oriented edges from  $i$  to  $j$  for every  $i < j$ . This assumption ensures that the graph is acyclic, simplifying the model estimation.

#### 2.1.2 | Semi-Markov model of bird migration

We consider that each bird trajectory is modelled as a semi-Markov model over a finite discrete time horizon  $H = \{0, 1, 2, \dots, T\}$ , and that the  $N$  bird trajectories are independent. The state of a trajectory at a given time can be one of the  $l$  sites, or the state 'death' which corresponds to a bird who dies before time  $T$ . Rigorously, to have a semi-Markov model, one should add the states corresponding to a bird flying towards a given site. Since flight durations are known and fixed, for sake of simplicity, we do not burden the model description with these extra states.

For bird  $n$  ( $1 \leq n \leq N$ ), the trajectory  $\pi_n$  can be summarised by the sequence of visited states and the time of arrival in the state:

$$\pi_n = \left( (i_0^n, t_0^n), (i_1^n, t_1^n), \dots, (i_{F_n}^n, t_{F_n}^n) \right). \quad (1)$$

In expression (1), trajectory  $\pi_n$  has  $F_n$  stages, bird  $n$  starts in site  $i_0^n$  at time  $t_0^n = 0$ , and  $t_k^n$  is the date of arrival of bird  $n$  at site  $i_k^n$ , for every  $1 \leq k \leq F_n$ . By convention,  $i_{F_n}^n$  is the last state occupied by bird  $n$ , that is, the bird entered state  $i_{F_n}^n$  at time  $t_{F_n}^n \leq T$  and is still in this state at time  $T$ .

If for some bird  $n$ , we have  $i_k^n = \text{'death'}$ , then  $t_k^n < T$  represents the date of death of bird  $n$ . In this case,  $\pi_n$  is stopped at  $t_k^n$ .

We assume that for every pair of sites  $(i, j)$  such that  $i < j \leq l$ , the flight duration between  $i$  and  $j$ ,  $f_{ij} \in \{1, 2, \dots\}$  is known. Under these hypotheses, expression (1) defines a unique bird trajectory.

Then, the semi-Markov model for a bird's trajectory is defined as follows:

- Transition probabilities between states. We define the  $l \times l$  matrix  $R$  of transition probabilities between states that are sites. The probability that any bird leaving site  $i$  at any given time goes to site  $j$  is  $R(i, j)$ . If  $i \geq j$  then  $R(i, j) = 0$ , so  $R$  is an upper triangular matrix. Note that, accounting for mortality, we may have, for any  $i < l$ ,

$$\sum_{j \in 1 \dots l} R(i, j) = \sum_{j=i+1}^l R(i, j) < 1.$$

The value  $\mu_i = 1 - \sum_{j \in 1 \dots l} R(i, j)$ , for  $i < l$  is the mortality probability in site  $i$  which is assumed known. For a bird leaving site  $i < l$ , the destination is thus selected according to a categorical distribution of parameters  $(R(i, i+1), \dots, R(i, l), \mu_i)$ . For the breeding site  $l$ , we assume that when a bird 'leaves' site  $l$ , it necessarily moves to death so  $\mu_l = 1$ . This assumption has no influence on the estimation, since the breeding node is the terminal node and we do not estimate sojourn time or transitions from it.

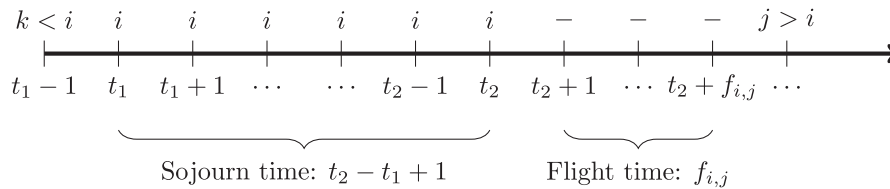
- Sojourn time. We assume that the sojourn time distribution in state  $i \leq l$  is a shifted Poisson distribution of parameter  $\lambda_i$ . The shift is equal to one to ensure that sojourn time is larger than 0 (as is done in the R package `MHSM`, O'Connell & Hojsgaard, 2011). Thus the probability that the sojourn time  $\tau_i$  in site  $i$  is equal to  $d$  is:

$$P_{\lambda_i}(\tau_i = d) = \frac{(\lambda_i)^{d-1}}{(d-1)!} e^{-\lambda_i}, \forall d = 1, 2, \dots$$

The sojourn time in state 'death' is infinite ('death' is an absorbing state). The definition for a sojourn of duration  $d$  is the following: if the state of bird  $n$ 's trajectory at time  $t$  is  $i$  for  $t = t_1, t_1 + 1, \dots, t_2$  and is not state  $i$  at  $t_1 - 1$  and at  $t_2 + 1$ , then  $d = t_2 - t_1 + 1$  (see Figure 1). Sojourn time distributions depend on the site, but are the same for each bird. Furthermore, we assume that sojourn times of two sites are independent.

- Initial distribution.  $\gamma_i^0$  is the distribution of the number of birds at site  $i$  at time zero. Rigorously, in the semi-Markov model framework, if we wanted to track the position of each individual bird, we would need to define a separate probability distribution for the initial position of each bird. However, since our model does not distinguish between birds, it is sufficient to summarise these individual distributions with a single distribution for site  $i$ ,  $\gamma_i^0$ . In this study we assume that the initial distribution is known.

These notions are formally defined in Supporting information S1.



**FIGURE 1** Sojourn time definition. The arrow represents time evolution, and the discrete times are indicated below the arrow. The state of the bird trajectory is indicated above the arrow, with  $k$ ,  $i$  and  $j$  being three distinct sites and  $-$  coding for a bird flying. In this example, the duration of sojourn in state  $i$  is  $t_2 - t_1 + 1$ .

### 2.1.3 | Observation model

Given the set of trajectories  $\Pi = \{\pi_1, \dots, \pi_n\}$ , we can determine  $N_i^t$ , the number of birds located in site  $i$  at time  $t$  (this variable is formally defined in Supporting Information S1).

Observations are observed counts  $O_i^t$  for a set  $\Omega \subseteq \{1 \dots I\} \times \{1 \dots T\}$  of observed site-times. We will consider that, conditional on the birds' trajectories, these counts are independent. Furthermore the distribution of  $O_i^t$  (for  $[i, t] \in \Omega$ ) conditional on  $\Pi$  is equal to the distribution of  $O_i^t$  conditional on  $N_i^t$ . We model it as a negative binomial distribution with parameters  $r_i^t$  and  $p$  where  $r_i^t = \delta N_i^t p / (1 - p)$  and  $\delta$  is the probability to report a bird (see Supporting Information S2 for details and also for other choices of observation model).

Since the observed counts may not be discrete (e.g. where they are averaged across several observers), observations are rounded to the nearest integer value.

The joint distribution of all the observations  $O = \{O_i^t\}_{(i,t) \in \Omega}$  given the trajectories is  $P(O|\Pi) = \prod_{(i,t) \in \Omega} \mathcal{NB}_{(r_i^t, p)}(O_i^t)$ , where  $\mathcal{NB}_{(r_i^t, p)}(\cdot)$  is the negative binomial distribution.

## 2.2 | Parameter interpretation and estimation

Let us denote  $\Lambda = (R, \{\lambda_i\}_{1 \leq i \leq I}, \delta)$  the set of parameters of the HSMM model that we would like to estimate. For the negative binomial observation model, we also estimate parameter  $p$  from data, but prior to the joint estimation of  $\Lambda$ , see Supporting Information S5.

Since we want to infer the most likely network of migration links between nodes, a parameter of interest is the matrix  $R$  of probabilities of transitions between sites. From this matrix, we can build a weighted migratory network where there is an edge from site  $i$  towards site  $j$  if  $R(i, j) > 0$ . The weight of the edge is  $R(i, j)$ . For a fixed  $i$ , the nonzero  $R(i, j)$ s provide the relative importance of the routes  $i \rightarrow j$ . The parameter  $\lambda_i$  indicates the expected duration that a bird stays at site  $i$ .

Estimating the model parameters is difficult for several reasons. First, this is a model with hidden data: neither the individual bird trajectories nor the real counts  $N_i^t$  are observed. Second, conditional on the observed counts, the  $N$  bird trajectories are no longer independent.

Direct optimisation of the likelihood is intractable, yet realisations of the  $O_i^t$  from the model are easy to simulate. Indeed, given parameters  $\Lambda$  we can first simulate each bird's trajectory, then compute the  $N_i^t$  and finally simulate each  $O_i^t$ . We designed two simulation-based methods to estimate the parameters, based on the Monte Carlo

expectation-maximisation method (MCEM, Andrieu et al., 2003) and the approximate Bayesian computation method (ABC, Csilléry et al., 2010) respectively. With MCEM, we obtain a pointwise estimate for each model parameter (frequentist approach) while with ABC we obtain an approximation of the posterior distribution of each parameter (Bayesian approach). The reason to design two estimation algorithms, from different approaches (i.e. frequentist and Bayesian) and with different optimisation criteria, is to help to diagnose the confidence we can have in the estimated parameters, that is, if the algorithms find different parameter values then this should prompt further investigation to understand the cause of the difference.

### 2.2.1 | Monte Carlo expectation-maximisation

The expectation-maximisation (EM) algorithm is an iterative algorithm that maximises the likelihood of the observed data when hidden variables preclude the use of direct maximisation of the likelihood. For our model, for a current value of the estimated parameters,  $\Lambda_{old}$  in the E-step, the conditional probabilities  $P_{\Lambda_{old}}(\Pi|O)$  are computed for all possible sets of trajectories  $\Pi$ . Then in the M-step, the parameter estimator is updated to  $\Lambda_{new}$ :

$$\Lambda_{new} = \operatorname{argmax}_{\Lambda} \sum_{\Pi} \log(P_{\Lambda}(\Pi, O)) P_{\Lambda_{old}}(\Pi|O),$$

where the sum is taken over every possible sets  $\Pi$  of  $N$  independent trajectories. So computing  $\Lambda_{new}$  requires the evaluation of the distribution  $P_{\Lambda_{old}}(\Pi|O)$  which is too complex (E-step). It is possible to approximate the updating formulas using Monte Carlo techniques (Andrieu et al., 2003; Levine & Casella, 2001) by drawing many samples from  $P_{\Lambda_{old}}(\Pi|O)$ . The corresponding updating formulas are equations 6 and 7 in Supporting information S2.

The challenging part of the MCEM approach is therefore to draw samples from  $P_{\Lambda_{old}}(\Pi|O)$ . To do this, we used a Metropolis-Hastings algorithm (Hastings, 1970). This approach, as well as a more complete presentation of the MCEM algorithm, is described in detail in Supporting Information S2.

### 2.2.2 | Approximate Bayesian computation

The idea of an ABC algorithm (Csilléry et al., 2010; Jabot et al., 2013) is to generate parameter values  $\Lambda$  from proposed prior distributions

(or in our case, a particle filter, since we use a more complex version of ABC; see Supporting Information S3), then to generate observations  $O_\Lambda$  for these values of the model parameters. If the simulated observations  $O_\Lambda$  are close to the true observation  $O$  then the parameter values  $\Lambda$  are accepted. The procedure is repeated a large number of times. The histogram of the set of accepted values is then used as an approximation of the true posterior distribution  $P(\Lambda|O)$ .

We used the Lenormand sequential sampling method of the EASYABC package in R (Jabot et al., 2015) to obtain the posterior distribution of every parameter of the model. We selected the set of all observed counts,  $O$ , as the summary statistics. Further details of the ABC algorithm, including the particle filtering algorithm for drawing candidate parameter values, are included in Supporting Information S3.

## 2.3 | Benchmarking

The performance of the MCEM and ABC algorithms was assessed by estimating the model parameters from data simulated from the HSMM model for several networks with known values of  $\Lambda$ . In these experiments the parameter  $p$  of the negative binomial distribution was not estimated but was fixed to its true value. Because the parameters of these benchmark networks are known a priori, we can test the performance of MCEM and ABC by comparing how well they recover the transition probabilities, the sojourn times and the reporting probability given different numbers of nodes and network structures.

The network structure tested during these experiments varied depending on the number of nodes (4–10 nodes) and the maximum number of outgoing edges per node (two to four outgoing edges/node, where outgoing edges are departure routes from a node). We used a range of [1, 3] for generating sojourn times. Five sets of parameter values (transition probabilities and sojourn times) were generated for each structure. The total population of birds was 10,000 and it was distributed equally over the set of nodes that have no incoming edge (source nodes). The number of parameters to estimate ranged from 6 to 20 depending on the problem structure. The total number of generated problems in this benchmark was 300. As well as varying the network structure, we tested the effect of missing observation data on parameter estimation by removing varying proportions of the observed data and re-estimating parameter values.

Performance was assessed quantitatively by comparing the log likelihood of estimated parameters and the mean absolute error (*meanAE*) of estimated parameters rescaled into [0, 1]. Computation of log likelihood as well as the *meanAE* are detailed in the Supporting Information S4. ABC provides an estimate of the posterior distribution of each parameter, so to compute *meanAE* we extracted point estimates from these distributions. Point estimates were represented using the mode of each distribution (the mean being less representative, in particular for nonsymmetric distributions). We compared several methods to estimate the mode of a distribution. Among them, the Lientz function (Lientz, 1972)

and the Venter method (Venter, 1967; both with bandwidth 0.2) returned similar results and led to the lowest *meanAE* values. We selected the Venter method because the bandwidth parameter is easier to interpret.

Wilcoxon tests (Wilcoxon, 1945) were performed on the *meanAE* and log likelihoods obtained for each of the benchmark networks. The tests compared the differences between the results of ABC and MCEM algorithms. The Wilcoxon method was used because we did not want to make any assumption on the distribution of the differences and the pairing option was chosen when we could focus on differences within benchmark problems. Common notations were used when displaying the  $p$  values computed by the test, that is:  $ns$  if  $p > 0.05$ , \* if  $p \leq 0.05$ , \*\* if  $p \leq 0.01$ , \*\*\* if  $p \leq 0.001$  and \*\*\*\* if  $p \leq 0.0001$ .

## 2.4 | Case study: Eastern curlews in the east Asian–Australasian flyway

We applied our model to infer the northward migration of the Eastern Curlew (*Numenius madagascariensis*) population in the East Asian–Australasian Flyway (EAAF). Eastern Curlews are the largest migratory shorebirds in the world, making an annual migration from their breeding grounds in Siberia and Kamchatka through east Asia to their predominantly Australian nonbreeding grounds, before returning to breed. Approximately 80% of the population is estimated to utilise the Yellow Sea during the northern migration Department of the Environment (2015), making the Yellow Sea a critically important stopover site for the species.

The global population of Eastern Curlews was estimated to be 32,000 birds in 2021 (Wetlands International, 2021). The population is declining at a rate of 81% over three generations, leading to the species being listed as Endangered globally (BirdLife International, 2017) and critically endangered in Australia (Department of the Environment, 2015). An identified priority information need is to better quantify the dependence of the species on key migratory staging sites (Garnett et al., 2011).

Individual Eastern Curlews are known to follow different routes on their northward and southward migrations (Minton, Jessop, et al., 2011). To demonstrate our approach, we focus on the northward migration, which is better understood (Minton, Jessop, et al., 2011). Birds depart the Australian nonbreeding grounds in late February and March, with more southerly birds departing and arriving at their destinations earlier. Most birds make a nonstop, long-distance flight to the southern parts of Japan, Korea and the Yellow Sea, in 3–4 weeks, arriving in late March or early April. Birds depart the Yellow Sea and Korean peninsula and arrive on their breeding ground during April and early May. The 10,000 km journey from the southerly Victorian nonbreeding grounds to the breeding grounds is completed in 6–8 weeks, while the shorter trip from the Southeastern (8000 km) and Northwestern (7500 km) Australian nonbreeding areas to the breeding grounds takes roughly 5–6 weeks (Minton, Jessop, et al., 2011). For our case study of northward



migration, we model the first 26 weeks of the year (i.e. 1 January–late June) with a weekly time step.

Although the major migration linkages are known from recaptures and resightings, as well as some satellite tracking and count data, little is known about the timing or the duration that curlews spend at stopover sites (Minton, Jessop, et al., 2011) or how the individual sightings data can be extrapolated to the population level.

In our case study we model the migration network using eight nodes representing the major known stopover regions for Eastern Curlews (see Table 1), connected by 12 edges (Figure 6). Nodes and edges are based on observed sightings and descriptions (Minton, Jessop, et al., 2011; Minton, Wahl, et al., 2011), an expert-derived network (Iwamura et al., 2013), distribution maps (BirdLife International, 2017) and eBird observations for Eastern Curlew. eBird sample data are incomplete and spatially biased (Strimas-Mackey et al., 2020), so to obtain estimates of total observed count ( $O_i^t$ ) within each node  $i$  at time  $t$ , we completed the following steps:

1. Drew approximate node boundaries based on the expert-derived network (Iwamura et al., 2013). The geographical extent of the approximate nodes was defined to capture the Internationally Important Sites designated on the basis of Eastern Curlew numbers (Bamford et al., 2008) and include as many eBird checklists as practicable.
2. Clipped the geographical extent to the intersection of the approximate node boundaries and the Birdlife species distribution maps (BirdLife International, 2017) to give a refined node area.
3. Overlaid these intersected areas with a hexagonal grid (cell size 100km<sup>2</sup>, which roughly coincides with the 10×10 km grid cell sizes used to extrapolate populations in Hansen et al., 2016).
4. Within each hexagon, computed the mean count observed in each hexagon in each week. This step aimed to reduce the impacts of double counting and spatially variable survey effort.
5. Obtained an extrapolated count estimate for the node by assigning the mean count per hexagon (from hexagons with observed records) to all hexagons. The sum of counts over all hexagons was assumed to be an estimate of  $O_i^t$ .

Weekly extrapolated count estimates were extracted for the first 26 weeks of 2018 and 2019. Initial node counts were assigned based

TABLE 1 Description of defined migration network nodes

| Node name | Description             |
|-----------|-------------------------|
| SAUS      | Southern Australia      |
| SE AUS    | Southeastern Australia  |
| NWAUS     | Northwestern Australia  |
| NAUS      | Northern Australia      |
| MSIA-IND  | Malaysia, Indonesia     |
| JPN-SK    | Japan, South Korea      |
| YS-NK     | Yellow Sea, North Korea |
| BREED     | Breeding                |

on expert estimates from Iwamura et al. (2013). See Supporting Information S5 for additional details regarding node and edge definition as well as how eBird count data were assigned to the nodes.

Flight durations between nodes (Supporting Information S5) were estimated using the distance between key aggregation sites in nodes and assuming a migration flight speed of 50km/h (ground speed) consistent with Driscoll and Ueta (2002) (estimated flight speed of 50km/h); and (Minton, Jessop, et al., 2011), Minton et al. (2013) (median tracked speed 50.2 km/h). Flight durations were rounded to the nearest week with a minimum assumed travel time of 1 week.

For our case study, we assume that mortality during the migration is zero. In the absence of mortality estimates during migration for Eastern Curlew, and most other species, it remains an open question whether the impacts of loss of staging habitat impact species directly during the migration, or indirectly through reduced breeding success or survival while at breeding or nonbreeding sites. However, if future analyses are able to determine mortality during migration, it would be simple to include this estimate in our analysis.

Parameter  $p$  of the negative binomial distribution was estimated directly from the data prior to the estimation of the other model parameters with MCEM or ABC (see Supporting Information S5).

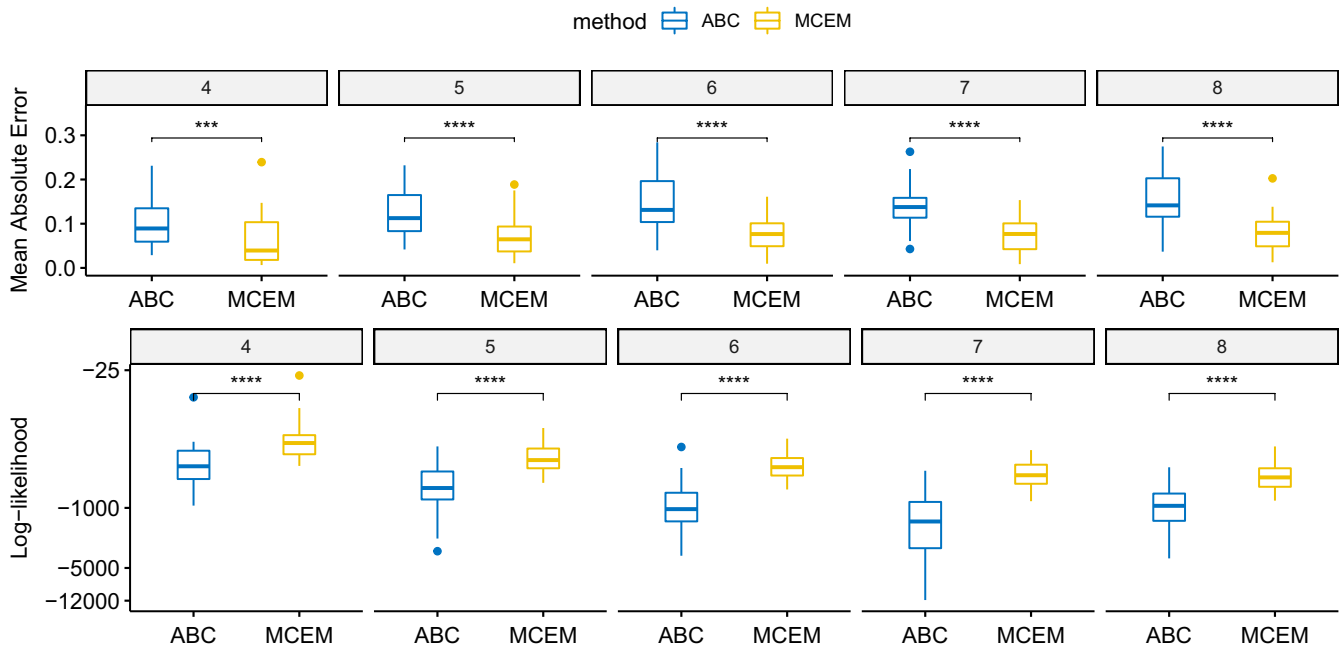
Our objective is to estimate the edge strength between nodes ( $R(i, j)$ ), providing estimates of population connectivity during migration, and the sojourn durations at each node. The uncertainties tested in the case study are the routes taken by birds migrating from Southern, Southeast and Northern Australia—specifically the proportions of internal migration along the eastern and northern Australian nodes and the relative proportions of birds using stopovers in the Yellow Sea compared to those using South Korea and Japan (Figure 6).

## 3 | RESULTS

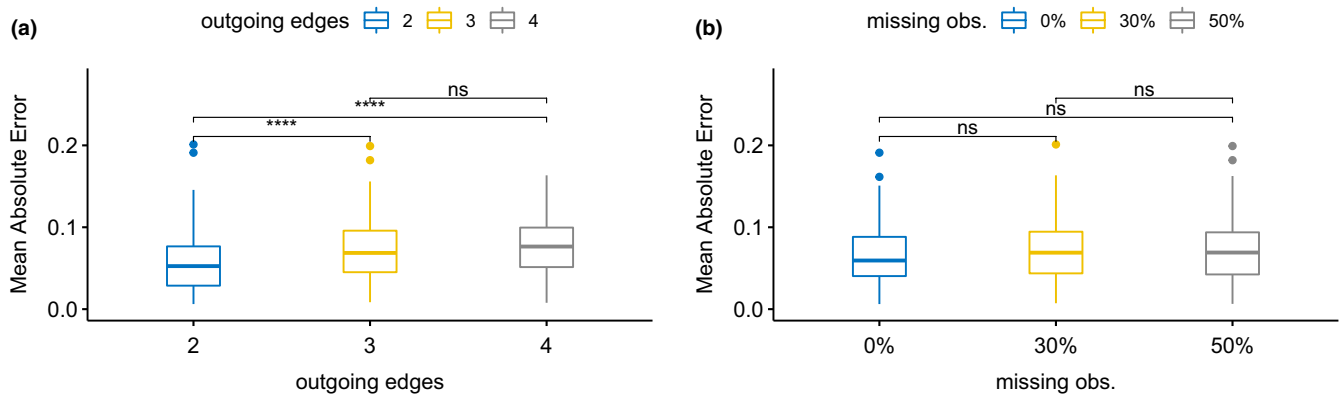
### 3.1 | Benchmarking the MCEM and ABC algorithms

Both algorithms performed well on the benchmark experiments. Across all benchmark problems, MCEM and ABC estimates of the parameters were associated with mean absolute errors of 0.07 and 0.13 respectively. Parameter estimates were well correlated with the true parameters, with a correlation of 0.84 and 0.65 respectively ( $p < 2.2e^{-16}$  for the Pearson's correlation tests). MCEM statistically performed better than ABC. Increasing the number of nodes (Figure 2) and the maximum number of outgoing edges (Figure 3a) increased the error of estimation, however, mean absolute errors remained reasonably low across all benchmark problems.

Missing observations did not substantially impact the quality of estimation (Figure 3b). Transition probabilities were estimated with lower error than sojourn mean times and appeared less sensitive to the number of sites (Figure 4). Error on the estimation of the observation parameter  $\delta$  was close to 0 regardless of the number of sites



**FIGURE 2** Mean absolute error (above), and log likelihood (below) of estimated parameters using ABC and MCEM, according to the number of sites. Mean absolute error is computed using the average error over all sojourn, transition and observation parameters. The statistical tests performed were paired Wilcoxon tests since the differences are based on the same benchmark problems. Significance test symbols \*\*\* and \*\*\*\* refer to a  $p$ -value less than 0.001 and 0.0001 respectively.



**FIGURE 3** Mean absolute error according to the maximal number of outgoing edges (a) and mean absolute error according to the percentage of missing observations (b). Mean absolute error was computed using the combined parameter estimates from the ABC and MCEM algorithms. The statistical tests performed were unpaired Wilcoxon tests since the differences are based on different benchmark problems. Significance test symbols *ns* and \*\*\*\* refer to a  $p$ -value greater than 0.05 and less than 0.0001 respectively.

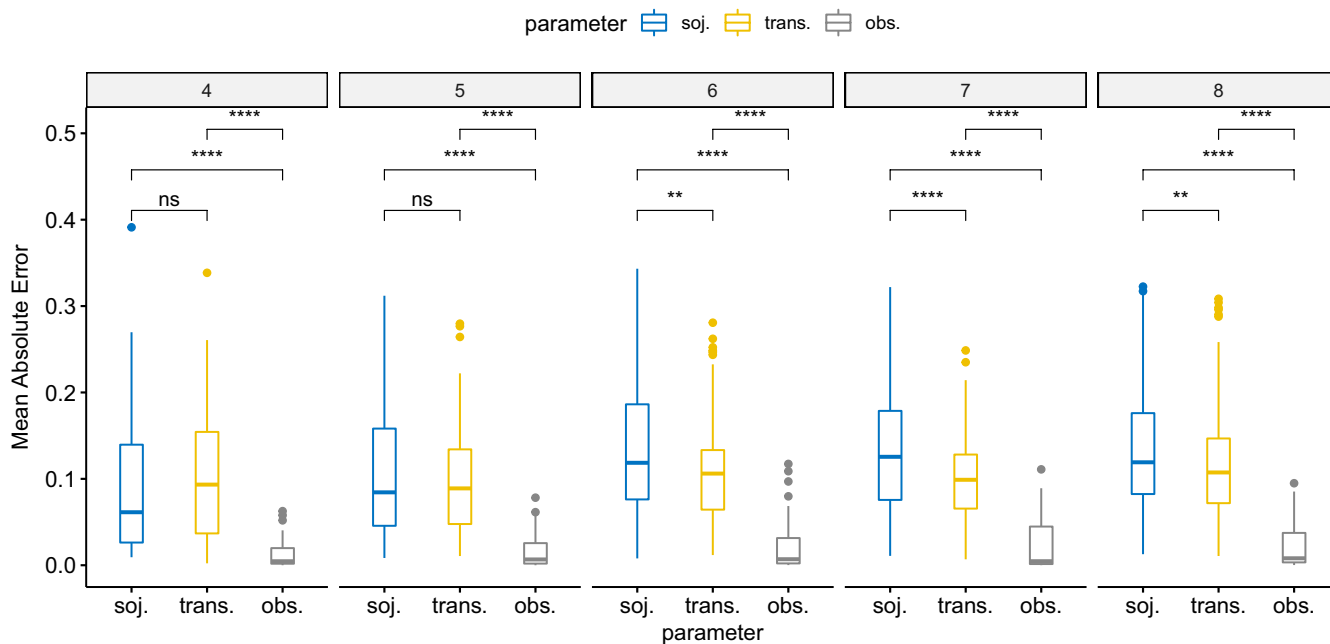
(Figure 4). An additional illustration of the estimation results on a benchmark problem is provided in Supporting information S6.

In terms of number of simulations, MCEM required a median number of 106,650 simulations and a median time of 0.5h to reach convergence per problem, taking into account the five optimisations using different initial values of the parameters. ABC was much more expensive with a median number of 245,000 simulations and a median time of 4h. ABC is expensive due to a high rejection rate of simulated observations. Rejections occur because there is a low probability that a sampled parameter set generates observations that are close to the true observations.

### 3.2 | Case study: The Eastern Curlew

MCEM and ABC estimates of parameters for the migratory network supported different hypotheses about the routes taken by Eastern Curlews (Figures 5 and 6; Venter mode parameter estimates are included in Supporting Information S7). ABC results estimated strong reliance on the Yellow Sea in both 2018 and 2019, with only small proportions visiting the Japan–South Korean node. Although the majority of birds flew directly to north Asia from their origin, ABC estimated that many birds staged in a more northerly Australian node before undertaking their migration, particularly in 2018.





**FIGURE 4** Mean absolute error of estimated sojourn mean time compared to mean absolute error of estimated transition probabilities and the estimated observation parameter  $\delta$  according to the number of sites. Mean absolute error was computed using the combined parameter estimates from the ABC and MCEM algorithms. The statistical tests performed were paired Wilcoxon tests since the differences are based on the same benchmark problems. Significance test symbols ns, \*\* and \*\*\*\* refer to a  $p$ -value greater than 0.05, less than 0.01 and less than 0.0001 respectively.

In 2018, MCEM estimated that Australian birds flew to the Yellow Sea, but in 2019 most birds instead flew to the South Korean–Japanese node. In MCEM, the amount of staging in a more northerly Australian node was much stronger than in ABC. Unlike in ABC, where many birds flew directly to the Yellow Sea from their origin node, in MCEM large majorities of birds (70%–90%) ‘hopped’ north to the Northern Australian node before undertaking their long flight to the Yellow Sea or Japan–South Korea.

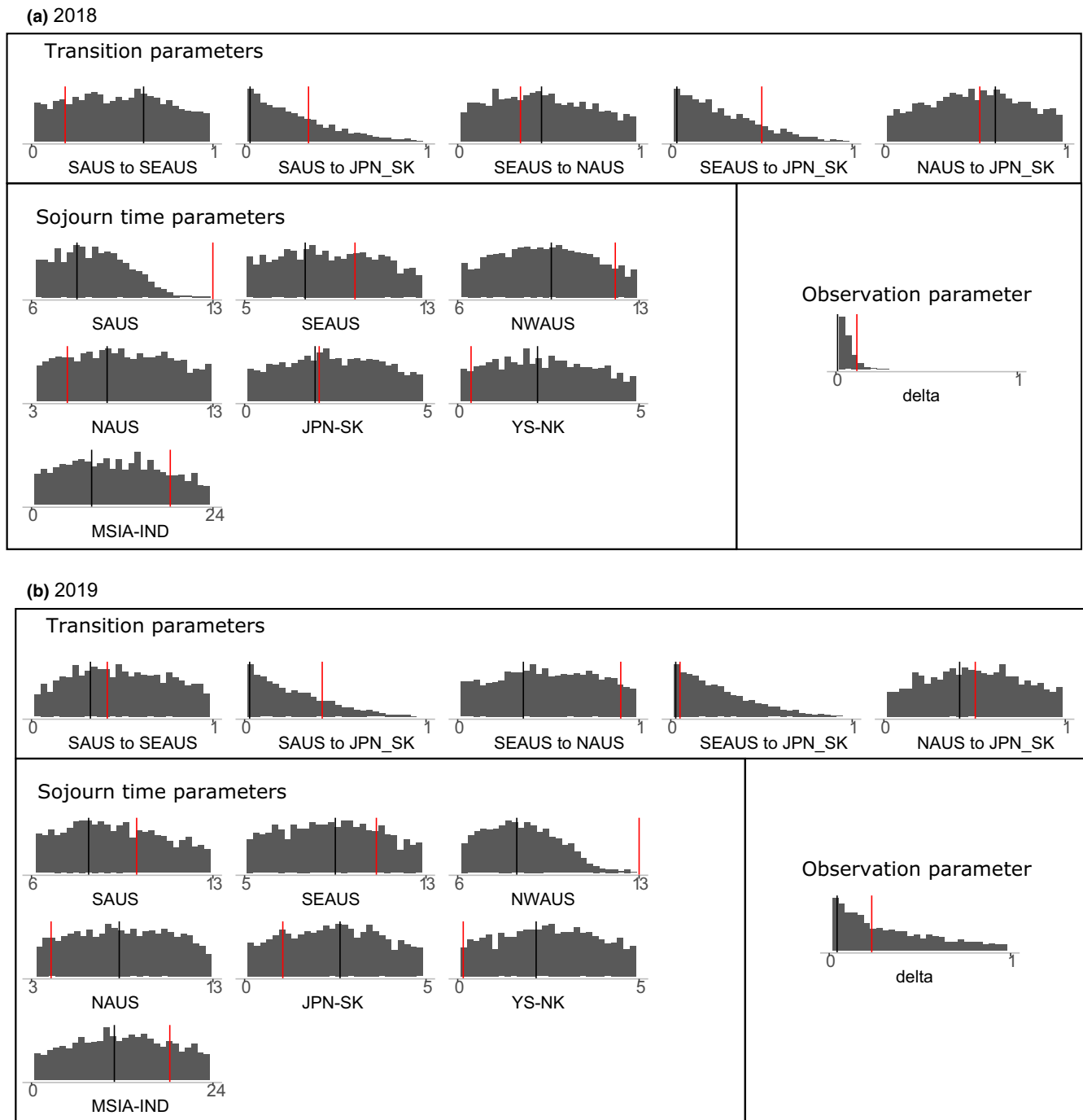
For both algorithms, sojourn time parameter estimates were stable between years. MCEM estimated sojourn times at Southern (13 weeks) and Northwestern Australia (13 weeks) that were longer than expected, suggesting a later start to migration than expected for the species. In contrast, ABC estimates for these two nodes were approximately 8 weeks for both nodes, matching the late February departure expected from observations. MCEM also predicted unusually short sojourn times (~0 weeks) for both the Yellow Sea and Japan–South Korean nodes (ABC predicted 2 weeks for both nodes).

To further highlight the differences between the two algorithms we simulated trajectories using FlywayNet with either MCEM or ABC estimators and computed the number  $N_i^t$  of birds at site  $i$  and time  $t$  from these trajectories (Figure 8). Reasonable departure times were estimated from the ABC simulated trajectories (approximately week 8; February/March). For MCEM, birds began departing Southeastern Australia slightly early (mean departure week 5), and left Southern Australia later than expected (mean departure week 12). Since the majority of birds originate in Southeastern Australia, this had the effect of causing the migration to be shifted earlier for MCEM. Consequently,

combined with very short estimated sojourn times at the Yellow Sea and South Korean–Japanese nodes (mean sojourn durations for both nodes <1 week), MCEM estimated trajectories had very early arrivals at the breeding node (first arrivals late January). With ABC estimators, trajectories were closer to observed Curlew behaviour (Figure 7): there was a peak departure for both Northwestern and Southern Australian nodes near the end of February and a peak in bird numbers at the Yellow Sea in April. The mean migration time for ABC was 7.0 weeks in 2018 and 6.2 weeks in 2019, close to the 6–8 weeks estimated in the literature (Minton, Jessop, et al., 2011); for MCEM it was 5.1 weeks and 5.7 weeks for 2018–19 respectively.

## 4 | DISCUSSION

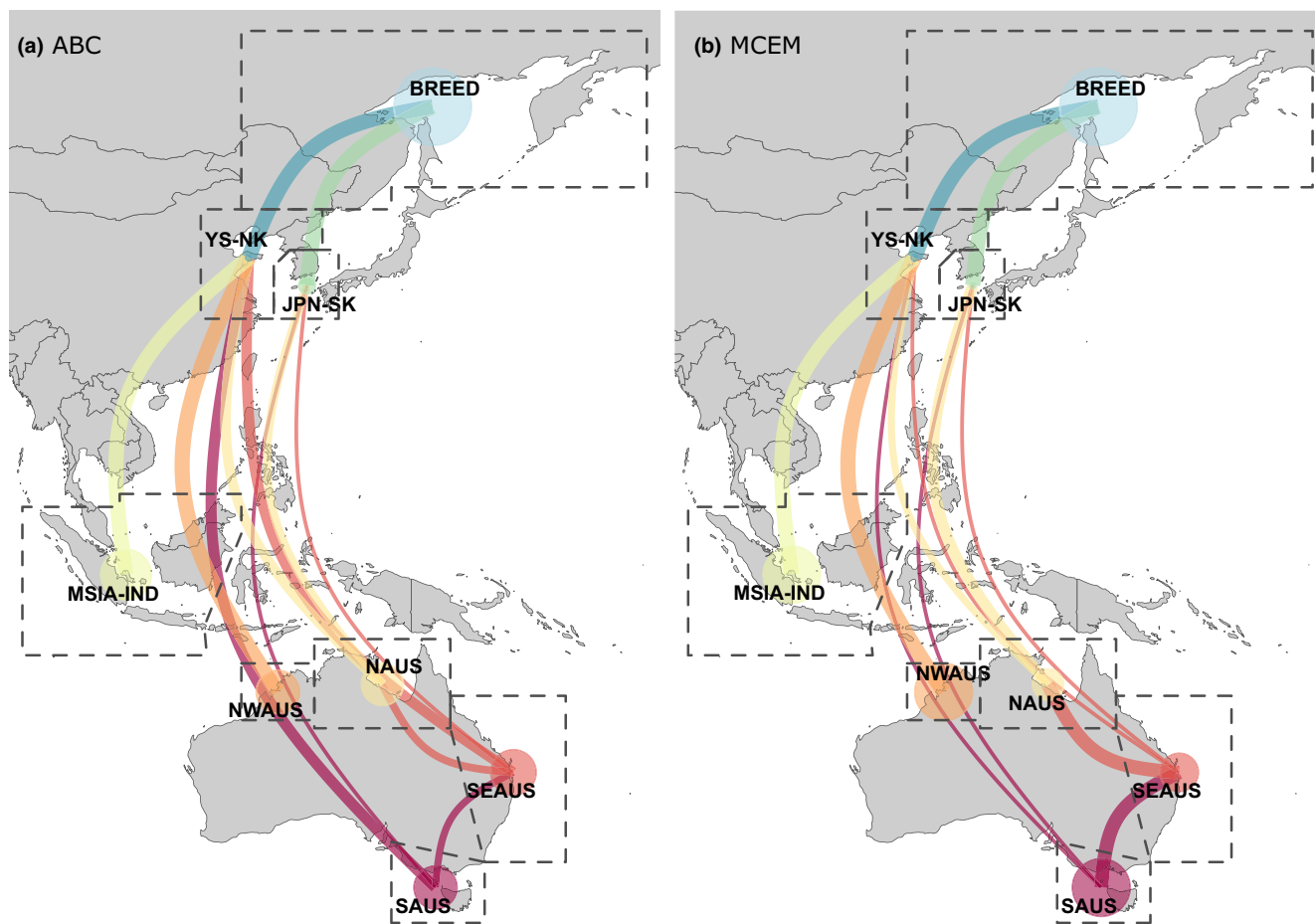
When tested on benchmark problems, both the ABC and MCEM algorithms performed well and recovered parameter values with good accuracy across various sized networks and numbers of connected neighbours. This suggests that, if our HSMM model assumptions hold and sufficiently accurate counts are available, our algorithms provide a powerful way to recover network structure (i.e. the relative importance of edges and the durations of stopovers). Unlike previous approaches (Kölzsch et al., 2018), using this approach enables us to account for error on count data and provides a flexible framework that can make inference with incomplete spatiotemporal count data. In contrast to previous studies which require high-resolution data which is difficult to obtain for species that are not included in formal monitoring programs



**FIGURE 5** The marginals of the posterior distribution for each parameter estimated with ABC for the Curlew problem using the (a) 2018 and (b) 2019 datasets; the black lines represent the Venter mode of the marginals of the posterior distribution. Red vertical lines represent the estimated parameters computed by MCEM.

(Jain & Dilkina, 2015; Kölzsch et al., 2018), our approach requires only a basic network structure and count information, which is more widely available than individual trajectories (e.g. from satellite tagging of birds), providing a useful complementary source of inference to traditional bird tracking studies. Count data are the default method of data collection for bird watchers globally, so harnessing this data source is a powerful way to make best use of a global citizen science network.

Although our algorithms worked well for connectivity estimation on benchmark problems, we obtained contrasting results when we attempted to estimate real Eastern Curlew networks from eBird data. In contrast to our benchmark testing results where MCEM had lower estimation error, the ABC results for the case study appeared to better match existing knowledge about curlew movements, particularly in terms of reproducing the dependence on the Yellow Sea rather than South Korea–Japan. This could be partly due to how



**FIGURE 6** Estimated Eastern Curlew networks using 2019 eBird records and (a) ABC and (b) MCEM algorithms. Edge widths depict the relative transition probabilities between nodes; node sizes represent relative sojourn time lengths. Dotted lines depict the node boundaries. Colours depict edges from the same origin.

the node boundaries are selected, however, the different estimates make it hard to conclude which algorithm best estimates the true migratory behaviour of Eastern Curlews.

We have two main hypotheses that could explain the differences between the algorithm estimations for Eastern Curlews. First, they have different objective functions: MCEM maximises the likelihood while ABC optimises a customised set of statistics (here we minimise a weighted sum squared error between observed and simulated data). We investigated this hypothesis by changing the ABC acceptance criterion to better match the MCEM likelihood and obtained results that were closer to the MCEM estimates (see Supporting Information S8). This suggested that some of the difference between algorithms may be due to different optimisation criteria, however, it does not suggest which results are closer to the real bird network dynamics.

Second, the eBird count data could be too noisy for the algorithms to reach a common estimation. Given that our benchmark performances were similar despite their different optimisation criteria, we believe that this is the most probable explanation for the discrepancy between algorithms. Although some nodes had substantial observed count data and we used the best available estimates of suitable range to develop weekly count estimates, geographical and

temporal coverage is variable in all nodes and our node abundance estimates had high week-to-week variability for all nodes (Figure 8, top panel). The variability was evident in the posterior distributions of the ABC results, which were relatively flat for several nodes (Figure 5). Furthermore, our model assumes that the initial population size at each node is known, but we drew this information from expert knowledge rather than count data. We pursued data from the International Waterbird Census (IWC, Delany, 2005), which is a highly promising dataset since it contains systematic count data recorded at the same time of year (January; which roughly corresponds to the beginning of our simulation period in our migration model). However, the IWC data have incomplete spatial coverage and variable survey effort in different areas, so additional research is needed to use IWC data to generate node abundance estimates that could be used in this study. Further research to better estimate node abundance from count data would likely greatly improve our estimation ability. Potentially useful methods may include smoothing the weekly observed data to infer observations at missing locations (e.g. Sheldon et al., 2007) or clustering to test the locations of the node boundaries (Jain & Dilkina, 2015). Although we deliberately tested how well we could estimate using very minimal data from eBird, another promising approach may be to incorporate additional

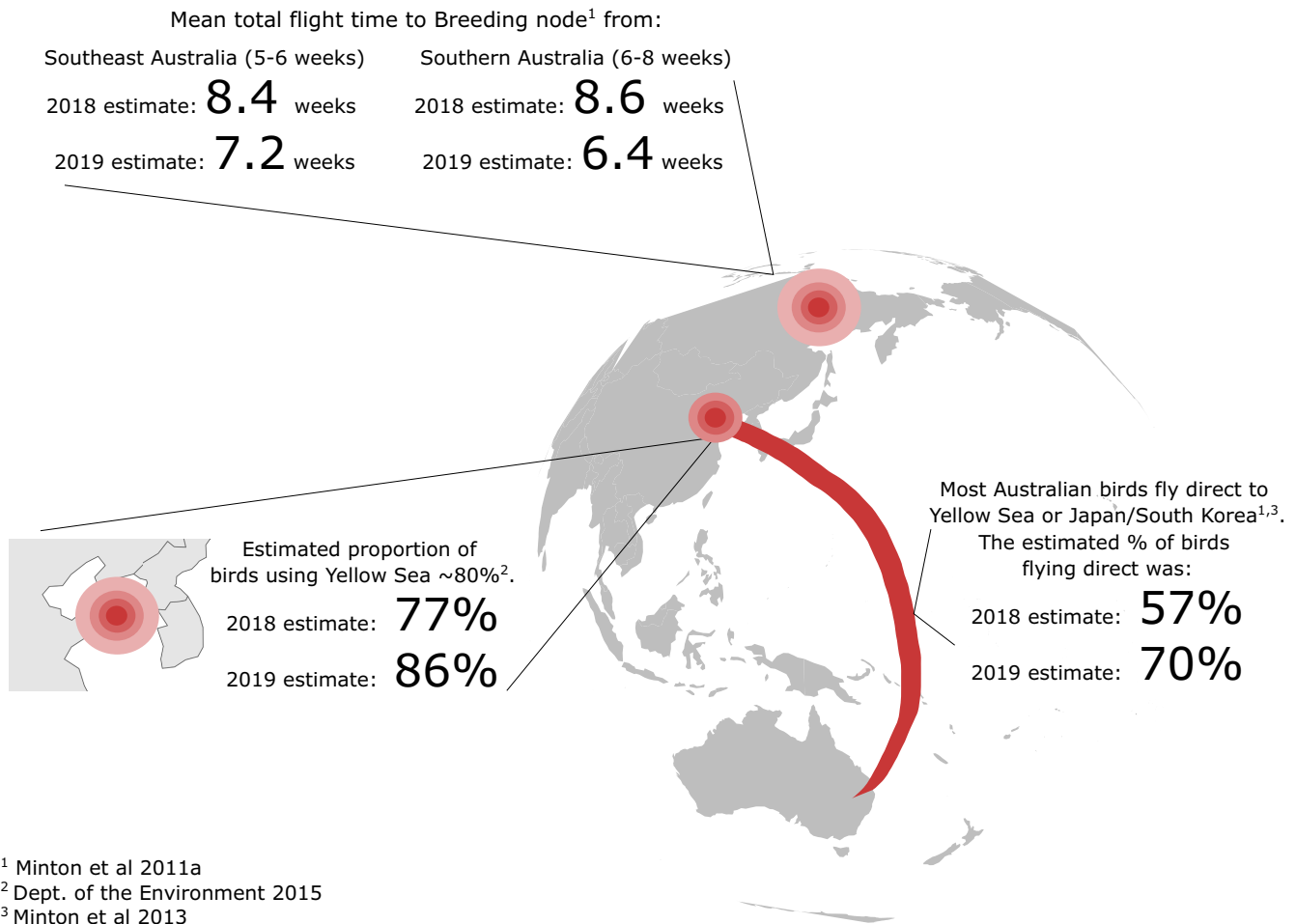


FIGURE 7 Comparison of estimated ABC results with values derived from the literature.

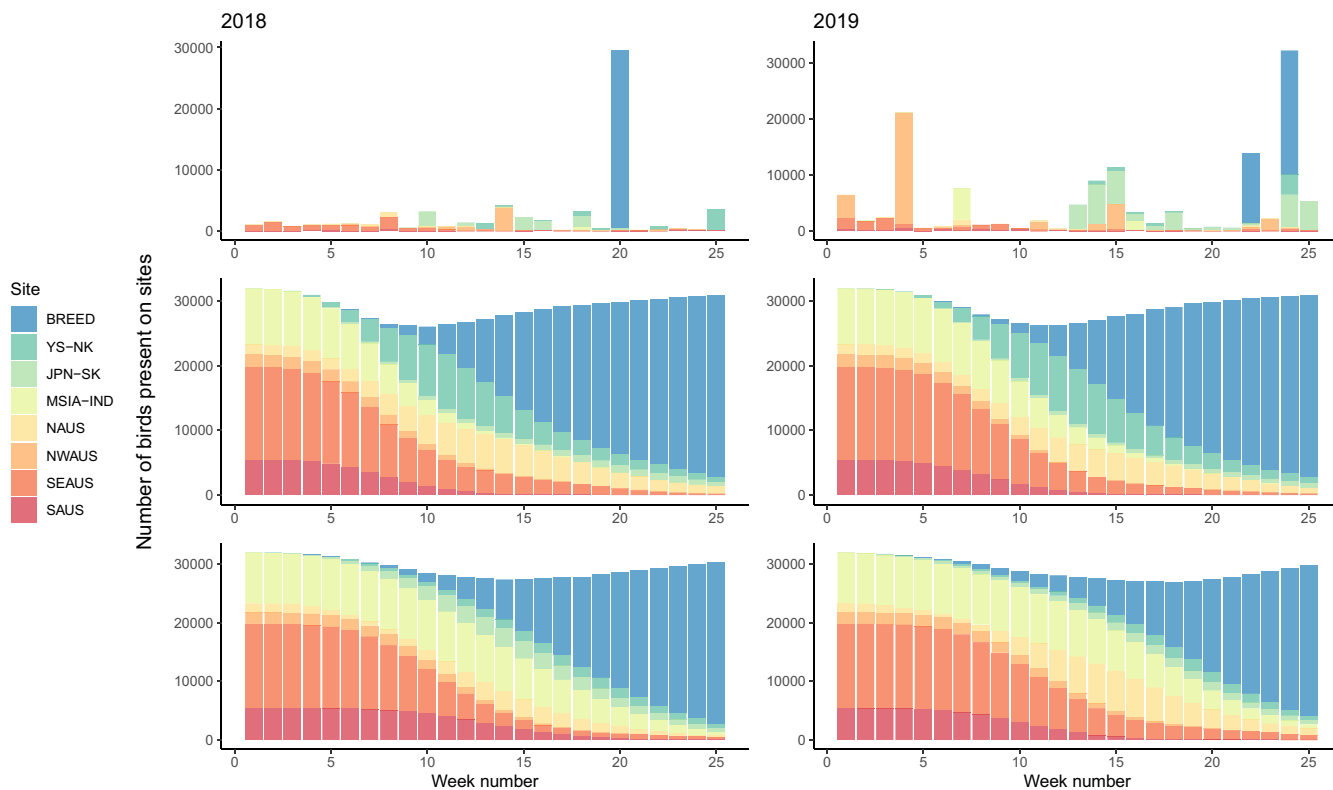
datasets, such as tracking and/or banding data, to guide the simulated trajectories. Formally incorporating environmental covariates such as habitat type, temperature or the results of species distribution modelling may provide additional information to improve abundance estimates. Including covariates would also be useful for predicting other parameters, most notably sojourn times, since birds make stopover decisions based on environmental conditions such as the time of year and wind conditions (Kölzsch et al., 2016). It may also be useful to try the method on other well-documented animal migratory data, particularly on species which are easier to track and count, such as ungulate migration (Convention on Migratory Species Secretariat, 2021; Sawyer et al., 2009), where observation errors may be lower than for shorebirds.

A key question for application is how to diagnose when the algorithms are performing well. Clearly, if the two algorithms estimate very different parameter values, users should seek to diagnose the cause of the difference. However, there may also be other useful indicators of reliability. For example, in the curlew problem, MCEM tended to seek the boundaries of its domain, suggesting that the likelihood surface is increasing monotonically (and therefore that it is unlikely that parameter estimates will be reliable). The posterior distributions estimated by ABC for some parameters were relatively

flat, which also acts as a simple check to test for parameters that are difficult to estimate. Where users find that MCEM estimates a majority of parameter values on the edge of their domain and ABC finds numerous 'flat' posterior distributions, we recommend reviewing the data quality rather than accepting parameter estimates.

For our curlew case study, the difference in the results between algorithms means that the findings should be interpreted with reference to the general movement patterns of the species. Expert knowledge, such as that used here (e.g. Figure 7), should be used to verify the predictions of the algorithms. The general migratory behaviour of most bird species is known, and if an algorithm does not reproduce this behaviour then it can be said to be performing poorly. Where the algorithm results do align with known movement behaviour, the tests outlined in the previous paragraph are a useful guide. In particular, for our curlew case study, we suggest that although the ABC results are encouraging, some parameters have flat marginal posterior distributions (Figure 5) and these should be further scrutinised before being used in applied conservation.

The migratory movement patterns of some bird species are poorly known. Our model accommodates this by requiring minimal input information, specifically: count data, the location boundaries of the nodes and the suspected connections between nodes. Count



**FIGURE 8** Simulated weekly trajectories for the Eastern Curlew in 2018 (left) and 2019 (right). Plots show (top) observed counts, (middle) trajectories simulated using ABC estimators and (bottom) trajectories simulated using MCEM estimators.

data are readily available for any species via eBird, so this should not be a limiting factor except where there are few observations of the species recorded. The locations of the nodes and the hypothesised connections between nodes can be obtained from a combination of eBird list locations, published studies and expert knowledge (see Supporting Information S5). For poorly known species where studies and expert knowledge are lacking, eBird data alone could be used to set the hypothesised network structure, although further research would be valuable to determine the most robust way to set node boundaries (e.g. clustering techniques).

Our model makes some assumptions that could be improved in future iterations. First, we assumed that the time spent in a node is the same regardless of the origin and destination (i.e. sojourn times are independent of the origin and destination). Where nodes are both a nonbreeding origin node as well as a stopover node (e.g. SEAUS in our model), this may confound sojourn times that occur between week zero and when the migration starts with sojourn times due to true stopovers during migration. If these sojourn times are considerably different, then sojourn estimation may be affected. Modelling conditional sojourn times is possible within our HSMM framework, however, it would increase the number of parameters at each node and increase the difficulty of estimation.

Second, we assumed that sojourn times at different sites are independent. Strictly, since we have an idea of the total duration of the migration, the sum of the sojourn times along a trajectory should not substantially exceed the expected total duration, so the

independence assumption may not hold. We expect that in practice the observed data will minimise the impacts of this assumption by enforcing average movement between nodes at reasonable times (i.e. that match the observed movement times), even without explicitly modelling dependence between nodes. There could be some trajectories that are overly long or short due to independent sojourn times, but these should be minimised by fitting to observed movements.

Third, we assumed that migratory birds progress in one direction (northward migration only) and that migration time between nodes was constant. In practice, some birds are known to abort migrations and return to their origin (Driscoll and Ueta (2002)), but for modelling reasons we assumed that this proportion was small. We also only modelled the northward migration; it would be theoretically trivial to model the full annual cycle of migration, however, doing so would increase the number of parameters that need to be estimated. It may be more practical to model northward and southward migrations separately as we have done here.

Fourth, eBird observations (and count data in general) will tend to be underestimates of the true population, since at any time it is unlikely that an observer will see and record all the birds in an area. This creates the possibility of systematic bias in the count data, which is not explicitly captured in our negative binomial observation model. However, there are other sources of error in the node abundance estimates, most notably the extrapolation process used to estimate counts in areas of the node where no

lists have been recorded. For the curlew study, the extrapolation process was likely to dwarf the errors caused by underestimation due to the area of extrapolation required, minimising the effects of systematic bias. However, for other studies where minimal extrapolation is required, further attention (and perhaps alternative observation distributions) may be required to deal with systematic under-counting.

In a situation where we are sure that the variable used for modelling the observed count is an underestimation of the real count, we should use a Binomial distribution,  $B(N_i, p_i)$ , with  $N_i$  the true count, and  $p_i$  the probability to see a bird (this formulation is included in Supporting Information S1). However, the drawback of the Binomial formulation is that we do not avoid overdispersion with this distribution. Since we are not certain that the node count estimates are underestimates, we use the negative Binomial to manage overdispersion.

The ABC algorithm estimates marginal modes for each parameter, but strictly speaking, the multivariate posterior mode is most comparable to the MCEM estimate. The multivariate mode was not used because it is more computationally expensive to generate (requires estimation of a multivariate kernel density function and optimisation of the density function) than the straightforward computation required to generate the Venter mode of the marginals. Other methods have also used the mode of the marginals to represent the ABC posterior (e.g. Nunes and Prangle, 2015). We are also confident that posterior is 'nice' enough to be summarised by the mode of marginals (see e.g. fig. 2 in Supporting Information S6), at least for the benchmark experiments.

Both algorithms became time-consuming to run as the networks became complex, particularly for the curlew case study. Runtime may limit performance on large networks, so it may be beneficial to investigate alternative methods to estimate the network connectivity. Variational EM in a frequentist approach (VEM Neal & Hinton, 2000) or Bayes expectation-maximisation (VBEM Beal, 2003) in a Bayesian approach may be interesting solutions for a trade-off between runtime and the quality of estimators. Instead of relying on simulations, variational approaches perform estimation by replacing the complex distribution (here the HSMM model) by a closer model in a family of tractable distributions. We are currently investigating whether VEM or VBEM can be used to solve our migratory network problem.

## 5 | CONCLUSIONS

Understanding how migratory populations move is crucial because it allows us to design conservation measures accordingly. Here we have developed a new way to estimate the connectivity of migratory populations based only on limited count data at discrete locations. The method accounts for observation error and predicts both migratory structure and sojourn times. Although information about migratory connectivity can be inferred from individual tracking studies, few studies have attempted to extrapolate individual behaviours to population-level movements. Our study complements existing

tracking work by providing a statistical model to exploit the most commonly collected form of bird data. Although questions remain about how best to estimate node abundance, our approach has tremendous promise because of the explosion in availability of citizen science count data through platforms such as eBird. As these datasets grow, existing geographical and temporal gaps in the datasets will be filled. As this happens, there will be increasing demand for algorithms that are sufficiently flexible to draw inference from unstructured data.

## AUTHOR CONTRIBUTIONS

Sam Nicol, Marie-Josée Cros, Nathalie Peyrard, Régis Sabbadin and Ronan Trépos contributed equally to this manuscript. Specifically, Nathalie Peyrard, Régis Sabbadin and Sam Nicol conceived and designed the HSMM model. Ronan Trépos implemented the MCEM and ABC algorithms and designed benchmarks. Marie-Josée Cros implemented the case study. Sam Nicol extracted the data for the curlew case study and wrote the draft manuscript. Richard Fuller and Bradley Woodworth provided advice on the curlew case study. All authors edited the manuscript and provided critical comments.

## ACKNOWLEDGEMENTS

The authors acknowledge the computational resources provided by RECORD platform (Bergez et al., 2013). S.N. was supported by a CSIRO Julius Career Award. N.P., R.S., M.J.C. and R.T. acknowledge the support of the French Agence Nationale de la Recherche (ANR), under grant ANR-21-CE40-005 (project HSMM-INCA).

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## PEER REVIEW

The peer review history for this article is available at <https://publons.com/publon/10.1111/2041-210X.14011>.

## DATA AVAILABILITY STATEMENT

The FLYWAYNET package (Trépos et al., 2022), including installation instructions and a vignette demonstrating a minimal example, is available from the INRAE git repository at: <https://doi.org/10.5281/zenodo.7156292>. Results and additional materials to run the study are available on figshare: <https://doi.org/10.6084/m9.figshare.16658185.v5>. Specifically, the figshare contains two files:

- 'Flywaynet\_experiments.zip' contains the results of the benchmarking and curlew experiments.
- 'CurlewCaseStudy\_GetWeeklyObservations.tar.gz' contains the scripts and data used to convert the raw eBird eastern curlew counts into weekly observation data, as outlined in Supporting Information S5.

Further information about each of the files and how to reproduce the results is contained in the readme files contained with the downloads. Scripts to re-run the benchmarking experiments and



the curlew case study are available from: [https://forgemia.inra.fr/birdnet/FlywayNet\\_experiments](https://forgemia.inra.fr/birdnet/FlywayNet_experiments). Note that re-running the scripts is time-consuming without substantial computing resources—if users want to run their own example or view the results, we instead recommend using the FLYWAYNET R package or viewing the figshare repository respectively.

## ORCID

Sam Nicol  <https://orcid.org/0000-0002-1160-7444>

Marie-Josée Cros  <https://orcid.org/0000-0002-6395-5563>

Nathalie Peyrard  <https://orcid.org/0000-0002-0356-1255>

Régis Sabbadin  <https://orcid.org/0000-0002-6286-1821>

Ronan Trépos  <https://orcid.org/0000-0002-3338-9337>

Richard A. Fuller  <https://orcid.org/0000-0001-9468-9678>

Bradley K. Woodworth  <https://orcid.org/0000-0002-4528-8250>

## REFERENCES

- Andrieu, C., Freitas, N. D., Doucet, A., & Jordan, M. (2003). An introduction to MCMC for machine learning. *Machine Learning*, 50, 5–43.
- Bamford, M., Watkins, D., Bancroft, W., Tischler, G., & Wahl, J. (2008). *Migratory shorebirds of the east Asian-Australasian flyway: Population estimates and internationally important sites*, Technical report. Wetlands International-Oceania.
- Beal, M. (2003). *Variational algorithms for approximate Bayesian inference*, Ph.d. thesis. Gatsby Computational Neuroscience Unit, University College London.
- Bergez, J.-E., Chabrier, P., Gary, C., Jeuffroy, M., Makowski, D., Quesnel, G., Ramat, E., Raynal, H., Rousse, N., Wallach, D., Debaeke, P., Durand, P., Duru, M., Dury, J., Faverdin, P., Gascuel-Oudou, C., & Garcia, F. (2013). An open platform to build, evaluate and simulate integrated models of farming and agro-ecosystems. *Environmental Modelling and Software*, 39(1), 39–49.
- BirdLife International. (2017). *Numenius madagascariensis* (amended version of 2016 assessment). the IUCN Red List of Threatened Species 2017: e.t22693199a118601473, Technical report. <https://doi.org/10.2305/iucn.uk.2017-3.rtls.t22693199a118601473.en>
- Clemens, R., Rogers, D. I., Hansen, B. D., Gosbell, K., Minton, C. D. T., Straw, P., Bamford, M., Woehler, E. J., Milton, D. A., Weston, M. A., Venables, B., Wellet, D., Hassell, C., Rutherford, B., Onton, K., Herrod, A., Studds, C. E., Choi, C.-Y., Dhanjal-Adams, K. L., ... Fuller, R. A. (2016). Continental-scale decreases in shorebird populations in Australia. *Emu - Austral Ornithology*, 116(2), 119–135.
- Cohen, E. B., Hostetler, J. A., Royle, J. A., & Marra, P. P. (2014). Estimating migratory connectivity of birds when re-encounter probabilities are heterogeneous. *Ecology and Evolution*, 4(9), 1659–1670.
- Convention on Migratory Species Secretariat. (2021). *Global initiative on ungulate migration*. Secretariat of the Convention on the Conservation of Migratory Species of Wild Animals (CMS). <https://www.cms.int/en/gium>
- Csilléry, K., Blum, M. G. B., Gaggiotti, O. E., & François, O. (2010). Approximate Bayesian computation (ABC) in practice. *Trends in Ecology & Evolution*, 25(7), 410–418.
- Delany, S. (2005). *Guidelines for participants in the international waterbird census (IWC)*, Technical report. Wetlands International.
- Department of the Environment. (2015). *Conservation advice Numenius madagascariensis eastern curlew* Technical report. Department of the Environment.
- Driscoll, P. V., & Ueta, M. (2002). The migration route and behaviour of eastern curlews *Numenius madagascariensis*. *Ibis*, 144(3), E119–E130.
- Fournier, A. M. V., Sullivan, A. R., Bump, J. K., Perkins, M., Shieldcastle, M. C., & King, S. L. (2017). Combining citizen science species distribution models and stable isotopes reveals migratory connectivity in the secretive Virginia rail. *Journal of Applied Ecology*, 54(2), 618–627.
- Garnett, S., Szabo, J., & Dutton, G. (2011). *The action plan for Australian birds 2010*, Birds Australia. CSIRO Publishing.
- Hallworth, M. T., Sillett, T. S., Van Wilgenburg, S. L., Hobson, K. A., & Marra, P. P. (2015). Migratory connectivity of a neotropical migratory songbird revealed by archival light-level geolocators. *Ecological Applications*, 25(2), 336–347.
- Hansen, B., Fuller, R., Watkins, D., Rogers, D., Clemens, R., Newman, M., Woehler, E., & Weller, D. (2016). *Revision of the east Asian-Australasia flyway population estimate for 37 listed migratory shorebird species*. Unpublished report for the Department of the Environment, Technical report. BirdLife Australia.
- Hastings, W. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1), 97–109.
- Iwamura, T., Possingham, H. P., Chadés, I., Minton, C., Murray, N. J., Rogers, D. I., Treml, E. A., & Fuller, R. A. (2013). Migratory connectivity magnifies the consequences of habitat loss from sea-level rise for shorebird populations. *Proceedings of the Royal Society B: Biological Sciences*, 280(1761), 20130325.
- Jabot, F., Faure, T., & Dumoulin, N. (2013). EasyABC: Performing efficient approximate Bayesian computation sampling schemes using R. *Methods in Ecology and Evolution*, 4(7), 684–687.
- Jabot, F., Faure, T., Dumoulin, N., & Albert, C. (2015). EasyABC: Efficient approximate Bayesian computation sampling schemes. R package version 1.5. <https://CRAN.R-project.org/package=EasyABC>
- Jain, N., & Dilkina, B. (2015). Coarse models for bird migrations using clustering and non-stationary Markov chains. In *AAAI workshop: Computational sustainability* (pp. 63–68). Association for Artificial Intelligence (AAAI). Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI-15). January 25–30, 2015, Austin, TX USA. <http://aaai.org/ocs/index.php/WS/AAAIW15/paper/view/10209>
- Joo, R., Bertrand, S., Tam, J., & Fablet, R. (2013). Hidden Markov models: The best models for forager movements? *PLoS ONE*, 8(8), e71246.
- Kölzsch, A., Kleyheeg, E., Kruckenberg, H., Kaatz, M., & Blasius, B. (2018). A periodic Markov model to formalize animal migration on a network. *Royal Society Open Science*, 5(6), 180438.
- Kölzsch, A., Müskens, G., Kruckenberg, H., Glazov, P., Weinzierl, R., Nolet, B., & Wikelski, M. (2016). Towards a new understanding of migration timing: Slower spring than autumn migration in geese reflects differing decision rules for stopover use and departure. *Oikos*, 125(10), 1496–1507.
- Levine, R. A., & Casella, G. (2001). Implementations of the Monte Carlo EM algorithm. *Journal of Computational and Graphical Statistics*, 10(3), 422–439.
- Lientz, B. (1972). Properties of modal intervals. *SIAM Journal on Applied Mathematics*, 23, 1–5.
- Marra, P. P., Cohen, E., Harrison, A.-L., Studds, C. E., & Webster, M. S. (2019). *Migratory connectivity* (2nd ed., pp. 455–461). Academic Press.
- Minton, C., Gosbell, K., Johns, P., Christie, M., Klaassen, M., Hassell, C., Boyle, A., Jessop, R., & Fox, J. (2013). New insights from geolocators deployed on waders in Australia. *Wader Study Group Bulletin*, 120, 37–46.
- Minton, C., Jessop, R., Collings, P., & Standen, R. (2011). The migration of eastern curlew *Numenius Madagascariensis* to and from Australia. *Stilt*, 59, 6–16.
- Minton, C., Wahl, J., Gibbs, H., Jessop, R., Hassell, C., & Boyle, A. (2011). Recoveries and flag sightings of waders which spend the non-breeding season in Australia. *Stilt*, 59, 17–43.
- Neal, R., & Hinton, G. (2000). A view of the EM algorithm that justifies incremental, sparse, and other variants. *Learning in Graphical Models*, 89, 355–368.

- Nunes, M. A., & Prangle, D. (2015). Abctools: An R package for tuning approximate Bayesian computation analyses. *The R Journal*, 7(2), 189–205.
- O'Connell, J., & Hojsgaard, S. (2011). Hidden semi Markov models for multiple observation sequences: The mhsmm package for R. *Journal of Statistical Software*, 39(4), 1–22.
- Rappole, J. H., & McDonald, M. V. (1994). Cause and effect in population declines of migratory birds. *The Auk*, 111(3), 652–660.
- Robbins, C. S., Sauer, J. R., Greenberg, R. S., & Droege, S. (1989). Population declines in north American birds that migrate to the neotropics. *Proceedings of the National Academy of Sciences of the United States of America*, 86(19), 7658–7662.
- Runge, C. A., Martin, T. G., Possingham, H. P., Willis, S. G., & Fuller, R. A. (2014). Conserving mobile species. *Frontiers in Ecology and the Environment*, 12(7), 395–402.
- Runge, C. A., Watson, J. E. M., Butchart, S. H. M., Hanson, J. O., Possingham, H. P., & Fuller, R. A. (2015). Protected areas and global conservation of migratory birds. *Science*, 350(6265), 1255–1258.
- Sawyer, H., Kauffman, M. J., Nielson, R. M., & Horne, J. S. (2009). Identifying and prioritizing ungulate migration routes for landscape-level conservation. *Ecological Applications*, 19(8), 2016–2025. <https://esajournals.onlinelibrary.wiley.com/doi/abs/10.1890/08-2034.1>
- Semmens, D. J., Diffendorfer, J. E., López-Hoffman, L., & Shapiro, C. D. (2011). Accounting for the ecosystem services of migratory species: Quantifying migration support and spatial subsidies. *Ecological Economics*, 70(12), 2236–2242.
- Sheldon, D., Saleh Elmohamed, M., & Kozen, D. (2007). Collective inference on Markov models for modeling bird migration. In J. C. Platt, D. Koller, Y. Singer, & S. T. Roweis (Eds.), *Advances in neural information processing systems 20, proceedings of the twenty-first annual conference on neural information processing systems, Vancouver, British Columbia, Canada, December 3–6, 2007* (pp. 1321–1328). Curran Associates, Inc.
- Sillett, T. S., Holmes, R. T., & Sherry, T. W. (2000). Impacts of a global climate cycle on population dynamics of a migratory songbird. *Science*, 288(5473), 2040–2042.
- Strimas-Mackey, M., Hochachka, W., Ruiz-Gutierrez, V., Robinson, O., Miller, E., Auer, T., Kelling, S., Fink, D., & Johnston, A. (2020). *Best practices for using eBird data*. Version 1.0. Cornell Lab of Ornithology. <https://cornelllabofornithology.github.io/ebird-best-practices/>
- Sullivan, B. L., Wood, C. L., Iliff, M. J., Bonney, R. E., Fink, D., & Kelling, S. (2009). eBird: A citizen-based bird observation network in the biological sciences. *Biological Conservation*, 142(10), 2282–2292.
- Taylor, C. M., & Norris, D. R. (2010). Population dynamics in migratory networks. *Theoretical Ecology*, 2, 65–73.
- Tonachella, N., Nastasi, A., Kaufman, G., Maldini, D., & Rankin, R. W. (2012). Predicting trends in humpback whale (*Megaptera novaeangliae*) abundance using citizen science. *Pacific Conservation Biology*, 18(4), 297–309. <https://www.publish.csiro.au/paper/PC120297>
- Trépos, R., Nicol, S., Cros, M., Peyrard, N., & Sabbadin, R. (2022). R package FlywayNet. *Zenodo*. <https://doi.org/10.5281/zenodo.7156292>
- Venter, J. H. (1967). On estimation of the mode. *The Annals of Mathematical Statistics*, 38(5), 1446–1455.
- Webster, M. S., Marra, P. P., Haig, S. M., Bensch, S., & Holmes, R. T. (2002). Links between worlds: Unraveling migratory connectivity. *Trends in Ecology & Evolution*, 17(2), 76–83.
- Wei, G. C. G., & Tanner, M. A. (1990). A Monte Carlo implementation of the EM algorithm and the poor man's data augmentation algorithms. *Journal of the American Statistical Association*, 85(411), 699–704.
- Wetlands International. (2021). IWC Online database. Wetlands International. <http://iwc.wetlands.org>.
- Wilcove, D. S., & Wikelski, M. (2008). Going, going, gone: Is animal migration disappearing. *PLoS Biology*, 6(7), e188.
- Wilcoxon, F. (1945). Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6), 80–83.
- Yu, S.-Z. (2010). Hidden semi-Markov models. *Artificial Intelligence*, 174, 215–243.

## SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

**How to cite this article:** Nicol, S., Cros, M.-J., Peyrard, N., Sabbadin, R., Trépos, R., Fuller, R. A., & Woodworth, B. K. (2023). FlywayNet: A hidden semi-Markov model for inferring the structure of migratory bird networks from count data. *Methods in Ecology and Evolution*, 14, 265–279. <https://doi.org/10.1111/2041-210X.14011>