**RESEARCH ARTICLE**

Methods in Ecology and Evolution

# Assessing the dynamics of natural populations by fitting individual-based models with approximate Bayesian computation

Jukka Sirén[1] 🆔 | Luc Lens[2] | Laurence Cousseau[2] | Otso Ovaskainen[1,3]

[1]Metapopulation Research Centre, Department of Biosciences, University of Helsinki, Helsinki, Finland

[2]Terrestrial Ecology Unit, Department of Biology, Ghent University, Ghent, Belgium

[3]Department of Biology, Centre for Biodiversity Dynamics, Norwegian University of Science and Technology, Trondheim, Norway

**Correspondence**
Jukka Sirén
Email: jukka.2.siren@aalto.fi

**Present address**
Jukka Sirén, Department of Computer Science, Aalto University, Espoo, Finland

**Funding information**
Norges Forskningsråd, Grant/Award Number: 223257; Suomen Akatemia, Grant/Award Number: 1273253 and 284601; Fonds Wetenschappelijk Onderzoek, Grant/Award Number: G014909 and G030813

Handling Editor: Michael Morrissey

**Abstract**

1. Individual-based models (IBMs) allow realistic and flexible modelling of ecological systems, but their parameterization with empirical data is statistically and computationally challenging. Approximate Bayesian computation (ABC) has been proposed as an efficient approach for inference with IBMs, but its applicability to data on natural populations has not been yet fully explored.

2. We construct an IBM for the metapopulation dynamics of a species inhabiting a fragmented patch network, and develop an ABC method for parameterization of the model. We consider several scenarios of data availability from count data to combination of mark-recapture and genetic data. We analyse both simulated and real data on white-starred robin (*Pogonocichla stellata*), a passerine bird living in montane forest environment in Kenya, and assess how the amount and type of data affect the estimates of model parameters and indicators of population state.

3. The indicators of the population state could be reliably estimated using the ABC method, but full parameterization was not achieved due to strong posterior correlations between model parameters. While the combination of the data types did not provide more accurate estimates for most of the indicators of population state or model parameters than the most informative data type (ringing data or genetic data) alone, the combined data allowed robust simultaneous estimation of all unknown quantities.

4. Our results show that ABC methods provide a powerful and flexible technique for parameterizing complex IBMs with multiple data sources, and assessing the dynamics of the population in a robust manner.

**KEYWORDS**

approximate Bayesian computation, individual-based models, metapopulation dynamics, multiple data sources, population dynamics

## 1 | INTRODUCTION

Ecological systems are often very complex, as they are influenced simultaneously by myriads of factors. Model-based analyses of empirical data have usually been based on drastic simplifications, partly because fitting complex models to data has been technically difficult. As an example, metapopulation studies have been typically based on population-level models, even if the data have been collected at individual level (Hanski, 1999). This can result in inefficient use of data, and prevent researchers to address questions that

wileyonlinelibrary.com/journal/mee3

*Methods Ecol Evol.* 2018;9:1286–1295.

require mechanistic (e.g. individual-based) understanding of how a system works.

One obstacle that makes it difficult to obtain a full picture of the factors affecting the dynamics of ecological systems is that many statistical methods have been developed for the analysis of a specific type of data. This complicates the utilization of multiple datasets, which could harbour complementary information about the system, as they need to be analysed separately. Integrated population modelling represents one approach for joint analysis of multiple datasets (Besbeas, Freeman, Morgan, & Catchpole, 2002; Schaub & Abadi, 2011), and similar techniques exist in other fields such as in fishery stock assessment (Maunder & Punt, 2013). However, the integrated population modelling approach has usually been restricted to computationally convenient formulations, such as matrix population models and Gaussian errors terms, which may lead to biased results if the assumptions behind them are not satisfied. Recently, Chandler and Clark (2014) introduced a spatially explicit integrated population model, which avoided many of the common simplifications, but the model was relatively simple yet still computationally expensive to fit.

As likelihood-based methods can be very difficult to apply to complex individual-based models (IBMs; see, for examples, Chandler & Clark, 2014; Harrison, Hanski, & Ovaskainen, 2011; Kattwinkel & Reichert, 2017), alternative approaches are needed. In ecology, pattern-oriented modelling (POM) has gained popularity as a strategy for building, selection and calibration of IBMs using simulations (Grimm & Railsback, 2011; Grimm et al., 2005). Pattern-oriented modelling is a protocol for building, evaluating and fitting IBMs based on multiple patterns that serve as filters. However, the approach is best suited for building and selection of models, but its suitability for estimating model parameters is questionable, because it is not based on any statistically rigorous foundation.

Approximate Bayesian computation (ABC) is a family of statistical methods that is based on utilizing model simulations in place of likelihood computations (Beaumont, Zhang, & Balding, 2002; Marin, Pudlo, Robert, & Ryder, 2012). Approximate Bayesian computation methods share the flexibility of POMs, but they are based on probability theory, allowing a rigorous treatment of uncertainty in estimation and prediction. Approximate Bayesian computation methods have been most widely used in population genetics, but lately they have been applied in other fields including ecology (Beaumont, 2010; Jabot & Bascompte, 2012; Jabot & Chave, 2011; Morales, Mermoz, Gowda, & Kitzberger, 2015). Approximate Bayesian computation methods have also started to gain popularity in parameterizing ecological IBMs (Chen et al., 2017; van der Vaart, Beaumont, Johnston, & Sibly, 2015; Zhang, Dennis, Landers, Bell, & Perry, 2017). However, the full potential of ABC as a method to parameterize IBMs based on ecological data on natural systems is yet to be explored. In particular, the utility of ABC in parametrizing IBMs with multiple heterogeneous data sources remains open.

The aim of the present work is to test the feasibility of ABC as a general tool for integrated population modelling. We ask if and how ABC methods allow one to parameterize an individual-based model of metapopulation dynamics based on either count data, mark-recapture data or genetic data, or a combination of these data types. In particular, we ask whether different data types carry complementary information that enables one to infer the model parameters and make predictions in a more accurate way than if using a single data type in isolation. We tackle these questions by formulating an IBM for the metapopulation dynamics of the white-starred robin, an Afrotropical passerine bird, living in the highly fragmented Taita Hills forest network in Kenya. We introduce a new post-processing strategy for ABC, which allows efficient analysis of multiple datasets with the same simulations by choosing locally optimal summaries separately for each dataset. We apply the method to simulated and real data on white-starred robin to test the effect of data availability on the accuracy of both parameter estimates as well as predictions of ecologically relevant indicators that characterize the state of the population and its sensitivity to perturbations.

## 2 | MATERIALS AND METHODS

### 2.1 | Empirical data on white-starred robin

As a case study, we consider data on the white-starred robin acquired from the Taita Hills area (SE Kenya, 03°20′S, 38°15′E). White-starred robin (*Pogonocichla stellata*, Vieillot, 1818) is a forest-dependent resident of montane forests across eastern to southern Africa (Keith, Urban, & Fry, 1992). The species is confined to indigenous forests, the extent of which has decreased by *c.* 50% between 1955 and 2004 in the study area, the remaining 12 fragments covering a total area of 470 ha (Pellikka, Ltjnen, Siljander, & Lens, 2009). During 1996–2009, 2,979 robins were trapped in 2,466 mist netting session that included multiple visits to each fragment and year. Variation in trapping effort was quantified by recording the duration and net length of each mist netting session. Missing values for duration (435 sessions) and net length (162 sessions) were imputed as averages over the recorded values. Genetic data on 5 microsatellite loci was acquired for 619 individuals, out of which sex was determined for 210 males and 88 females. Detailed information about the study area, the study species and the data are provided in Appendix S1.

### 2.2 | An individual-based model of the white-starred robin metapopulation

Based on prior knowledge of the life-history and ecology of white-starred robin (Githiru, 2003; Githiru & Lens, 2006b; Keith et al., 1992), we built an individual-based model (IBM) to simulate the population dynamics of the species in the Taita Hills forest fragment network. Here we provide a brief overview of the model, a more detailed description following the ODD (Overview, Design concepts, Details) protocol (Grimm et al., 2006, 2010) is available in Appendix S2.

The spatial structure of the model consists of $P = 12$ forest patches surrounded by matrix of unsuitable habitat. Each patch $i$ is characterized by its area $A_i$ and distances $d_{ij}$ to other patches $j$. The patch areas $A_i$ are presented in Appendix S1. The individuals are characterized by their sex and genotype on $l$ microsatellite loci. A new-born bird is

**TABLE 1** Model parameters $\theta_i$ and their true values $\theta_{i,T}$ used to generate simulated data

| $i$ | $\theta_i$ | $\theta_{i,T}$ | Explanation |
|---|---|---|---|
| 1 | $q$ | 1.6 | Mean number of territories per hectare in the patches |
| 2 | $\zeta_d$ | −8.18 | Intercept for male mortality |
| 3 | $\zeta_s$ | 0.25 | Difference between female and male mortality on logit scale |
| 4 | $\nu_f$ | −9.21 | Intercept for daily floater emigration probability |
| 5 | $\nu_a$ | −0.2 | Effect of the patch area on emigration probability |
| 6 | $\alpha$ | 0.12 | Migration distance |
| 7 | $p^J$ | 0.29 | Mating success and survival of fledgling phase |
| 8 | $\nu_j$ | −2.2 | Intercept for juvenile emigration probability |
| 9 | $\mu$ | 0.001 | Per generation mutation probability |
| 10 | $\eta_1$ | −4 | Intercept for observation probability of territorial birds |
| 11 | $\eta_2$ | 0.5 | Effect of sampling intensity on the observation probability |
| 12 | $\eta_3$ | −0.5 | Effect of patch area on the observation probability |
| 13 | $\eta_4$ | 0.5 | Difference in observation probability between floaters and others on logit scale |

assumed to be in the juvenile state for the first 2 years, after which it becomes an adult.

Each patch $i$ is assumed to host $K_i$ territories, which can be occupied by a pair of birds (a male and a female). The number of territories is modelled as $K_i \sim \text{Poisson}(qA_i)$, where $A_i$ is the area of the patch and $q$ is the average density of territories. An adult bird can either hold a territory or be a floater. In the beginning of each year, vacant territories are filled randomly among the floaters in each patch. If not enough floaters are present, then the slots remain empty. The birds do not leave the territories once they have occupied them.

Movements proceed in daily time steps. Each day, adult floaters and juveniles are assumed to emigrate ($E$) from the patch $i$ with a probability $p_i^E = \text{logit}^{-1}(\nu_f + \nu_a A_i^*)$, where the $A_i^*$ are log-transformed and zero-centred patch areas, and the parameters $\nu_f$ and $\nu_a$ model the mean emigration probability and its dependency on the patch area, respectively. A bird emigrating from patch $i$ is assumed to immigrate ($I$) to patch $j$ with probability $p_{ij}^I = \exp(-\alpha d_{ij}) \sum_{k \neq i} \exp(-\alpha d_{ik})$, where the parameter $\alpha$ models the distance dependency of migration.

Mortality of individuals is modelled on daily basis. Each day individual $j$ is assumed to die ($D$) with probability $p_j^D = \text{logit}^{-1}(\zeta_d + \zeta_s S_j)$, where $S_j = 1$ if individual $j$ is a female and $S_j = 0$ if individual $j$ is a male. The parameters $\zeta_d$ and $\zeta_s$ model the daily death probability of males and the difference between female and male death probabilities, respectively.

During the breeding season each pair holding a territory is assumed to produce juveniles ($J$), each of which is a male or female with

equal probability, and the number of which follows the binomial distribution $\text{Binomial}(2, p^J)$. The parameter $p^J$ models both mating success and juvenile survival until after fledging. Immediately after fledging, each juvenile is assumed to emigrate ($EJ$) from its natal patch $i$ to another patch with probability $p_i^{EJ} = \text{logit}^{-1}(\nu_j + \nu_a A_i^*)$. The parameter $\nu_j$ models the mean emigration probability for juveniles, and the effect of patch area $\nu_a$ is assumed to be the same as for regular emigration $E$. The target patch is chosen according to the same distribution with probabilities $P_{ij}^I$ as for regular dispersal. The genotypes of the juveniles are constructed from the parental genotypes according to Mendelian laws, assuming a one-step mutation for each allele with probability $\mu$ under a stepping stone mutation model.

## 2.3 | Observation model for mist netting sessions

The probability of observing individual $j$ in patch $i$ during mist netting session $k$ was modelled as $p_{ijk}^O = \text{logit}^{-1}(\eta_1 + \eta_2 L_k + \eta_3 A_i^* + \eta_4 F_j)$. Here, $L_k$ is the zero-centred logarithm of the sampling effort $h_k l_k$, where $h_k$ is the duration and $l_k$ is the length of the net in mist netting session $k$, and $F_j = 1$ if the individual $j$ is a floater or juvenile and $F_j = 0$ if the individual is holding a territory. The parameter $\eta_2$ models the effect of sampling intensity, $\eta_3$ the effect of patch area and $\eta_4$ the difference in observation probability between juveniles and floaters, and territorial birds on the observation probability, and $\eta_1$ is the intercept.

## 2.4 | Scenarios of data availability

To examine the effect of amount and type of data, we considered three types of data: (1) Count data $C$ includes the number of captured individuals at each session, separately for juveniles and adults; (2) Ring data $R$ includes the identities and ages (juvenile or adult) of all captured individuals at each session; (3) Genetic data $G(l, a)$ includes the genotypes for proportion $a$ of all individuals at $l$ loci, and the sex for either all individuals (if $a = 1$) or for half of the genotyped individuals (if $a < 1$). For genetic data and sex, we assumed no observation error. When fitting the model to simulated data (see below), we considered six scenarios for the type and amount of available data: $C$, $R$, $G(5, 0.2)$, $R + G(5, 0.2)$, $G(20, 1)$ and $R + G(20, 1)$, where + indicates combination of two data types. When fitting the model to the real data on white-starred robin, we considered four scenarios: $C$, $R$, $G(5, 0.2)$ and $R + G(5, 0.2)$, where the last scenario corresponds to the availability of real data.

## 2.5 | Statistical inference by approximate Bayesian computation (ABC)

We used ABC to approximate the posterior distribution of the model parameters listed in Table 1 (see van der Vaart et al., 2015; for introduction in ABC parameterization of IBMs). As described in full detail in the Appendix S3, we sampled a candidate parameter vector from the prior distribution, and simulated the model for a total of 264 years, out of which the first 250 years were ignored as a transient and the remaining 14 years were used to generate

a pseudo-dataset. The pseudo-datasets were compared to the observed dataset using a large number of raw summary statistics (up to 359 depending on the scenario for the type and amount of data), such as means and standard deviations of ecological summaries (e.g. number of birds recorded in a session or the number of birds that moved between patches) and genetic summaries (e.g. number of distinct alleles or homozygosity). We transformed the raw summary statistics to lower dimensional summaries with partial least squares (Wegmann, Leuenberger, & Excoffier, 2009) to improve the statistical and computational efficiency of the method. We used a two-stage approach, where in the first stage we applied the transformation to the raw summaries of all pseudo-datasets, whereas in the second step we restricted the transformation to a smaller number of pseudo-datasets, which were closest to the real dataset based on the first step. The number of pseudo-datasets retained after the first stage, as well as the dimensionalities of the transformed summaries in both stages were optimized to produce most accurate estimates. We used a total of $N = 100, 000$ simulated pseudo-datasets, and approximated the posterior distributions by applying local linear regression (Beaumont et al., 2002) to those $N_a = 100$ samples for which the pseudo-dataset was closest to the real dataset in terms of the transformed summary statistics. All distances were calculated as Euclidean distances between summary statistics.

We tested the performance of the ABC algorithm and the influence of the amount and type of data (see the six scenarios described above) by applying it to 100 datasets generated by simulating the individual-based model. We assumed the parameter values were centred on values shown in Table 1, but with small amount of variation between simulated datasets (see Appendix S3 for more details and for the choice of the parameter values). Conditional on the data $Y$, we measured the accuracy of the estimation procedure for each model parameter $\theta_i$ with root mean squared error (RMSE)

$$\mathrm{RMSE}(\theta_i | Y) = \sqrt{\frac{1}{N_a} \sum_{j=1}^{N_a} (\theta_{i,j} - \theta_{i,T})^2}, \tag{1}$$

where the $\theta_{i,j}$ are the accepted values and $\theta_{i,T}$ is the true value used for simulating the test dataset. We normalized RMSE as NRMSE$(\theta_i | Y)$ = RMSE$(\theta_i | Y)$/RMSE$(\theta_i)$, where RMSE$(\theta_i)$ is the RMSE under the prior distribution.

## 2.6 | Indicators of population state

We used six indicators of population state to characterize the state and dynamics of the bird population and how it responds to perturbation: (1) Population size $M$ before reproductive season; (2) Average number $D$ of dispersal events per year and per individual; (3) The proportion $O$ of individuals with at least one of the parents originating from another patch; (4) Global $F_{ST}$ statistics computed using the method of Weir and Cockerham (1984) with the MATLAB package of Strauss (2015), with individuals being assigned to their natal patches; (5) The population size after 10 years of

deforestation, measured as a proportion $L$ of the baseline population size $M$. Deforestation was introduced by decreasing the areas of patches by half and removing half of the territories in each patch; (6) The population size after 10 years of nest predation, measured as proportion $P$ of the baseline population size $M$. With nest predation, all juveniles in a nest were assumed to die with probability 0.15.

The indicators $M$, $D$, $F_{ST}$ and $O$ describe the average state of the system during the 14 year sampling period, while the indicators $L$ and $P$ predict what would happen if the system was perturbed at the end of the study period, and the response of the population is measured 10 years after the perturbation.

## 2.7 | Effect of model misspecification

Analysing empirical data with a structurally misspecified model can lead to compromised quantification of uncertainty, and potentially even to meaningless results (Freedman, 2006). Therefore, while a real generating model for the empirical data is not known and the possible misspecification is difficult to assess, it is important to understand how the model and associated inference method would behave under misspecification. We tested our method on this by analysing separate test datasets generated with additional environmental stochasticity. We added annual random effect with variance $\tau$ to the logit-scaled mortality, and generated 100 test datasets with both low ($\tau = 0.2$) and high ($\tau = 0.6$) level of variation, using same distributions for parameter values as for the original 100 test datasets. The new test datasets were then analysed similarly as the empirical data. See Appendix S3 for more details.

## 2.8 | Comparison of post-processing techniques

We compared the introduced two-step post-processing technique to two other dimension reduction techniques commonly used with ABC: regression-based estimation of posterior mean of Fearnhead and Prangle (2012) and standard PLS of Wegmann et al. (2009). We evaluated the performance of the techniques by calculating NRMSE values for the indicators and model parameters for the 100 simulated datasets under scenario $G(5, 0.2)$ (see Appendix S3 for more details).

## 3 | RESULTS

## 3.1 | Simulated data

The results from the simulated data show that using the full data (ringing data and full genetic data) yielded the most accurate estimates, both for all indicators of population state (Figure 1) as well as for all model parameters (Figure 2). However, the combination of the data types did not provide more accurate estimates for any of the indicators of population state or for most of the model parameters than the most informative data type (ringing data or genetic data) alone. As the sole exception, male mortality ($\zeta_d$) and female mortality ($\zeta_d + \zeta_s$) were estimated more accurately with the combined data than either
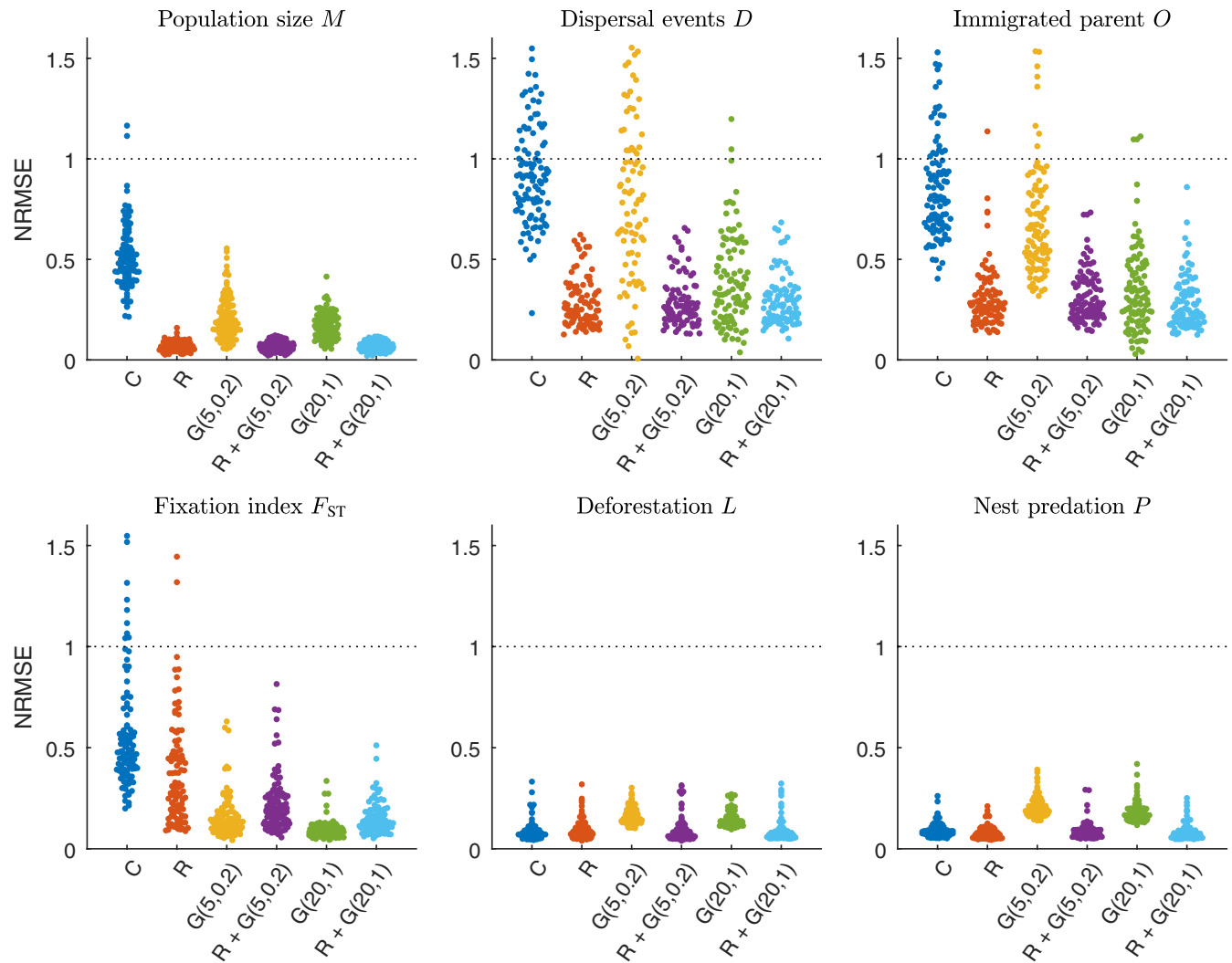
**FIGURE 1** The performance of the Approximate Bayesian computation algorithm in estimating indicators of population state based on fitting the model to simulated data with known parameter values (Table 1). The plots show the normalized root mean squared errors (NRMSE) of the indicators of population state for different scenarios of data availability, evaluated over 100 replicate test datasets. The scenarios are as follows: count data (C), ringing data (R), sparse genetic data (G(5,0.2)), full genetic data (20,1) and combination of ringing and genetic data (R + G)

ringing or genetic data alone (Figure 2, Appendix S4: Figure S1). Thus, the main utility of combining the data types into a single analysis was not in providing better estimates per se, but in providing the best estimates simultaneously for all indicators of population state and model parameters.

All data types provided information about population size, and the effect of deforestation and nest predation on it, as well as population genetic structure $F_{ST}$: for these indicators, the normalized mean squared errors were clearly smaller than one for all data types (Figure 1). In contrast, count data or genetic data with 5 loci alone did not provide information about the number of dispersal events $D$ nor about the proportion of individuals whose parents had dispersed from natal patch $O$ (Figure 1). Genetic data with 20 loci provided almost as accurate information as ringing data about the dispersal indicators $D$ and $O$. Ringing data provided most accurate estimates for all other indicators than $F_{ST}$, for which genetic data expectedly provided the most accurate estimate (Figure 1). Compared to the success of estimating

the indicators of population state, the estimation of the model parameters turned out to be difficult: for many of the model parameters the normalized root mean squared errors overlapped with one, suggesting that the data provided essentially no information over the prior (Figure 2). Failure to estimate the model parameters was at least partially explained by strong posterior correlations between some of the parameters (Appendix S4: Figure S2). Among the parameters that could be identified, ringing data provided the most accurate information for the estimation of emigration rate and its dependence on the patch area ($\nu_f$, $\nu_a$), as well as the observation probability and its dependency on sampling intensity and patch area ($\eta_1$, $\eta_2$, $\eta_3$), whereas genetic data provided the most accurate information for the estimation of mutation rate ($m$), and on the difference between female and male mortality ($\zeta_s$) (Figure 2).

The number of individuals in the data was found to have a significant positive effect on the accuracy of the estimates for all of the indicators and for nine of the 13 model parameters in the 100 test
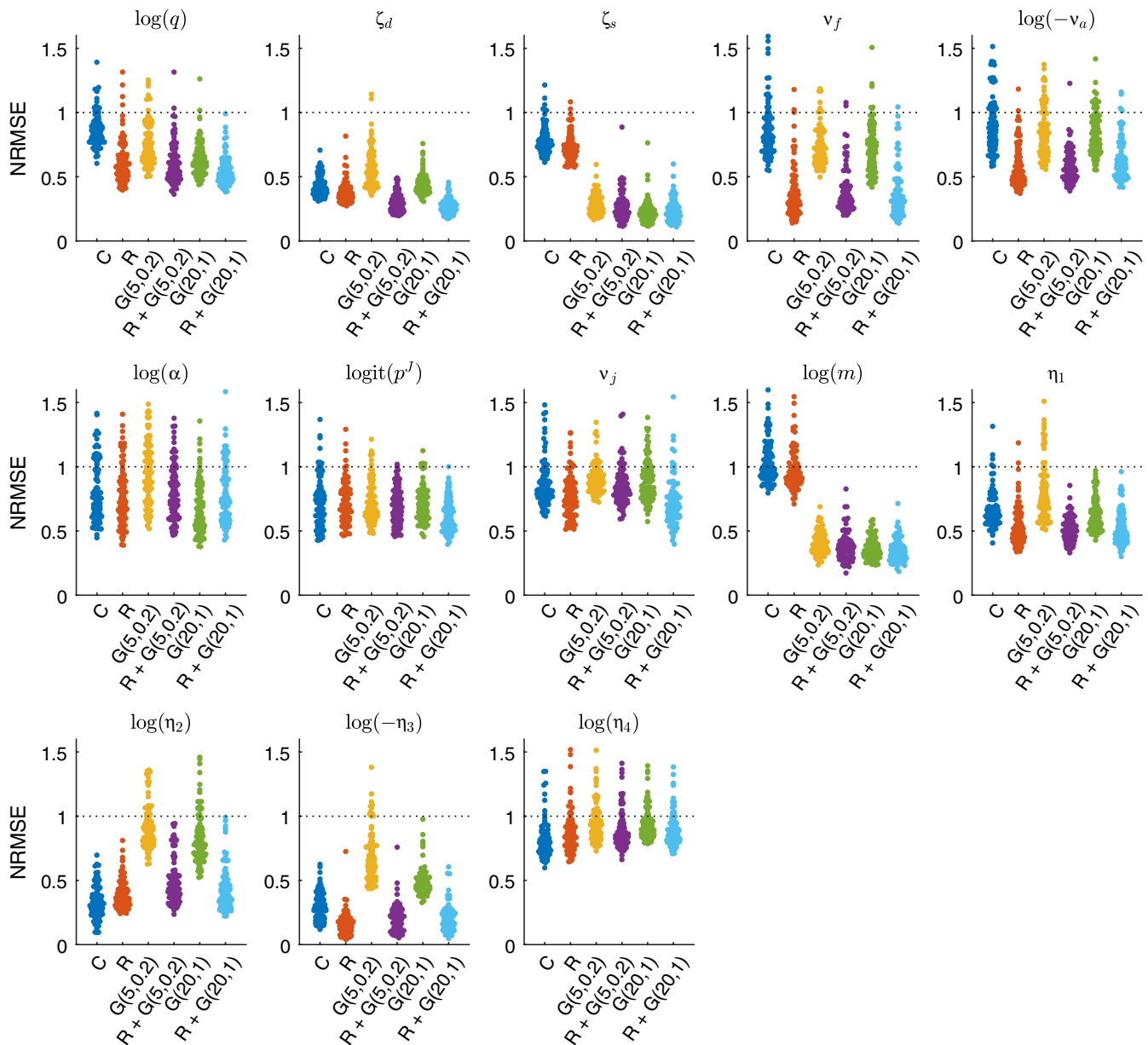
**FIGURE 2** The performance of the Approximate Bayesian computation algorithm in estimating model parameters based on fitting the model to simulated data with known parameter values (Table 1). The plots show the normalized root mean squared errors (NRMSE) of model parameters for different scenarios of data availability, evaluated over 100 replicate test datasets. The scenarios are: count data (C), ringing data (R), sparse genetic data (G(5,0.2)), full genetic data (20,1) and combination of ringing and genetic data (R + G). The function of each model parameter is explained in Table 1

datasets (Appendix S4: Figures S3 and S4). Results showing RMSE values for model parameters and indicators in natural units, and scatter plots of true and estimated values for indicators are shown in Appendix S4: Figures S5–S7.

## 3.2 | White-starred robin

As with the simulated data, also with real data, the combination of ringing and genetic data produced generically the tightest posterior distributions (Figure 3). Based on the posterior mean estimates derived from the combined data, the average metapopulation size was *c*. 4,200 individuals,

*c*. 3.3% of the individuals dispersed to another patch each year, *c*. 16% of the individuals had one of their parents originating from another patch, and the $F_{ST}$ of the metapopulation was *c*. 0.13. Concerning the sensitivity of the system to perturbations, population size was predicted to drop to *c*. 68% of its present value in the deforestation scenario and to *c*. 75% of its present value in the nest predation scenario. The posterior distributions of the model parameters and the posterior correlations between them are shown in Figures S8 and S9 of Appendix S4, respectively. Based on distributions of selected summary statistics for the accepted simulations, we did not find evidence for lack of model fit for the real data on white-starred robin (Appendix S4: Figure S10).
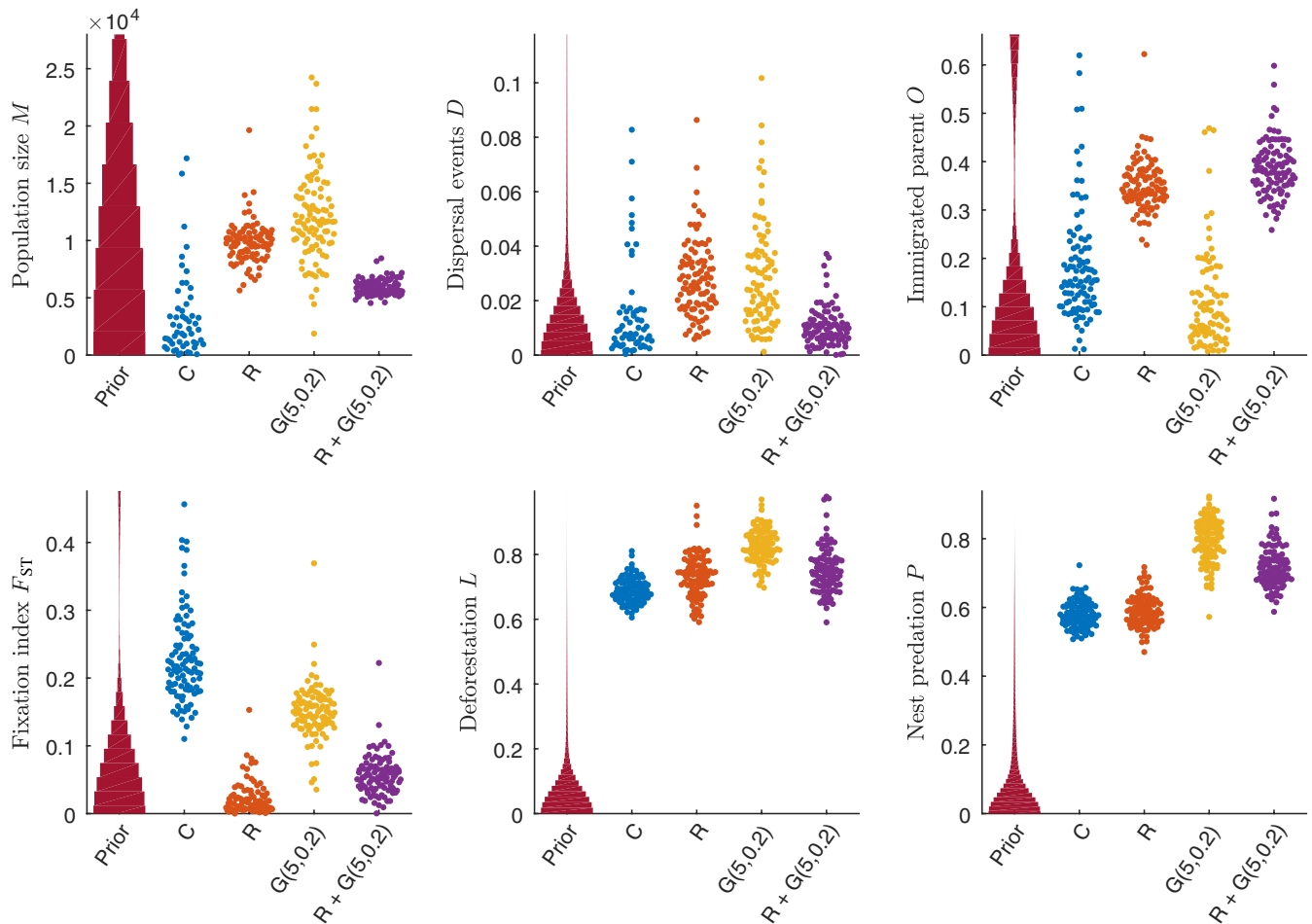
**FIGURE 3** Posterior distributions of the indicators of population state estimated by the Approximate Bayesian computation algorithm for the white-starred robin data for different scenarios of data availability. The plots show the prior predictive and posterior distributions for each indicator. The range of the prior predictive distribution has been cut from above to aid visualization. The scenarios shown are: count data (C), ringing data (R), sparse genetic data (G(5,0.2)) and combination of ringing and genetic data (R + G(5,0.2))

## 3.3 | The effect of model misspecification

The analyses of test datasets generated with additional annual variation in mortality resulted in slightly less precise estimation of indicators of population state (Figure 4). The difference was clear only for deforestation $D$ and nest predation $P$ with the high annual variation in mortality ($\tau = 0.6$). This difference was likely caused by the higher annual variation in population size as shown in Figure 4. For the other indicators and for the low annual variation ($\tau = 0.2$), the difference to standard model with no annual variation was negligible (Appendix S4: Figures S11–S16).

## 3.4 | Comparison of post-processing techniques

The two-step post-processing technique was found to perform significantly better than the regression-based technique and standard PLS for estimating the indicators of population state (Appendix S4: Figure S17). Model parameters were similarly most accurately estimated with the two-step post-processing, but the differences between techniques were not as great (Appendix S4: Figure S18).

## 4 | DISCUSSION

In this work, we developed an ABC approach for statistically and computationally efficient parameterization of IBMs based on combination of ecological and genetic data. Our results show that while the accurate estimation of the large number of primary parameters involved in an IBM may not be feasible even with high availability of data, the parameterized model can still provide informative estimates of ecologically relevant indicators of population state. Our results further show that combining different data types into a single analysis increases the accuracy of the estimates compared to separate analyses of different data types. While individual data types can provide equally accurate estimates for individual parameters or indicators of population state as the combined data, the utility of combining the data is in obtaining most reliable estimates for all unknown quantities simultaneously.

The reason why the ABC approach failed to estimate the individual model parameters accurately yet succeeded in predicting the indicators of population state is that similar patterns of population dynamics, and hence similar indicators, could be created with different combinations of the parameter values. This is evident in the strong correlations
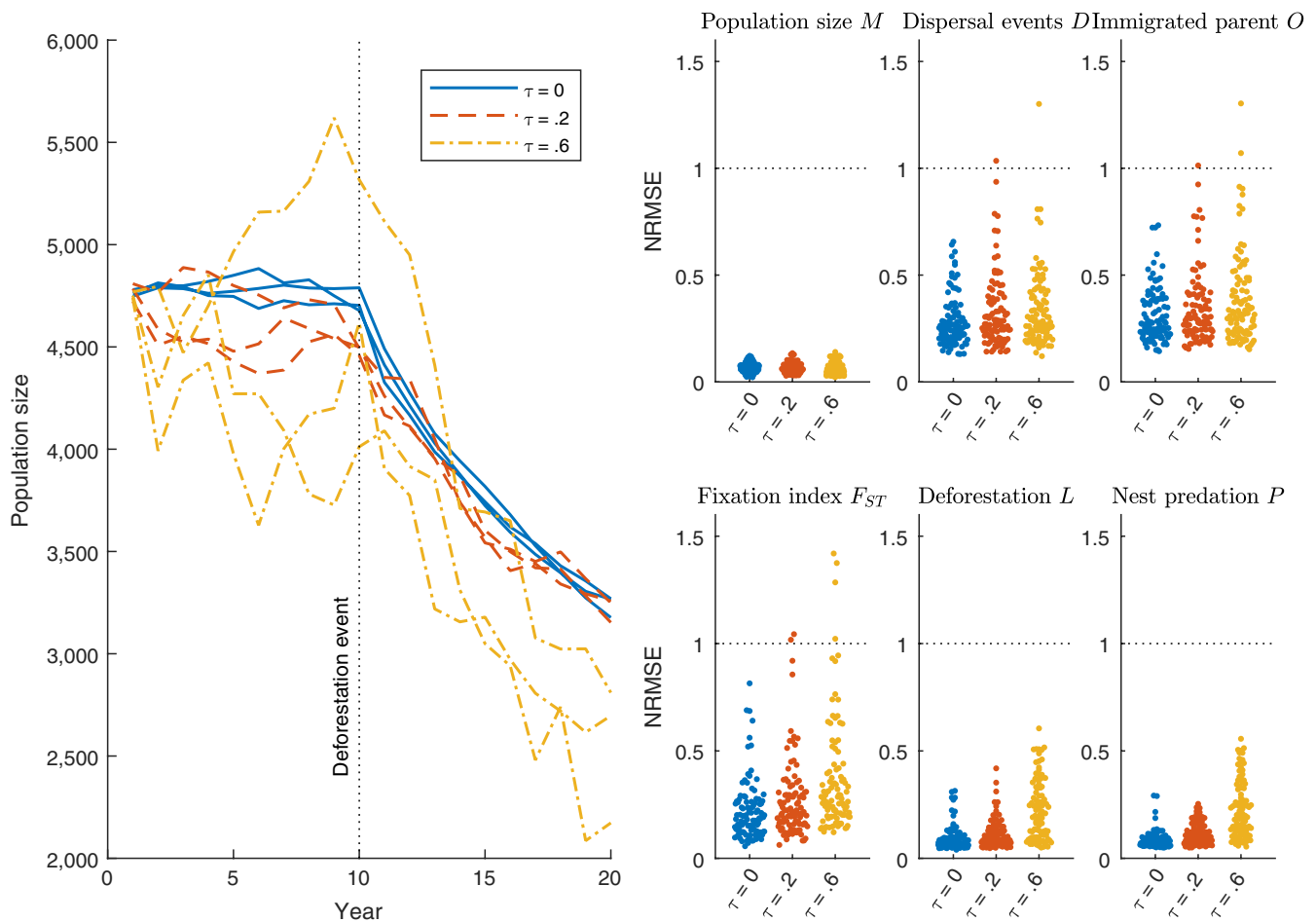
**FIGURE 4** The influence of model misspecification on the ecological inference. Left: Population size in several model simulations started from the same initial state with no (τ = 0), low (τ = 0.2) and high (τ = 0.6) annual variation in mortality. After the first 10 years of simulation, the patch areas are halved similarly as in the deforestation indicator $D$. Right: The performance of the Approximate Bayesian computation algorithm in estimating indicators of population state based on fitting the model with no annual variation in mortality to simulated data generated with no (τ = 0), low (τ = 0.2) and high (τ = 0.6) annual variation. The plots show the normalized root mean squared errors (NRMSE) of the indicators of population state for different values of τ, evaluated over 100 replicate test datasets under scenario $R + G(5, 0.2)$ (ringing data and sparse genetic data)

among the model parameters in the joint posterior distribution approximated with the ABC method, and these correlations among parameters restrict the ranges of possible outcomes for predictions. More accurate estimates of the model parameters could be obtained by using informative prior distributions motivated by the biology of the bird over the wide non-informative prior distributions used here. However, accurate estimation of the model parameters with the white-starred robin data was not the primary aim of the study.

The parameter estimates based on the combination of ringing and genetic data were in line with previously published results. In our results, the median number of territories per hectare was 1.30 (95% CI: 0.63–4.07) corresponding to a home range of 0.77 hectare. In Githiru, Lens, and Bennun (2007), mean home range size was estimated at 0.7 ha (SD: 0.33). Annual dispersal probability from multistate model controlling for dispersal distances was estimated as 0.022 in Lens, Van Dongen, Norris, Githiru, and Matthysen (2002) against 0.0223 in our results. We estimated males and females to have an annual survival probability of 0.95 and 0.55, respectively, while in Githiru and Lens

(2006a) estimates from Cormack–Jolly–Seber model were 0.83 and 0.43, respectively. However, parameters from Cormack–Jolly–Seber model are likely to be biased downward since they estimate apparent and not true survival (Schaub & Royle, 2014). In our results, the meta-population $F_{ST}$ of 0.13 was greater than the 0.034 estimated over the period 1996–1999 in Galbusera, Githiru, Lens, and Matthysen (2004), but the difference could be related to the differences in the definitions of the measures between the studies.

While there is an extensive literature on ABC algorithms, we needed to develop a methodological extension over the published algorithms in order to make the method computationally feasible for parameterizing the IBM considered here. One advantage of the rejection sampling-based ABC that we used is that it enabled us to use the same set of $N = 100,000$ model simulations for the analysis of all the 101 test datasets (100 simulated and 1 empirical) and 6 scenarios of data availability. Another advantage that motivated the two-step approach was that the second step allowed us to construct summary statistics that were efficient locally around the test dataset. While our approach was computationally

feasible for parameterizing the IBM considered here, further computational efficiency could be obtained by incorporating elements from other kinds of ABC algorithms, such as Sequential Monte Carlo or Markov chain Monte Carlo (Beaumont, 2010), or alternative likelihood free inference techniques, such as BOLFI (Gutmann & Corander, 2016).

Our case study on white-starred robin demonstrated that ABC parameterization of IBMs can be a powerful technique for population viability analyses, as it allows one to infer many kinds of ecologically relevant indicators of population state that could be otherwise difficult to estimate in a robust manner. The likelihood free computation allows for a more flexible model specification compared to standard integrated population modelling techniques that are usually restricted to a small class of data types and models (Schaub & Abadi, 2011). This flexibility creates possibility of utilizing pre-existing, heterogeneous data sources, which might provide important information on the system, but whose integration with standard analysis techniques might be difficult to achieve. While separate analyses of the data sources could provide similar accuracy in individual parameters relating to the dynamics of the population, such as dispersal or survival rate, the main benefit of the integration comes from prediction of more complex phenomena that combine many different processes.

Our estimates of the population state indicators account for spatio-temporal variation in sampling effort, as they combine an explicit observation model with the biological process model. They also account both for process uncertainty as well as for parameter uncertainty, as they were derived from simulations of the stochastic IBM, with parameters sampled from the approximate posterior distribution. However, they ignore structural model uncertainty as they are based on an IBM the structure of which was kept fixed. While the model structure was based on relatively high amount of prior information on the life-history of the focal species, it is clearly at best a rough approximation of the underlying reality. Our analyses of test datasets with misspecified models show that violation of the modelling assumptions could lead to more noisy results, but they also indicate that minor violations might not have significant impact on prediction of indicators of population state. Validation of model structure in the context of parameterizing IBMs with ABC algorithms presents a major challenge, and formal approaches with ABC face considerable difficulties (Robert, Cornuet, Marin, & Pillai, 2011). Whether an IBM developed for understanding a biological system matches the underlying reality closely enough to produce accurate predictions is a central question, for which definite answers are not available. Model building involves multiple choices and compromises made under uncertainty about the system, and these relate to the quality of the predictions made with the model in an largely unknown manner. Therefore, more research is needed to develop robust and general methods for structural model validation, which could provide alternative to methods based on more subjective assessment of model fit.

## ACKNOWLEDGEMENTS

## AUTHORS' CONTRIBUTIONS

J.S. and O.O. conceived the ideas and designed methodology; L.C. and L.L. collected the data; J.S. analysed the data; J.S. and O.O. led the writing of the manuscript. All authors contributed critically to the drafts and gave final approval for publication.

## DATA ACCESSIBILITY

The Matlab source for the individual-based model and ABC methods, as well as the White-starred robin data is available from the Dryad Digital Repository https://doi.org/10.5061/dryad.851jr (Sirén, Lens, Cousseau, & Ovaskainen, 2017).

## ORCID

*Jukka Sirén* iD http://orcid.org/0000-0002-2680-0597

## REFERENCES

Beaumont, M. A. (2010). Approximate Bayesian computation in evolution and ecology. *Annual Review of Ecology, Evolution, and Systematics*, *41*, 379–406. https://doi.org/10.1146/annurev-ecolsys-102209-144621

Beaumont, M. A., Zhang, W., & Balding, D. J. (2002). Approximate Bayesian computation in population genetics. *Genetics*, *162*, 2025–2035.

Besbeas, P., Freeman, S. N., Morgan, B. J. T., & Catchpole, E. A. (2002). Integrating mark-recapture-recovery and census data to estimate animal abundance and demographic parameters. *Biometrics*, *58*, 540–547. https://doi.org/10.1111/j.0006-341X.2002.00540.x

Chandler, R. B., & Clark, J. D. (2014). Spatially explicit integrated population models. *Methods in Ecology and Evolution*, *5*, 1351–1360. https://doi.org/10.1111/2041-210X.12153

Chen, C. M., Drovandi, C., Keith, J., Anthony, K., Caley, M., & Mengersen, K. (2017). Bayesian semi-individual based model with approximate Bayesian computation for parameters calibration: Modelling crown-of-thorns populations on the great barrier reef. *Ecological Modelling*, *364*, 113–123. https://doi.org/10.1016/j.ecolmodel.2017.09.006

Fearnhead, P., & Prangle, D. (2012). Constructing summary statistics for approximate Bayesian computation: Semi-automatic approximate Bayesian computation. *Journal of the Royal Statistical Society Series B (Methodological)*, *74*, 1–28.

Freedman, D. A. (2006). On the so-called huber sandwich estimator and robust standard errors. *The American Statistician*, *60*, 299–302. https://doi.org/10.1198/000313006X152207

Galbusera, P., Githiru, M., Lens, L., & Matthysen, E. (2004). Genetic equilibrium despite habitat fragmentation in an afrotropical bird. *Molecular Ecology*, *13*, 1409–1421. https://doi.org/10.1111/j.1365-294X.2004.02175.x

Githiru, M. (2003). *Endemic forest birds of the Taita Hills: Using a model species to understand the effects of habitat fragmentation on small populations*. DPhil thesis, University of Oxford.

Githiru, M., & Lens, L. (2006a). Annual survival and turnover rates of an afrotropical robin in a fragmented forest. *Biodiversity & Conservation*, *15*, 3315–3327. https://doi.org/10.1007/s10531-005-1213-6

Githiru, M., & Lens, L. (2006b). Demography of an afrotropical passerine in a highly fragmented landscape. *Animal Conservation*, 9, 21–27. https://doi.org/10.1111/j.1469-1795.2005.00002.x

Githiru, M., Lens, L., & Bennun, L. (2007). Ranging behaviour and habitat use by an afrotropical songbird in a fragmented landscape. *African Journal of Ecology*, 45, 581–589. https://doi.org/10.1111/j.1365-2028.2007.00772.x

Grimm, V., Berger, U., Bastiansen, F., Eliassen, S., Ginot, V., Giske, J., … DeAngelis, D. L. (2006). A standard protocol for describing individual-based and agent-based models. *Ecological Modelling*, 198, 115–126. https://doi.org/10.1016/j.ecolmodel.2006.04.023

Grimm, V., Berger, U., DeAngelis, D. L., Polhill, J. G., Giske, J., & Railsback, S. F. (2010). The ODD protocol: A review and first update. *Ecological Modelling*, 221, 2760–2768. https://doi.org/10.1016/j.ecolmodel.2010.08.019

Grimm, V., & Railsback, S. F. (2011). Pattern-oriented modelling: A multi-scope for predictive systems ecology. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 367, 298–310.

Grimm, V., Revilla, E., Berger, U., Jeltsch, F., Mooij, W. M., Railsback, S. F., … DeAngelis, D. L. (2005). Pattern-oriented modeling of agent-based complex systems: Lessons from ecology. *Science*, 310, 987–991. https://doi.org/10.1126/science.1116681

Gutmann, M. U., & Corander, J. (2016). Bayesian optimization for likelihood-free inference of simulator-based statistical models. *The Journal of Machine Learning Research*, 17, 4256–4302.

Hanski, I. (1999). *Metapopulation ecology*. Oxford, UK: Oxford University Press.

Harrison, P. J., Hanski, I., & Ovaskainen, O. (2011). Bayesian state-space modeling of metapopulation dynamics in the glanville fritillary butterfly. *Ecological Monographs*, 81, 581–598. https://doi.org/10.1890/11-0192.1

Jabot, F., & Bascompte, J. (2012). Bitrophic interactions shape biodiversity in space. *Proceedings of the National Academy of Sciences*, 109, 4521–4526. https://doi.org/10.1073/pnas.1107004109

Jabot, F., & Chave, J. (2011). Analyzing tropical forest tree species abundance distributions using a nonneutral model and through approximate bayesian inference. *The American Naturalist*, 178, E37–E47. https://doi.org/10.1086/660829

Kattwinkel, M., & Reichert, P. (2017). Bayesian parameter inference for individual-based models using a Particle Markov Chain Monte Carlo method. *Environmental Modelling & Software*, 87, 110–119. https://doi.org/10.1016/j.envsoft.2016.11.001

Keith, S., Urban, E. K., & Fry, C. H. (Eds.) (1992). *The birds of Africa*. New York, NY: Academic Press Inc Ltd..

Lens, L., Van Dongen, S., Norris, K., Githiru, M., & Matthysen, E. (2002). Avian persistence in fragmented rainforest. *Science*, 298, 1236–1238. https://doi.org/10.1126/science.1075664

Marin, J. M., Pudlo, P., Robert, C., & Ryder, R. (2012). Approximate Bayesian computational methods. *Statistics and Computing*, 22, 1167–1180. https://doi.org/10.1007/s11222-011-9288-2

Maunder, M. N., & Punt, A. E. (2013). A review of integrated analysis in fisheries stock assessment. *Fisheries Research*, 142, 61–74. https://doi.org/10.1016/j.fishres.2012.07.025

Morales, J. M., Mermoz, M., Gowda, J. H., & Kitzberger, T. (2015). A stochastic fire spread model for north Patagonia based on fire occurrence maps. *Ecological Modelling*, 300, 73–80. https://doi.org/10.1016/j.ecolmodel.2015.01.004

Pellikka, P. K., Ltjnen, M., Siljander, M., & Lens, L. (2009). Airborne remote sensing of spatiotemporal change (1955–2004) in indigenous and exotic forest cover in the Taita Hills, Kenya. *International Journal of Applied Earth Observation and Geoinformation*, 11, 221–232. https://doi.org/10.1016/j.jag.2009.02.002

Robert, C. P., Cornuet, J. M., Marin, J. M., & Pillai, N. S. (2011). Lack of confidence in approximate Bayesian computation model choice. *Proceedings of the National Academy of Sciences*, 108, 15112–15117. https://doi.org/10.1073/pnas.1102900108

Schaub, M., & Abadi, F. (2011). Integrated population models: A novel analysis framework for deeper insights into population dynamics. *Journal of Ornithology*, 152, 227–237. https://doi.org/10.1007/s10336-010-0632-7

Schaub, M., & Royle, J. A. (2014). Estimating true instead of apparent survival using spatial Cormack–Jolly–Seber models. *Methods in Ecology and Evolution*, 5, 1316–1326. https://doi.org/10.1111/2041-210X.12134

Sirén, J., Lens, L., Cousseau, L., & Ovaskainen, O. (2017). Data from: Assessing the dynamics of natural populations by fitting individual-based models with approximate Bayesian computation. *Dryad Digital Repository*, https://doi.org/10.5061/dryad.851jr

Strauss, R. E. (2015). Matlab statistical functions. Retrieved from http://www.faculty.biol.ttu.edu/Strauss/Matlab/matlab.htm

van der Vaart, E., Beaumont, M. A., Johnston, A. S., & Sibly, R. M. (2015). Calibration and evaluation of individual-based models using approximate Bayesian computation. *Ecological Modelling*, 312, 182–190. https://doi.org/10.1016/j.ecolmodel.2015.05.020

Wegmann, D., Leuenberger, C., & Excoffier, L. (2009). Efficient approximate Bayesian computation coupled with Markov chain Monte Carlo without likelihood. *Genetics*, 182, 1207–1218. https://doi.org/10.1534/genetics.109.102509

Weir, B. S., & Cockerham, C. C. (1984). Estimating F-statistics for the analysis of population structure. *Evolution*, 38, 1358–1370.

Zhang, J., Dennis, T. E., Landers, T. J., Bell, E., & Perry, G. L. (2017). Linking individual-based and statistical inferential models in movement ecology: A case study with black petrels (*Procellaria parkinsoni*). *Ecological Modelling*, 360, 425–436. https://doi.org/10.1016/j.ecolmodel.2017.07.017

## SUPPORTING INFORMATION

Additional Supporting Information may be found online in the supporting information tab for this article.

---

**How to cite this article:** Sirén J, Lens L, Cousseau L, Ovaskainen O. Assessing the dynamics of natural populations by fitting individual-based models with approximate Bayesian computation. *Methods Ecol Evol*. 2018;9:1286–1295. https://doi.org/10.1111/2041-210X.12964