

BIRDIE: A data pipeline to inform wetland and waterbird conservation at multiple scales

Francisco Cervantes^{1,2*}, Res Altwegg¹, Francis Strobbe³, Andrew Skowno², Vernon Visser¹, Michael Brooks⁴, Yvan Stojanov³, Doug Harebottle⁵, Nancy Job²

¹ Centre for Statistics in Ecology, the Environment and Conservation, University of Cape Town, Cape Town, South Africa

² South African National Biodiversity Institute, Kirstenbosch Research Centre, Cape Town, South Africa

³ Operational Directorate Natural Environment, Royal Belgian Institute of Natural Sciences, Brussels, Belgium

⁴ FitzPatrick Institute of African Ornithology, University of Cape Town, Cape Town, South Africa

⁵ Risk and Vulnerability Science Centre, Sol Plaatje University, Kimberley, South Africa

Correspondence*:

Francisco Cervantes

f.cervantesperalta@gmail.com

2 ABSTRACT

3 With increasing pressure on ecosystems, we have seen substantial efforts in the collection
4 and management of ecological data, over the last decade. This is especially true for freshwater
5 ecosystems, which are among the most impacted by human activity and yet they have been
6 lagging behind in terms of data availability, compared to other terrestrial ecosystems. Despite
7 improvements in data collection and accessibility, raw data still need to be analysed and
8 interrogated to support conservation programmes and management decisions, a process that can
9 be complex and time consuming. The South African Biodiversity Data Pipeline for Wetlands and
10 Waterbirds (BIRDIE) aims to help fast and efficient information uptake, bridging the gap between
11 raw ecological datasets and final users. BIRDIE is a full data pipeline that takes up raw data, and
12 through a series of processing steps, estimates indicators related to abundance, distribution and
13 diversity of waterbirds, keeping track of their associated uncertainty. It focuses on the assessment
14 of waterbird populations in South Africa and uses two citizen-science bird monitoring datasets,
15 namely: the African Bird Atlas Project and the Coordinated Waterbird Counts. In addition, a
16 suite of environmental layers help contextualise waterbird population indicators, and link these to
17 the ecological condition of the supporting wetlands. Data processing is conveniently organised
18 in modules that can be run independently, and include tasks, such as: occupancy modelling,
19 state-space modelling, and computation of indicators at multiple temporal and spatial scales. Both
20 data and indicators are accessible to end users through an online portal and through web services.
21 Envisioned users of BIRDIE include government officials, conservation managers, researchers
22 and the general public. Acknowledging that conservation programmes run at multiple scales, from
23 site management to international agreements, we have developed a granular framework in which

waterbird population indicators are estimated at small scales, and then these are aggregated to compute similar indicators at broader scales. The online portal is designed to provide spatial and temporal visualisation of the indicators using maps and time series, to help contextualization. This paper describes the structure of the BIRDIE pipeline and the technical features underpinning its components.

Keywords: Biodiversity informatics, Citizen science, Data pipeline, Waterbirds, Wetlands, Species distribution, Species Abundance, Diversity

INTRODUCTION

Freshwater ecosystems are among the most productive, biodiverse, and most efficient at capturing and storing carbon (Convention on Wetlands, 2021). Unfortunately, they are also among the most impacted by human activity (Convention on Wetlands, 2021; Skowno et al., 2019), and climate change will likely increasingly exacerbate the pressure on freshwater resources. This is particularly true for the African continent, home to some of the largest wetlands, which not only host a wealth of freshwater species, but are also key in supporting human communities (Stephenson et al., 2020). Such critical issues have fuelled unprecedented efforts to collect and mobilise freshwater biodiversity data (Dallas et al., 2021; Wetzel et al., 2015). While we must strive to keep monitoring programmes that deliver data funded and alive, it is clear that data on their own are not enough (MacFadyen et al., 2022). If we are to take effective action to stop ecosystem degradation, it is important that data are analysed to extract indicators that are meaningful for decision- and policy-making (Jetz et al., 2019; P. Stephenson et al., 2017).

Bridging the gap between biodiversity data sets and conservation action requires, not only collecting, but also analysing and interpreting the data (Pollock et al., 2020). Furthermore, with continuous data collection, we need to implement workflows that update indicators and support decisions in a timely fashion (MacFadyen et al., 2022; Yenni et al., 2019). Automated data pipelines allow us to keep datasets updated and free of errors (Yenni et al., 2019), make model-based forecasts, and evaluate previous forecasts in light of new data (White et al., 2019). These modern and automated data workflows require multidisciplinary skills in ecology, statistics, data science, and software development, but their end products should ideally be free, accessible and easy to interpret (P. Stephenson et al., 2017). It would also be desirable that they integrate multiple datasets and environmental layers to produce a holistic understanding of biodiversity structure and function (MacFadyen et al., 2022).

South Africa is leading the African continent in terms of biodiversity data availability (Barnard et al., 2017), with successful citizen-science programmes such as the Southern African Bird Atlas Project (Brooks et al., 2022), and biodiversity data platforms, such as the Biodiversity Advisor (SANBI, 2023) or the Freshwater Biodiversity Information System (FBIS, Dallas et al., 2021). In contrast, dashboards and tools that facilitate the timely uptake of information and unlock the utility of current data are still limited. Here, we describe a data pipeline that implements a workflow of wetland- and waterbird-related biodiversity data, the South African Biodiversity Data Pipeline for Wetlands and Waterbirds (BIRDIE). At its core, BIRDIE utilises two long-term citizen-science programmes that have collected waterbird data in South Africa for more than two decades, and are still active: the Southern African Bird Atlas Project (SABAP; Brooks et al., 2022) and the Coordinated Waterbird Counts (CWAC; FIAO, 2022). These datasets have well defined protocols supported by an established community of citizen scientists, and they provide information about: 1) bird species occurrence, with visits to a grid of pentads (5' x 5' grid cells) across Southern Africa (SABAP), and 2) bird abundance, with waterbird counts taken twice a year at ca. 730 sites across South Africa (CWAC). These data provide the necessary information to investigate changes in species abundance

and distribution, which are considered the minimum set of variables necessary to study changes in species populations (Pereira et al. 2013, Jetz et al. 2020). Apart from waterbird data, BIRDIE uses and serves ancillary environmental data for contextualising the aforementioned variables associated with waterbird populations, and also for describing the state of the wetlands that support them.

The BIRDIE data pipeline was designed to run periodically (yearly in principle), and automatically (but supervised), to keep information updated. It is meant to inform the general public about the state of wetlands and waterbird populations, but also to support the work of particular end users. A first group of users includes authorities that need to report on the state of wetlands (e.g., Ramsar Convention), or migratory waterbirds (e.g. African-Eurasian Migratory Waterbird Agreement) as required by international agreements, for national level processes (e.g., National Biodiversity Assessment) and for provincial or municipal level regulation (e.g. to set hunting quotas for waterfowl). A second group of users includes site managers and other stakeholders who need to make a range of decisions specific to certain wetlands (e.g., related to wetland condition and extent) and the waterbirds supported by these habitats. Such policy-linked outputs are used to inform environmental strategies, identify priorities for protection and for sustainable use of biodiversity, and guide land-use management.

FRAMEWORK AND TARGET USERS

Indicators on the state of biodiversity have been adopted by a range of multilateral environmental agreements including the United Nations Convention on Biological Diversity (CBD, 2022) and Sustainable Development Goals (SDGs; UN, 2022). New indicators are under development and established processes, such as the International Union for the Conservation of Nature (IUCN, 2022) species red-listing efforts, are receiving renewed attention (Han et al., 2017). With these indicators come various global and national initiatives and targets for reducing rates of biodiversity loss (Mace et al., 2018). Essential Biodiversity Variables (EBVs) have been conceptualised and developed to help standardise and improve interoperability of biodiversity data for monitoring (Pereira et al., 2013).

Within this framework, BIRDIE gives support to both national and international programs contributing information about the state of waterbird populations in South Africa, with a view to expand to the Southern Africa region. We focus primarily on species population EBVs, with the assessment of waterbird abundance, distribution and diversity, and changes of these over time (Jetz et al., 2019; Kissling et al., 2018). Species diversity falls under the community composition EBVs rather than species populations.

At the international scale, South Africa is signatory to the Ramsar Convention (Convention on Wetlands, 2021), hosts 28 Wetlands of International Importance, and needs to produce reports on the state of sites every three years. Change in condition and extent of wetland habitat, as well as those Ramsar Criteria invoked when listing the different sites as internationally important are all core reporting requirements. These include changes in overall abundance and distribution of waterbirds at these sites, with special attention to threatened and migratory species.

National reports must also be compiled for the Agreement on the Conservation of African-Eurasian Migratory Waterbirds (AEWA; UNEP, 2022), another international agreement, framed under the Convention on Migratory Species, and focused on protecting migratory waterbirds and their habitats. Reporting for this agreement requires the full suite of EBVs outlined earlier in this section, now specific to migratory species, namely, abundance, distribution, diversity, and changes in these over time. Within South Africa, there is alignment and cooperation between Ramsar and AEWA, and as such, reporting on the change in wetland extent and condition is also relevant.

At the national level, South Africa produces a National Biodiversity Assessment every four years, which constitutes the main reporting tool of the state of biodiversity in the country, and informs policy and conservation strategies (Skowno et al., 2019). At the same time, there are regular efforts to address the conservation status of South African species within the IUCN Red-List framework. Changes in abundance and distribution of species are key in these assessments to track and report on population trends, and shifts in species ranges and community diversity.

Keeping these main reporting channels in mind, BIRDIE also intends to support local management actions and basic research. Site-scale wetland monitoring is severely limited in South Africa, lagging far behind monitoring of other aquatic ecosystems such as rivers and estuaries. Managers ideally need to report on the state of the wetland (e.g., wetland condition, flux in surface water extent) as well as the species that the wetland supports, including species of special concern. Local waterbird and wetland information can facilitate the development of site-specific management actions and management plans, and support permitting decisions. At the same time, linking the local manager inputs and feedback into the data pipeline closes the gap between large-scale assessments and local data collection. The data pipeline also allows citizen scientists to more actively interact with the data they have collected, and to see it taken up into the statistical analyses and data visualisations.

INPUT DATA

South Africa needs to monitor its biodiversity to inform decisions that affect the environment, and to fulfil reporting obligations linked to international conventions. Fortunately, we have a number of long-running citizen science projects that help monitor waterbird populations throughout the country. At its core, the project leverages two bird-related datasets: the Coordinated Waterbird Counts (CWAC) and the South African Bird Atlas Project (SABAP2), which is part of the larger African Bird Atlas Project (ABAP). These datasets have well established citizen scientist support and offer information about: 1) bird abundance, with waterbird counts taken twice a year at 731 water bodies across Southern Africa (mostly South Africa) since 1992, and 2) species occurrence, with visits to a grid of pentads (5' x 5' grid cells) initiated in 2007 and covering several African countries.

The Coordinated Waterbird Counts project provides regular counts of all waterbirds at just over 700 sites throughout South Africa. The project was launched in 1992 and since then, it has accumulated a long time series for many sites. However, not all sites have been monitored since the start of the project, some regions are better represented than others, and not all sites have been monitored continuously (Figure 3). Waterbird species have diverse habitat requirements and life histories; some use the same sites year-round, whereas others are migratory. To capture this diversity, CWAC counts are carried out twice per year: one in mid-summer and one in mid-winter. These data are noisy because it is difficult to count waterbirds precisely, but with appropriate statistical analyses, they can reveal long-term temporal trends and seasonal fluctuations in waterbird populations. Waterbird population sizes and trends are important reporting parameters for international agreements such as the Ramsar Convention, or AEWA.

ABAP offers occurrence, rather than abundance data. In ABAP, volunteers collect checklists of all birds observed over a grid of pentads (5' x 5' minute grid) covering different African countries. We are currently restricting our analysis to South Africa, and therefore we are using the SABAP2 component of ABAP (Figure 4). However, in the future we would like to expand our functionality to cover other countries contributing data to ABAP, such as Kenya or Nigeria. Under the SABAP2 protocol, which started in 2007, observers need to spend at least two hours of intensive birding at a pentad and are asked to visit as many habitats within it as possible. They can add new species for up to five days. SABAP2 currently has ca. 17

million records, and > 2 million records are added per year. The structured sampling protocol, together with the spatial and temporal extent of SABAP2 allow us to examine how bird distributions are changing over time, although statistical modelling is required to account for imperfect detection and sampling biases (Figure 4).

There are a variety of other data sources that BIRDIE uses for adding environmental information into its analytical workflows. Most of these data sources are conveniently accessed through Google Earth Engine, such as TerraClimate (Abatzoglou et al., 2018), the JRC surface water dataset (Pekel et al., 2016), MODIS Vegetation Indices (Didan, Kamel, 2015) and Digital Elevation Models (DEM, (Yamazaki et al., 2017). Other data not yet available on Google Earth Engine, such as the National Wetland Map (SANBI, in prep.) are managed independently.

INDICATORS AND STATISTICAL METHODS

Capturing good quality raw data is a fundamental first step to monitor the state of biodiversity. However, capturing error-free comprehensive datasets over large temporal and spatial scales requires huge efforts that are unrealistic in most situations. In fact, the very nature of natural systems often prevents them from being observed perfectly and entirely. Therefore, some level of statistical analysis is required to estimate the state of the system of interest, and separate it from observational artefacts introduced by the observation process used for capturing the data (Gimenez et al., 2014; King, 2014; Yoccoz et al., 2001). The BIRDIE pipeline broadly uses two types of models: i) occupancy models (Altwegg and Nichols, 2019; MacKenzie et al., 2002) to estimate the probability of species being present at the different SABAP2 pentads, and ii) state-space models (Buckland et al., 2004; Newman et al., 2014) to estimate the number of individuals at the sites monitored by the CWAC programme. Contrary to raw observations (counts and detection/non-detection of a species), model-based estimates (abundance and occupancy probabilities) allow us to quantify uncertainty.

The variety of end-user needs requires a pipeline that provides waterbird population indicators at multiple spatial and temporal scales. Therefore, in addition to estimating basic occupancy and abundance at small scales (i.e., individual site/pentad), the BIRDIE pipeline produces other high-level indicators obtained by aggregation (Table 1). The idea is to follow a process whereby raw data are used to inform models that estimate indicators at the smallest temporal and spatial scales possible, and then to aggregate these estimates at larger scales, as required. For example, species abundance can be estimated for a set of regularly monitored wetlands in South Africa, and these site-specific estimates can then be combined to calculate an abundance index for all sites as a group. We can follow this procedure to estimate abundance and occupancy probabilities at national, regional and local levels, as well as for specific groups of wetlands (e.g., designated Ramsar sites, estuaries or artificial sites (e.g., dams).

The main indicators computed by the BIRDIE pipeline for waterbird species are:

- Population size: estimated for CWAC sites in two seasons per year. For each species, only those wetlands with at least a ten-year coverage between 1997 and 2021 are analysed.
- Occurrence: estimated for ABAP pentads on an annual basis.
- Diversity: the simplest and most easily understood metric is species richness. Species richness can be calculated based on the occupancy analysis, by summing occupancy probabilities of all species potentially present in each pentad, to estimate the expected number of species present.
- Important records: sightings of rarities, invasive species. Although this information does not require any statistical processing, it does make particular records more visible.

In addition to estimates of static indicators, the pipeline also estimates their associated dynamics, such as: changes in abundance, occupancy probabilities and diversity. The temporal reference for these dynamics can also vary ranging from a single season to multiple years (typically ca. 5 years, for short-term changes, and ca. 15 years for long-term changes). It is important that uncertainty is correctly propagated when aggregating, and also when estimating dynamic indicators. We work in a Bayesian framework and use the posterior distribution of occupancy probabilities and species abundance to define indicators at the various scales. Working with full posterior distributions allows us to conveniently keep track of the uncertainty in the estimates used as building blocks for other derived indicators.

Delineating species distributions

Occupancy models are fitted to detection/non-detection data from SABAP2 to delineate the distribution of waterbird species and its dynamics over time. Within the SABAP2 framework, observers visit pentads and make a list of the bird species detected during the visit. We assume that observers identify species correctly and only list species observed (the rigorous vetting process of SABAP2 data justifies this assumption), but non-detections may be caused by either species not being present in the pentad or by observers not being able to detect them, when present. Therefore, occupancy models describe two processes simultaneously: i) the underlying occupancy of the sites (pentads), and ii) the observation process whereby species present might or might not be observed.

More precisely, we define z_{jt} to be the true occupancy of site j in year t , which can be 1 (if species present) or 0 (if species absent) and has distribution:

$$z_{jt}|\psi_{jt} \sim \text{Bernoulli}(\psi_{jt})$$

where ψ_{jt} is the occupancy probability at site j and year t . The logit transformation of ψ_{jt} can be modelled as a linear combination of covariates and smooth functions of covariates, such that:

$$\text{logit}(\psi_{jt}) = \mathbf{x}_{jt}^T \boldsymbol{\beta} + \sum_{k=1}^K f_k(u_{jk})$$

where u_{jk} is a smooth function of the covariate u_k , which is defined as

$$f_k(u_{jk}) = \sum_{l=1}^L \mathbf{B}_l(u_{jkl}) \gamma_{jkl}$$

where the smooth function f is represented by a set of L basis functions \mathbf{B}_l evaluated at the value of the covariates associated with site j at year t .

Then, the probability of detection of a species that is present in site j on visit i is denoted by p_{ij} . Following the same logic as for the probability of occupancy, the logit transformation of p is modelled as a linear combination of covariates and smooth functions:

$$\text{logit}(p_{ij}) = \mathbf{w}_{ij}^T \boldsymbol{\alpha} + \sum_{h=1}^H f_h(v_{ih})$$

,

218 Finally, the likelihood of observation y_{ij} can be written as:

$$y_{ij}|z_{jt}, p_{ij} \sim \text{Bernoulli}(z_{jt}p_{ij})$$

219 Spatial, spatio-temporal, and unstructured random effects can be specified for either occupancy or
220 detection probabilities to account for variation across sites, observers and visits, not accounted for by the
221 covariates incorporated in the models.

222 Each checklist is treated as an independent survey, but occupancy is assessed on a yearly basis. This
223 means that if a species is detected in any one survey it is considered present that year. Therefore, missing a
224 species because it left the site is considered part of the observation process and not the occupancy process.
225 Migratory birds, for example, are considered present at a site even if they are only there for part of the year.

226 We are fitting single-season occupancy models without spatial random effects to most species. However,
227 all models incorporate random effects to account for pentad- and observer-specific detection probabilities.
228 If model diagnostics indicate poor model fit (see model diagnostics section below), we try incorporating
229 spatial random effects for occupancy probabilities with an exponential decay function. Currently, we fit
230 the models in R (R Core Team, 2022), in a Bayesian framework using the package spOccupancy (Doser
231 et al., 2022), and running three MCMC chains for 20,000 iterations, with a thinning interval of 20. We
232 use non-informative priors for all parameters when no information from other years is available, but we
233 incorporate information obtained from other model fits if available, by centering the priors on the closest
234 model's posterior means. However, it is important noticing that modelling details may differ among species
235 and may be updated in the future.

236 **Estimating abundance and population trends**

237 State-space models (Buckland et al., 2004; Newman et al., 2014) are used to describe and understand
238 dynamic systems that may not be perfectly observed. Within this framework, we consider waterbird
239 abundance to be a process that evolves over time, and which we observe during visits to CWAC sites.
240 However, counts conducted by observers are distorted by imperfect detection that translates into counting
241 errors. By counting repeatedly over time, and assuming that abundance evolves smoothly over time
242 compared to observation error, we can disentangle these two processes.

243 We consider that the observed counts (y_i) at sampling occasion i (generally there were two sampling
244 occasions per year, one in mid-summer and one in mid-winter), at any given site, arise from a Poisson(λ_i)
245 distribution

$$y_i \sim \text{Poisson}(\lambda_i)$$

246 And we model the log of the intensity λ_i as:

$$\log(\lambda_i) \sim N(\mu_i, \sigma^2)$$

247 where μ_i is the mean abundance of waterbirds present at a site on sampling occasion i and σ^2 is the
248 corresponding variance of the observers counting error, both in the log scale. Therefore, counts depend
249 both on the number of waterbirds present on site, and on errors in the counts of these birds.

250 To model changes in waterbird abundance between the two-seasons of year t , we define s_t to be the
251 summer abundance and w_t the winter abundance. Note that there might be multiple counts in a single year

and season, but the underlying true abundance is considered to stay constant in any given year and season (for clarity, note also that while sampling occasions were indexed by i , years are indexed by t). Thus, the expected (log) abundance for any given count can be written as

$$\mu_i = s_t \text{summer} + w_t \text{winter}$$

where ‘summer’ is an indicator variable that takes on the value 1 in summer and 0 in winter, and ‘winter’ is the opposite. We then define abundance dynamics as:

$$s_t = s_{t-1} + \beta_t$$

$$w_t = s_t + \xi_t$$

where β_t corresponds to the change in summer abundance from year $t-1$ to year t , and ξ_t is the difference between summer and winter abundance, both in the log scale. If exponentiated, these parameters can be interpreted as the rate of change in the population and the winter-to-summer ratio of the population, respectively.

We impose relatively smooth changes in abundance by defining autocorrelation in and terms over time. In addition, we define relationships between the rate of change in the population and environmental covariates. These relationships facilitate the estimation of abundance for those years in which counts are missing, and it is particularly useful to contain uncertainty in long periods with missing data between counts. Thus, we set

$$\beta_t = \phi \beta_{t-1} + \eta_{t-1} + \zeta_{t-1}$$

$$\xi_t = \xi_{t-1} + \epsilon_{t-1}$$

where ϕ lies between zero and one, and it defines an autoregressive term on β_{t-1} ; η_t captures the effect of covariates in the expected change in abundance, and can be expanded to $\gamma^T U$, where U is a matrix of covariate values and γ a vector of coefficients; ζ_t and ϵ_t are random variables that represent change in abundance change, and change in winter to summer ratio, respectively.

We mentioned at the beginning that this model applies to each monitored site. However, we have multiple sites, and counts are often missing for some seasons or even full years. To facilitate the estimation of abundance with missing data, we borrow information from sites with counts, by defining a hierarchical structure such that:

$$\zeta_{tj} \sim N(0, \sigma_{\zeta t}^2)$$

$$\epsilon_{tj} \sim N(0, \sigma_{\epsilon t}^2)$$

so that random changes at any site and year come from a common distribution of changes across all sites for that year. These distributions are normal with variances $\sigma_{\zeta t}^2$ and $\sigma_{\epsilon t}^2$ for changes in abundance and winter to summer ratio, respectively.

We fit these models in R (R Core Team, 2022) with the additional functionality provided by JAGS (Plummer, 2003) using the package jagsUI (Kellner, 2021). We work on a Bayesian framework, using non-informative priors, and running three chains for 10,000 iterations each. Similar to the occupancy models, these are the details of the models we are working with at the moment, and they are intended to give an idea of the type of model we are using. These models might be updated in the future and the updated modelling details will be published on the BIRDIE website.

Data and model diagnostics

The pipeline needs to run for a multitude of species, with different ecological requirements and geographical distributions. Therefore, finding a model that suits all species is challenging. Not only may a model not be a good fit for a particular species, but the algorithms used for fitting the model may fail to converge due to characteristics of the data.

In a first control stage, we have defined some minimum requirements data should meet to enter the model-fitting process. Species that have been observed in five or less pentads in a year are considered to not have enough data to inform an occupancy model. Similarly, we chose only those CWAC sites that have been counted at least ten times between 1993 and 2021, to fit state-space models. Otherwise, data tend to be too sparse to assess trends in abundance reliably. These thresholds are rather arbitrary, based on our own preliminary experiences, and they are considered to be the minimum requirements for models to converge successfully. However, meeting these requirements does not guarantee model convergence or a good fit. To keep track of potential issues arising during model fitting, and to improve the algorithms of the pipeline, each time the pipeline runs it generates several reports that are later examined.

To decide whether any occupancy or state-space model converges, we calculate the Gelman-Rubin (Rhat) diagnostic (Gelman et al., 2014) for each estimated parameter. These diagnostics are then tabulated and stored for future revision. Any Rhat value above 1.1 or below 0.9 is considered to represent lack of convergence. Distinctive characteristics of the models with convergence issues are explored and addressed on a case by case basis, after the pipeline has finished running.

In addition to convergence, we assess goodness of fit using posterior predictive checks (Doser et al., 2022; Gelman et al., 2014). This procedure compares some quantity of interest calculated using pseudo-data simulated from the model posterior distribution, with that same quantity calculated from the observed data. In a well-fitting model we would expect real and synthetic data to produce similar values. For occupancy models, we produce simulated detection/non-detection data for each site, species and year and compute the expected number of detections out of as many visits as there were in the data. We compare the results of the simulations with the observed number of detections recorded in the data using a Chi-square test. For state-space models we follow a similar procedure, but instead of simulating detection/non-detection data for one year, we simulate count data for summer and winter, and aggregate these in a single annual count. Results from the goodness of fit Chi-square tests are also tabulated and stored for revision. Significant deviations detected with these tests are addressed for each case individually.

Due to the computational burden of the pipeline, it is not possible to run multiple models for each species, site and year, to perform model selection. Therefore, model selection is performed on a sample of species, selected to have representation of common and scarce taxa, but that are otherwise selected arbitrarily. Our general approach has been to include a rich set of variables that we believe can explain

the main environmental gradients within our geographical range, without paying too much attention to multi-collinearity and overfitting. We are therefore cautious about making causal inference or predictions out of the range of the data, and so should be other users.

SYSTEMS AND TECHNOLOGY

This section describes the servers, database and dashboard technology. We present the process whereby the data is processed along the pipeline and transformed into indicators that are presented to the final user. There are multiple formats in which the information is presented to be able to accommodate the needs of downstream systems. BIRDIE's data, code and outputs are stored and run on three main systems (Figure 5): the Africa Bird Data servers, and the two BIRDIE servers (servers A and B).

The Africa Bird Data servers are hosted at the FitzPatrick Institute for African Ornithology (FIAO), University of Cape Town, and contain the CWAC and ABAP databases. They also serve these data through an Application Programming Interface (API).

BIRDIE's server A is the access point of the final user to the information generated by BIRDIE. This information is stored in a data mart. At its core, the data mart is a MySQL database (version 8.0.27), a widely used, open source, relational database management system. Its main objective is to store BIRDIE's outputs in a way that provides easy and flexible access to the final user. At the same time, the structure of the database ensures that the inputs and outputs conform to a given standard, and creates security back-ups for the stored data. The main mechanism BIRDIE uses to present data to the user is through OpenAPI web services (OpenAPI Specification, Version 3.1.0), which was designed to provide a standard interface for documenting and exposing APIs. The public web services offered by the OpenAPI give users the flexibility to access and download data from the database in multiple ways, without being constrained by specific functionality of a web application. This technology facilitates the integration of BIRDIE's outputs into other workflows. For the user that is interested in readily accessing the information through a dashboard, we have deployed a web application, written in HTML5, CSS, and the most common and popular JavaScript libraries, including OpenLayers (<https://openlayers.org/>) and Plotly (<https://plotly.com/>). Among other elements (see section 6), the web application features a map viewer, based on mviewer (<https://mviewer.netlify.app/en/>), a free and open-source cartographic application, that has an easy-to-use and intuitive interface.

If we thought of server A as being the face of the pipeline, server B would be the brain. All the functionality in this server revolves around statistical modelling. This server connects with the Africa Bird Data servers to obtain CWAC and ABAP data, and with external systems, such as Google Earth Engine or SANBI servers to obtain environmental information. It then runs the main analytical modules of the pipeline, where occupancy and state-space models are fitted. The analytical workflows are supported by an Intel Xeon Dual 8 core, with 64 GB RAM and an 8 TB hard drive. The model outputs are made available to server A, where they are incorporated into the data mart, used to compute derived high-level indicators by aggregation (see section 4), and prepared to be presented to the final users.

WEB APPLICATION

BIRDIE's web application offers four main menus that provide access to the pipeline outputs in different ways (see Figure 6):

1. An exploration map. Through this menu the user is free to explore the different indicators BIRDIE computes using a map. This spatial framework can be configured to display information layers, such

- as occupancy probabilities for ABAP pentads or waterbird abundance at CWAC sites. Users can also zoom in and out to find the scale that best fits their needs.
2. Site and species summaries, are detailed reports elaborated for users specifically interested in sites or species, rather than in general exploration. These summaries can also be accessed indirectly through the exploration map.
 3. Reporting tools. We mentioned in section 2 that BIRDIE was developed to support reports for national and international conservation programmes. In this menu, users interested in elaborating, or accessing the information underpinning these reports, will find this information conveniently packed in programme-specific summaries.
 4. Web services. Through this menu users can access BIRDIE's API and retrieve its outputs in the most flexible way. It is through BIRDIE's API that all maps and plots in the web application are produced. By accessing this functionality directly, users can download the data themselves and incorporate them into their own workflows.

DISCUSSION

Data on biodiversity and related environmental drivers are collected at increasingly faster rates. Although these data can be accessed to support decisions at various levels, it can be difficult for decision makers to extract relevant information in a timely fashion (MacFadyen et al., 2022; P. J. Stephenson et al., 2017). Here, we introduce BIRDIE, the South African Biodiversity Data Pipeline for Wetlands and Waterbirds; a data pipeline that aims to provide the information needed for making evidence-based decisions on wetlands and waterbirds in southern Africa. Target users of BIRDIE include government and public entities that need to report on the status of wetlands and waterbirds, as well as site managers, and the general public (e.g., birdwatchers). Apart from data availability and accessibility, obstacles for using biodiversity data in decision-making include (MacFadyen et al., 2022; P. J. Stephenson et al., 2017): lack of analysis and interpretation, lack of technical accessibility with excessive use of jargon, and timely use of data.

All data used by BIRDIE are freely available, so one of the core objectives of BIRDIE is to facilitate information uptake by statistically analysing these data and filtering out observational artefacts introduced during data collection. Uneven sampling efforts, imperfect detection and missing data are all examples of how data collection methods can affect data, and if undealt with, mislead decision making. Furthermore, statistical models also provide measures of uncertainty in their estimates, which must be clearly communicated to the stakeholders. With all their benefits, these statistical analyses require technical knowledge and are time-consuming. Therefore, having their outputs pre-computed and readily available could dramatically increase the impact of the data.

BIRDIE relies heavily on citizen science projects such as ABAP and CWAC, which poses clear challenges in terms of uneven efforts and imperfect detection, but also adds the advantage of having the support of a large community of observers that provides a continuous and steady flow of data. These data inputs allow us to run the pipeline periodically to keep the indicators updated and timely. Although we would like to update our indicators more often, for now we only update once per year due to the computational requirements of the pipeline, and certain characteristics of the data (e.g., CWAC counts are conducted only twice a year).

Our web application was designed to present indicators that are easily accessed, visualised and interpreted, avoiding unnecessary jargon. At the same time, and for the interested user, we have allocated some space for clearly explaining the analytical routines used in all the analyses in dedicated sections. In BIRDIE, we try to follow the Findable, Accessible, Interoperable, Reusable (FAIR) principles (Wilkinson et al.,

2016), making all processes reproducible and transparent. All the code used by the pipeline is public, freely available and based on open-source software.

Whenever possible we allow the user to display multiple pieces of information in the same screen, including some environmental data that we consider relevant, to facilitate information contextualization. For certain specific users, namely those interested in the main national and international reporting channels (e.g., Ramsar, AEWA, NBA), we prepare dedicated reports that combine the information that suits their needs. The contents of these reports have been defined during stakeholder engagement, and are intended to be dynamically adapted to new necessities.

The BIRDIE data pipeline consists of several fundamental building blocks. The first block, which we call the data source layer (Figure 5) hosts and curates the raw data. Vetting and data cleaning also happens in this block. The second block, the analysis layer, analyses the data and estimates the fundamental quantities of interest, like abundance and occurrence of each species at each wetland. The third block consists of a data mart where the outputs of the analysis are stored and indicators are aggregated or disaggregated to multiple scales. The final block serves the information to the user via a web application. The modular structure of BIRDIE enables us to maintain and update individual parts independently. For example, we could replace the current statistical routines with more efficient ones without changing the other parts of the pipeline. Or we could add new indicators to the data mart layer without needing to change the statistical routines that produce the underlying components.

Running models automatically and periodically for multiple species requires pre-defining models and using similar models for all species. Therefore there is a trade-off between having accurate models that fit individual species well and having a pipeline that works well for all species in general. Users should keep in mind this compromise, and think of BIRDIE's outputs as useful approximations rather than precise estimates. We recommend designing bespoke models for those species for which precision is required.

In the next phase, we intend to develop more profound links between waterbird population indicators and wetlands. Waterbirds are often regarded as good indicators of wetland biodiversity and condition. However, this assumption is rarely proven empirically, and it is apparent that it needs careful consideration on a case by case basis (Amat and Green, 2010). With advances in the accessibility to biodiversity data, we are now in a better position to investigate whether these claims hold, and if so, under which conditions. Data portals such as GBIF.org and in South Africa, the Freshwater Biodiversity Information System (FBIS), and SANBI's biodiversity data portal. SANBI has also been developing the wetland condition map that provides information about the ecological status of wetlands, looking at elements such as water quality and vegetation. These additional sources of information could help us understand how waterbird occurrence, abundance and diversity relates to the general ecological condition of the hosting wetlands.

We will also extend BIRDIE's functionality to cover other African countries with similar available data, such as Kenya and Nigeria that also use the ABAP protocol. There is also a wealth of information that BIRDIE has not yet used, such as eBird or iNaturalist, that could improve the outputs of the pipeline. However, integrating data sources with different sampling designs, coverages and biases is not trivial. Data integration is a very active topic in the field of statistical ecology. Approaches to combining data range from pooling multiple data sources together disregarding their different assumptions and biases, to much more accurate integrated models in which characteristics of each data source are explicitly accounted for (Fletcher et al., 2019). Although at the expense of increased model complexity, with the application of newly-developed statistical methods for data integration, we can now explore how different species interrelate, and inform more effective and efficient conservation actions.

CONFLICT OF INTEREST

449 The authors declare that the research was conducted in the absence of any commercial or financial
450 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

451 All authors contributed to conception and design of the project. NJ is the project director; FC, RA, and
452 VV developed the analyses of the pipeline; NJ, AS and DH work on reporting and indicators; FS and YS,
453 designed and implement the data mart and web application; MB manages the citizen science database. FC
454 lead the writing of the manuscript with contribution from all authors, who also revised, read, and approved
455 the submitted version.

FUNDING

456 This project is funded by the JRS Biodiversity Foundation, grant number 60908.

ACKNOWLEDGMENTS

457 We are really grateful to other members of the BIRDIE team without whom this project would not be
458 viable: Sediqa Khatieb, Monica Klass, and Carol Poole. We are also grateful for the support of the JRS
459 Biodiversity Foundation, and to the many interested users that have engaged and shared useful insights
460 with us. Finally, we would like to recognize the tremendous contribution of all the citizen scientists that
461 devote their time and effort to collect the valuable data that we use.

REFERENCES

- 462 Abatzoglou, J.T., Dobrowski, S.Z., Parks, S.A., Hegewisch, K.C., 2018. TerraClimate, a high-resolution
463 global dataset of monthly climate and climatic water balance from 1958–2015. *Sci. Data* 5, 170191.
464 <https://doi.org/10.1038/sdata.2017.191>
- 465 Altwegg, R., Nichols, J.D., 2019. Occupancy models for citizen-science data. *Methods Ecol. Evol.* 10,
466 8–21. <https://doi.org/10.1111/2041-210X.13090>
- 467 Barnard, P., Altwegg, R., Ebrahim, I., Underhill, L.G., 2017. Early warning systems for biodiversity in
468 southern Africa – How much can citizen science mitigate imperfect data? *Biol. Conserv.* 208, 183–188.
469 <https://doi.org/10.1016/j.biocon.2016.09.011>
- 470 Brooks, M., Rose, S., Altwegg, R., Lee, A.T., Nel, H., Ottosson, U., Retief, E., Reynolds, C., Ryan, P.G.,
471 Shema, S., Tende, T., Underhill, L.G., Thomson, R.L., 2022. The African Bird Atlas Project: a description
472 of the project and BirdMap data-collection protocol. *Ostrich* 1–10. [https://doi.org/10.2989/](https://doi.org/10.2989/00306525.2022.2125097)
473 [00306525.2022.2125097](https://doi.org/10.2989/00306525.2022.2125097)
- 474 Buckland, S.T., Newman, K.B., Thomas, L., Koesters, N.B., 2004. State-space models for the
475 dynamics of wild animal populations. *Ecol. Model.* 171, 157–175. [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.ecolmodel.2003.08.002)
476 [ecolmodel.2003.08.002](https://doi.org/10.1016/j.ecolmodel.2003.08.002)
- 477 CBD, 2022. Convention on Biological Diversity [WWW Document]. *Conv. Biol. Divers.* URL [https :](https://www.cbd.int/)
478 [/www.cbd.int/](https://www.cbd.int/) (accessed 12.22.22).
- 479 Convention on Wetlands, 2021. *Global Wetland Outlook: Special Edition 2021*. Secretariat of the
480 Convention on Wetlands, Gland, Switzerland.

- 481 Dallas, H., Shelton, J., Sutton, T., Tri Cuptura, D., Kajee, M., Job, N., 2021. The Freshwater Biodiversity
482 Information System (FBIS) – mobilising data for evaluating long-term change in South African rivers. *Afr.*
483 *J. Aquat. Sci.* 1–16. <https://doi.org/10.2989/16085914.2021.1982672>
- 484 Didan, Kamel, 2015. MOD13A2 MODIS/Terra Vegetation Indices 16-Day L3 Global 1km SIN Grid
485 V006. <https://doi.org/10.5067/MODIS/MOD13A2.006>
- 486 Doser, J.W., Finley, A.O., Kéry, M., Zipkin, E.F., 2022. spOccupancy: An R package for single-
487 species, multi-species, and integrated spatial occupancy models. *Methods Ecol. Evol.* 13, 1670–1678.
488 <https://doi.org/10.1111/2041-210X.13897>
- 489 FIAO, F.I. of A.O., 2022. CWAC: Coordinated Waterbird Counts [WWW Document]. URL [https://](https://cwac.birdmap.africa/index.php)
490 cwac.birdmap.africa/index.php (accessed 12.21.22).
- 491 Gelman, A., Carlin, J., Stern, H., Dunson, D., Vehtari, A., Rubin, D., 2014. Bayesian Data Analysis.
492 CRC Press, Taylor and Francis Group, Boca Raton, FL.
- 493 Gimenez, O., Buckland, S.T., Morgan, B.J.T., Bez, N., Bertrand, S., Choquet, R., Dray, S., Etienne, M.-P.,
494 Fewster, R., Gosselin, F., Mérigot, B., Monestiez, P., Morales, J.M., Mortier, F., Munoz, F., Ovaskainen,
495 O., Pavoine, S., Pradel, R., Schurr, F.M., Thomas, L., Thuiller, W., Trenkel, V., de Valpine, P., Rexstad,
496 E., 2014. Statistical ecology comes of age. *Biol. Lett.* 10, 20140698. [https://doi.org/10.1098/](https://doi.org/10.1098/rsbl.2014.0698)
497 [rsbl.2014.0698](https://doi.org/10.1098/rsbl.2014.0698)
- 498 Han, X., Josse, C., Young, B.E., Smyth, R.L., Hamilton, H.H., Bowles-Newark, N., 2017. Monitoring
499 national conservation progress with indicators derived from global and national datasets. *Biol. Conserv.*
500 213, 325–334. <https://doi.org/10.1016/j.biocon.2016.08.023>
- 501 IUCN, 2022. International Union for the Conservation of Nature [WWW Document]. IUCN. URL
502 <https://www.iucn.org/content/home-page> (accessed 12.22.22).
- 503 Jetz, W., McGeoch, M.A., Guralnick, R., Ferrier, S., Beck, J., Costello, M.J., Fernandez, M., Geller,
504 G.N., Keil, P., Merow, C., Meyer, C., Muller-Karger, F.E., Pereira, H.M., Regan, E.C., Schmeller, D.S.,
505 Turak, E., 2019. Essential biodiversity variables for mapping and monitoring species populations. *Nat.*
506 *Ecol. Evol.* 3, 539–551. <https://doi.org/10.1038/s41559-019-0826-1>
- 507 Kellner, K., 2021. jagsUI: A Wrapper Around “rjags” to Streamline “JAGS” Analyses.
- 508 King, R., 2014. Statistical Ecology. *Annu. Rev. Stat. Its Appl.* 1, 401–426. [https://doi.org/10.](https://doi.org/10.1146/annurev-statistics-022513-115633)
509 [1146/annurev-statistics-022513-115633](https://doi.org/10.1146/annurev-statistics-022513-115633)
- 510 Kissling, W.D., Ahumada, J.A., Bowser, A., Fernandez, M., Fernández, N., García, E.A., Guralnick, R.P.,
511 Isaac, N.J.B., Kelling, S., Los, W., McRae, L., Mihoub, J., Obst, M., Santamaria, M., Skidmore, A.K.,
512 Williams, K.J., Agosti, D., Amariles, D.,
- 513 Arvanitidis, C., Bastin, L., De Leo, F., Egloff, W., Elith, J., Hobern, D., Martin, D., Pereira, H.M., Pesole,
514 G., Peterseil, J., Saarenmaa, H., Schigel, D., Schmeller, D.S., Segata, N., Turak, E., Uhlir, P.F., Wee,
515 B., Hardisty, A.R., 2018. Building essential biodiversity variables (EBV s) of species distribution and
516 abundance at a global scale. *Biol. Rev.* 93, 600–625. <https://doi.org/10.1111/brv.12359>
- 517 Mace, G.M., Barrett, M., Burgess, N.D., Cornell, S.E., Freeman, R., Grooten, M., Purvis, A., 2018.
518 Aiming higher to bend the curve of biodiversity loss. *Nat. Sustain.* 1, 448–451. [https://doi.org/](https://doi.org/10.1038/s41893-018-0130-0)
519 [10.1038/s41893-018-0130-0](https://doi.org/10.1038/s41893-018-0130-0)

- MacFadyen, S., Allsopp, N., Altwegg, R., Archibald, S., Botha, J., Bradshaw, K., Carruthers, J., De Klerk, H., de Vos, A., Distiller, G., Foord, S., Freitag-Ronaldson, S., Gibbs, R., Hamer, M., Landi, P., MacFadyen, D., Manuel, J., Midgley, G., Moncrieff, G., Munch, Z., Mutanga, O., Sershen, Nenguda, R., Ngwenya, M., Parker, D., Peel, M., Power, J., Pretorius, J., Ramdhani, S., Robertson, M., Rushworth, I., Skowno, A., Slingsby, J., Turner, A., Visser, V., Van Wageningen, G., Hui, C., 2022. Drowning in data, thirsty for information and starved for understanding: A biodiversity information hub for cooperative environmental monitoring in South Africa. *Biol. Conserv.* 274, 109736. <https://doi.org/10.1016/j.biocon.2022.109736>
- MacKenzie, D.I., Nichols, J.D., Lachman, G.B., Droege, S., Andrew Royle, J., Langtimm, C.A., 2002. Estimating site occupancy rates when detection probabilities are less than one. *Ecology* 83, 2248–2255. [https://doi.org/10.1890/0012-9658\(2002\)083%5B2248:ESORWD%5D2.0.CO;2](https://doi.org/10.1890/0012-9658(2002)083%5B2248:ESORWD%5D2.0.CO;2)
- Newman, K.B., Buckland, S.T., Morgan, B.J.T., King, R., Borchers, D.L., Cole, D.J., Besbeas, P., Gimenez, O., Thomas, L., 2014. *Modelling Population Dynamics: model formulation, fitting and assessment using state-space methods*. Springer, New York, NY.
- Pekel, J.-F., Cottam, A., Gorelick, N., Belward, A.S., 2016. High-resolution mapping of global surface water and its long-term changes. *Nature* 540, 418–422. <https://doi.org/10.1038/nature20584>
- Pereira, H.M., Ferrier, S., Walters, M., Geller, G.N., Jongman, R.H.G., Scholes, R.J., Bruford, M.W., Brummitt, N., Butchart, S.H.M., Cardoso, A.C., Coops, N.C., Dulloo, E., Faith, D.P., Freyhof, J., Gregory, R.D., Heip, C., Höft, R., Hurtt, G., Jetz, W., Karp, D.S., McGeoch, M.A., Obura, D., Onoda, Y., Pettorelli, N., Reyers, B., Sayre, R., Scharlemann, J.P.W., Stuart, S.N., Turak, E., Walpole, M., Wegmann, M., 2013. Essential Biodiversity Variables. *Science* 339, 277–278. <https://doi.org/10.1126/science.1229931>
- Plummer, M., 2003. JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling, in: *Proceedings of the 3rd International Workshop on Distributed Statistical Computing*.
- Pollock, L.J., O'Connor, L.M.J., Mokany, K., Rosauer, D.F., Talluto, M.V., Thuiller, W., 2020. Protecting Biodiversity (in All Its Complexity): New Models and Methods. *Trends Ecol. Evol.* 35, 1119–1128. <https://doi.org/10.1016/j.tree.2020.08.015>
- R Core Team, 2022. *R: A Language and Environment for Statistical Computing*.
- SANBI, South African National Biodiversity Institute, 2023. Biodiversity Advisor [WWW Document]. URL <http://biodiversityadvisor.sanbi.org/> (accessed 12.21.22).
- SANBI, South African National Biodiversity Institute, in prep. National Wetland Map version 6.
- Skowno, A., Poole, C.J., Raimondo, D.C., Sink, K.J., Van Deventer, H., Van Niekerk, L., Harris, L.R., Smith-Adao, L.B., Tolley, K.A., Zengeya, T.A., Foden, W.B., Midgley, G.F., Driver, A., 2019. National biodiversity assessment 2018: the status of South Africa's ecosystems and biodiversity: synthesis report. South African National Biodiversity Institute, Department of Environment, Forestry and Fisheries, Pretoria.
- Stephenson, P., Brooks, T.M., Butchart, S.H., Fegraus, E., Geller, G.N., Hoft, R., Hutton, J., Kingston, N., Long, B., McRae, L., 2017. Priorities for big biodiversity data. *Front. Ecol. Environ.* 15, 124–125. <https://doi.org/10.1002/fee.1473>
- Stephenson, P.J., Bowles-Newark, N., Regan, E., Stanwell-Smith, D., Diagana, M., Höft, R., Abarchi, H., Abrahamse, T., Akello, C., Allison, H., Banki, O., Batieno, B., Dieme, S., Domingos, A., Galt, R., Githaiga,

- 561 C.W., Guindo, A.B., Hafashimana, D.L.N., Hirsch, T., Hobern, D., Kaaya, J., Kaggwa, R., Kalembe, M.M.,
 562 Linjouom, I., Manaka, B., Mbwambo, Z., Musasa, M., Okoree, E., Rwetsiba, A., Siam, A.B., Thiombiano,
 563 A., 2017. Unblocking the flow of biodiversity data for decision-making in Africa. *Biol. Conserv.* 213,
 564 335–340. <https://doi.org/10.1016/j.biocon.2016.09.003>
- 565 Stephenson, P.J., Ntiamoa-Baidu, Y., Simaika, J.P., 2020. The Use of Traditional and Modern Tools for
 566 Monitoring Wetlands Biodiversity in Africa: Challenges and Opportunities. *Front. Environ. Sci.* 8, 1–12.
 567 <https://doi.org/10.3389/fenvs.2020.00061>
- 568 UN, United Nations, 2022. Sustainable Development Goals [WWW Document]. URL [https://sdgs.](https://sdgs.un.org/)
 569 [un.org/](https://sdgs.un.org/) (accessed 12.22.22).
- 570 UNEP, United Nations Environmental Programme, 2022. AEWA: Agreement on the Conservation of
 571 African-Eurasian Migratory Waterbirds [WWW Document]. URL <https://www.unep-aewa.org/>
 572 (accessed 12.22.22).
- 573 Wetzel, F.T., Saarenmaa, H., Regan, E., Martin, C.S., Mergen, P., Smirnova, L., Tuama, É.Ó., García
 574 Camacho, F.A., Hoffmann, A., Vohland, K., Häuser, C.L., 2015. The roles and contributions of Biodiversity
 575 Observation Networks (BONs) in better tracking progress to 2020 biodiversity targets: a European case
 576 study. *Biodiversity* 16, 137–149. <https://doi.org/10.1080/14888386.2015.1075902>
- 577 White, E.P., Yenni, G.M., Taylor, S.D., Christensen, E.M., Bledsoe, E.K., Simonis, J.L., Ernest, S.K.M.,
 578 2019. Developing an automated iterative near-term forecasting system for an ecological study. *Methods*
 579 *Ecol. Evol.* 10, 332–344. <https://doi.org/10.1111/2041-210X.13104>
- 580 Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M., Baak, A., Blomberg, N.,
 581 Boiten, J.-W., da Silva Santos, L.B., Bourne, P.E., Bouwman, J., Brookes, A.J., Clark, T., Crosas, M.,
 582 Dillo, I., Dumon, O., Edmunds, S., Evelo, C.T., Finkers, R., Gonzalez-Beltran, A., Gray, A.J.G., Groth, P.,
 583 Goble, C., Grethe, J.S., Heringa, J., 't Hoen, P.A.C., Hooft, R., Kuhn, T., Kok, R., Kok, J., Lusher, S.J.,
 584 Martone, M.E., Mons, A., Packer, A.L., Persson, B., Rocca-Serra, P., Roos, M., van Schaik, R., Sansone,
 585 S.-A., Schultes, E., Sengstag, T., Slater, T., Strawn, G., Swertz, M.A., Thompson, M., van der Lei, J.,
 586 van Mulligen, E., Velterop, J., Waagmeester, A., Wittenburg, P., Wolstencroft, K., Zhao, J., Mons, B.,
 587 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* 3, 160018.
 588 <https://doi.org/10.1038/sdata.2016.18>
- 589 Yamazaki, D., Ikeshima, D., Tawatari, R., Yamaguchi, T., O'Loughlin, F., Neal, J.C., Sampson, C.C.,
 590 Kanae, S., Bates, P.D., 2017. A high-accuracy map of global terrain elevations: Accurate Global Terrain
 591 Elevation map. *Geophys. Res. Lett.* 44, 5844–5853. <https://doi.org/10.1002/2017GL072874>
- 592 Yenni, G.M., Christensen, E.M., Bledsoe, E.K., Supp, S.R., Diaz, R.M., White, E.P., Ernest, S.K.M.,
 593 2019. Developing a modern data workflow for regularly updated data. *PLOS Biol.* 17, e3000125. <https://doi.org/10.1371/journal.pbio.3000125>
- 595 Yoccoz, N.G., Nichols, J.D., Boulinier, T., 2001. Monitoring of biological diversity in space and time.
 596 *Trends Ecol. Evol.* 16, 446–453. [https://doi.org/10.1016/S0169-5347\(01\)02205-4](https://doi.org/10.1016/S0169-5347(01)02205-4)

FIGURES AND CAPTIONS

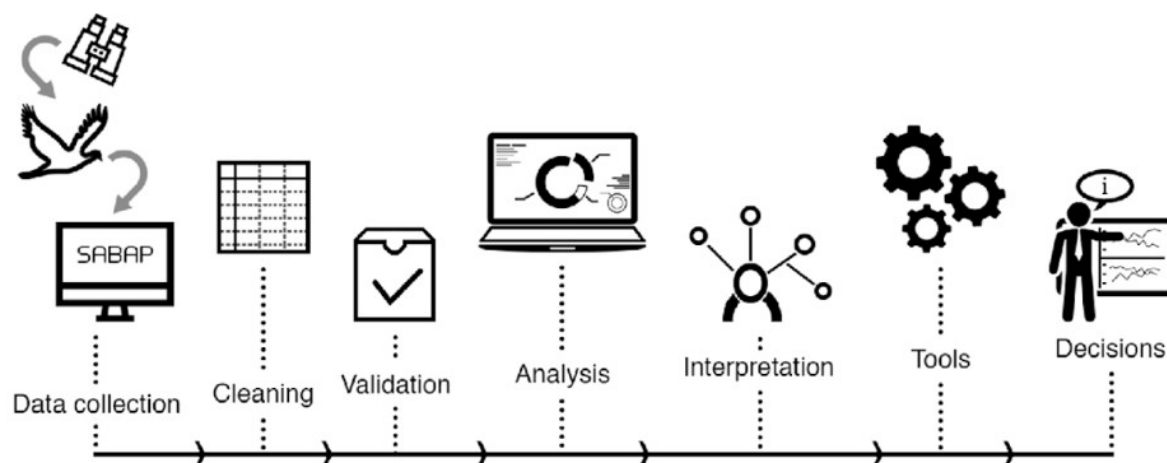


Figure 1. Basic workflow of the BIRDIE pipeline covering all steps from data collection, to analysis and presentation of digested, decision-ready indicators.

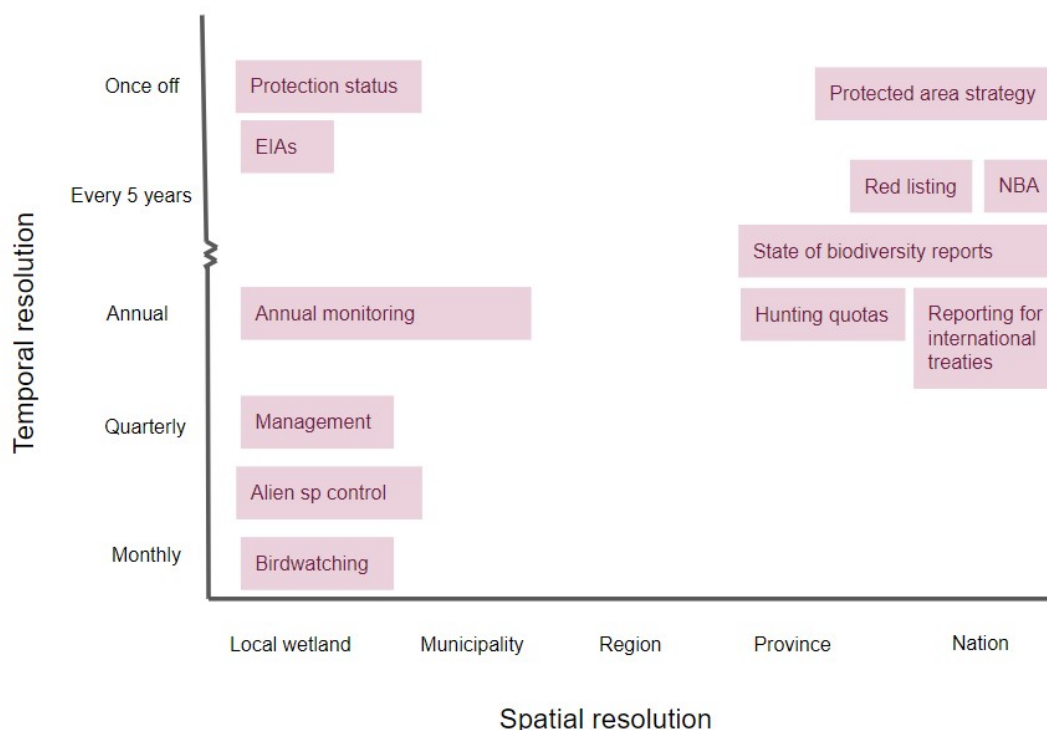


Figure 2. Main users targeted by BIRDIE in relation to their spatial and temporal assessment scales. At present, we focus on computing indicators at an annual temporal resolution or coarser. Finer resolutions are typically based on access to raw data.

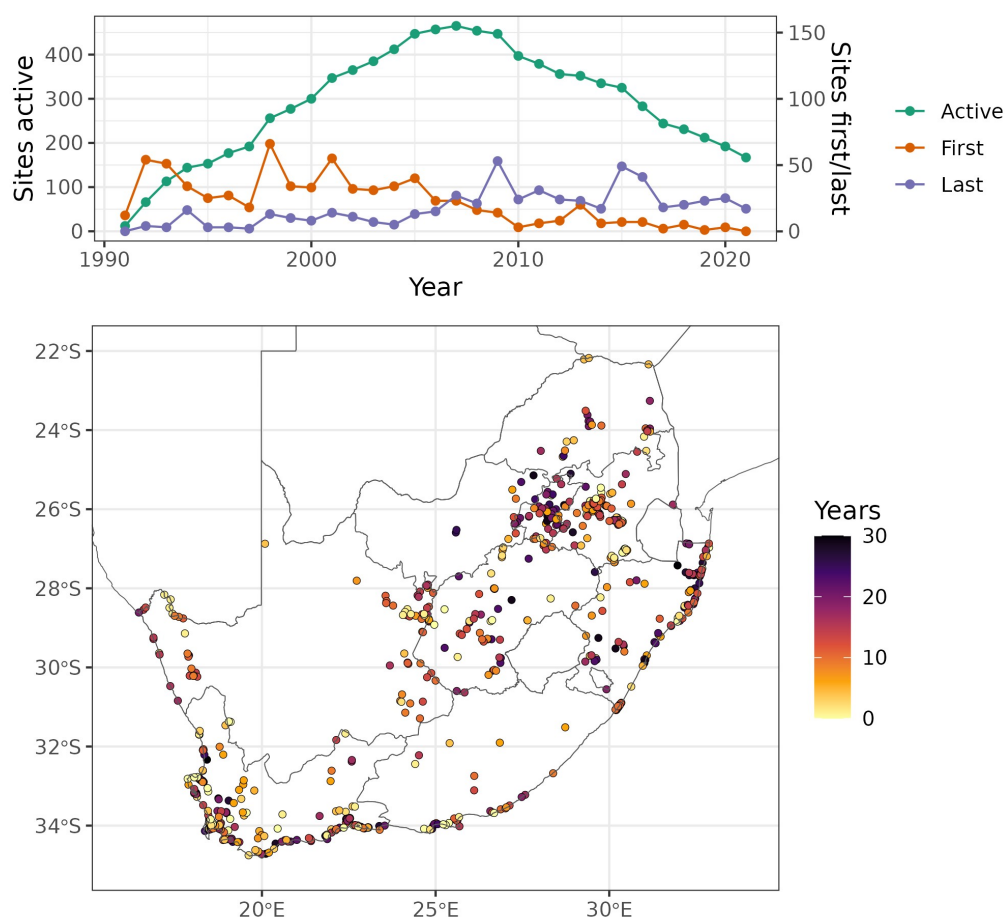


Figure 3. The graph shows, the number of CWAC sites active (green), firstly counted (red) and last counted (purple), per year, between 1991 and 2021. Note that some of the sites that were last counted before 2021, might be counted again in the future. In the map, the spatial location of CWAC sites in South Africa. The colour gradient represent the duration of the period the site was counted for.

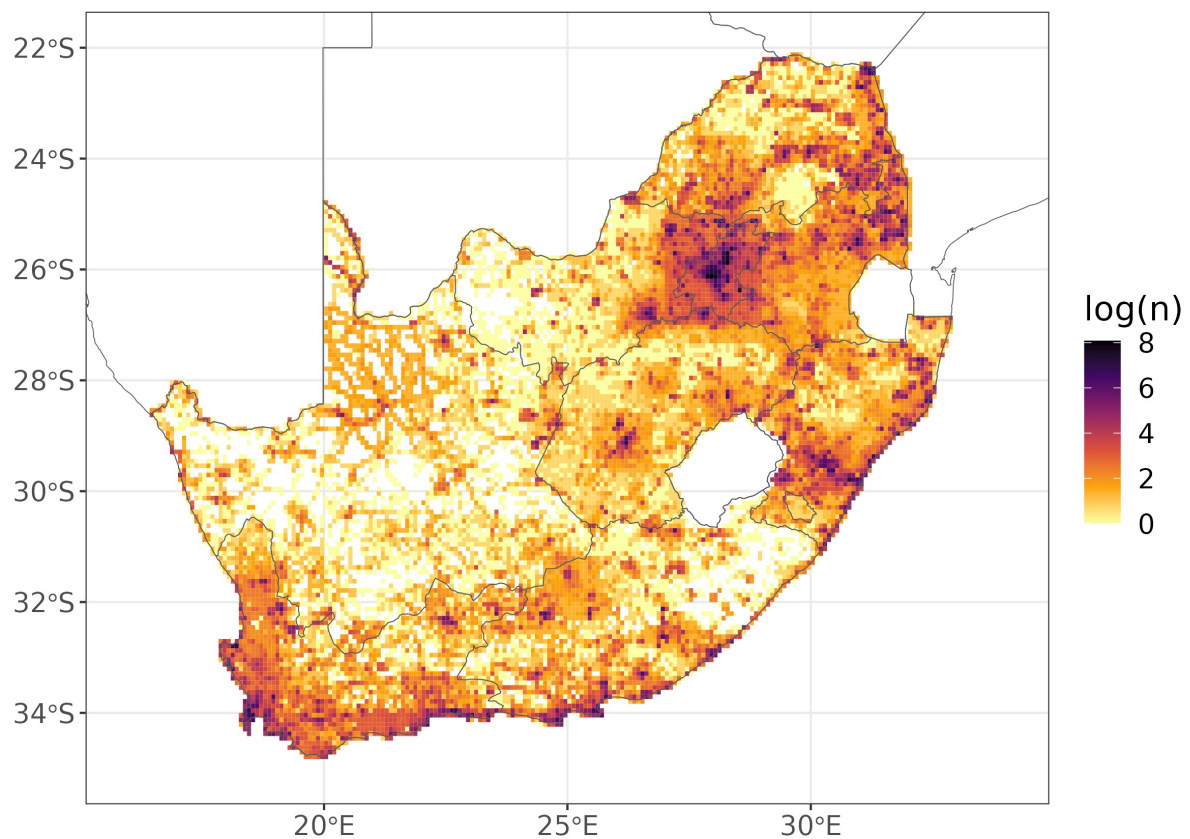


Figure 4. Number of SABAP2 cards recorded for the South African pentads between 2008-2021, in logarithmic scale. We can see how areas close to large cities in the Western Cape and Gauteng provinces, accumulate larger efforts. We can also appreciate sampling biased towards roads, particularly in the northwest of the country.

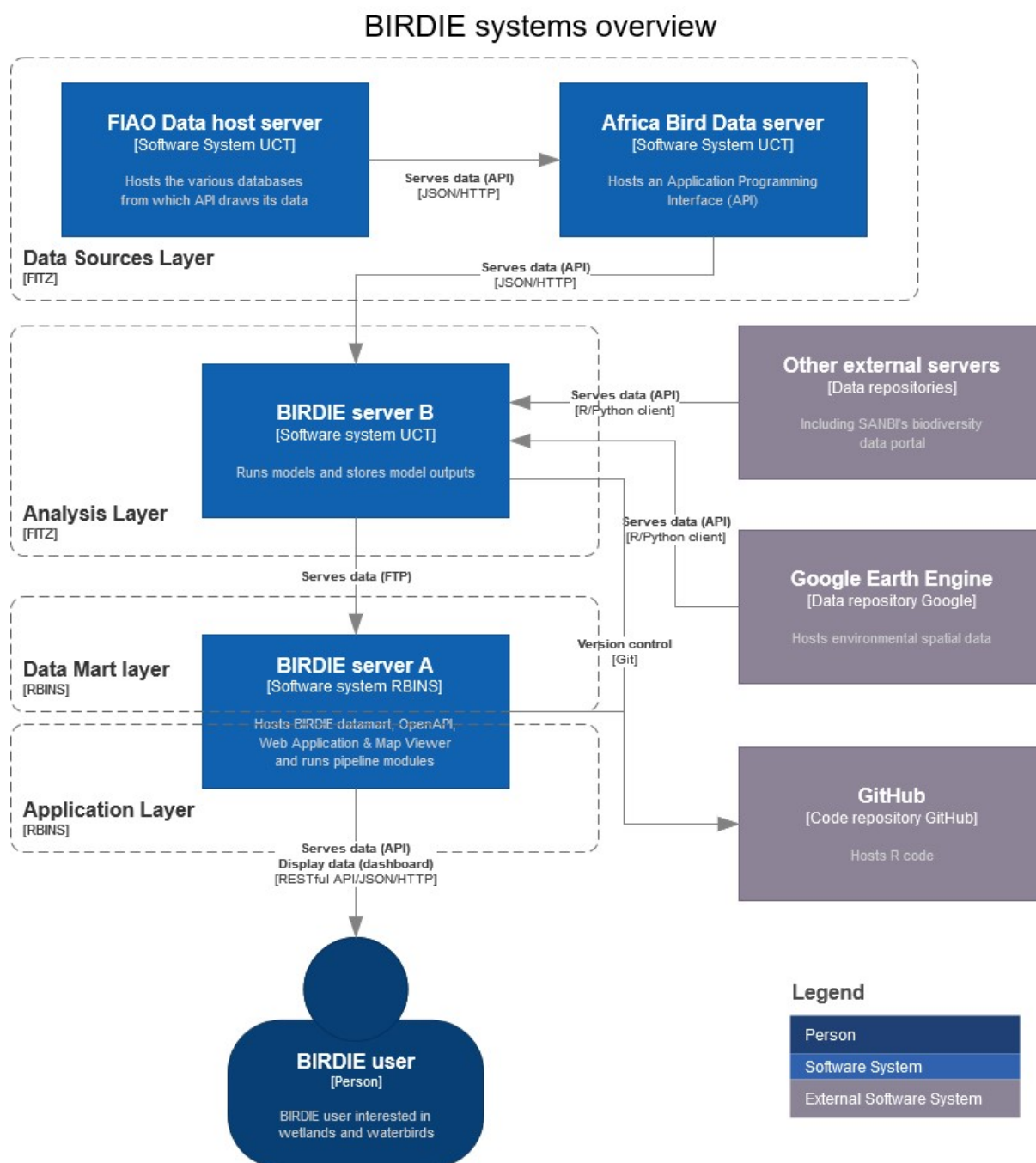
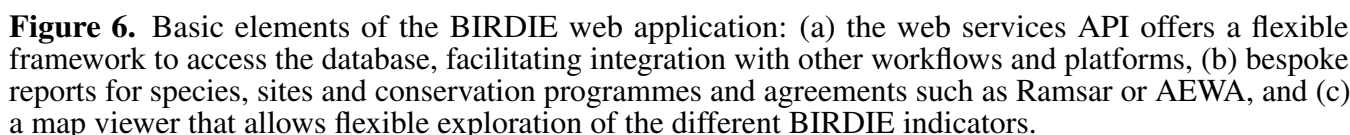


Figure 5. Overview of BIRDIE’s server architecture. Data flows from CWAC, ABAP and other external servers into BIRDIE server B to be processed and analysed by the R modules, then these outputs move into the data mart in BIRDIE server A, which is the gateway for the dashboard and the final users.



TABLES

Table 1. Main indicators produced by the BIRDIE pipeline for waterbird species. For each indicator, we show the inputs, which can be databases (Coordinated Waterbird Counts - CWAC and the African Bird Atlas Project - ABAP), or other indicators; models used to compute the indicator (state-space model -SSM, and occupancy model - Occupancy) or whether it was computed by aggregating other lower-level indicators; the smaller spatial scale of assessment; and the smaller temporal scale of assessment. Annual changes in all of these indicators are also computed, and other indicators will be added over time as needed.

Indicator	Input	Model	Spatial scale	Temporal scale
Abundance	CWAC	SSM	CWAC site	2 seasons/year
Diversity	ABAP	Occupancy	Pentad	Annual
Extent of occurrence	Occurrence	Aggregated	National	Annual
Area of occupancy	Occurrence	Aggregated	National	Annual
Population size	Abundance	Aggregated	National	2 seasons/year
Pop. proportion on site	Abundance	Aggregated	CWAC site/national	2 seasons/year
Waterbird Conservation Value	Abundance	Aggregated	CWAC site/national	2 seasons/year
Number of sites	Abu./occu.	Aggregated	National	Annual