

Evaluating the method reproducibility of deep learning models in the biodiversity domain

Waqas Ahmed^{1,5}, Vamsi Krishna Kommineni^{1,3,4,5}, Birgitta König-Ries^{1,2,4}, Jitendra Gaikwad¹, Luiz Gadelha¹, and Sheeba Samuel^{1,2}

¹Heinz Nixdorf Chair for Distributed Information Systems, Friedrich Schiller University Jena, Ernst-Abbe-Platz 2, Jena, 07743, Thuringia, Germany

²Michael Stifel Center Jena, Friedrich Schiller University Jena, Leutragraben 1, Jena, 07743, Thuringia, Germany

³Department of Functional Biogeography, Max Planck Institute for Biogeochemistry, Hans-Knoell-Str. 10, Jena, 07745, Thuringia, Germany

⁴German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Puschstraße 4, Leipzig, 04103, Saxony, Germany

Corresponding author:

Waqas Ahmed¹

Email address: ahmed.waqas@uni-jena.de

ABSTRACT

Artificial Intelligence (AI) is revolutionizing biodiversity research by enabling advanced data analysis, species identification, and habitats monitoring, thereby enhancing conservation efforts. Ensuring reproducibility in AI-driven biodiversity research is crucial for fostering transparency, verifying results, and promoting the credibility of ecological findings. This study investigates the reproducibility of deep learning (DL) methods within the biodiversity domain. We design a methodology for evaluating the reproducibility of biodiversity-related publications that employ DL techniques across three stages. We define ten variables essential for method reproducibility, divided into four categories: resource requirements, methodological information, uncontrolled randomness, and statistical considerations. These categories subsequently serve as the basis for defining different levels of reproducibility. We manually extract the availability of these variables from a curated dataset comprising 61 publications identified using the keywords provided by biodiversity experts. Our study shows that the dataset is shared in 47% of the publications; however, a significant number of the publications lack comprehensive information on deep learning methods, including details regarding randomness.

INTRODUCTION

In recent years, deep learning methods have been increasingly applied to understand complex ecological systems, particularly in the field of biodiversity. These methods have the potential to process and analyze large amounts of biological data rapidly, leading to significant insights. For instance, August et al. (2020) demonstrated how AI image classifiers can create new biodiversity datasets from social media imagery, highlighting the spatial and taxonomic biases that can influence ecological inference. Similarly, DL models have been utilized to analyze camera trap images for wildlife monitoring, enabling researchers to identify species and infer ecological patterns and processes with high accuracy in Tabak et al. (2019). However, there is a growing concern about the reproducibility, transparency, and trustworthiness of research findings produced using deep learning methods in this domain (Feng et al., 2019; GPAI, 2022).

Reproducibility is essential in scientific research as it allows researchers to validate and advance methods and results, ensuring the reliability of scientific claims (Samuel and König-Ries, 2021). Goodman et al. (2016) define research reproducibility in three aspects: methods reproducibility, results reproducibility, and inferential reproducibility. In this paper, we focus on methods reproducibility, which involves the ability to exactly replicate the experimental and computational processes using the same data and tools to

achieve the same outcomes. This is especially important in biodiversity conservation, where decisions can directly impact ecological health. For example, Norouzzadeh et al. (2018) demonstrated the application of deep learning for automatically identifying and counting wildlife in camera trap images, which is vital for monitoring species populations and informing conservation strategies. However, failure to reproduce these methods accurately could lead to erroneous population estimates, thereby compromising conservation efforts and resource allocation. Similarly, in the study by Rovero et al. (2013), the deployment of camera traps for wildlife monitoring highlighted the need for reproducible methods to ensure consistent data collection and analysis across different geographical locations. Inconsistent methods could result in data that are not comparable, leading to flawed conclusions and ineffective conservation measures. Moreover, the inability to reproduce results due to methodological inconsistencies can prevent the detection of errors that might hide underlying biases or issues in the data. For instance, Christin et al. (2019) discuss how the reproducibility of deep learning models in ecological research is challenged by the complexity and heterogeneity of biodiversity data, which includes interactions between variables, missing values, and non-linearity. Ensuring reproducibility helps to uncover these complexities and improve the reliability of ecological models used for decision-making.

A deep learning pipeline is a structured sequence of processes in training and deploying a DL model (El-Amir and Hamdy, 2020). The pipeline typically begins with data collection and preprocessing, involving tasks such as data cleaning, normalization, and transformation. Following data preprocessing, the next stage consists of designing and selecting an appropriate deep learning architecture, considering factors like model complexity and the nature of the problem. Subsequently, model training takes place, where the chosen architecture is trained on the preprocessed data using optimization algorithms and specific hyperparameter configurations. After training, the model is evaluated and fine-tuned using various performance metrics to ensure its effectiveness in solving the targeted problem. The best performing model is run on the test data. To ensure an unbiased evaluation of the model's predictive performance. Finally, the trained model is deployed for real-world applications or further refinement.

To ensure the reproducibility of the deep learning pipeline, comprehensive documentation is crucial at each stage. This includes detailed records of the data collection steps (including providing persistent identifier for each data point), data preprocessing steps, such as the specific data transformation techniques applied and any data augmentation strategies employed. Additionally, it is vital to document the specifics of the chosen deep learning architecture, including the exact configurations and versions of the neural network layers utilized. Detailed notes on the hyperparameter values selected during the model training phase are essential, as well as the training convergence criteria and the optimization algorithms employed. Proper documentation of the evaluation metrics used and the testing dataset ensures the reproducibility of the model's performance assessment. Finally, maintaining a record of the software libraries, hardware, frameworks, and versions utilized throughout the pipeline aids in replicating the experimental setup accurately.

In this paper, we aim to shed light on the current state of reproducibility of deep learning methods applied to biodiversity research. We conducted a systematic literature review to identify publications that use deep learning techniques in biodiversity research using keywords provided by biodiversity experts (Abdelmageed et al., 2022). We define various reproducibility-related variables inspired by the current literature. We then curated a dataset of 61 publications and manually extracted reproducibility-related variables, such as the availability of datasets and code. We also analyzed advanced reproducibility variables, such as the specific deep learning methods employed, models, and hyperparameters. Our findings suggest that the reproducibility of deep learning methods in biodiversity research is generally low. However, there is a growing trend towards improving reproducibility, with more and more publications making their datasets and code available. Our study will contribute to the ongoing discourse on the reproducibility of deep learning methods in biodiversity research and help to improve the credibility and impact of these methods in this vital field.

In the following sections, we provide a detailed description of our study. We start with an overview of the state-of-the-art ("Related Work"). We provide the methodology of our study ("Methodology") We describe our results and discuss the implications of our work ("Discussion"). Finally, we summarize the key aspects of our study and provide future directions of our research ("Conclusion").

RELATED WORK

Method reproducibility, as mentioned in Goodman et al. (2016), is an important aspect of progress in any scientific field. Raff (2019) provides some insight into reproducibility in ML. The author studied the reproducibility of 255 ML papers published from 1984 to 2017 and tried to correlate the reproducibility of the papers with their characteristics. According to his observation, the main obstacle to reproducibility of results is insufficient explanation of the method and implementation details when the source code is not provided. However, some of the defined attributes were very subjective, such as the algorithmic difficulty of the work or the readability of the work. While Raff (2019) provides insight into ML reproducibility in general, there is also research specifically related to the field of biodiversity (Feng et al., 2019; Schnitzer and Carson, 2016). They pointed out that reproducibility is problematic because erroneous data is widespread and there is a lack of empirical studies to verify previous research findings. Several authors have pointed out the need for better data management and reporting practices to ensure the reproducibility of research results (Stark, 2018; Waide et al., 2017; Samuel et al., 2021). In Waide et al. (2017), the authors discuss the challenges associated with managing different data sets in ecological research, while Leonelli (2018) further argues that reproducibility should be a criterion for assessing research quality. However, setting standards for assessing reproducibility itself is an important issue for research. Gundersen et al. (2022) in his latest study defined 22 variables that could be categorized according to different degrees of reproducibility. Three categories were mentioned: Reproducible Data, Reproducible Experiment and Reproducible Method. Gundersen et al. (2022) used these variables to evaluate reproducibility support on 13 different open source machine learning platforms. Similarly, Heil et al. (2021) proposed 3 different standards for reproducibility in machine learning that have applications in the life sciences. Bronze, silver or gold standards are defined according to the availability of data, models and code. Tatman et al. (2018) went a step further and analyzed ML papers from ICML and NeurIPS conferences and distinguished levels of reproducibility based on the resources provided with the paper. To achieve the desired result, they recommended some practical steps such as providing code and data in an executable environment where all libraries and dependencies are linked. This allows code or experiments to run smoothly on a new machine. This recommendation is in line with the gold standard provided in (Heil et al., 2021) that authors should produce reproducible results with the execution of a single command. In addition to the general recommendations, Pineau et al. (2021) proposed standard practices and activities to improve reproducibility in the AI community. Some of these are the reproducibility programme at the NeurIPS conference and the ML reproducibility checklist, which provides guidelines to authors before submitting to conferences. Inspired by Pineau et al. (2021), some of the major AI conferences (AAAI and IJCAI) have introduced similar reproducibility checklists. We have considered these guidelines and also previous work in Gundersen et al. (2022) to develop 10 reproducibility variables. However, the work in Gundersen et al. (2022) is designed to assess the reproducibility of different ML platforms, whereas we aim to assess the reproducibility of biodiversity research. Our variables also include the aspect of uncontrolled randomness and the statistical information necessary for reproducibility.

METHODOLOGY

In this section, we will discuss the steps of our work analyzing the reproducibility of research papers in the biodiversity domain.

Identification:

To assess the reproducibility of methods in biodiversity research, we first tried to obtain unbiased and relevant publications. To do this, we adopted keywords from Abdelmageed et al. (2022), which were provided by biodiversity experts and were also used to develop a corpus in the field of biodiversity for large language models (Abdelmageed et al., 2022). These keywords are: "biodivers*" OR "genetic diversity" OR "*omic diversity" OR "phylogenetic diversity" OR "population diversity" OR "species diversity" OR "ecosystem diversity" OR "functional diversity" OR "microbial diversity" OR "soil diversity" AND "Deep Learning".

With this search query, Google Scholar suggested more than 8000 articles in a period from 2015 to 2021. However, we only selected the first 100 results for our analysis, as manual processing of such a large number of articles was not possible in a limited time frame and we believe that the characteristics of these papers are statistically not different from those of the complete set. We acknowledge that the first

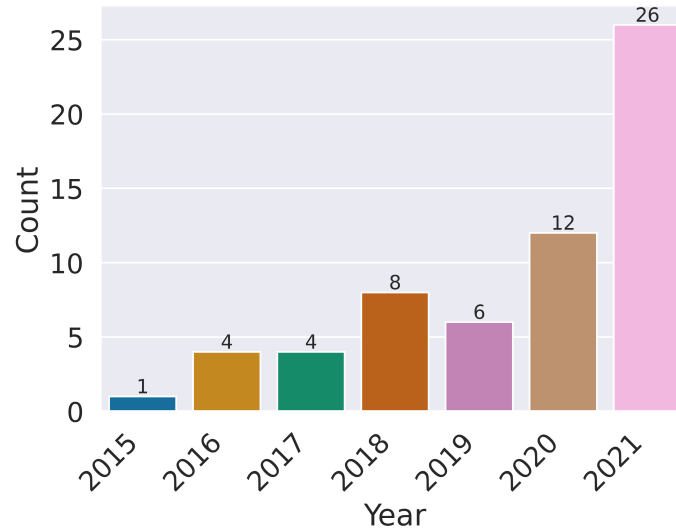


Figure 1. Number of publications considered for collecting the binary responses as per the reproducibility criteria based on year

100 articles selected for our analysis might not represent an entirely unbiased sample of the search results. However, our selection covers a broad range of publishers and publication years, thus capturing a diverse cross-section of the available literature. Additionally, we manually reviewed each selected publication to extract relevant information on ten defined variables, dedicating approximately 40 minutes per paper. This thorough manual curation aimed to enhance the reliability of our dataset. This manually curated dataset will serve as a benchmark for developing an automated system to extract variable information, enabling future analyses to encompass a larger number of articles while maintaining methodological rigor and reproducibility.

Screening:

Before analyzing the individual articles for reproducibility, we found that some articles could not be considered further for the following reasons: 1) duplication, 2) only an abstract was provided or the full article was not accessible, and 3) some articles were not empirical studies and therefore did not include experiments. After considering all these limitations, the number of articles was reduced to 61. Figure 1 shows the distribution of articles over the period from 2015 to 2021, with most articles being from 2020 and 2021. Figure 2 shows the publisher information for the considered articles, which are distributed quite randomly across more than ten publishers.

Selection of reproducibility variables:

There is no standard for assessing the reproducibility of DL biodiversity methods. A common practice for analyzing the reproducibility of research articles is to rerun the experiments using the same methodology as the original author. However, this requires a lot of time and computing resources. Also, sometimes it is not possible to work with the same resources such as hardware and software. Inspired by (Gundersen et al., 2022), we have formulated a set of 10 indicators that can be used as proxies for the probability that a result is reproducible. Instead of repeating each experiment, we look for these variables or factors that are considered important determinants of the reproducibility. Using the literature and the reproducibility checklists at conferences such as NeurIPS and AAAI (Pineau et al., 2021), we can divide these variables into four categories. Resource Information (ReI) details the availability of datasets, source code, open-source tools, and proprietary model details crucial for reproduction efforts. Methodological Information (MI) captures the specifics of software and hardware used, as well as a high-level overview of the deep learning methods employed. Randomness Information (RaI) addresses aspects of unpredictability in computational processes, ensuring they are documented for consistency. Statistical Information (SI) focuses on the rigor of result analysis, advocating for multiple metrics and averaging techniques for reliable

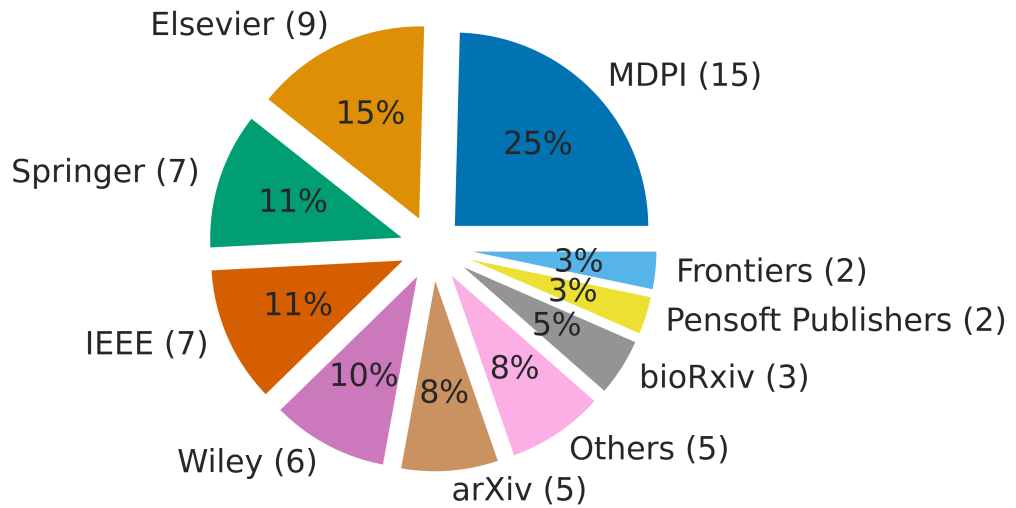


Figure 2. Publisher information for the 61 publications selected for collecting the binary responses as per the reproducibility criteria

evaluation. The comprehensive details that define these aspects of reproducibility are systematically itemized in Table 1

Reproducibility check:

For all selected papers, we first manually checked the variables for each paper. To reduce the degree of subjectivity, we (two of the authors) began by independently recording the binary responses (availability or non-availability). In the initial phase of our assessment, we encountered notable discrepancies due to ambiguous interpretations of the definitions of each variable, leading to inconsistent binary responses. This was quantitatively evidenced by an average Cohen's Kappa value (Cohen, 1960) of 0.54, indicating a moderate level of agreement and highlighting the initial inconsistencies between annotators. To address these issues and improve the reliability of our analysis, we undertook a review and clarification of the variable definitions. In the end, we obtained the same binary responses for each paper, as shown in Table 2. These responses formed the basis for analysing reproducibility. We verified whether the paper covered the functionality of each previously defined variable. After completing the binary responses for each paper, we quantified the reproducibility in 5 levels. The idea behind the different levels is to give independent researchers an insight into the chances of obtaining accurate results independently. The higher the level of reproducibility, the more variables are covered and the greater the chances of reproducing the result. These levels are set considering the time and computational effort required for reproducibility when certain variables are not present. However, the criteria for the basic level (Level 1) must be met to be reproducible according to the definition we refer to Goodman et al. (2016). Figure 3 shows the different levels of reproducibility together with the categories covered by each level. Level 1 should at least cover all variables defined in ReI. Level 2 should cover all variables defined in ReI and MI. Level 3 comprises the variables ReI, MI and SI, while Level 4 comprises the variables ReI, MI and RaI and these two levels do not build on each other. Finally, the highest level of reproducibility combines all categories at Level 5.

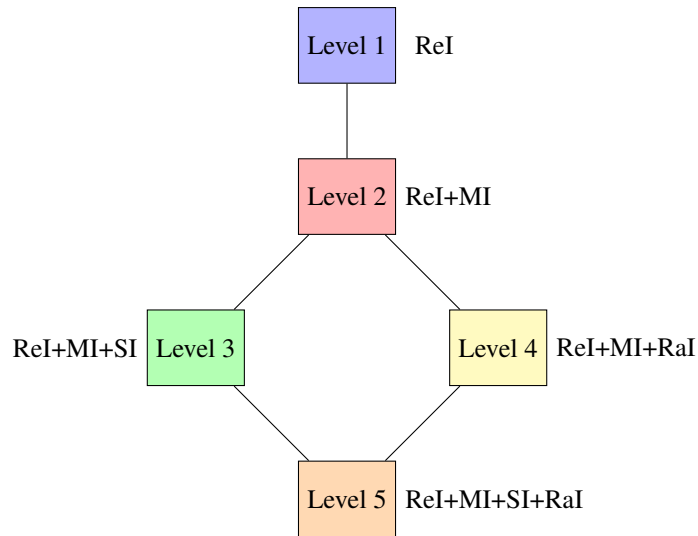


Figure 3. Number of levels along with the categories covered for respective levels

Table 1. Definition of the various reproducibility variables used in this research paper

Category	Variable	Variable short name	Variable Description
Resource Information (ReI)	Dataset	V1	Availability of dataset with persistent identifiers.
	Source Code	V2	Availability of source code to recreate the experiment (e.g. GitHub, GitLab, Zenodo).
	Open source frameworks or environment	V3	Availability of open source software and hardware tools required to reproduce the work (e.g., docker container, virtual environments).
	Model architecture	V4	Availability of a deep learning model architecture or the accessibility of its internal working.
Methodological Information (MI)	Software and Hardware Specification	V5	Availability of information related to the type of hardware used (e.g., GPU), its specifications, and the version of software and libraries used.
	Methods	V6	Availability of high-level information about the deep learning pipeline including the pre-processing and post-processing steps. The other parameters of the pipeline are assessed as separate variables.
	Hyper-parameters	V7	Availability of hyperparameters used to train the model (e.g., number of epochs, learning rate, optimizer, etc.).
Randomness Information (RaI)	Randomness	R	Availability of information about weight initialization, data shuffling, data augmentation, data train-test split and cuDNN GPU library.

Continued on next page

Table 1 – Continued from previous page

Category	Variable	Variable short name	Variable Description
Statistical Information (SI)	Averaging result	S1	Availability of multiple model training and averaging them out rather than selecting only the highest value.
	Evaluation metrics	S2	Availability of more than one evaluated metric (e.g., R^2 score along with Root Mean Square Error (RMSE) or Mean Absolute Error (MAE) or Loss of the model).

RESULTS

As discussed in Section , we divided the ten reproducibility variables we had defined into four categories 1) Resource Information 2) Methodological Information 3) Randomness Information, and 4) Statistical Information. These four categories were further categorized into five different levels of reproducibility, allowing us to assess the reproducibility level of the publications more comprehensively.

Variable level information

Figure 4 contains four plots, one for each category of information. Within each category, specific reproducibility variables are depicted with binary counts, illustrated by two overlapping bars: one denoted by 'Yes' in a brick orange color and the other by 'No' in blue, along with the corresponding number of publications. In total, we made 610 individual judgements for the ten variables, 353 of these were positive i.e the respective information was provided, while this was not the case in 257 cases resulting in negative judgements. Detailed descriptions of all responses can be found at both the category and individual variable levels.

In the Resources Information category, 29 publications provided a dataset, 16 indicated the code repositories, 58 used open-source frameworks or environments, and 57 mentioned the model architectures (Figure 4a).

For the Methodological Information category, around a quarter of the publications included details about the hardware and the software (libraries) they used, all the publications explained methods that were used to build a machine learning pipeline and 40 publications provided the basic hyperparameters (Figure 4b).

The number of publications that used a random seed in all possible ways (weight initialization, data shuffling, data augmentation, data train-test split and cuDNN GPU library) in their code is 3 (Figure 4c). Regarding statistical considerations, 21 publications provided the average result of multiple model trainings, 54 publications evaluated their models with more than one evaluation metric. (Figure 4d).

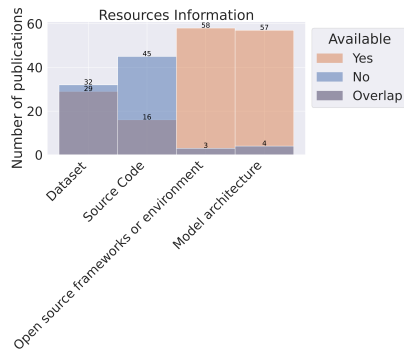
Categorical level information

We have harmonized the binary responses of reproducibility variables for each categorical level in such a way that all the individual reproducibility variables in a category must be available (Yes) to mark that the specific categorical information is available (Yes). If one of the individual reproducibility variables is unavailable (No), the specific categorical information is unavailable (No). Mathematically, $V1 \& V2 \& V3 \& V4 == Yes$ is used to denote the availability of resource information (Yes).

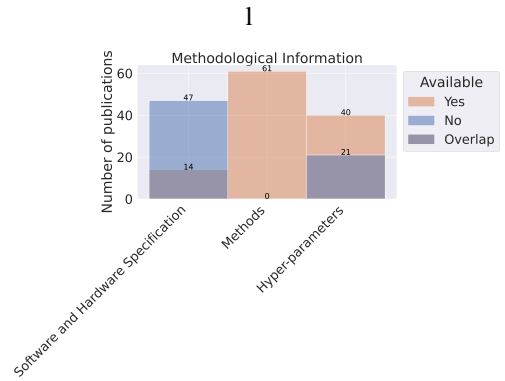
Figure 5 depicts the distribution of publications that satisfy the defined categories. The number of publications that satisfy the Resources, Methodological, Randomness, and Statistical information are 10, 13, 3, and 20, respectively.

Reproducibility levels of publications

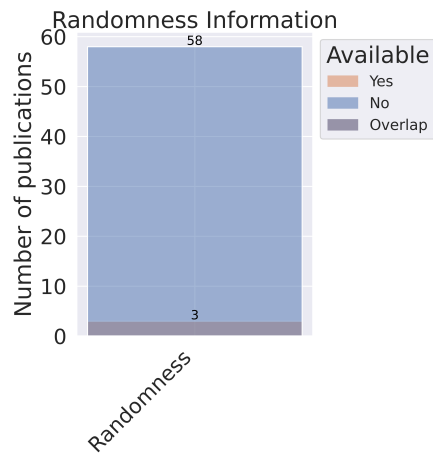
Per the definitions of the reproducibility levels described in Section , Levels 1 and 5 are the lowest and highest, respectively. According to Figure 6, only one publication fulfils the highest reproducibility level, while ten publications meet the lowest.



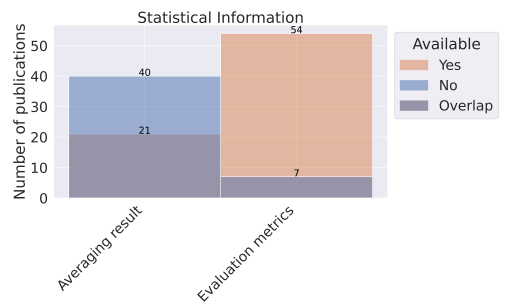
(a) Binary responses of variables for the category 'Resources'



(b) Binary responses of variables for the category 'Methodological information'



(c) Binary responses of variables for the category 'Randomness'



(d) Binary responses of variables for the category 'Statistical consideration'

Figure 4. Binary responses of the considered reproducible variables in four categories 1) Resources 2) Methodological information 3) Randomness 4) Statistical consideration for selected research publication

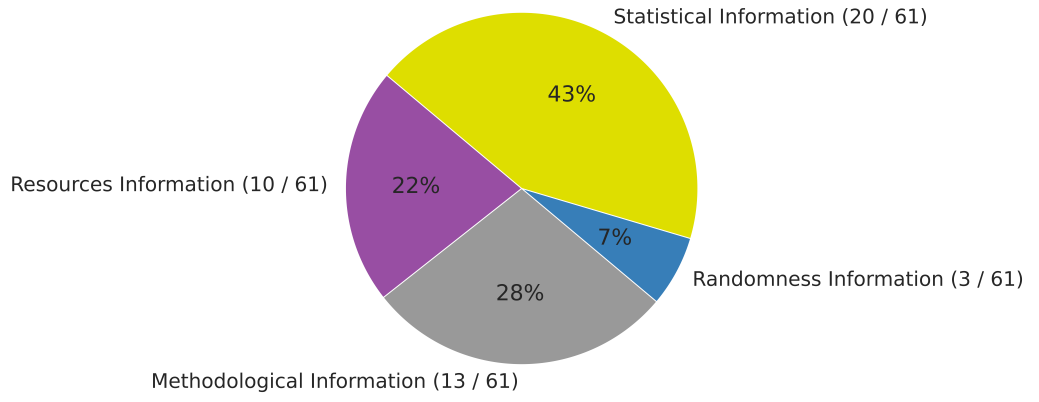


Figure 5. Distribution of selected research publications along four categories

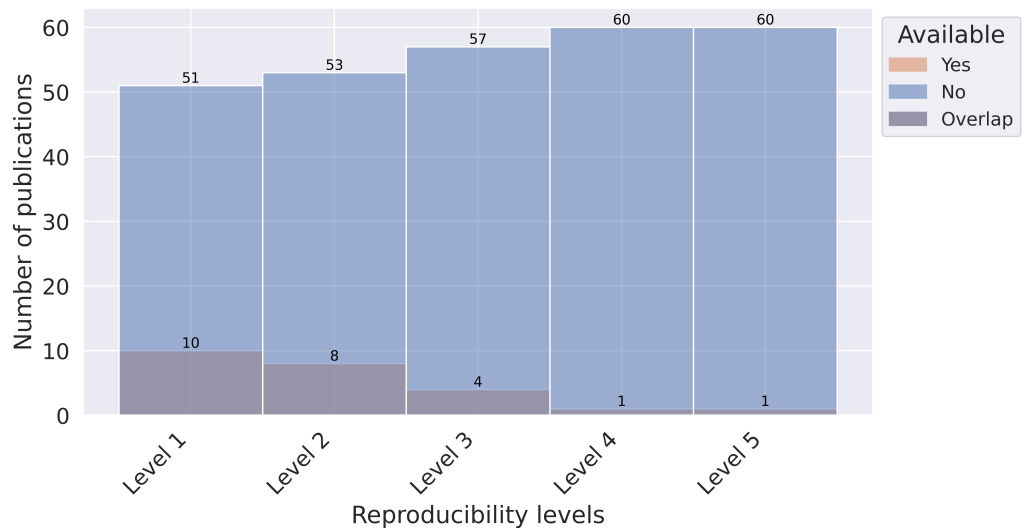
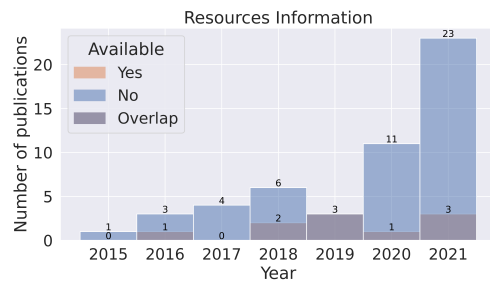
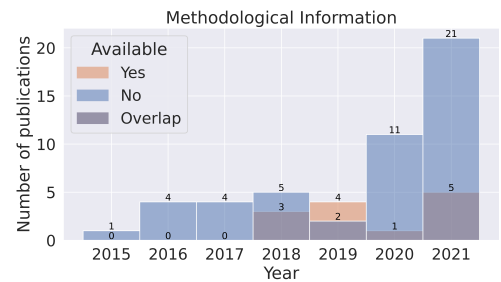


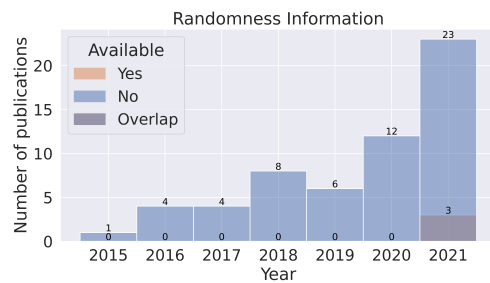
Figure 6. Bar plot indicating the number of publications satisfying the different levels of reproducibility



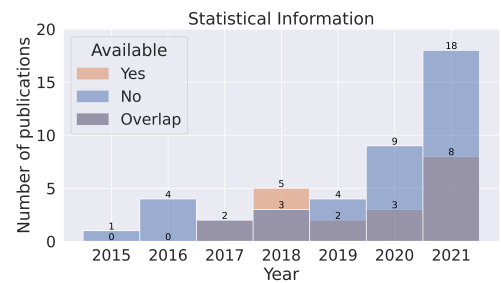
(a) Number of papers meeting criteria for the category 'Resources' according to year



(b) Number of papers meeting criteria for the category 'Methodological information' according to year



(c) Number of papers meeting criteria for the category 'Randomness' according to year



(d) Number of papers meeting criteria for the category 'Statistical consideration' according to year

Figure 7. Number of papers meeting criteria for the four categories 1) Resources 2) Methodological information 3) Randomness 4) Statistical consideration for selected research publications by year

Reproducibility status of publications by year

The number of publications that meet the defined reproducibility criteria follows a linear trend with respect to year (Figure 7).

DISCUSSION

In the Biodiversity field, deep-learning methods are becoming part of many studies that run large-scale experiments. These gave us the opportunity to orchestrate and extract the binary responses of 10 reproducibility variables from 61 publications (Table 1 and Table 2). As a result, we recorded 610 total responses: 353 were positive, and 257 were negative. All the positive responses were dominated by four variables (230 responses): 1) Open source frameworks or environment, 2) Model architecture, 3) Methods, and 4) Evaluation metrics and the negative responses were dominated by only one variable: Randomness.

Most of the publications that employed deep learning models use open-source frameworks or environments like Tensorflow and PyTorch with the programming language Python/R, and they also provide model architectures either as a figure/table in the publication or described in the text with respective citations. Some of the publications used licensed programming languages like Matlab, comparatively, it is negligible. We looked for high-level information on the whole deep-learning pipeline in methods, all the publications provided compact information.

Most of the publications use more than one metric to evaluate their models, for example, when it comes to regression tasks, they use R^2 score along with Root Mean Square Error (RMSE) or Mean Absolute Error (MAE) or Loss of the model, etc. The information from the publications is in compliance with reproducible workflow guidelines for the positively dominated variables.

We found datasets only in 29 publications; in the other 32 publications, there was no tangential information about the dataset. There were some cases where the authors linked to some data-providing websites to find the data that was used in the publication, but those websites will change over time and finding the exact data that was used in the publication is not possible without providing persistent identifiers of the respective data points.

Source code is one of the fundamental variables of reproducibility. However, a little more than 25 % of the publications provided their source code.

In 14 publications, authors have provided specifications about both software and hardware. Without specific information about the hardware and software, the reproducibility results will change because the random generators work differently with different hardware and software changes with each version.

Hyperparameters are the values that are chosen to control the model learning process. This means with each set of hyperparameters, the model will provide a varied range of results. However, information about basic hyperparameters (epochs, learning rate, optimizer and loss function) was missing from 21 publications. Reproducing the results of these 21 publications is unfeasible due to the missing essential hyperparameter information.

Averaging results is also an important aspect of reproducibility, after each training process, results will change slightly because of the certain random initializations through the deep learning pipeline. However, in our study, 40 publications didn't report the multiple training results.

We opted for manually updating the variables by going through each publication and extracting the required information as a binary response because we did not find a system or technique that could automatically extract the required variable information from a publication. Since we are extracting the information about variables from publications manually, it is only possible to work with a small dataset, which is also the limitation of this study.

Due to recent developments with Large Language Models (LLMs), we are considering extracting the reproducible variables information from publications using LLMs from the year 2022 onwards (Ahmed et al., 2023; Kommineni et al., 2024). This will allow us to implement our analysis on large-scale publications.

CONCLUSION

In this paper, we presented our pipeline for assessing the reproducibility of deep learning methods in biodiversity research. Inspired by the current state of the art, we established a comprehensive set of ten variables, categorized into four distinct groups, to effectively quantify the reproducibility of DL empirical research. Based on the defined categories, we documented the availability of each variable across 61

selected publications over the period from 2015 to 2021. From the total 610 responses for the 10 variables, 57.9% show the availability of the variables in the publication, while the remaining 42.1% are primarily characterized by the absence of randomness-related information. The highest and lowest reproducibility levels are satisfied by only one and ten publications, respectively. Given the use of deep learning to advance biodiversity research, improving reproducibility of the DL methods is crucial. Considering the limitations of the manual approach and the relatively small dataset analyzed until 2021, our future endeavors will focus on implementing a semi-automatic approach that leverages Large Language models for extracting information on reproducible variables from publications.

DATA AVAILABILITY STATEMENT

The data and the code used to extract and analyse the reproducibility information of Deep Learning methods from publications in the Biodiversity domain is publicly available: <https://github.com/fusion-jena/Reproduce-DLmethods-Biodiv>

Appendices

Table 2. Binary responses of different reproducibility variables. 'y' denotes the presence of variable information, while 'n' signifies the absence of variable information.

Publication	V1	V2	V3	V4	V5	V6	V7	R	S1	S2
Klein et al. (2015)	n	n	y	n	n	y	n	n	n	n
Khalighifar et al. (2021)	n	n	y	y	n	y	n	n	n	y
Choe et al. (2021)	n	n	y	y	y	y	y	n	y	y
Mahmood et al. (2016)	y	n	y	y	n	y	n	n	n	y
Younis et al. (2020)	y	y	y	y	y	y	y	n	y	y
Schwartz and Alfaro (2021)	y	y	y	y	y	y	y	n	y	y
Potamitis (2016)	n	n	y	y	n	y	n	n	n	n
Chen et al. (2020)	n	n	y	y	n	y	n	n	n	y
Fujisawa et al. (2021)	n	n	y	y	n	y	n	n	y	y
Weinstein (2018)	n	y	y	y	y	y	y	n	n	y
Chalmers et al. (2021)	n	n	y	y	y	y	n	n	n	y
Guirado et al. (2019)	y	n	y	y	y	y	y	n	y	y
Zualkernan et al. (2020)	n	n	y	y	n	y	y	n	n	y
Villon et al. (2018)	n	n	y	y	n	y	y	n	y	y
Fairbrass et al. (2019)	y	y	y	y	y	y	y	n	n	y
Weinstein et al. (2019)	y	y	y	y	y	y	y	n	n	y
Hu et al. (2020)	n	n	y	y	n	y	n	n	n	y
Alshahrani et al. (2021)	y	n	y	y	n	y	y	n	y	y
Villon et al. (2020)	n	n	y	y	n	y	n	n	y	y
Botella et al. (2018)	y	n	y	y	n	y	y	n	y	y
Salamon et al. (2017)	y	n	y	y	n	y	y	n	y	y
Mac Aodha et al. (2018)	y	y	y	y	y	y	y	n	y	y

Continued on next page

Table 2 – Continued from previous page

Publication	V1	V2	V3	V4	V5	V6	V7	R	S1	S2
Schindler and Steinhage (2021)	n	n	y	y	n	y	y	n	n	y
Rakshit et al. (2018)	n	n	y	y	n	y	y	n	n	y
Bjerge et al. (2021)	n	y	y	y	n	y	y	y	y	y
Guirado et al. (2017)	n	n	y	y	n	y	n	n	n	y
Hussein et al. (2021a)	n	n	y	y	n	y	y	n	n	y
Rammer and Seidl (2019)	y	y	y	y	y	y	y	n	n	y
Anand et al. (2021)	n	n	n	y	n	y	n	n	n	y
Zizka et al. (2021)	y	y	y	y	n	y	y	y	n	n
Miele et al. (2020)	n	n	y	y	n	y	n	n	y	n
Demertzis et al. (2018)	n	n	y	y	n	y	n	n	y	y
Malerba et al. (2021)	y	n	y	y	n	y	y	n	n	n
Huang and Basanta (2021)	n	n	y	y	n	y	y	n	y	y
Campos-Taberner et al. (2020)	y	n	n	y	n	y	n	n	n	y
Rousset et al. (2021)	y	n	y	y	n	y	y	n	n	y
López-Jiménez et al. (2019)	y	n	y	y	n	y	y	n	n	y
Schuettpelz et al. (2017)	n	y	y	y	n	y	y	n	n	y
Schiller et al. (2021)	y	y	y	y	y	y	y	y	y	y
Heredia (2017)	n	y	y	y	n	y	y	n	y	y
Martins et al. (2021)	y	n	y	y	y	y	y	n	n	y
Browning et al. (2018)	y	n	y	y	n	y	n	n	y	y
Guirado et al. (2020)	y	n	y	y	n	y	y	n	n	y
Ayhan et al. (2020)	n	n	y	y	n	y	y	n	n	y
Loddo et al. (2021)	n	n	y	y	n	y	n	n	y	y
Jamil et al. (2021)	y	n	y	y	n	y	y	n	n	y

Continued on next page

Table 2 – Continued from previous page

Publication	V1	V2	V3	V4	V5	V6	V7	R	S1	S2
Neves et al. (2021)	n	n	y	y	n	y	n	n	n	y
Jin et al. (2021)	y	n	y	n	n	y	n	n	n	y
Mohanty et al. (2016)	y	y	y	y	n	y	y	n	n	y
Xie et al. (2019)	y	n	y	y	n	y	y	n	y	y
Tian et al. (2020)	n	n	n	n	n	y	n	n	y	y
Gimenez et al. (2021)	n	y	y	y	n	y	y	n	n	y
Ortega Adarme et al. (2020)	y	n	y	y	n	y	y	n	n	y
Boer and Vos (2018)	y	y	y	y	y	y	y	n	n	n
Dunker et al. (2021)	n	n	y	y	n	y	n	n	n	y
Villon et al. (2016)	n	n	y	y	n	y	n	n	n	y
Dyrmann et al. (2021)	n	n	y	y	n	y	y	n	n	y
Arruda et al. (2021)	y	n	y	n	n	y	y	n	n	n
Kislov and Korznikov (2020)	y	n	y	y	n	y	y	n	n	y
Hussein et al. (2021b)	y	n	y	y	n	y	y	n	n	y
Becker et al. (2021)	n	y	y	y	y	y	y	n	n	y

REFERENCES

- Abdelmageed, N., Löffler, F., Feddoul, L., Algergawy, A., Samuel, S., Gaikwad, J., Kazem, A., and König-Ries, B. (2022). Biodivnere: Gold standard corpora for named entity recognition and relation extraction in the biodiversity domain. *Biodiversity Data Journal*, 10.
- Ahmed, W., Kommineni, V. K., Koenig-ries, B., and Samuel, S. (2023). How reproducible are the results gained with the help of deep learning methods in biodiversity research? *Biodiversity Information Science and Standards*, 7.
- Alshahrani, H. M., Al-Wesabi, F. N., Al Duhayyim, M., Nemri, N., Kadry, S., and Alqaralleh, B. A. Y. (2021). An automated deep learning based satellite imagery analysis for ecology management. *Ecol. Inform.*, 66(101452):101452.
- Anand, A., Pandey, M. K., Srivastava, P. K., Gupta, A., and Khan, M. L. (2021). Integrating multi-sensors data for species distribution mapping using deep learning and envelope models. *Remote Sens. (Basel)*, 13(16):3284.
- Arruda, V. L. S., Piontekowski, V. J., Alencar, A., Pereira, R. S., and Matricardi, E. A. T. (2021). An alternative approach for mapping burn scars using landsat imagery, google earth engine, and deep learning in the brazilian savanna. *Remote Sens. Appl. Soc. Environ.*, 22(100472):100472.
- August, T. A., Pescott, O. L., Joly, A., and Bonnet, P. (2020). AI naturalists might hold the key to unlocking biodiversity data in social media imagery. *Patterns*, 1(7):100116.

- Ayhan, B., Kwan, C., Budavari, B., Kwan, L., Lu, Y., Perez, D., Li, J., Skarlatos, D., and Vlachos, M. (2020). Vegetation detection using deep learning and conventional methods. *Remote Sens. (Basel)*, 12(15):2502.
- Becker, A., Russo, S., Puliti, S., Lang, N., Schindler, K., and Wegner, J. D. (2021). Country-wide retrieval of forest structure from optical and SAR satellite imagery with deep ensembles.
- Bjerge, K., Nielsen, J. B., Sepstrup, M. V., Helsing-Nielsen, F., and Høye, T. T. (2021). An automated light trap to monitor moths (lepidoptera) using computer vision-based tracking and deep learning. *Sensors (Basel)*, 21(2):343.
- Boer, M. J. and Vos, R. A. (2018). Taxonomic classification of ants (formicidae) from images using deep learning. *bioRxiv*, page 407452.
- Botella, C., Joly, A., Bonnet, P., Monestiez, P., and Munoz, F. (2018). A deep learning approach to species distribution modelling. In *Multimedia Tools and Applications for Environmental & Biodiversity Informatics*, pages 169–199. Springer International Publishing, Cham.
- Browning, E., Bolton, M., Owen, E., Shoji, A., Guilford, T., and Freeman, R. (2018). Predicting animal behaviour using deep learning: GPS data alone accurately predict diving in seabirds. *Methods Ecol. Evol.*, 9(3):681–692.
- Campos-Taberner, M., García-Haro, F. J., Martínez, B., Izquierdo-Verdiguier, E., Atzberger, C., Camps-Valls, G., and Gilabert, M. A. (2020). Understanding deep learning in land use classification based on sentinel-2 time series. *Sci. Rep.*, 10(1):17188.
- Chalmers, C., Fergus, P., Wich, S., and Longmore, S. N. (2021). Modelling animal biodiversity using acoustic monitoring and deep learning. In *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE.
- Chen, X., Zhao, J., Chen, Y.-H., Zhou, W., and Hughes, A. C. (2020). Automatic standardized processing and identification of tropical bat calls using deep learning approaches. *Biol. Conserv.*, 241(108269):108269.
- Choe, H., Chi, J., and Thorne, J. H. (2021). Mapping potential plant species richness over large areas with deep learning, MODIS, and species distribution models. *Remote Sens. (Basel)*, 13(13):2490.
- Christin, S., Hervet, É., and Lecomte, N. (2019). Applications for deep learning in ecology. *Methods in Ecology and Evolution*, 10(10):1632–1644.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
- Demertzis, K., Iliadis, L. S., and Anezakis, V.-D. (2018). Extreme deep learning in biosecurity: the case of machine hearing for marine species identification. *J. Inf. Telecommun.*, 2(4):492–510.
- Dunker, S., Motivans, E., Rakosy, D., Boho, D., Mäder, P., Hornick, T., and Knight, T. M. (2021). Pollen analysis using multispectral imaging flow cytometry and deep learning. *New Phytol.*, 229(1):593–606.
- Dyrmann, M., Mortensen, A. K., Linneberg, L., Høye, T. T., and Bjerge, K. (2021). Camera assisted roadside monitoring for invasive alien plant species using deep learning. *Sensors (Basel)*, 21(18):6126.
- El-Amir, H. and Hamdy, M. (2020). Deep learning pipeline. *Apress: Berkeley, CA, USA*.
- Fairbrass, A. J., Firman, M., Williams, C., Brostow, G. J., Titheridge, H., and Jones, K. E. (2019). CityNet—Deep learning tools for urban ecoacoustic assessment. *Methods Ecol. Evol.*, 10(2):186–197.
- Feng, X., Park, D. S., Walker, C., Peterson, A. T., Merow, C., and Papeş, M. (2019). A checklist for maximizing reproducibility of ecological niche models. *Nature Ecology & Evolution*, 3(10):1382–1395.
- Fujisawa, T., Nogueras, V., Meramveliotakis, E., Papadopoulou, A., and Vogler, A. P. (2021). Image-based taxonomic classification of bulk biodiversity samples using deep learning and domain adaptation.
- Gimenez, O., Kervellec, M., Fanjul, J.-B., Chaine, A., Marescot, L., Bollet, Y., and Duchamp, C. (2021). Trade-off between deep learning for species identification and inference about predator-prey co-occurrence: Reproducible R workflow integrating models in computer vision and ecological statistics.
- Goodman, S. N., Fanelli, D., and Ioannidis, J. P. A. (2016). What does research reproducibility mean? *Science Translational Medicine*, 8(341):341ps12–341ps12.
- GPAI (2022). Biodiversity and artificial intelligence, opportunities and recommendations report.
- Guirado, E., Alcaraz-Segura, D., Cabello, J., Puertas-Ruiz, S., Herrera, F., and Tabik, S. (2020). Tree cover estimation in global drylands from space using deep learning. *Remote Sens. (Basel)*, 12(3):343.
- Guirado, E., Tabik, S., Alcaraz-Segura, D., Cabello, J., and Herrera, F. (2017). Deep-Learning convolutional neural networks for scattered shrub detection with google earth imagery.

- Guirado, E., Tabik, S., Rivas, M. L., Alcaraz-Segura, D., and Herrera, F. (2019). Whale counting in satellite and aerial images with deep learning. *Sci. Rep.*, 9(1):14259.
- Gundersen, O. E., Shamsaliei, S., and Isdahl, R. J. (2022). Do machine learning platforms provide out-of-the-box reproducibility? *Future Generation Computer Systems*, 126:34–47.
- Heil, B. J., Hoffman, M. M., Markowitz, F., Lee, S.-I., Greene, C. S., and Hicks, S. C. (2021). Reproducibility standards for machine learning in the life sciences. *Nature Methods*, 18(10):1132–1135.
- Heredia, I. (2017). Large-scale plant classification with deep neural networks. In *Proceedings of the Computing Frontiers Conference*, New York, NY, USA. ACM.
- Hu, J., Huang, W., Su, Y., Liu, Y., and Xiao, P. (2020). BatNet++: A robust deep learning-based predicting models for calls recognition. In *2020 5th International Conference on Smart Grid and Electrical Automation (ICSGEA)*. IEEE.
- Huang, Y.-P. and Basanta, H. (2021). Recognition of endemic bird species using deep learning models. *IEEE Access*, 9:102975–102984.
- Hussein, B. R., Malik, O. A., Ong, W.-H., and Slik, J. W. F. (2021a). Automated extraction of phenotypic leaf traits of individual intact herbarium leaves from herbarium specimen images using deep learning based semantic segmentation. *Sensors (Basel)*, 21(13):4549.
- Hussein, B. R., Malik, O. A., Ong, W.-H., and Slik, J. W. F. (2021b). Reconstruction of damaged herbarium leaves using deep learning techniques for improving classification accuracy. *Ecol. Inform.*, 61(101243):101243.
- Jamil, S., Rahman, M., and Haider, A. (2021). Bag of features (BoF) based deep learning framework for bleached corals detection. *Big Data Cogn. Comput.*, 5(4):53.
- Jin, L., Yu, J., Yuan, X., and Du, X. (2021). Fish classification using DNA barcode sequences through deep learning method. *Symmetry (Basel)*, 13(9):1599.
- Khalighifar, A., Brown, R. M., Goyes Vallejos, J., and Peterson, A. T. (2021). Deep learning improves acoustic biodiversity monitoring and new candidate forest frog species identification (genus *platymantis*) in the Philippines. *Biodivers. Conserv.*, 30(3):643–657.
- Kislov, D. E. and Korznikov, K. A. (2020). Automatic windthrow detection using very-high-resolution satellite imagery and deep learning. *Remote Sens. (Basel)*, 12(7):1145.
- Klein, D. J., McKown, M. W., and Tershy, B. R. (2015). Deep learning for large scale biodiversity monitoring. In *Bloomberg Data for Good Exchange Conference*.
- Kommineni, V. K., König-Ries, B., and Samuel, S. (2024). From human experts to machines: An llm supported approach to ontology and knowledge graph construction. *arXiv preprint arXiv:2403.08345*.
- Leonelli, S. (2018). Rethinking reproducibility as a criterion for research quality. In *Including a symposium on Mary Morgan: curiosity, imagination, and surprise*, volume 36, pages 129–146. Emerald Publishing Limited.
- Loddo, A., Loddo, M., and Di Ruberto, C. (2021). A novel deep learning based approach for seed image classification and retrieval. *Comput. Electron. Agric.*, 187(106269):106269.
- López-Jiménez, E., Vasquez-Gomez, J. I., Sanchez-Acevedo, M. A., Herrera-Lozada, J. C., and Uriarte-Arcia, A. V. (2019). Columnar cactus recognition in aerial images using a deep learning approach. *Ecol. Inform.*, 52:131–138.
- Mac Aodha, O., Gibb, R., Barlow, K. E., Browning, E., Firman, M., Freeman, R., Harder, B., Kinsey, L., Mead, G. R., Newson, S. E., Pandourski, I., Parsons, S., Russ, J., Szodoray-Paradi, A., Szodoray-Paradi, F., Tilova, E., Girolami, M., Brostow, G., and Jones, K. E. (2018). Bat detective—deep learning tools for bat acoustic signal detection. *PLoS Comput. Biol.*, 14(3):e1005995.
- Mahmood, A., Bennamoun, M., An, S., Sohel, F., Boussaid, F., Hovey, R., Kendrick, G., and Fisher, R. B. (2016). Automatic annotation of coral reefs using deep learning. In *OCEANS 2016 MTS/IEEE Monterey*. IEEE.
- Malerba, M. E., Wright, N., and Macreadie, P. I. (2021). A continental-scale assessment of density, size, distribution and historical trends of farm dams using deep learning convolutional neural networks. *Remote Sens. (Basel)*, 13(2):319.
- Martins, J. A. C., Nogueira, K., Osco, L. P., Gomes, F. D. G., Furuya, D. E. G., Gonçalves, W. N., Sant’Ana, D. A., Ramos, A. P. M., Liesenberg, V., dos Santos, J. A., de Oliveira, P. T. S., and Junior, J. M. (2021). Semantic segmentation of tree-canopy in urban environment with pixel-wise deep learning. *Remote Sens. (Basel)*, 13(16):3054.
- Miele, V., Dussert, G., Cucchi, T., and Renaud, S. (2020). Deep learning for species identification of

modern and fossil rodent molars.

- Mohanty, S. P., Hughes, D. P., and Salathé, M. (2016). Using deep learning for image-based plant disease detection. *Front. Plant Sci.*, 7:1419.
- Neves, A. K., Körting, T. S., Fonseca, L. M. G., Soares, A. R., Girolamo-Neto, C. D., and Heipke, C. (2021). Hierarchical mapping of brazilian savanna (cerrado) physiognomies based on deep learning. *J. Appl. Remote Sens.*, 15(04).
- Norouzzadeh, M. S., Nguyen, A., Kosmala, M., Swanson, A., Palmer, M. S., Packer, C., and Clune, J. (2018). Automatically identifying, counting, and describing wild animals in camera-trap images with deep learning. *Proceedings of the National Academy of Sciences*, 115(25):E5716–E5725.
- Ortega Adarme, M., Queiroz Feitosa, R., Nigri Happ, P., Aparecido De Almeida, C., and Rodrigues Gomes, A. (2020). Evaluation of deep learning techniques for deforestation detection in the brazilian amazon and cerrado biomes from remote sensing imagery. *Remote Sens. (Basel)*, 12(6):910.
- Pineau, J., Vincent-Lamarre, P., Sinha, K., Larivière, V., Beygelzimer, A., d'Alché Buc, F., Fox, E., and Larochelle, H. (2021). Improving reproducibility in machine learning research (a report from the neurips 2019 reproducibility program). *The Journal of Machine Learning Research*, 22(1):7459–7478.
- Potamitis, I. (2016). Deep learning for detection of bird vocalisations. *arXiv preprint arXiv:1609.08408*.
- Raff, E. (2019). A step toward quantifying independently reproducible machine learning research. *Advances in Neural Information Processing Systems*, 32.
- Rakshit, S., Debnath, S., and Mondal, D. (2018). Identifying land patterns from satellite imagery in amazon rainforest using deep learning.
- Rammer, W. and Seidl, R. (2019). Harnessing deep learning in ecology: An example predicting bark beetle outbreaks. *Front. Plant Sci.*, 10:1327.
- Rousset, G., Despinoy, M., Schindler, K., and Mangeas, M. (2021). Assessment of deep learning techniques for land use land cover classification in southern new caledonia. *Remote Sens. (Basel)*, 13(12):2257.
- Rovero, F., Zimmermann, F., Berzi, D., and Meek, P. (2013). "which camera trap type and how many do i need?" a review of camera features and study designs for a range of wildlife research applications. *Hystrix*.
- Salamon, J., Bello, J. P., Farnsworth, A., and Kelling, S. (2017). Fusing shallow and deep learning for bioacoustic bird species classification. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE.
- Samuel, S. and König-Ries, B. (2021). Understanding experiments and research practices for reproducibility: an exploratory study. *PeerJ*, 9:e11140.
- Samuel, S., Löffler, F., and König-Ries, B. (2021). Machine learning pipelines: Provenance, reproducibility and FAIR data principles. In Glavic, B., Braganholo, V., and Koop, D., editors, *Provenance and Annotation of Data and Processes - 8th and 9th International Provenance and Annotation Workshop, IPAW 2020 + IPAW 2021, Virtual Event, July 19-22, 2021, Proceedings*, volume 12839 of *Lecture Notes in Computer Science*, pages 226–230. Springer.
- Schiller, C., Schmidlein, S., Boonman, C., Moreno-Martínez, A., and Kattenborn, T. (2021). Deep learning and citizen science enable automated plant trait predictions from photographs. *Sci. Rep.*, 11(1):16395.
- Schindler, F. and Steinhage, V. (2021). Identification of animals and recognition of their actions in wildlife videos using deep learning techniques. *Ecol. Inform.*, 61(101215):101215.
- Schnitzer, S. A. and Carson, W. P. (2016). Would ecology fail the repeatability test? *BioScience*, 66(2):98–99.
- Schuettpelz, E., Frandsen, P., Dikow, R., Brown, A., Orli, S., Peters, M., Metallo, A., Funk, V., and Dorr, L. (2017). Applications of deep convolutional neural networks to digitized natural history collections. *Biodivers. Data J.*, 5:e21139.
- Schwartz, S. T. and Alfaro, M. E. (2021). *Sashimi*: A toolkit for facilitating high-throughput organismal image segmentation using deep learning. *Methods Ecol. Evol.*, 12(12):2341–2354.
- Stark, P. B. (2018). Before reproducibility must come preproducibility. *Nature*, 557(7706):613–614.
- Tabak, M. A., Norouzzadeh, M. S., Wolfson, D. W., Sweeney, S. J., VerCauteren, K. C., Snow, N. P., Halseth, J. M., Di Salvo, P. A., Lewis, J. S., White, M. D., et al. (2019). Machine learning to classify animal species in camera trap images: Applications in ecology. *Methods in Ecology and Evolution*, 10(4):585–590.

- Tatman, R., VanderPlas, J., and Dane, S. (2018). A practical taxonomy of reproducibility for machine learning research.
- Tian, J., Wang, L., Yin, D., Li, X., Diao, C., Gong, H., Shi, C., Menenti, M., Ge, Y., Nie, S., Ou, Y., Song, X., and Liu, X. (2020). Development of spectral-phenological features for deep learning to understand spartina alterniflora invasion. *Remote Sens. Environ.*, 242(111745):111745.
- Villon, S., Chaumont, M., Subsol, G., Villéger, S., Claverie, T., and Mouillot, D. (2016). Coral reef fish detection and recognition in underwater videos by supervised machine learning: Comparison between deep learning and HOG+SVM methods. In *Advanced Concepts for Intelligent Vision Systems*, Lecture notes in computer science, pages 160–171. Springer International Publishing, Cham.
- Villon, S., Mouillot, D., Chaumont, M., Darling, E. S., Subsol, G., Claverie, T., and Villéger, S. (2018). A deep learning method for accurate and fast identification of coral reef fishes in underwater images. *Ecol. Inform.*, 48:238–244.
- Villon, S., Mouillot, D., Chaumont, M., Subsol, G., Claverie, T., and Villéger, S. (2020). A new method to control error rates in automated species identification with deep learning algorithms. *Sci. Rep.*, 10(1):10972.
- Waide, R. B., Brunt, J. W., and Servilla, M. S. (2017). Demystifying the landscape of ecological data repositories in the United States. *BioScience*, 67(12):1044–1051.
- Weinstein, B. G. (2018). Scene-specific convolutional neural networks for video-based biodiversity detection. *Methods Ecol. Evol.*, 9(6):1435–1441.
- Weinstein, B. G., Marconi, S., Bohlman, S., Zare, A., and White, E. (2019). Individual tree-crown detection in RGB imagery using semi-supervised deep learning neural networks. *Remote Sens. (Basel)*, 11(11):1309.
- Xie, J., Hu, K., Zhu, M., Yu, J., and Zhu, Q. (2019). Investigation of different CNN-based models for improved bird sound classification. *IEEE Access*, 7:175353–175361.
- Younis, S., Schmidt, M., Weiland, C., Dressler, S., Seeger, B., and Hickler, T. (2020). Detection and annotation of plant organs from digitised herbarium scans using deep learning. *Biodivers. Data J.*, 8:e57090.
- Zizka, A., Silvestro, D., Vitt, P., and Knight, T. M. (2021). Automated conservation assessment of the orchid family with deep learning. *Conserv. Biol.*, 35(3):897–908.
- Zualkernan, I. A., Dhou, S., Judas, J., Sajun, A. R., Gomez, B. R., Hussain, L. A., and Sakhini, D. (2020). Towards an IoT-based deep learning architecture for camera trap image classification. In *2020 IEEE Global Conference on Artificial Intelligence and Internet of Things (GCAIoT)*. IEEE.