

gibbonNetR: an R Package for the Use of Convolutional Neural Networks for Automated Detection of Acoustic Data

Dena Jane Clink¹ and Abdul Hamid Ahmad²

1 K. Lisa Yang Center for Conservation Bioacoustics, Cornell Lab of Ornithology, Cornell University, Ithaca, New York, United States 2 Institute for Tropical Biology and Conservation, Universiti Malaysia Sabah (UMS), Kota Kinabalu, Sabah, Malaysia

DOI: 10.21105/joss.07250

Software

- Review I
- Repository 🖒
- Archive 🗗

Editor: Fabian-Robert Stöter ♂ Reviewers:

@Desjongu

Østeffilazerte

Submitted: 28 May 2024 Published: 07 June 2025

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License (CC BY 4.0).

Summary

Automated detection of acoustic signals is crucial for effectively monitoring vocal animals and their habitats across large spatial and temporal scales. Recent advances in deep learning have made high-performance automated detection approaches accessible to more practitioners. However, there are few deep learning approaches that can be implemented natively in R. The 'torch for R' ecosystem has made the use of convolutional neural networks (CNNs) accessible for R users. Here, we provide an R package and workflow to use CNNs for automated detection and classification of acoustics signals from passive acoustic monitoring data. We provide examples using data collected in Sabah, Malaysia. The package provides functions to create spectrogram images from labeled data, compare the performance of different CNN architectures, deploy trained models over directories of sound files, and extract embeddings from trained models. The R programming language remains one of the most commonly used languages among ecologists, and we hope that this package makes deep learning approaches more accessible to this audience. In addition, these models can serve as important benchmarks for future automated detection work.

Statement of need

Passive acoustic monitoring

We are in a biodiversity crisis, and there is a great need for the ability to rapidly assess biodiversity in order to understand and mitigate anthropogenic impacts. One approach that can be especially effective for monitoring of sound-producing yet cryptic animals is the use of passive acoustic monitoring (Gibb et al., 2018), a technique that relies on autonomous acoustic recording units. PAM allows researchers to monitor acoustically active animals and their habitats at temporal and spatial scales that are impossible to achieve using only human observers. Interest in use of PAM in terrestrial environments has increased substantially in recent years (Sugai et al., 2019), due to the reduced price of autonomous recording units and improved battery life and data storage capabilities. However, the use of PAM often leads to the collection of terabytes of data that is time- and cost-prohibitive to analyze manually.

Automated detection

Automated detection for PAM data refers to identifying the start and stop time of signals of interest within a longer sound recording (Stowell, 2022). Some of the early non-deep learning approaches for the automated detection of acoustic signals in terrestrial PAM data include binary point matching (Katz et al., 2016), spectrogram cross-correlation (Balantic &

Clink, & Ahmad. (2025). gibbonNetR: an R Package for the Use of Convolutional Neural Networks for Automated Detection of Acoustic Data. 1 *Journal of Open Source Software*, 10(110), 7250. https://doi.org/10.21105/joss.07250.



Donovan, 2020), or the use of a band-limited energy detector and subsequent classifier, such as support vector machine (Clink et al., 2023; Kalan et al., 2015). Recent advances in deep learning have revolutionized image and speech recognition (LeCun et al., 2015), with important cross-over for the analysis of PAM data. Traditional approaches to machine learning relied heavily on feature engineering, since early machine learning algorithms required a reduced set of representative features that were manually chosen by researchers, such as features estimated from the spectrogram.

Deep learning does not require feature engineering (Stevens et al., 2020), as the algorithms include a step that identifies relevant features from the input. This can lead to faster development time and increased ability to represent complex patterns typically seen in image and acoustic data. Convolutional neural networks (CNNs) — one of the most widely used deep learning algorithms—are useful for processing data that have a 'grid-like topology', such as image data that can be considered a 2-dimensional grid of pixels (Goodfellow et al., 2016). The 'convolutional' layer learns the feature representations of the inputs; these convolutional layers consist of a set of filters, which are two-dimensional matrices of numbers, and the primary parameter is the number of filters (Gu et al., 2018). If training data are scarce, overfitting may occur as representations of images tend to be large with many variables (LeCun et al., 1995).

Transfer learning

Training deep learning models generally requires a large amount of training data and substantial computing resources. Transfer learning is an approach wherein the architecture of a pre-trained CNN (which is generally trained on a very large dataset) is applied to a new classification problem. For example, CNNs trained on the ImageNet dataset of > 1 million images (Deng et al., 2009) such as ResNet have been applied to automated detection/classification of primate and bird species from PAM data (Dufourq et al., 2022; Ruan et al., 2022). Generally, very few practitioners train a CNN from scratch, and there are two common approaches for transfer learning. The first option is to use the CNN as a feature extractor, and train only the last classification layer. The second option is known as 'fine-tuning', where instead of initializing a neural network with random weights, the initialization is done using the pre-trained network. Using these pre-trained weights is valuable because the model has already learned useful feature representations (Takhirov, 2021). Both approaches require substantially less computing power than training from scratch. The functions in the 'gibbonNetR' package allow users to train models using both types of transfer learning.

State of the field

The two most popular open-source programming languages are R and Python (Scavetta & Angelov, 2021). Python has surpassed R in terms of overall popularity, but R remains an important language for the life sciences (Lawlor et al., 2022). 'Keras' (Chollet & others, 2015), 'PyTorch' (Paszke et al., 2019) and 'Tensorflow' (Martín Abadi et al., 2015) are some of the more popular neural network libraries; these libraries were all initially developed for the Python programming language. One of the earliest implementations of automated detection using R was the 'monitoR' package, which included functions for template detection, which identifies signals of interest in a certain frequency range above specified energy thresholds (Araya-Salas & Smith-Vidaurre, 2017). The 'gibbonR' package combined energy-based detection with traditional machine learning classification (Clink & Klinck, 2019).

Until recently, deep learning implementations in R relied on the 'reticulate' package, which served as an interface to Python (Ushey et al., 2022). For example, early packages for automated detection using deep learning in R used the 'reticulate' package Silva et al. (2022). However, the recent release of the 'torch for R' ecosystem provides a framework based on 'PyTorch' that runs natively in R and has no dependency on Python (Falbel, 2023). Running natively in R means more straightforward installation, and higher accessibility for users of



the R programming environment. Keydana (2023) provides tutorials for image and audio classification in the 'torch for R' ecosystem, and the functionality in 'gibbonNetR' relies heavily on these tutorials. Variations of the transfer learning approaches included in this package have already been implemented in Python (Dufourq et al., 2022). Recent advances have used embeddings from audio classification models trained on bird songs for new classification problems, and in many cases, these embeddings led to better performance than general audio or image datasets (Ghani et al., 2023).

Overview

The package 'gibbonNetR' provides functions to create spectrogram images using the 'seewave' package (Sueur et al., 2008), and train and deploy six CNN architectures: AlexNet (Krizhevsky et al., 2017), VGG16, VGG19 (Simonyan & Zisserman, 2014), ResNet18, ResNet50, and ResNet152 (He et al., 2016)) trained on the ImageNet dataset (Deng et al., 2009). This package has been used for automated detection of gunshots (Vu et al., 2024) and the calls of two gibbon species (Clink, Kim, et al., 2024; Clink, Cross-Jaya, et al., 2024). The package also has functions to evaluate model performance, deploy the highest-performing model over a directory of sound files, and extract embeddings from trained models to visualize acoustic data. We provide an example dataset that consists of labelled vocalizations of the loud calls of four vertebrates (see detailed description below) from Danum Valley Conservation Area, Sabah, Malaysia (Clink & Hamid Ahmad, 2024). Detailed usage instructions for 'gibbonNetR' can be found on the 'gibbonNetR' documentation site

Data summary

We include sound files and spectrogram images of five sound classes: great argus pheasant (*Argusianus argus*) long calls (Clink et al., 2021), helmeted hornbills (*Rhinoplax vigil*), and rhinoceros hornbills (*Buceros rhinoceros*) (Kennedy et al., 2023), female gibbons (*Hylobates funereus*) and a catch-all "noise" category. The data come from two separate PAM arrays in Danum Valley Conservation Area, Sabah, Malaysia. The training and validation data come from a wide array of Swift autonomous recording units placed on ~750 m spacing (Clink et al., 2023), and the test data come from a different, smaller array (~250 m spacing) within the same area. We used a band-limited energy detector to identify signals that were 3-sec or longer duration within the 400-1600 Hz range, and then a single observer (DJC) manually sorted the detections into their respective categories (Clink et al., 2023).

Preparing training, validation, and test data

The package currently uses spectrogram images (Figure 1) to train and evaluate CNN model performance, and we include a function that can be used to create spectrogram images from Waveform Audio File Format (.wav) files. The .wav files should be organized into separate folders, with each folder named according to the class label of the files it contains. We highly recommend that your test data come from a different recording time and/or location to better understand the generalizability of the models (Stowell, 2022).





Figure 1: Spectrograms of training clips for CNNs. The x- and y-axis labels are not included when training the model, so are not shown here. The duration of the clips is variable, and the frequency range is 0.4-1.6 kHz.

Model training

The package currently allows for the training of six different CNN architectures ('alexnet', 'vgg16', 'vgg19', 'resnet18', 'resnet50', or 'resnet152'), and the user can specify if they want to freeze the feature extraction layers or not. There is also the option to train a binary or multi-class classifer.

Evaluate model performance

We can compare the performance of different CNN architectures (Figure 2). Using the 'get_best_performance()' function, we can evaluate the performance of different model architectures on the test dataset for the specified class. We can calculate the best F1, precision, and recall using the 'caret' package (Kuhn, 2008), and the area under the ROC (receiver operating characteristic) curve using the 'ROCR' package (Sing et al., 2005), which is a threshold or confidence independent metric that evaluates the classifier's ability to discriminate between positive and negative classes.

```
PerformanceOutput <- get_best_performance(
    performancetables.dir =
        performancetables.dir,
    class = 'female.gibbon',
    model.type = "multi",
    Thresh.val = 0
)</pre>
```

PerformanceOutput\$f1_plot





Results for female.gibbon class

Figure 2: Evaluating performance of pre-trained CNNs.

Extract embeddings

Embeddings from deep learning models can be used as features in unsupervised approaches, with promising results for call repertoires (Best et al., 2023) and individual identity (Lakdari et al., 2024). This package contains a function to use pre-trained CNNs to extract embeddings, where the trained model path, along with test data location and target class are specified. Depending on the research question, this output could be used to visualize true and false positives from automated detection, or to explore differences in call types or potential number of individuals in the dataset.

We can plot the unsupervised clustering results

In Figure 3, the top plot is a Uniform Manifold Approximation and Projection (UMAP) where each point represents one call, and the colors indicate the original class label. The bottom plot is the same UMAP plot, but with points colored based on cluster assignment by the 'hdbscan' algorithm (Hahsler et al., 2019).



True labels



Figure 3: UMAP plot of embeddings from test data set colored by actual label (top) and unsupervised cluster assignment (bottom).

Explore the unsupervised clustering results

For a specified target class we can calculate the Normalize Mutual Information score, which provides a value between 0 and 1, indicating the match between cluster labels and actual labels. We also create a confusion matrix using the 'caret' package (Kuhn, 2008), which returns the results when we use the unsupervised clustering algorithm function 'hdbscan' (Hahsler et al., 2019) to match the target class to the cluster with the largest number of observations of that particular class.

Future directions

There have been huge advances in the fields of deep learning and automated detection for PAM data in recent years. The approach presented in this package is one of the first to use the 'torch for R' ecosystem and to employ automated detection using deep learning natively in R. More recent approaches that use models that are explicitly trained on bioacoustics data, such as BirdNET (Ghani et al., 2023), have been introduced. There is a huge need in the field of bioacoustics to do benchmarking, wherein different model architectures and performance are compared across diverse datasets. The methods presented here can provide important benchmarks for future work and for understanding how and if deep learning advances improve performance over more traditional methods. In addition, this package provides a comprehensive suite of tools for processing, analyzing, and visualizing acoustic data, providing robust support for tasks such as automated detection, feature extraction, classification, and data visualization, which are critical for conservation work using PAM. The R package is available on Github, where issues can be opened.

Clink, & Ahmad. (2025). gibbonNetR: an R Package for the Use of Convolutional Neural Networks for Automated Detection of Acoustic Data. 6 *Journal of Open Source Software*, *10*(110), 7250. https://doi.org/10.21105/joss.07250.



Ethical statement

The research presented here adhered to all local and international laws. Institutional approval was provided by Cornell University (IACUC 2017–0098). Sabah Biodiversity Centre and the Danum Valley Management Committee provided permission for the collection of acoustic recordings.

Acknowledgments

We would like to thank the Sabah Biodiversity Centre and Danum Valley Conservation Area for granting us permission to conduct research. We are incredibly grateful for the detailed comments provided by Steffi LaZerte and Camille Desjonquères, which substantially improved the package and documentation.

References

- Araya-Salas, M., & Smith-Vidaurre, G. (2017). warbleR: An r package to streamline analysis of animal acoustic signals. *Methods in Ecology and Evolution*, 8(2), 184–191. https: //doi.org/10.1111/2041-210X.12624
- Balantic, C., & Donovan, T. (2020). AMMonitor: Remote monitoring of biodiversity in an adaptive framework with r. *Methods in Ecology and Evolution*, 11(7), 869877. https: //doi.org/10.1111/2041-210X.13397
- Best, P., Paris, S., Glotin, H., & Marxer, R. (2023). Deep audio embeddings for vocalisation clustering. PLOS ONE, 18(7), 1–18. https://doi.org/10.1371/journal.pone.0283396
- Chollet, F., & others. (2015). Keras. https://doi.org/10.1163/1574-9347_bnp_e612900
- Clink, D. J., Cross-Jaya, H., Kim, J., Ahmad, A. H., Hong, M., Sala, R., Birot, H., Agger, C., Vu, T. T., Thi, H. N., Chi, T. N., & Klinck, H. (2024). Benchmarking for the automated detection and classification of southern yellow-cheeked crested gibbon calls from passive acoustic monitoring data. *bioRxiv*. https://doi.org/10.1101/2024.08.17.608420
- Clink, D. J., Groves, T., Ahmad, A. H., & Klinck, H. (2021). Not by the light of the moon: Investigating circadian rhythms and environmental predictors of calling in bornean great argus. *Plos One*, 16(2), e0246564. https://doi.org/10.1371/journal.pone.0246564
- Clink, D. J., & Hamid Ahmad, A. (2024). A labelled dataset of the loud calls of four vertebrates collected using passive acoustic monitoring in malaysian borneo. https://doi.org/10.5281/ zenodo.14213067
- Clink, D. J., Kier, I., Ahmad, A. H., & Klinck, H. (2023). A workflow for the automated detection and classification of female gibbon calls from long-term acoustic recordings. *Frontiers in Ecology and Evolution*, 11. https://doi.org/10.3389/fevo.2023.1071640
- Clink, D. J., Kim, J., Cross-Jaya, H., Ahmad, A. H., Hong, M., Sala, R., Birot, H., Agger, C., Vu, T. T., Thi, H. N., & others. (2024). Automated detection of gibbon calls from passive acoustic monitoring data using convolutional neural networks in the" torch for r" ecosystem. arXiv Preprint arXiv:2407.09976. https://doi.org/10.48550/arXiv.2407.09976
- Clink, D. J., & Klinck, H. (2019). gibbonR: An r package for the detection and classification of acoustic signals. arXiv Preprint arXiv:1906.02572. https://doi.org/10.48550/arXiv.1906. 02572
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. 248255. https://doi.org/10.1109/cvpr.2009.5206848

Dufourq, E., Batist, C., Foquet, R., & Durbach, I. (2022). Passive acoustic monitoring of



animal populations with transfer learning. *Ecological Informatics*, 70, 101688. https://doi.org/10.1016/j.ecoinf.2022.101688

- Falbel, D. (2023). Luz: Higher level 'API' for 'torch'. https://doi.org/10.32614/CRAN. package.luz
- Ghani, B., Denton, T., Kahl, S., & Klinck, H. (2023). Global birdsong embeddings enable superior transfer learning for bioacoustic classification. *Scientific Reports*, 13(1), 22876. https://doi.org/10.1038/s41598-023-49989-z
- Gibb, R., Browning, E., Glover-Kapfer, P., & Jones, K. E. (2018). Emerging opportunities and challenges for passive acoustics in ecological assessment and monitoring. *Methods in Ecology and Evolution*. https://doi.org/10.1111/2041-210X.13101
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep learning. MIT Press.
- Gu, J., Wang, Z., Kuen, J., Ma, L., Shahroudy, A., Shuai, B., Liu, T., Wang, X., Wang, G., Cai, J., & others. (2018). Recent advances in convolutional neural networks. *Pattern Recognition*, 77, 354377. https://doi.org/10.1016/j.patcog.2017.10.013
- Hahsler, M., Piekenbrock, M., & Doran, D. (2019). dbscan: Fast density-based clustering with R. Journal of Statistical Software, 91(1), 1–30. https://doi.org/10.18637/jss.v091.i01
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. 770778. https://doi.org/10.1109/cvpr.2016.90
- Kalan, A. K., Mundry, R., Wagner, O. J. J., Heinicke, S., Boesch, C., & Kühl, H. S. (2015). Towards the automated detection and occupancy estimation of primates using passive acoustic monitoring. *Ecological Indicators*, 54(July 2015), 217226. https://doi.org/10. 1016/j.ecolind.2015.02.023
- Katz, J., Hafner, S. D., & Donovan, T. (2016). Assessment of error rates in acoustic monitoring with the r package monitoR. *Bioacoustics*, 25(2), 177196. https://doi.org/10. 1080/09524622.2015.1133320
- Kennedy, A. G., Ahmad, A. H., Klinck, H., Johnson, L. M., & Clink, D. J. (2023). Evidence for acoustic niche partitioning depends on the temporal scale in two sympatric bornean hornbill species. *Biotropica*, 55(2), 517–528. https://doi.org/10.1111/btp.13205
- Keydana, S. (2023). Deep learning and scientific computing with r torch. CRC Press. https://doi.org/10.1201/9781003275923
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2017). Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6), 8490. https://doi. org/10.1145/3065386
- Kuhn, M. (2008). Caret package. Journal of Statistical Software, 28(5), 126. https://doi.org/10.18637/jss.v028.i05
- Lakdari, M. W., Ahmad, A. H., Sethi, S., Bohn, G. A., & Clink, D. J. (2024). Mel-frequency cepstral coefficients outperform embeddings from pre-trained convolutional neural networks under noisy conditions for discrimination tasks of individual gibbons. *Ecological Informatics*, 80, 102457. https://doi.org/10.1016/j.ecoinf.2023.102457
- Lawlor, J., Banville, F., Forero-Muñoz, N.-R., Hébert, K., Martínez-Lanfranco, J. A., Rogy, P., & MacDonald, A. A. M. (2022). Ten simple rules for teaching yourself R. *PLOS Computational Biology*, 18(9), e1010372. https://doi.org/10.1371/journal.pcbi.1010372
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, *521*(7553), 436–444. https://doi.org/10.1038/nature14539
- LeCun, Y., Bengio, Y., & others. (1995). Convolutional networks for images, speech, and time series. In *The handbook of brain theory and neural networks* (Vol. 3361, p. 1995).



- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Jia, Y., Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, ... Xiaoqiang Zheng. (2015). *TensorFlow: Large-scale machine learning* on heterogeneous systems. https://doi.org/10.48550/arXiv.1605.08695
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., ... Chintala, S. (2019). *PyTorch: An imperative style, high-performance deep learning library* (p. 80248035). Curran Associates, Inc. https://doi.org/10.48550/arXiv.1912.01703
- Ruan, W., Wu, K., Chen, Q., & Zhang, C. (2022). ResNet-based bio-acoustics presence detection technology of hainan gibbon calls. *Applied Acoustics*, 198, 108939. https: //doi.org/10.1016/j.apacoust.2022.108939
- Ruff, Z. J., Lesmeister, D. B., Appel, C. L., & Sullivan, C. M. (2021). Workflow and convolutional neural network for automated identification of animal sounds. *Ecological Indicators*, 124, 107419. https://doi.org/10.1016/j.ecolind.2021.107419
- Scavetta, R. J., & Angelov, B. (2021). Python and r for the modern data scientist. O'Reilly Media, Inc. https://doi.org/10.18637/jss.v103.b02
- Silva, B., Mestre, F., Barreiro, S., Alves, P. J., & Herrera, J. M. (2022). soundClass: An automatic sound classification tool for biodiversity monitoring using machine learning. *Methods in Ecology and Evolution*. https://doi.org/10.1111/2041-210X.13964
- Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. arXiv Preprint arXiv:1409.1556. https://doi.org/10.48550/arXiv.1409.1556
- Sing, T., Sander, O., Beerenwinkel, N., & Lengauer, T. (2005). ROCR: Visualizing classifier performance in r. *Bioinformatics*, 21(20), 7881. https://doi.org/10.1093/bioinformatics/ bti623
- Stevens, E., Antiga, L., & Viehmann, T. (2020). Deep Learning with PyTorch. Simon; Schuster.
- Stowell, D. (2022). Computational bioacoustics with deep learning: a review and roadmap. *PeerJ*, 10, e13152. https://doi.org/10.7717/peerj.13152
- Sueur, J., Aubin, T., & Simonis, C. (2008). Seewave: A free modular tool for sound analysis and synthesis. *Bioacoustics*, 18, 213–226. https://doi.org/10.1080/09524622.2008.9753600
- Sugai, L. S. M., Silva, T. S. F., Ribeiro, J. W., & Llusia, D. (2019). Terrestrial passive acoustic monitoring: Review and perspectives. *BioScience*, 69(1), 1525. https://doi.org/10.1093/ biosci/biy147
- Takhirov, Z. (2021). Quantized transfer learning tutorial. https://pytorch.org/tutorials/ intermediate/quantized_transfer_learning_tutorial.html
- Ushey, K., Allaire, J. J., & Tang, Y. (2022). *Reticulate: Interface to 'python'*. https://doi.org/10.32614/CRAN.package.reticulat
- Vu, T. T., Phan, D. V., Le, T. S., & Clink, D. J. (2024). Investigating hunting in a protected area in southeast asia using passive acoustic monitoring with mobile smartphones and transfer learning. *Ecological Indicators*. https://doi.org/10.1016/j.ecolind.2024.112501