
GIBBONNETR: AN R PACKAGE FOR THE USE OF CONVOLUTIONAL NEURAL NETWORKS AND TRANSFER LEARNING ON ACOUSTIC DATA

A PREPRINT

Dena Jane Clink

K. Lisa Yang Center for Conservation Bioacoustics, Cornell Lab of Ornithology
Cornell University
Ithaca, New York, United States
dena.clink@cornell.edu

Abdul Hamid Ahmad

Institute for Tropical Biology and Conservation
Universiti Malaysia Sabah (UMS)
Kota Kinabalu, Sabah, Malaysia
abdulhamidahmad.sabah@gmail.com

July 13, 2024

Abstract

1 Automated detection of acoustic signals is crucial for effective monitoring of vocal animals
2 and their habitats across large spatial and temporal scales. Recent advances in deep learning
3 have made high performance automated detection approaches more accessible to more
4 practitioners. However, there are few deep learning approaches that can be implemented
5 natively in R. The ‘torch for R’ ecosystem has made the use of transfer learning with
6 convolutional neural networks accessible for R users. Here we provide an R package and
7 workflow to use transfer learning for the automated detection of acoustics signals from
8 passive acoustic monitoring (PAM) data collected in Sabah, Malaysia. The package provides
9 functions to create spectrogram images from PAM data, compare the performance of different
10 pre-trained CNN architectures trained on the ImageNet dataset, deploy trained models over
11 directories of sound files, and extract embeddings from trained models. The R programming
12 language remains one of the most commonly used languages among ecologists, and we
13 hope that this package makes deep learning approaches more accessible to this audience.
14 In addition, these models can serve as important benchmarks for more state-of-the-art
15 approaches.

16 **Keywords** deep learning · passive acoustic monitoring · gibbon · automated detection

17 1 Statement of need

18 1.1 *Passive acoustic monitoring*

19 We are in a biodiversity crisis, and there is a great need for the ability to rapidly assess biodiversity in order to
20 understand and mitigate anthropogenic impacts. One approach that can be especially effective for monitoring
21 of vocal yet cryptic animals is the use of passive acoustic monitoring (Gibb et al. 2018), a technique that
22 relies on autonomous acoustic recording units. PAM allows researchers to monitor vocal animals and their
23 habitats at temporal and spatial scales that are impossible to achieve using only human observers. Interest in
24 use of PAM in terrestrial environments has increased substantially in recent years (Sugai et al. 2019), due to

25 reduced price of the recording units and improved battery life and data storage capabilities. However, the
 26 use of PAM often leads to the collection of terabytes of data that is time- and cost-prohibitive to analyze
 27 manually.

28 1.2 Automated detection

29 Some of the early non-deep learning approaches for the automated detection of acoustic signals in terrestrial
 30 PAM data include binary point matching (Katz, Hafner, and Donovan 2016), spectrogram cross-correlation
 31 (Balantic and Donovan 2020), or the use of a band- limited energy detector and subsequent classifier, such
 32 as support vector machine (Clink et al. 2023; Kalan et al. 2015). Recent advances in deep learning have
 33 revolutionized image and speech recognition (LeCun, Bengio, and Hinton 2015), with important cross-over for
 34 the analysis of PAM data. Traditional approaches to machine learning relied heavily on feature engineering,
 35 as early machine learning algorithms required a reduced set of representative features, such as features
 36 estimated from the spectrogram. Deep learning does not require feature engineering (Stevens, Antiga, and
 37 Viehmann 2020). Convolutional neural networks (CNNs) — one of the most widely used deep learning
 38 algorithms—are useful for processing data that have a ‘grid-like topology’, such as image data that can be
 39 considered a 2-dimensional grid of pixels (Goodfellow, Bengio, and Courville 2016). The ‘convolutional’ layer
 40 learns the feature representations of the inputs; these convolutional layers consist of a set of filters which are
 41 basically two-dimensional matrices of numbers and the primary parameter is the number of filters (Gu et al.
 42 2018). Therefore, with CNN’s there is no feature engineering required. However, if training data are scarce,
 43 overfitting may occur as representations of images tend to be large with many variables (LeCun, Bengio, and
 44 others 1995).

45 2 Transfer learning?

46 Training deep learning models generally requires a large amount of training data and substantial computing
 47 resources, which can be hard to obtain with PAM data. Transfer learning is an approach wherein the
 48 architecture of a pretrained CNN (which is generally trained on a very large dataset) is applied to a new
 49 classification problem. For example, CNNs trained on the ImageNet dataset of > 1 million images (Deng
 50 et al. 2009) such as ResNet have been applied to automated detection/classification of primate and bird
 51 species from PAM data (Dufourq et al. 2022; Ruan et al. 2022). At the most basic level, transfer learning in
 52 computer vision applications retains the feature extraction or embedding layers, and modifies the last few
 53 classification layers to be trained for a new classification task (Dufourq et al. 2022). Recent advances have
 54 used embeddings from audio classification models trained on bird songs for new classification problems, and
 55 in most cases these embeddings led to better performance than general audio or image datasets (Ghani et al.
 56 2023).

57 3 State of the field

58 The two most popular open-source programming languages are R and Python (Scavetta and Angelov 2021).
 59 Python has surpassed R in terms of overall popularity, but R remains an important language for the life sciences
 60 (Lawlor et al. 2022). ‘Keras’ (Chollet and others 2015), ‘PyTorch’ (Paszke et al. 2019) and ‘Tensorflow’
 61 (Martín Abadi et al. 2015) are some of the more popular neural network libraries; these libraries were all
 62 initially developed for the Python programming language. Until recently, deep learning implementations in R
 63 relied on the ‘reticulate’ package which served as an interface to Python (Ushey, Allaire, and Tang 2022).
 64 Previous implementations of automated detection using deep learning in R relied on the ‘reticulate’ package
 65 Silva et al. (2022). However, the recent release of the ‘torch for R’ ecosystem provides a framework based on
 66 ‘PyTorch’ that runs natively in R and has no dependency on Python (Falbel 2023). Running natively in R
 67 means more straightforward installation, and higher accessibility for users of the R programming environment.
 68 Keydana (2023) provides tutorials for transfer learning in the ‘torch for R’ ecosystem, and the functions in
 69 ‘gibbonNetR’ rely heavily on these tutorials. The transfer learning approaches included in this package have
 70 already been implemented in Python (Dufourq et al. 2022).

71 4 Usage

72 4.1 Overview

73 This package provides functions to create spectrogram images, use transfer learning for six pretrained CNN
 74 architectures: AlexNet (Krizhevsky, Sutskever, and Hinton 2017) , VGG16, VGG19 (Simonyan and Zisserman
 75 2014), ResNet18, ResNet50, and ResNet152 (He et al. 2016)) trained on the ImageNet dataset (Deng et al.
 76 2009). The package also has functions to evaluate model performance, deploy the highest performing model
 77 over a directory of sound files, and extract embeddings from trained models to visualize acoustic data. We
 78 provide an example dataset that consists of labelled vocalizations of the loud calls of four vertebrates from
 79 Danum Valley Conservation Area, Sabah, Malaysia.

80 4.2 Data summary

81 We include spectrogram images of five classes: great argus pheasant (*Argusianus argus*) long calls (Clink et
 82 al. 2021), helmeted hornbills (*Rhinoplax vigil*), and rhinoceros hornbills (*Buceros rhinoceros*) (Kennedy et
 83 al. 2023), female gibbons (*Hylobates funereus*) and a catch-all “noise” category. The data come from two
 84 separate PAM arrays in Danum Valley Conservation Area, Sabah, Malaysia. The training and validation
 85 data come from a wide array of Swift autonomous recording units placed on ~750 m spacing (Clink et al.
 86 2023), and the test data come from a different, smaller array (~250 m spacing) within the same area. We
 87 used a band-limited energy detector to identify signals that were 3-sec or longer duration within the 400-1600
 88 Hz range, and then a single observer (DJC) manually sorted the detections into their respective categories
 89 (Clink et al. 2023).

90 4.3 Spectrogram images

91 The package currently uses spectrogram images (Figure 1), and ‘gibbonNetR’ includes a function that can be
 92 used to create spectrogram images from .wav files. The ‘splits’ option will assign the specified proportion
 93 of clips to either training, validation, or test folders. We highly recommend that your test data come from
 94 a different recording time and/or location to better understand the generalizability of the models (Stowell
 95 2022).

```

library(gibbonNetR)
# Generate spectrogram images for training and validation.
# Test data should come from a different source.
gibbonNetR::spectrogram_images(
  trainingBasePath = trainingBasePath,
  outputBasePath = outputBasePath,
  minfreq.khz = 0.4,
  maxfreq.khz = 1.6,
  splits = c(0.7, 0.3, 0), # Assign proportion to training, validation, or test folders
  new.sampleratehz = 'NA'
)

```

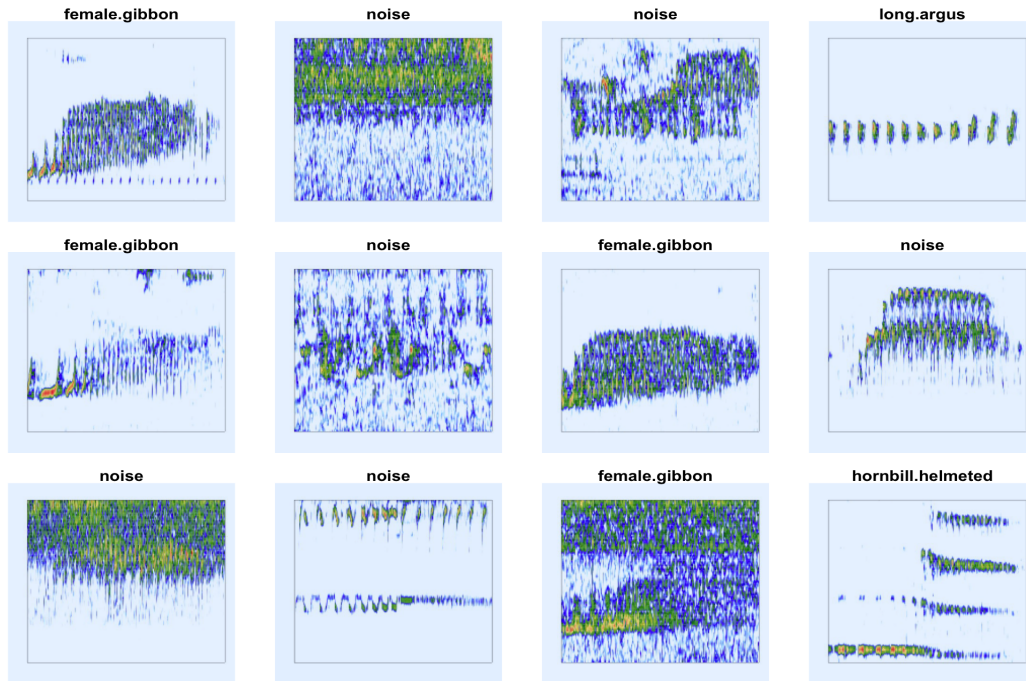


Figure 1: Spectrograms of training clips for CNNs

96 4.4 Model training

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

```
# Location of spectrogram images for training
input.data.path <- 'data/examples/'

# Location of spectrogram images for testing
test.data.path <- 'data/examples/test/'

# User specified training data label for metadata
trainingfolder.short <- 'danummulticlassexample'

# We can specify the number of epochs to train here
epoch.iterations <- c(20)

# Function to train a multi-class CNN
gibbonNetR::train_CNN_multi(input.data.path=input.data.path,
                             architecture = 'resnet18',
                             learning_rate = 0.001,
                             class_weights = rep( (1/5), 5),
                             test.data=test.data.path,
                             unfreeze.param = TRUE,
                             epoch.iterations=epoch.iterations,
                             save.model= TRUE,
                             early.stop = "yes",
                             output.base.path = "model_output/",
                             trainingfolder=trainingfolder.short,
                             noise.category = "noise")
```

100 4.5 Evaluate model performance

101 We can compare the performance of different CNN architectures (Figure 2). Using the ‘get_best_performance’
 102 function we can evaluate the performance of different model architectures on the test dataset for the specified
 103 class.

```
PerformanceOutput <- gibbonNetR::get_best_performance(performancetables.dir=performancetables.dir,
                                                    class='female.gibbon',
                                                    model.type = "multi",Thresh.val=0)
PerformanceOutput$f1_plot
```

Results for female.gibbon class

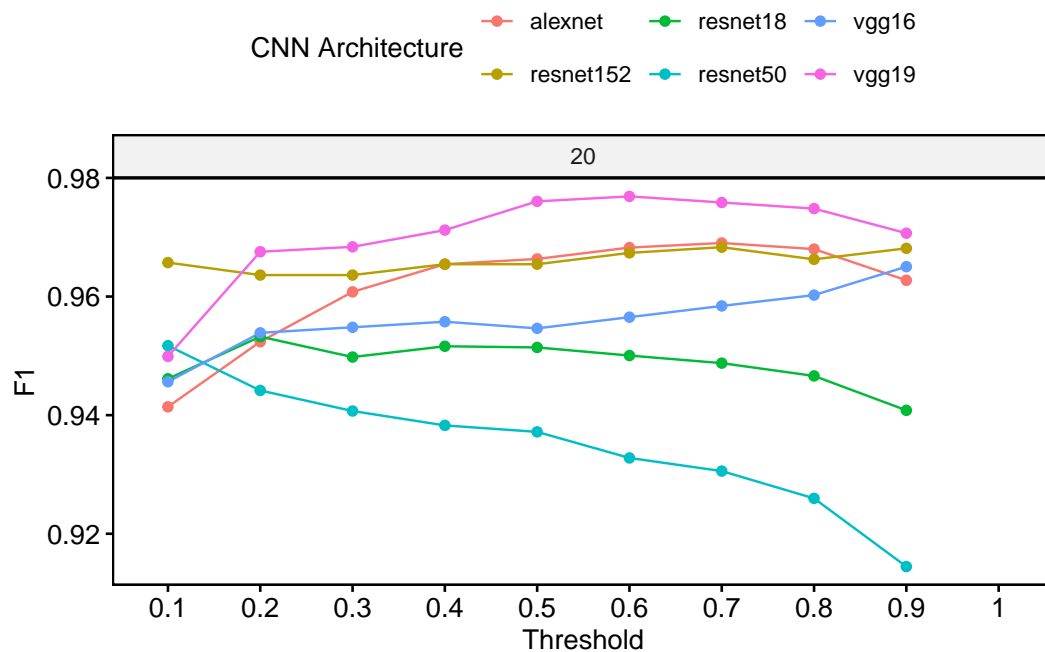


Figure 2: Evaluating performance of pretrained CNNs

104 4.6 Extract embeddings

105 Embeddings from deep learning models can be used as features in unsupervised approaches, with promising
 106 results for call repertoires (Best et al. 2023) and individual identity (Lakdari et al. 2024). This package
 107 contains the function ‘extract_embeddings’ to use pretrained CNNs to extract embeddings, where the trained
 108 model path, along with test data location and target class are specified.

```
result <- gibbonNetR::extract_embeddings(test_input="data/examples/test/",
                                       model_path=ModelPath,
                                       target_class = "female.gibbon")
```

109 4.7 We can plot the unsupervised clustering results

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

```
result$EmbeddingsCombined
```

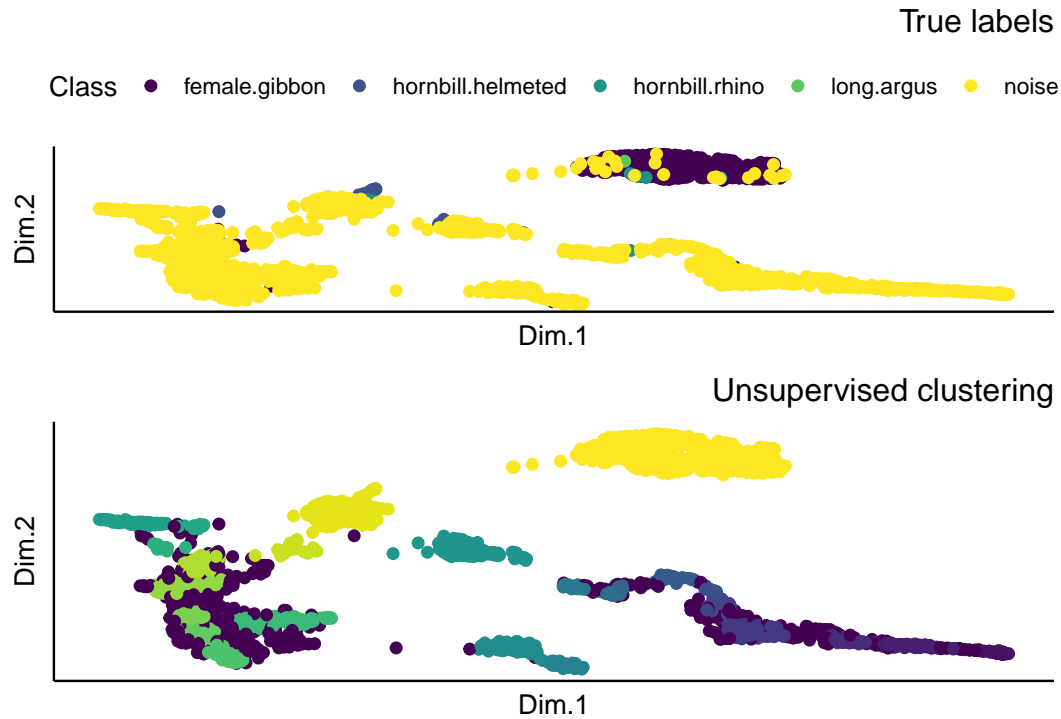


Figure 3: UMAP plot of embeddings from test data set colored by actual label (top) and unsupervised cluster assignment (bottom)

114 4.7.1 We can explore the unsupervised clustering results

115 Here we can see the Normalize Mutual Information score, which provides a value between 0 and 1, indicating
116 the match between cluster labels and actual labels.

```
result$NMI
#> [1] 0.7317309
```

117 The confusion matrix is made using the ‘caret’ package (Kuhn 2008) returns the results when we use ‘hdbscan’
118 (Hahsler, Piekenbrock, and Doran 2019) to match the target class to the cluster with the largest number of
119 observations of that particular class.

```
result$ConfusionMatrix
#>      Sensitivity      Specificity      Pos Pred Value
#>      0.9110672      0.9835556      0.9257028
#>      Neg Pred Value      Precision      Recall
#>      0.9800709      0.9257028      0.9110672
#>      F1      Prevalence      Detection Rate
#>      0.9183267      0.1835994      0.1672714
#>      Detection Prevalence      Balanced Accuracy
#>      0.1806967      0.9473114
```

120 5 Future directions

121 There have been huge advances in the fields of deep learning and automated detection approaches for PAM
122 data in recent years. The approach presented in this package (as of 2024) is no longer state of the art,

123 although it was just recently developed. More recent approaches use transfer learning from models that are
 124 explicitly trained on bioacoustics data, such as BirdNET (Ghani et al. 2023). However, there is a huge need
 125 in the field of bioacoustics to do benchmarking, wherein different model architectures and performance are
 126 compared across diverse datasets. Therefore, the methods presented here can provide important benchmarks
 127 for future work, and for understanding how and if advances improve performance over more traditional
 128 methods. The R package is available on Github, where issues can be opened.

129 6 Acknowledgments

130 We would like to thank the Sabah Biodiversity Centre and Danum Valley Conservation Area for granting us
 131 permission to conduct research.

132 References

- 133 Balantic, Cathleen, and Therese Donovan. 2020. “AMMonitor: Remote Monitoring of Biodiversity in an
 134 Adaptive Framework with r.” *Methods in Ecology and Evolution* 11 (7): 869877.
- 135 Best, Paul, Sébastien Paris, Hervé Glotin, and Ricard Marxer. 2023. “Deep Audio Embeddings for Vocalisation
 136 Clustering.” *PLOS ONE* 18 (7): 1–18. <https://doi.org/10.1371/journal.pone.0283396>.
- 137 Chollet, François, and others. 2015. “Keras.” <https://keras.io>.
- 138 Clink, Dena J., Tom Groves, Abdul Hamid Ahmad, and Holger Klinck. 2021. “Not by the Light of the Moon:
 139 Investigating Circadian Rhythms and Environmental Predictors of Calling in Bornean Great Argus.” *Plos*
 140 *One* 16 (2): e0246564.
- 141 Clink, Dena J., Isabel Kier, Abdul Hamid Ahmad, and Holger Klinck. 2023. “A Workflow for the Automated
 142 Detection and Classification of Female Gibbon Calls from Long-Term Acoustic Recordings.” *Frontiers in*
 143 *Ecology and Evolution* 11. <https://www.frontiersin.org/articles/10.3389/fevo.2023.1071640>.
- 144 Deng, Jia, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. “Imagenet: A Large-Scale
 145 Hierarchical Image Database.” In, 248255. Ieee.
- 146 Dufourq, Emmanuel, Carly Batist, Ruben Foquet, and Ian Durbach. 2022. “Passive Acoustic Monitoring of
 147 Animal Populations with Transfer Learning.” *Ecological Informatics* 70: 101688. <https://doi.org/https://doi.org/10.1016/j.ecoinf.2022.101688>.
- 148 Falbel, Daniel. 2023. *Luz: Higher Level 'API' for 'Torch'*. <https://CRAN.R-project.org/package=luz>.
- 149 Ghani, Burooj, Tom Denton, Stefan Kahl, and Holger Klinck. 2023. “Global Birdsong Embeddings Enable
 150 Superior Transfer Learning for Bioacoustic Classification.” *Scientific Reports* 13 (1): 22876.
- 151 Gibb, Rory, Ella Browning, Paul Glover-Kapfer, and Kate E. Jones. 2018. “Emerging Opportunities and
 152 Challenges for Passive Acoustics in Ecological Assessment and Monitoring.” *Methods in Ecology and*
 153 *Evolution*, October. <https://doi.org/10.1111/2041-210X.13101>.
- 154 Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- 155 Gu, Jiuxiang, Zhenhua Wang, Jason Kuen, Lianyang Ma, Amir Shahroudy, Bing Shuai, Ting Liu, et al. 2018.
 156 “Recent Advances in Convolutional Neural Networks.” *Pattern Recognition* 77: 354377.
- 157 Hahsler, Michael, Matthew Piekenbrock, and Derek Doran. 2019. “dbscan: Fast Density-Based Clustering
 158 with R.” *Journal of Statistical Software* 91 (1): 1–30. <https://doi.org/10.18637/jss.v091.i01>.
- 159 He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. “Deep Residual Learning for Image
 160 Recognition.” In, 770778.
- 161 Kalan, Ammie K., Roger Mundry, Oliver J J Wagner, Stefanie Heinicke, Christophe Boesch, and Hjalmar
 162 S. Köhl. 2015. “Towards the Automated Detection and Occupancy Estimation of Primates Using
 163 Passive Acoustic Monitoring.” *Ecological Indicators* 54 (July 2015): 217226. <https://doi.org/10.1016/j.ecolind.2015.02.023>.
- 164 Katz, Jonathan, Sasha D Hafner, and Therese Donovan. 2016. “Assessment of Error Rates in Acoustic
 165 Monitoring with the r Package monitoR.” *Bioacoustics* 25 (2): 177196.
- 166 Kennedy, Amy G, Abdul Hamid Ahmad, Holger Klinck, Lynn M Johnson, and Dena J Clink. 2023. “Evidence
 167 for Acoustic Niche Partitioning Depends on the Temporal Scale in Two Sympatric Bornean Hornbill
 168 Species.” *Biotropica* 55 (2): 517–28.
- 169 Keydana, Sigrid. 2023. *Deep Learning and Scientific Computing with r Torch*. CRC Press.
- 170 Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E Hinton. 2017. “Imagenet Classification with Deep
 171 Convolutional Neural Networks.” *Communications of the ACM* 60 (6): 8490.
- 172 Kuhn, Max. 2008. “Caret Package.” *Journal of Statistical Software* 28 (5): 126.
- 173 Lakdari, Mohamed Walid, Abdul Hamid Ahmad, Sarab Sethi, Gabriel A Bohn, and Dena J Clink. 2024.
 174 “Mel-Frequency Cepstral Coefficients Outperform Embeddings from Pre-Trained Convolutional Neural
 175
 176

- 177 Networks Under Noisy Conditions for Discrimination Tasks of Individual Gibbons.” *Ecological Informatics*
 178 80: 102457.
- 179 Lawlor, Jake, Francis Banville, Norma-Rocio Forero-Muñoz, Katherine Hébert, Juan Andrés Martínez-
 180 Lanfranco, Pierre Rogy, and A. Andrew M. MacDonald. 2022. “Ten Simple Rules for Teaching Yourself R.”
 181 *PLOS Computational Biology* 18 (9): e1010372. <https://doi.org/10.1371/journal.pcbi.1010372>.
- 182 LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. 2015. “Deep Learning.” *Nature* 521 (7553): 436–44.
 183 <https://doi.org/10.1038/nature14539>.
- 184 LeCun, Yann, Yoshua Bengio, and others. 1995. “Convolutional Networks for Images, Speech, and Time
 185 Series.” *The Handbook of Brain Theory and Neural Networks* 3361 (10): 1995.
- 186 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado,
 187 et al. 2015. “TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems.” [https://www.
 188 tensorflow.org/](https://www.tensorflow.org/).
- 189 Paszke, Adam, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
 190 Trevor Killeen, et al. 2019. “PyTorch: An Imperative Style, High-Performance Deep Learn-
 191 ing Library.” In, 80248035. Curran Associates, Inc. [http://papers.neurips.cc/paper/
 192 9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 193 Ruan, Wenda, Keyi Wu, Qingchun Chen, and Chengyun Zhang. 2022. “ResNet-Based Bio-Acoustics Presence
 194 Detection Technology of Hainan Gibbon Calls.” *Applied Acoustics* 198: 108939. <https://doi.org/https://doi.org/10.1016/j.apacoust.2022.108939>.
- 196 Ruff, Zachary J., Damon B. Lesmeister, Cara L. Appel, and Christopher M. Sullivan. 2021. “Workflow and
 197 Convolutional Neural Network for Automated Identification of Animal Sounds.” *Ecological Indicators* 124
 198 (May): 107419. <https://doi.org/10.1016/j.ecolind.2021.107419>.
- 199 Scavetta, Rick J, and Boyan Angelov. 2021. *Python and r for the Modern Data Scientist*. O’Reilly Media,
 200 Inc.
- 201 Silva, Bruno, Frederico Mestre, Sílvia Barreiro, Pedro J Alves, and José M Herrera. 2022. “soundClass: An
 202 Automatic Sound Classification Tool for Biodiversity Monitoring Using Machine Learning.” *Methods in
 203 Ecology and Evolution*.
- 204 Simonyan, Karen, and Andrew Zisserman. 2014. “Very Deep Convolutional Networks for Large-Scale Image
 205 Recognition.” *arXiv Preprint arXiv:1409.1556*.
- 206 Stevens, Eli, Luca Antiga, and Thomas Viehmann. 2020. *Deep Learning with PyTorch*. Simon; Schuster.
- 207 Stowell, Dan. 2022. “Computational Bioacoustics with Deep Learning: A Review and Roadmap.” *PeerJ* 10
 208 (March): e13152. <https://doi.org/10.7717/peerj.13152>.
- 209 Sugai, Larissa Sayuri Moreira, Thiago Sanna Freire Silva, José Wagner Ribeiro, and Diego Llusia. 2019.
 210 “Terrestrial Passive Acoustic Monitoring: Review and Perspectives.” *BioScience* 69 (1): 1525. <https://doi.org/10.1093/biosci/biy147>.
- 211
- 212 Ushey, Kevin, J. J. Allaire, and Yuan Tang. 2022. *Reticulate: Interface to ‘Python’*.