

---

# An Expert-Aligned Toolbox for Explainable AI in Animal Communication

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Explainable AI (XAI) remains underdeveloped in bioacoustics, despite the growing  
2 reliance on high-performance black-box models. We evaluate the explainability of  
3 state-of-the-art models for capuchin monkey individual identification and introduce  
4 new methods to make bioacoustic classifiers more interpretable. Our approach com-  
5 bines participatory evaluation with domain experts through a web-based interface,  
6 with quantitative metrics that assess alignment between saliency maps and expert  
7 annotations. Specifically, we report metrics on ranking quality, spatial overlap and  
8 distributional similarity. Each metric is computed under complementary feature  
9 importance formulations. To facilitate annotation, we introduce a web interface  
10 for pixel-level spectrogram labeling with interactive, mask-exclusive audio play-  
11 back, allowing experts to listen separately to masked foreground or background  
12 regions and optional semi-automated segmentation. Together, these tools provide a  
13 reproducible framework for benchmarking explainability in bioacoustic models,  
14 advancing toward more transparent, collaborative, and biologically meaningful AI  
15 for animal communication.

## 16 1 Introduction

17 **Motivation.** There is no established toolbox for evaluating how bioacoustic models attend to  
18 meaningful spectro-temporal features at the pixel level, nor a participatory framework that allows  
19 domain experts to validate these models through time–frequency spectrogram annotations. This gap  
20 hinders scientific discovery and cross-disciplinary trust.

21 In order to dig deeper into the unknown semantic spaces of other species, AI models need to be  
22 carefully designed with appropriate architectures that enable effective domain expert input. Never-  
23 theless, these two priorities conflict with each other as the trend is to develop increasingly complex  
24 models [15] for which internal representations are not interpretable by design but represent the state  
25 of the art in classification performance [19]. In bioacoustics, explainability has been addressed only  
26 sparingly; recent XAI studies have explored custom models [18], while most researchers continue to  
27 rely on black-box pretrained models [2] due to their powerful transfer learning capabilities. These  
28 models achieve high performance, but their learned representations are difficult to interpret, limiting  
29 biological insight, responsible use, and trust from cross-domain expertise.

30 Capuchin monkeys produce over 27 distinct call types and exhibit cultural evolution and complex  
31 social cognition, making them an excellent model for studying animal communication [6]. Recent  
32 advances in joint cross-species embedding models have unlocked superior classification performance  
33 of caller identity for this species, offering new opportunities for remote monitoring and analysis [20].  
34 However, to uncover which features of their rich vocal repertoire convey individual identity, new  
35 methodological approaches are needed. XAI provides one such avenue, motivating this study and  
36 offering a strong test case for advancing interpretable methods in bioacoustics. While recent work has

emphasized system scalability [3] and the use of large language models for cross-modal representation learning [15], little attention has been given to model explainability. Our work addresses this gap by introducing an XAI toolbox and participatory framework to interpret black-box bioacoustic models in collaboration with domain experts.

Our work contributes (1) an XAI toolbox for bioacoustic models, combining spectrogram feature importance maps with simple, interpretable saliency evaluation metrics adapted from computer vision to compare model attention against expert annotations, and (2) a web-based annotation platform that enables pixel-level, participatory validation of model attention.

## 2 Background and Related Work

**Explainable AI.** Explainability can play an invaluable role in scientific exploration by identifying and refining target phenomena, motivating hypotheses and guiding inquiry [21]. On the other hand, explainability has become a central concern in AI as models grow in size and complexity, and society increasingly questions the consequences of their inner workings, with XAI techniques as a deciding factor for user trust and adoption [17]. In animal communication studies, ground-truth labels are usually tied to observed behavioral states and contexts hypothesized to motivate specific signals, and statistical models have long been used to test such hypotheses about semantics and linguistics phenomena. As the field shifts from using simple statistical descriptions to complex AI models, explainability becomes crucial for it to situate algorithmic insights within the rich ecological and evolutionary context of biological signals.

Compared to language and vision, interpretability in audio models, and especially in bioacoustics, is far less developed. Most work in audio has focused on acoustic event detection [11, 14, 9]. However, for bioacoustic applications explainability remains scarce. Models often operate as black-box classifiers, offering little insight into what acoustic features drive decisions. This lack of interpretability not only limits scientific understanding of animal communication but also hampers trust in deployed systems for conservation and ecological monitoring. Addressing this gap requires adapting or developing interpretability frameworks that are sensitive to the unique structure and semantics of acoustic signals in biological contexts.

More recently, interpretability has started to gain traction. Heinrich et al. [7] proposed incorporating interpretability directly into the model architecture, demonstrating how a network can learn prototypical patterns for bird species. In contrast, Silva et al. [18] focus on post-hoc analysis of trained models, using SHAP to interpret learned features. Our work follows this post-hoc perspective, as we find it more practical to study interpretability after the model has already been trained.

**Participatory Design and Human–AI Collaboration.** Explainability is not only a technical concern but also a design principle for effective human–AI collaboration. According to established guidelines [1], systems should support transparency, provide rationales, and enable meaningful human control. In wildlife monitoring, where technologies often interact with communities in overseas territories, justice-oriented design principles are equally important to strengthen governance, community agency, and cultural appropriateness [12, 13].

## 3 Approach

**Multi-Grid Spectrogram Occlusion.** We generate explanations using a *multi-grid spectrogram occlusion* procedure.<sup>1</sup>

The method systematically masks local spectro-temporal regions of the input and measures the change in model confidence. For a waveform  $x$  sampled at 48 kHz, we compute its spectrogram and partition it into grids with fixed cell sizes of 75 ms in time ( $t_w$ ) and 3 kHz in frequency ( $f_w$ ). To avoid aliasing explanations to a single grid alignment, we generate multiple translated grids by shifting the partition

<sup>1</sup>We use the term *saliency map* to refer to explanations that highlight influential input regions [16, 10, 5]. In our case, saliency maps are obtained through occlusion and visualized as spectrogram heatmaps. Regions of high saliency are interpreted as *feature importance*, i.e., spectro-temporal components most critical for model predictions.

82 along time and frequency (by  $\Delta t$ ,  $\Delta f$ ). This increases effective resolution when aggregating results,  
83 similar to adaptive strategies in other domains [4].

84 Each perturbed input  $\tilde{x}_i$  is produced by occluding a single cell. In our implementation, occlusion  
85 is applied directly in the time domain: the band-limited signal corresponding to the selected time-  
86 frequency window is extracted with zero-phase filtering, tapered with short Tukey ramps, and set  
87 to silence. This approach ensures that only the target region is modified while the remainder of the  
88 waveform remains undistorted.

89 The trained classifier for acoustic individual identification is then applied to perturbed inputs. For  
90 each occlusion we obtain class probabilities  $p(y | \tilde{x}_i)$ , to be compared with the unperturbed prediction  
91  $p(y | x)$ . Feature importance is quantified in two complementary ways: (i) *distributional change* via  
92 Jensen–Shannon divergence (JS Div) between the two predictive distributions, and (ii) *label-specific*  
93 *change* via the difference in cross-entropy with respect to the true label ( $\Delta\text{CE}$ ). Aggregating these  
94 values across all grids yields a prediction-drop heatmap aligned to the original spectrogram, providing  
95 a saliency map with higher resolution.

96 While our experiments employ silence-based occlusion, the procedure is fully parameterizable.  
97 Window sizes, translation steps, and masking strategies (e.g., pink noise or band-limited noise) can  
98 be adapted to the requirements of other tasks or domains.

99 **Web-Based Annotation Toolbox.** Pixel-level annotation is well established in computer vision,  
100 but the analogous task of segmenting spectrograms into time and frequency bins remains largely  
101 absent in bioacoustics. Our lightweight web interface, built on Flutter and Firebase, sequentially  
102 serves spectrogram–audio pairs to annotators. Users can draw masks manually or provide foreground  
103 and background points that trigger AI-assisted segmentation with Meta’s Segment Anything (SAM)  
104 [8], refining suggestions with tools such as eraser, brush, and opacity controls. To aid validation,  
105 the interface also supports playback of masked foreground and background audio at variable speeds  
106 (e.g.,  $0.3\times$  for capuchin calls). Each completed mask is stored as a binary map with metadata for  
107 subsequent evaluation. The design prioritizes ease of use and remote collaboration, consistent with  
108 principles of human–AI interaction and participatory bioacoustics [1, 13]. The source code will be  
109 released upon publication.

110 **Evaluation Metrics.** We evaluate the alignment between saliency maps and expert annotations  
111 using overlap-, ranking-, and correlation-based metrics. First, we report the Area Under the Precision-  
112 Recall Curve (AUPRC), which measures how well saliency values discriminate between annotated  
113 and non-annotated pixels. To quantify spatial overlap, we compute the Intersection-over-Union (IoU)  
114 and Coverage at a fixed threshold of 0.2, capturing how much of the annotated region is recovered  
115 and how precisely it is localized. In addition, we measure distributional similarity with Pearson  
116 correlation between continuous saliency maps and binary annotation masks [10]. Each of these  
117 metrics is computed under the two complementary importance formulations mentioned above: JS  
118 Div and  $\Delta\text{CE}$ .

## 119 4 Results

120 **Qualitative: spectrogram saliency heatmaps.** Vocal production in primates arises from me-  
121 chanical and physiological processes that generate distinctive acoustic patterns, enabling individual  
122 recognition. AI explainability techniques should be capable of revealing this information by high-  
123 lighting salient spectro-temporal regions. As shown in Fig. 1, the Whisper–Perch MRMR joint  
124 embedding model isolates specific portions of the spectrogram that may be informative, even when  
125 their biological relevance is not yet established. Expert annotations on frequency–time bins (spectro-  
126 gram pixels) can highlight known salient features such as the call itself (manual source separation).  
127 While essential, these annotations capture only what humans can interpret from the calls, whereas  
128 models, either individually or through combinations such as MRMR joint embeddings [20], can  
129 surface complementary perspectives. In turn, these saliency maps can themselves become analyzable  
130 objects, supporting statistical methods to test hypotheses about which acoustic features may carry  
131 semantic or individual identity cues.

132 **Quantitative: expert annotation alignment.** Across all metrics, Perch aligns most closely with  
133 expert annotations, achieving the highest scores in ranking, spatial overlap, and correlation (Table 1).

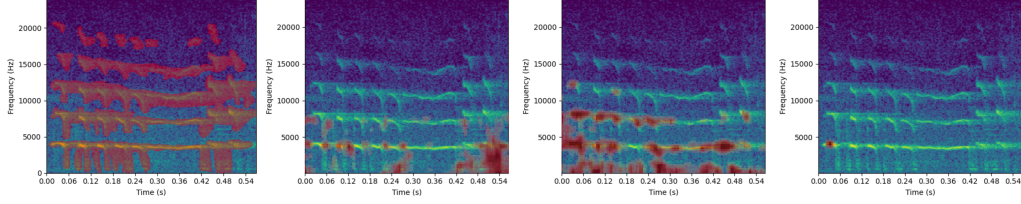


Figure 1: Qualitative evaluation of feature importance for acoustic individual identification in capuchin monkeys using saliency maps ( $\Delta CE$ ). Heatmaps (overlaid on spectrograms) highlight spectro-temporal regions that most influence model decisions about caller identity. From left to right: Annotated Mask, Whisper, Perch, and Whisper-Perch Joint Embedding. Note the variation in saliency across models, and how the Whisper-Perch plot could motivate hypothesis testing if the pattern is consistent.

Model	Importance	AUPRC	IoU <sub>0.2</sub>	Coverage <sub>0.2</sub>	Pearson Corr.
Perch-Whisper MRMR	JS Div.	$0.576 \pm 0.128$	$0.013 \pm 0.011$	$0.014 \pm 0.011$	$0.214 \pm 0.082$
	$\Delta CE$	$0.469 \pm 0.145$	$0.028 \pm 0.018$	$0.029 \pm 0.019$	$0.203 \pm 0.084$
Google Perch 2	JS Div.	<b><math>0.605 \pm 0.140</math></b>	<b><math>0.103 \pm 0.038</math></b>	<b><math>0.110 \pm 0.045</math></b>	<b><math>0.397 \pm 0.097</math></b>
	$\Delta CE$	$0.571 \pm 0.135$	$0.050 \pm 0.025$	$0.052 \pm 0.026$	$0.325 \pm 0.101$
Whisper Large V3	JS Div.	$0.372 \pm 0.172$	$0.033 \pm 0.020$	$0.037 \pm 0.022$	$0.055 \pm 0.122$
	$\Delta CE$	$0.386 \pm 0.161$	$0.030 \pm 0.020$	$0.032 \pm 0.021$	$0.055 \pm 0.102$
Baseline	JS Div	$0.347 \pm 0.148$	$0.026 \pm 0.011$	$0.029 \pm 0.011$	$0.001 \pm 0.047$
	$\Delta CE$	$0.352 \pm 0.148$	$0.017 \pm 0.010$	$0.018 \pm 0.011$	$0.000 \pm 0.047$

Table 1: Performance comparison of model architectures across expert alignment evaluation metrics using different feature importance methods.

134 Compared to Whisper, it improves overlap by about threefold and correlation by about sevenfold.  
135 Between importance formulations, JS Div generally outperforms  $\Delta CE$ , with Perch showing the  
136 clearest advantage. Although the Perch-Whisper MRMR achieves the best individual classification  
137 accuracy [20], it lags behind Perch in explainability. These results suggest that models can excel at  
138 classification while still diverging from human-recognizable cues, highlighting the importance of  
139 explainability as a complementary evaluation dimension.

## 140 5 Conclusion.

141 Human annotators can reliably identify the presence of calls on spectrograms, but we lack precise  
142 knowledge of which acoustic features are truly decisive for individual identity. As a result, the  
143 reported metrics measure alignment with human intuition rather than absolute ground truth, and  
144 should be interpreted with caution. This limitation reinforces the need for combining qualitative  
145 and quantitative XAI: saliency maps not only benchmark model interpretability but also serve as  
146 scientific tools, enabling researchers to ask new questions and generate hypotheses (Fig. 1) about  
147 animal communication.

148 **Limitations and Future Work.** The present work focuses on post-hoc saliency evaluation; integrat-  
149 ing interpretability directly into model architectures remains an open challenge. Broader validation  
150 across species and ecological contexts is also needed. Future work should explore statistical analyses  
151 on saliency maps themselves to test hypotheses on signals conveying meaning, and extend partic-  
152 ipatory platforms to more diverse user groups, and examine how explainable models can support  
153 decision-making in conservation practice.

## References

- [1] Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N Bennett, Kori Inkpen, et al. Guidelines for human-ai interaction. In *Proceedings of the 2019 chi conference on human factors in computing systems*, pages 1–13, 2019.
- [2] Jules Cauzinille, Benoit Favre, Ricard Marxer, and Arnaud Rey. Applying machine learning to primate bioacoustics: Review and perspectives. *American Journal of Primatology*, 86(10): e23666, 2024.
- [3] Vincent Dumoulin, Otilia Stretcu, Jenny Hamer, Lauren Harrell, Rob Laber, Hugo Larochelle, Bart van Merriënboer, Amanda Navine, Patrick Hart, Ben Williams, et al. The search for squawk: Agile modeling in bioacoustics. *arXiv preprint arXiv:2505.03071*, 2025.
- [4] Daniel Fink, Theodoros Damoulas, and Jaimin Dave. Adaptive spatio-temporal exploratory models: Hemisphere-wide species distributions from massively crowdsourced ebird data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 27, pages 1284–1290, 2013.
- [5] Tristan Gomez, Thomas Fréour, and Harold Mouchère. Metrics for saliency map evaluation of deep learning explanation methods. In *International Conference on Pattern Recognition and Artificial Intelligence*, pages 84–95. Springer, 2022.
- [6] Julie J Gros-Louis, Susan E Perry, Claudia Fichtel, Eva Wikberg, Hannah Gilkenson, Susan Wofsy, and Alex Fuentes. Vocal repertoire of cebus capucinus: acoustic structure, context, and usage. *International Journal of Primatology*, 29(3):641–670, 2008.
- [7] René Heinrich, Lukas Rauch, Bernhard Sick, and Christoph Scholz. Audioprotonet: An interpretable deep learning model for bird sound classification, 2024. URL <https://arxiv.org/abs/2404.10420>.
- [8] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023.
- [9] Holger Klinck, Maggie, Sohier Dane, Stefan Kahl, Tom Denton, and Vijay Ramesh. Birdclef 2024. <https://kaggle.com/competitions/birdclef-2024>, 2024. Kaggle.
- [10] Olivier Le Meur and Thierry Baccino. Methods for comparing scanpaths and saliency maps: strengths and weaknesses. *Behavior research methods*, 45(1):251–266, 2013.
- [11] Jinhua Liang, Inês Nolasco, Burooj Ghani, Huy Phan, Emmanouil Benetos, and Dan Stowell. Mind the domain gap: A systematic analysis on bioacoustic sound event detection. *2024 32nd European Signal Processing Conference (EUSIPCO)*, pages 1257–1261, 2024. URL <https://api.semanticscholar.org/CorpusID:268723684>.
- [12] Joycelyn Longdon. Environmental data justice. *The Lancet Planetary Health*, 4(11):e510–e511, 2020.
- [13] Joycelyn Longdon, Michelle Westerlaken, Alan F Blackwell, Jennifer Gabrys, Benjamin Ossom, Adham Ashton-Butt, and Emmanuel Acheampong. Justice-oriented design listening: Participatory ecoacoustics with a ghanaian forest community. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI ’24, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400703300. doi: 10.1145/3613904.3643044. URL <https://doi.org/10.1145/3613904.3643044>.
- [14] Inês Nolasco, Burooj Ghani, Shubhr Singh, Ester Vidaña-Vila, Helen Whitehead, Emily Grout, Michael G. Emmerson, Frants Havmand Jensen, Ivan Kiskin, Joe Morford, Ariana Strandburg-Peshkin, Lisa F. Gill, Hanna Pamula, Vincent Lostanlen, and Dan Stowell. Few-shot bioacoustic event detection at the dcase 2023 challenge. *ArXiv*, abs/2306.09223, 2023. URL <https://api.semanticscholar.org/CorpusID:260472804>.

- 203 [15] David Robinson, Marius Miron, Masato Hagiwara, Benno Weck, Sara Keen, Milad Alizadeh,  
204 Gagan Narula, Matthieu Geist, and Olivier Pietquin. Naturelm-audio: an audio-language  
205 foundation model for bioacoustics. *arXiv preprint arXiv:2411.07186*, 2024.
- 206 [16] Ramprasaath R. Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi  
207 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-  
208 based localization. *International Journal of Computer Vision*, 128:336 – 359, 2016. URL  
209 <https://api.semanticscholar.org/CorpusID:15019293>.
- 210 [17] Donghee Shin. The effects of explainability and causability on perception, trust, and acceptance:  
211 Implications for explainable ai. *International journal of human-computer studies*, 146:102551,  
212 2021.
- 213 [18] Larissa De Andrade Silva, Juan G Colonna, Bernardo B Gatto, and João Marcelo Protázio.  
214 Impacts of anthropogenic noise on the house wren’s song: An xai approach to bioacoustic  
215 insights. In *2025 IEEE Symposium on Trustworthy, Explainable and Responsible Computational*  
216 *Intelligence (CITREx)*, pages 1–7. IEEE, 2025.
- 217 [19] Bart van Merriënboer, Vincent Dumoulin, Jenny Hamer, Lauren Harrell, Andrea Burns, and  
218 Tom Denton. Perch 2.0: The bittern lesson for bioacoustics. *arXiv preprint arXiv:2508.04665*,  
219 2025.
- 220 [20] Álvaro Vega-Hidalgo, Artem Abzaliev, Thore Bergman, and Rada Mihalcea. Acoustic individ-  
221 ual identification of white-faced capuchin monkeys using joint multi-species embeddings. In  
222 Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors,  
223 *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*  
224 *(Volume 2: Short Papers)*, pages 645–659, Vienna, Austria, July 2025. Association for Com-  
225 putational Linguistics. ISBN 979-8-89176-252-7. doi: 10.18653/v1/2025.acl-short.51. URL  
226 <https://aclanthology.org/2025.acl-short.51/>.
- 227 [21] Carlos Zednik and Hannes Boelsen. Scientific exploration and explainable artificial intelligence.  
228 *Minds and Machines*, 32(1):219–239, 2022.