

Characterizing soundscapes across diverse ecosystems using a universal acoustic feature set

Sarab S. Sethi^{a,b,c,1}[®], Nick S. Jones^a, Ben D. Fulcher^d[®], Lorenzo Picinali^b[®], Dena Jane Clink^e[®], Holger Klinck^e[®], C. David L. Orme^c[®], Peter H. Wrege^e[®], and Robert M. Ewers^c

^aDepartment of Mathematics, Imperial College London, London, SW7 2AZ, United Kingdom; ^bDyson School of Design Engineering, Imperial College London, London, SW7 2AZ, United Kingdom; ^cDepartment of Life Sciences, Imperial College London, Ascot, SL5 7PY, United Kingdom; ^dSchool of Physics, University of Sydney, Sydney, NSW 2006, Australia; and ^eCenter for Conservation Bioacoustics, Cornell Lab of Ornithology, Cornell University, Ithaca, NY 14850

Edited by Simon A. Levin, Princeton University, Princeton, NJ, and approved June 10, 2020 (received for review March 12, 2020)

Natural habitats are being impacted by human pressures at an alarming rate. Monitoring these ecosystem-level changes often requires labor-intensive surveys that are unable to detect rapid or unanticipated environmental changes. Here we have developed a generalizable, data-driven solution to this challenge using ecoacoustic data. We exploited a convolutional neural network to embed soundscapes from a variety of ecosystems into a common acoustic space. In both supervised and unsupervised modes, this allowed us to accurately guantify variation in habitat guality across space and in biodiversity through time. On the scale of seconds, we learned a typical soundscape model that allowed automatic identification of anomalous sounds in playback experiments, providing a potential route for real-time automated detection of irregular environmental behavior including illegal logging and hunting. Our highly generalizable approach, and the common set of features, will enable scientists to unlock previously hidden insights from acoustic data and offers promise as a backbone technology for global collaborative autonomous ecosystem monitoring efforts.

machine learning | acoustic | soundscape | monitoring | ecology

With advances in sensor technology and wireless networks, automated passive monitoring is growing in fields such as healthcare (1), construction (2), surveillance (3), and manufacturing (4) as a scalable route to gain continuous insights into the behavior of complex systems. A particularly salient example of this is in ecology, where, due to accelerating global change (5), we urgently need to track changes in ecosystem health, accurately and in real time, in order to detect and respond to threats (6, 7). Traditional ecological field survey methods are poorly suited to this challenge: they tend to be slow, labor intensive, and narrowly focused and are often susceptible to observer bias (8). Using automated monitoring to provide scalable, rapid, and consistent data on ecosystem health seems an ideal solution (9, 10), yet progress in implementing such solutions has been slow. Existing automated systems tend to retain a narrow biotic or temporal focus and do not transfer well to novel ecosystems or threats (11, 12).

We present an innovative framework for automated ecosystem monitoring using eco-acoustic data (Fig. 1). We used a pretrained general-purpose audio classification convolutional neural net (CNN) (13, 14) to generate acoustic features and discovered that these are powerful ecological indicators that are highly descriptive across spatial, temporal, and ecological scales. We were able to discern acoustic differences among ecosystems, detect spatial variation in habitat quality, and track temporal biodiversity dynamics through days and seasons with accuracies surpassing that possible using conventional hand-crafted ecoacoustic indices. We extended this approach to demonstrate efficient exploration of large monitoring datasets, and the unsupervised detection of anomalous environmental sounds, providing a potential route for real-time automated detection of illegal logging and hunting behavior. Our approach avoids two pitfalls of previous algorithmic assessments of eco-acoustic data (15). First, we do not require supervised machine-learning techniques to detect (16, 17) or identify (18, 19) acoustic events indicating the presence of threats or species. Supervised methods use annotated training datasets to describe target audio exemplars. This approach can yield high accuracy (20), but is narrowly focused on the training datasets used, can be subverted [e.g., in the case of illegal activity detection (21)], requires investment in laborious data annotation, and frequently transfers poorly from one setting to another (22).

Second, we do not depend on specific hand-crafted ecoacoustic indices. Such indices share our approach of aggregating information across a whole audio sample (23)—a soundscape—but differ in their approach of identifying a small number of specific features [e.g., entropy of the audio waveform (24)] rather than a machine-learned, general acoustic fingerprint. Again, these indices can predict key ecological indicators in local contexts (25–27), but they often fail to discriminate even large ecological gradients (28, 29) and behave unpredictably when transferred to new environments (30).

Significance

Human pressures are causing natural ecosystems to change at an unprecedented rate. Understanding these changes is important (e.g., to inform policy decisions), but we are hampered by the slow, labor-intensive nature of traditional ecological surveys. In this study, we show that automated analysis of the sounds of an ecosystem—its soundscape—enables rapid and scalable ecological monitoring. We used a neural network to calculate fingerprints of soundscapes from a variety of ecosystems. From these acoustic fingerprints we could accurately predict habitat quality and biodiversity across multiple scales and automatically identify anomalous sounds such as gunshots and chainsaws. Crucially, our approach generalized well across ecosystems, offering promise as a backbone technology for global monitoring efforts.

The authors declare no competing interest.

This article is a PNAS Direct Submission.

Copyright © 2020 the Author(s). Published by PNAS.

Published under the PNAS license.

Data deposition: Code to reproduce results and figures from this study is available on Zenodo at https://doi.org/10.5281/zenodo.3530203, and the associated data can be found at https://doi.org/10.5281/zenodo.3530206.

¹To whom correspondence may be addressed. Email: s.sethi16@imperial.ac.uk.

This article contains supporting information online at https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2004702117/-/DCSupplemental.

APPLIED MATHEMATICS

Author contributions: S.S.S., N.S.J., B.D.F., L.P., and R.M.E. designed research; S.S.S., N.S.J., B.D.F., L.P., D.J.C., H.K., C.D.L.O., P.H.W., and R.M.E. performed research; S.S.S. analyzed data; and S.S.S., N.S.J., B.D.F., L.P., D.J.C., H.K., C.D.L.O., P.H.W., and R.M.E. wrote the paper.



Fig. 1. A common framework for monitoring ecosystems autonomously using soundscape data. (A) We embed eco-acoustic data in a high-dimensional feature space using a CNN. Remarkably, this common embedding means that we can both (B) draw out ecological insights into ecosystem health across multiple temporal and spatial scales, and (C) effectively identify anomalous sounds in an unsupervised manner.

A lack of transferability is characteristic of approaches that use site-specific calibration or training, where high local accuracy is achieved at the cost of generality (31). Lack of generalizability is a critical failure for monitoring applications, where rapid deployment is essential and the nature of both threats and responses cannot always be known in advance. Threats can be immediate, such as logging or hunting (32), or play out over longer timescales, such as the invasion of a new species (33) or climate change (34), and may drive unpredictable ecological responses (35). The remarkable efficacy of our feature set provides a general solution to these complex methodological challenges. The same acoustic features are highly descriptive across spatial and temporal scales and are capable of reliably detecting anomalous events and behavior across a diverse set of ecosystems.

A Common Feature Embedding Yields Multiscale Ecological Insight

We collected a wide range of acoustic data from the following ecosystems: protected temperate broadleaf forests in both Ithaca, New York, and Abel Tasman National Park, New Zealand; protected lowland rainforests in Sulawesi, Indonesia; protected and logged lowland rainforest in and surrounding Nouabalé-Ndoki National Park, Republic of Congo; and lowland rainforests across a gradient of habitat degradation in Sabah, Malaysia. These five study sites span temperate, tropical, managed, and protected forest ecosystems, allowing us to test the transferability of our approach. In total we analyzed over 2,750 h of audio, collected using a variety of devices including Audio-Moths (36), Tascam recorders, Cornell Lab Swifts, and custom setups using commercial microphones (Materials and Methods). We then embedded each 0.96-s sample of eco-acoustic data in a 128-dimensional feature space using a CNN pretrained on Google's AudioSet dataset (13, 14).

AudioSet is a collection of human-labeled sound clips, organized in an expanding ontology of audio events, which contains over 2 million short audio samples drawn from a wide range of sources appearing on YouTube. Although a small amount of eco-acoustic data is present, the vast majority of audio clips are unrelated to natural soundscapes (13), with the largest classes consisting of music, human speech, and machine noise. No ecological acoustic datasets provide labeled data on a similar magnitude to AudioSet, and when detecting "unknown unknowns" it is in fact desirable to have a feature space that is able to efficiently capture characteristics of nonsoundscape-specific audio. The resulting acoustic features are therefore both very general and of high resolution, placing each audio sample in high-dimensional feature space that is unlikely to show ecosystem-specific bias.

We first investigated whether this feature embedding revealed expected ecological, spatial, and temporal structure across our eco-acoustic datasets. Short audio samples are highly stochastic, so we averaged the learned acoustic features over 5 consecutive minutes. We were able to clearly differentiate eco-acoustic data from different ecosystems (Fig. 2A). Furthermore, samples from the same location clustered strongly, even when different recording techniques and equipment were used, and audio samples from similar ecosystems were more closely located in audio feature space (SI Appendix, Fig. S1). Within sampling locations, the acoustic features captured ecological structure appropriate to the spatial and temporal scale of recordings. Data recorded across a gradient of logging disturbance in Sabah (37) reflected independent assessment of habitat quality based on the quantity of above-ground biomass (AGB), except for sites near rivers where the background sound of water dominated the audio (Fig. 2B). Monthly recordings across 3 y (2016 to 2019) from Ithaca captured eco-acoustic trajectories describing consistent seasonal changes in community composition driven by migratory fluxes of birds (Fig. 2C). Similarly, daily recordings in Sabah strongly discriminated between the dawn and dusk choruses in the tropical rainforest of Malaysia, with large discontinuities at 05:00 and 17:00 h, respectively, that reflected diurnal turnover in the identity of vocalizing species (Fig. 2D). The same acoustic features also revealed diurnal patterns in data from the four other ecosystems used in this study (SI Appendix, Fig. S2). These results show that we are able to capture complex hierarchical structure in ecosystem dynamics using a common eco-acoustic embedding, with no modification required when moving across spatial and temporal scales.

While unsupervised approaches can thus be used to qualitatively visualize and explore ecosystem data in our feature space, a core aim of autonomous monitoring systems is to directly predict ecosystem health and to be able to do so longitudinally over long time periods. We showed that the same general acoustic features (derived from the pretrained CNN) were well suited to this problem by performing a series of classification tasks. Classifications were performed using a random forest classifier in the full feature space, and we compared the performance [measured by F1 score (39)] with a feature space made



Fig. 2. Embedding eco-acoustic data in a common, highly descriptive feature space yields ecological insight across spatial and temporal scales. (*A*) Seven ecoacoustic datasets from five countries are embedded in the same acoustic feature space, in which different ecosystems are distinguished. Features were robust to different recording technologies used in Sabah (Tascam, Audiomoth) and Ithaca (Swift, custom microphone) (*Materials and Methods*). (*B*) Tropical forest areas in Sabah that differ in habitat quality (measured by above-ground biomass, $log_{1o}[t-ha^{-1}]$) cluster in the same acoustic feature space. (C) Three years of soundscape data from a temperate forest in Ithaca reveals a clear seasonal cycle. (*D*) One month of acoustic data from a logged tropical forest site in Sabah shows a repeating diurnal pattern. In all panels, UMAP (38) was used to visualize a 2D embedding from the full 128-dimensional acoustic feature space, and centroids of classes are denoted by larger points.

up from five existing eco-acoustic indices (EAI) often used to assess ecosystem health (*Materials and Methods*). Our approach provided markedly more accurate predictions of biodiversity and habitat quality metrics in both temperate (avian richness; CNN: 88% versus EAI: 59%; Fig. 3A) and tropical (AGB; CNN 94% versus EAI 62%; Fig. 3B) landscapes. Importantly, our predictions of avian richness did not require individual identification of species within the soundscape—a process only possible given vast amounts of manually labeled, species-specific data. General acoustic features also allowed more accurate predictions of temporal variables at both seasonal (months within temperate soundscapes; CNN 86% versus EAI 42%; Fig. 3C) and daily (hours within tropical soundscapes; CNN 68% versus EAI 31%; Fig. 3D) timescales. These results pave the way for automated eco-acoustic monitoring to detect environmental changes over long time scales. For example, the loss of tree biomass from logging over a period of months, annual shifts in the seasonal phenology of bird communities (40), and the gradual increase of forest biomass through decades of forest recovery or restoration (41) may all be accurately tracked through time using this analysis framework.

A Common Feature Space Allows Effective Unsupervised Anomaly Detection and Eco-Acoustic Data Summarization

Given the huge volumes of audio data that are rapidly collected from autonomous monitoring networks, it is important to create automated summaries of these data that highlight the most typical or anomalous sounds at a given site—a task that is not possible given current approaches to eco-acoustic monitoring. In particular, the task of unsupervised anomaly detection is critical



Fig. 3. General acoustic features allow accurate classification of the degree of ecosystem degradation and position in diurnal and seasonal cycles. We performed a multiclass classification task using a 20% test set to assess the predictive power of the general acoustic features on a range of spatial and temporal scales of eco-acoustic data. For each task we measured the F1 score for each of the classes and compared the results using general acoustic features derived from a pretrained CNN (red) to a baseline made up of standard eco-acoustic indices regularly used in eco-acoustics (blue) (*Materials and Methods*). In *A*, we were able to accurately predict a measure of biodiversity (avian richness, species per hour) from a temperate forest site in Ithaca. (*B*) We were also able to predict habitat quality (as measured by above-ground biomass, $\log_{10}[t \cdot ha^{-1}]$) across a landscape degradation gradient in tropical Malaysia with high accuracy, with the exception of sites near rivers. In *C* and *D*, we show how temporal cyclicity on the scale of months and hours, respectively, can be predicted using the same acoustic feature set.

in real-time warning systems which need to automatically warn of unpredictable rapid changes to the environment or illegal activities such as logging and hunting (32). Our solution to both the problems of efficient data summarization and unsupervised anomaly detection involves performing density estimation in our general acoustic feature space.

We developed a site-specific anomaly-scoring algorithm using a Gaussian mixture model (GMM) fit to 5 full days of acoustic features from a given recording location. Here we used the original 0.96-s, 128-dimensional features, which best captured transient acoustic events. We then explored the most typical and anomalous sounds from a logged tropical forest in Sabah, Malaysia, to demonstrate how this approach allows efficient exploration of large amounts of data (Fig. 4A). High-probability, or typical, sounds corresponded to distinct background noise profiles, driven primarily by insect and frog vocalizations, which varied in composition throughout the day, and regular abiotic sounds such as rainfall. Low probability, or anomalous, sounds included sensor malfunctions, anthropogenic sounds (e.g., speech), distinctive species calls that were heard rarely during the recording period (e.g., gibbon trills), or unusually loud events (e.g., a cicada immediately adjacent to the microphone) (Fig. 4A). Exploring the data in this way, we were able to acquire a high-level, rounded summary of a 120-h (432,000-s) period of acoustic monitoring by listening to just 10 s of the most typical sounds and 12 s of anomalies (Audio File S1).

Real-time detection of human activities such as illegal logging and hunting is a particularly pressing problem in protected areas (32). One approach is to train supervised classifiers to search for sounds such as chainsaws (42) or gunshots (43). However, not only do these classifiers require specific training datasets, but also they can easily go out of date or be subverted [e.g., by using a different gun (21)]. We carried out calibrated playback experiments to test the efficacy of our unsupervised density estimation approach for detecting novel acoustic events without prior training. We used a speaker to play sounds including chainsaws, gunshots, chopping, lorries, and speech at distances of 1, 10, 25, 50, and 100 m from an acoustic recorder within the habitat (Fig. 4B) (for logistical reasons we were unable to use real chainsaws, guns, etc.). We then replicated this experiment across 10 sites from the land degradation gradient in Sabah, Malaysia. All sounds were scored as strongly anomalous at 1 m, but differed in how the score declined with distance. Chainsaws, gunshots, and, to a lesser extent, chopping all scored highly at distances of up to 25 m of forest from the recorder, but were not audible over background noise at greater distances (Fig. 4C). In contrast, lorries and speech were reliably detected only within about 10 m of the recorder. Detection ranges in real-world settings will be larger as our playback experiments were unable to fully replicate the sound pressure levels of events such as gunshots (Materials and Methods). The same playback experiment also detected chainsaw and gunshot sounds in a temperate setting in Ithaca with no modification to the algorithm (SI Appendix, Fig. S3), suggesting that this approach to automated anomaly detection is transferable among vastly differing ecosystems.

The Future of Automated Environmental Monitoring

We have shown how state-of-the-art machine-learning techniques can be used to draw out detailed information on the natural environment via its soundscape. Using a common learned feature

nloaded at Cornell University Library on July 8, 2020



Fig. 4. Density estimation in acoustic feature space allows unsupervised detection of anomalous sounds. (A) A projection of the GMM fit to five full days of data from one logged tropical forest site in Sabah, Malaysia. Principal component analysis was used to project the GMM centers and covariances from 128 to 2 dimensions for purposes of visualization, and shaded areas correspond to 2 SDs from each GMM center. Points close to the centers are typical background sounds and thus are given low anomaly scores (i, ii, and iii: ambient noise at different times of day; iv: light rain in a largely silent forest) (Audio File S2). Conversely, very unusual sounds are in low-density regions of acoustic space and are given high anomaly scores (1: human

embedding, derived from a large dataset of nonecosystem audio data, we were able to monitor diverse ecosystems on a wide variety of spatial and temporal scales and to predict biological metrics of ecosystem health with much higher accuracies than was previously possible from eco-acoustic data. Furthermore, we used the same feature-based approach to concisely summarize huge volumes of data and to identify anomalous events occurring in large datasets over long time periods in an unsupervised manner. Our approach offers a bridge from unpredictable handcrafted eco-acoustic indices and highly taxonomically specific detection-classification models to a truly generalizable approach to soundscape analysis. Although in this paper we have focused on monitoring of tropical and temperate forests, future work could employ learned features to analyze eco-acoustic data from grasslands, wetlands, or marine or freshwater ecosystems (23). Additionally, the same approach can easily be generalized to other fields employing acoustic analysis, for example, in healthcare (1), construction (2), surveillance (3), or manufacturing (4). Pairing these new computational methods with networked acoustic recording platforms (44, 45) offers promise as a general framework on which to base larger efforts at standardized, autonomous system monitoring.

Materials and Methods

Audio Data Collection. Audio data were collected from a wide variety of locations using different sampling protocols in this study.

In Sabah, Malaysia, two datasets using different recording devices contained data across an ecological gradient encompassing primary forest, logged forest, cleared forest, and oil palm sites (37) collected between February 2018 and June 2019. In the Tascam dataset, audio was recorded as 20-min sound files at 44.1 kHz using a Tascam DR-05 recorder mounted at chest height on a tripod at 14 sites. One 20-min file was recorded per hour at each site, and a total of 27 h 40 min was recorded. In the Audiomoth dataset, version 1.0.0. devices (36) were used. Audio was recorded continuously in consecutive 5-min sound files at 16 kHz. Audiomoths were secured to trees at chest height across 17 sites (14 overlapping with the Tascam dataset). A total of 748 h of audio was recorded.

Two datasets were recorded from Sapsucker Woods, Ithaca, New York, using the following methodologies. The first dataset was recorded from a single location, continuously over 3 y, between January 2016 and December 2019 (inclusive) using a Gras-41AC precision microphone and audio-digitized through a Barix Instreamer ADC at 48 kHz. A total of 797 h of audio was collected. The second dataset contains 24 h of audio from 13 May 2017 and was recorded using 30 Swift acoustic recorders across an area of 220 acres. Audio was recorded continuously in consecutive 1-h files at 48 kHz, and recorders were attached to trees at eye height. A total of 638 h of audio was recorded.

In New Zealand, audio was recorded using semiautonomous recorders from the New Zealand Department of Conservation from 8 to 20 December 2016. Ten units were deployed in the Abel Tasman National Park, with 5 on the mainland and 5 on Adele Island. Audio was recorded continuously in consecutive 15-min files at 32 kHz. Recorders were attached to trees at eye height. A total of 240 h of audio was recorded.

In Sulawesi, audio was recorded using Swift acoustic recorders with a standard condenser microphone in Tangkoko National Park, a protected

talking; 2: vocalizing gibbon in background; 3: sensor malfunction; 4: loud insect near microphone) (Audio File S2). (B) We used playback experiments to test the sensitivity of the anomaly score to novel acoustic events, illustrated here by chainsaw sounds. Spectrograms are shown for audio recorded from a fixed location when the anomalous audio file (original) was played from a speaker at a variety of distances. Blue and green represent time-frequency patches of low and high volume, respectively, and the black line is the anomaly score for each 0.96 s of audio. (C) We investigated the sensitivity of the algorithm to a variety of anomalous sounds typical of illegal activity (chainsaws, gunshots, chopping, lorries, talking). Anomaly scores were averaged across 10 sites from a logged tropical forest landscape in Sabah and vary with distance of playback. Dashed lines show where averaged anomaly scores entered the top 0.1 and 0.01%, respectively, of all 449,280 0.96-s audio clips that were used to fit the probability density function.

Sethi et al.

lowland tropical forest area. Data were recorded from four recording locations within the park during August 2018. Audio was recorded continuously in consecutive 40-min files at 48 kHz. Recorders were set at 1 m height from ground level. A total of 64 h of data was recorded.

In the Republic of Congo, audio was recorded using Swift acoustic recorders with a standard condenser microphone from 10 sites in and surrounding Nouabalé-Ndoki National Park between December 2017 and July 2018. Audio was recorded continuously in consecutive 24-h files at 8 kHz. Habitat types spanned mixed forest and *Gilbertiodendron* spp. from within a protected area, areas within a 6-y-old logging concession, and within active logging concessions. Recorders were set at 7 to 10 m from ground level, suspended below tree limbs. A total of 238 h 20 min of audio was recorded.

Acoustic Feature Embedding. Each 0.96-s chunk of eco-acoustic audio was first resampled to 16 kHz using a Kaiser window, and a log-scaled Mel-frequency spectrogram was generated (96 temporal frames, 64 frequency bands). Each audio sample was then passed through a CNN from Google's AudioSet project (13, 14) to generate a 128-dimensional embedding of the audio.

The architecture of the particular CNN that we used, VGGish, was based upon Configuration A of the VGG image classification model with 11 weight layers (46). VGGish was trained by Google to perform general-purpose audio classification using a preliminary version of the YouTube-8M dataset (47). Once trained, the final layer was removed from the network, leaving a 128-dimensional acoustic feature embedding as the CNN output. In this study, we used a Tensorflow implementation of VGGish provided at https:// github.com/tensorflow/models/tree/master/research/audioset/vggish.

Data from Nouabalé-Ndoki, Republic of Congo was recorded at 8 kHz, and then up-sampled to 16 kHz to enable its input to the CNN. While many animals produce sounds with fundamental frequencies below the original Nyquist limit of 4 kHz, it should be noted that audio from other datasets contained full spectrum information up to at least 8 kHz when features were calculated.

The CNN that we used takes a Mel-scaled spectrogram of 0.96 s of duration at a Nyquist frequency of 8 kHz as an input. Insects and bats in particular produce sounds reaching well into the ultrasonics (23) which contain important ecological information but will be missed by this embedding, although their presence may be indicated by species vocalizing under the 8 kHz Nyquist limit. Additionally, the features may be biased toward stationary signals occurring over longer durations, as very short acoustic events could be smoothed out by the window size of the CNN. To achieve a similar embedding that includes information from higher frequencies and can receive variable length inputs, one could train a new model. However, to completely retrain the model would require acquiring an extremely large dataset (the YouTube-8M dataset used by Hershey et al. (14) contains over 350,000 h of audio), and therefore a hybrid transfer learning approach would likely be more appropriate.

As a baseline comparison we created a similar embedding using a selection of standard metrics used extensively in the soundscape ecology literature. These were Sueur's α -index (24), temporal entropy (24), spectral entropy (24), Acoustic Diversity Index (48), and Acoustic Complexity Index (49). Each of the above features was computed over 1-s windows of audio and concatenated to create a five-dimensional feature vector. This is referred to as a compound index in standard eco-acoustic studies (25).

For the multiclass classification problems, for prediction of biodiversity, and to create the visualizations in Fig. 2, we averaged acoustic feature vectors over consecutive 5-min periods to account for the high stochasticity of short audio samples.

Dimensionality Reduction. To produce Fig. 2, we used a uniform manifold learning technique (UMAP) (38) to embed the 128-dimensional acoustic features into a two-dimensional (2D) space. For the global comparison (Fig. 2A) there was a large sample size imbalance among the datasets. To ensure that the dimensionality reduction was not biased, we randomly subsampled 27 h 40 min of data from each dataset before running the UMAP algorithm, and then all points were reprojected into 2D based on this embedding.

Multiclass Classification. We performed multiclass classification using a random forest classifier (50) with 100 trees on acoustic features averaged over 5 min. We used a five-fold cross-validation procedure in which data were split into stratified training and test sets using an 80:20 ratio. F1 score was chosen to report classifier accuracy as it integrates information regarding both precision and recall (39). The balanced accuracy of the classifier on the test set was reported as an average F1 score for each class to account for sample-size imbalances among classes.

Quantifying Biodiversity and Habitat Quality. In Ithaca, New York, between 25 February and 31 August 2017 near-continuous recordings were made using Swift recorders across 30 sites through the Sapsucker Woods area at a sample rate of 48 kHz. For each 1-h period of each day during this period, we randomly selected 1 of the 30 sites in which to quantify biodiversity within the audio recording. For the chosen site and hour combination, a 1-h audio clip was manually annotated to identify all avifaunal species vocalizing. Avian richness at each site was taken to be the total number of distinct species detected in the recordings. Finally, values of avian richness were normalized by sampling effort for all sites. Annotations were made using the Raven Pro software (51).

For each of the 17 sites across a logged tropical forest ecosystem in Sabah, Malaysia, we estimated AGB ($\log_{10}[t\cdotha^{-1}]$). Raw AGB values across the landscape were taken from Pfeifer et al.'s estimates based on ground surveys of the same study site (52). Pfeifer et al. (52) identified a number of 25- \times 25-m plots across the SAFE (Stability of Altered Forest Ecosystems) project landscape. Within each plot, tree diameter and height were recorded, and an allometric equation was applied to derive an estimate for AGB at that location. For each of our recording sites we averaged AGB from all plots surveyed by Pfeifer et al. (52) within 1 km of the recorder. This allowed us to gain a broader picture of ecosystem health, as acoustic data integrates information over larger spatial scales than the 25- \times 25-m plots used for the original AGB estimates.

While both avian richness and AGB were derived as numerical variables, we grouped sites into equidistant bins in both cases and treated them as categorical variables for the purposes of the multiclass classification task.

Anomaly Score Definition and Density Estimation. We used a GMM with 10 components and diagonal covariance matrices to fit a probability density function to 5 d of acoustic features from each site (449,280 clips of 0.96 s per site). Acoustic features were calculated at the 0.96-s resolution with no averaging over longer time windows in the full 128-dimensional feature space. We tested for improvements to the method by estimating the probability distribution using the following: 1) additional GMM components, 2) nondiagonal covariance matrices, and 3) a Dirichlet-process Bayesian GMM (53). Each of these modifications delivered only small advantages (with respect to the ability to identify synthetic anomalies) despite considerable increases in computational complexity. Accordingly, here we report the results of a 10-component GMM with diagonal covariance matrices in the 128-feature space.

The anomaly score of each 0.96-s audio clip was defined as the negative log likelihood of its acoustic feature vector, given the probability density function for the site at which the audio was recorded.

We used the GMM as a data exploration tool to pull out the most anomalous and typical sounds over a 5-d period in a logged forest site in Sabah, Malaysia (Fig. 4 A and B). To characterize the most typical sounds of the soundscape, we found the audio clips from the 5-d period which were closest (Euclidean distance) to each of the 10 GMM components in the feature space. To find a small set of distinct anomalous sounds, we first clustered the 50 most anomalous audio clips using affinity propagation clustering (54), which returns a variable number of clusters. Then, from each of the clusters we picked the clip which had the maximum anomaly score as a representative for the final list of anomalies.

In Fig. 4A, we show a 2D representation of a 128-dimensional acoustic feature space in which the GMM-derived probability density function is depicted from a logged tropical forest site in Sabah, Malaysia. Dimensionality reduction was performed by applying principal component analysis (PCA) to the 5 d of 0.96-s audio clips used to fit the GMM. Anomalous points and the centers and covariances of each of the GMM components were projected into 2D using the same embedding, and shaded areas represent two SDs from each of the centers. PCA was used over other nonlinear dimensionality reduction techniques to enable straightforward visualization of the probability density function.

Anomaly Playback Experiments. Three variants from the following five categories of sounds were used for the anomaly playback experiments: chainsaws, gunshots, lorries, chopping, and talking. All sounds were played in WAV format on a Behringer Europort HPA40 Handheld PA System, and the audio files and speaker together were calibrated to the following sound pressure levels (SPL) at 1 m (chainsaws: 110 dB SPL; gunshots: 110 dB SPL; lorries: 90 dB SPL; chopping: 90 dB SPL; talking: 65 dB SPL). All 15 playback sounds were played while holding the speaker at hip height facing an Audiomoth recording device affixed to a tree at chest height. This was repeated at distances of 1, 10, 25, 50, and 100 m.

Real-world SPL levels are higher for chainsaws and gunshots, but we were unable to reproduce sound pressure levels above 110 dB SPL with the speaker used. For this reason, we expect the detection distances of real events to be larger than reported here. For example, conservatively assuming spherical sound absorption, a real gunshot sound (~150 dB SPL at 1 m) will have traveled 100 m by the time it is attenuated to 110 dB SPL (the value used in our playback experiments).

Data and Materials Availability. Code to reproduce results and figures from this study is available on Zenodo at https://doi.org/10.5281/zenodo.3907296 (55), and the associated data can be found at https://doi.org/10.5281/zenodo.3530206 (56).

ACKNOWLEDGMENTS. We thank Till Hoffman for his input in selecting the audio features and the field staff and organizations that enabled the data collection from all our sites: Sabah: Jani Sleutel, Nursyamin Zulkifli, Adi Shabrani, Dr. Henry Bernard, and SAFE Project; Ithaca: Ray Mack, Ben Thomas, and Cornell Lab of Ornithology; Congo: Phael Malonga, Frelcia Bambi, and Elephant Listening Project/Wildlife Conservation Society; New Zealand: Mike Ogle, New Zealand Department of Conservation; and Sulawesi: Indonesian Ministry of Research. Data from Sulawesi were collected under permit number

- 1. M. Patino et al., Accuracy of acoustic respiration rate monitoring in pediatric patients. *Paediatr. Anaesth.* 23, 1166–1173 (2013).
- D. W. Cullington, D. MacNeil, P. Paulson, J. Elliott, Continuous acoustic monitoring of grouted post-tensioned concrete bridges. NDT Int. 34, 95–105 (2001).
- A. Harma, M. F. McKinney, J. Skowronek, "Automatic surveillance of the acoustic activity in our living environment" in 2005 IEEE International Conference on Multimedia and Expo (IEEE, Amsterdam, 2005), p. 4.
- L. E. Atlas, G. D. Bernard, S. B. Narayanan, Applications of time-frequency analysis to signals from manufacturing and machine monitoring sensors. *Proc. IEEE* 84, 1319–1329 (1996).
- P. M. Vitousek, Beyond global warming: Ecology and global change. *Ecology* 75, 1861–1876 (1994).
- 6. D. J. Rapport, What constitutes ecosystem health? Perspect. Biol. Med. 33, 120-132 (1989).
- D. J. Rapport, R. Costanza, A. J. McMichael, Assessing ecosystem health. *Trends Ecol. Evol.* (*Amst.*) 13, 397–402 (1998).
- M. C. Fitzpatrick, E. L. Preisser, A. M. Ellison, J. S. Elkinton, Observer bias and the detection of low-density populations. *Ecol. Appl.* **19**, 1673–1679 (2009).
- 9. S. E. Hampton et al., Big data and the future of ecology. Front. Ecol. Environ. 11, 156–162 (2013).
- 10. P. A. Soranno, D. S. Schimel, Macrosystems ecology: Big data, big ecology. Front. Ecol. Environ. 12, 3 (2014).
- F. Baret, S. Buis, "Estimating canopy characteristics from remote sensing observations: Review of methods and associated problems" in Advances in Land Remote Sensing: System, Modeling, Inversion and Application, S. Liang, Ed. (Springer Netherlands, 2008), pp. 173–201.
- R. Sollmann, A. Mohamed, H. Samejima, A. Wilting, Risky business or simple solution: Relative abundance indices from camera-trapping. *Biol. Conserv.* 159, 405–412 (2013).
- J. F. Gemmeke et al, Audio Set: An ontology and human-labeled dataset for audio events. Google Al. https://ai.google/research/pubs/pub45857. Accessed 1 November 2019.
- S. Hershey et al., "CNN architectures for large-scale audio classification" in 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), (IEEE, New Orleans, LA, 2017), pp. 131–135.
- R. Gibb, E. Browning, P. Glover-Kapfer, K. E. Jones, Emerging opportunities and challenges for passive acoustics in ecological assessment and monitoring. *Methods Ecol. Evol.* **10**, 169–185 (2019).
- C. J. C. Bravo, R. Á. Berríos, T. M. Aide, Species-specific audio detection: A comparison of three template-based detection algorithms using random forests. *PeerJ Comput. Sci.* 3, e113 (2017).
- D. Stowell, M. Wood, Y. Stylianou, H. Glotin, Bird detection in audio: A survey and a challenge. arXiv, 1608.03417 (11 August 2016).
- D. Stowell, E. Benetos, L. F. Gill, On-bird sound recordings: Automatic acoustic recognition of activities and contexts. arXiv, 1612.05489 (16 December 2016).
- M. Towsey, B. Planitz, A. Nantes, J. Wimmer, P. Roe, A toolbox for animal call recognition. *Bioacoustics* 21, 107–125 (2012).
- T. M. Aide et al., Real-time bioacoustics monitoring and automated species identification. PeerJ 1, e103 (2013).
- R. C. Maher, "Acoustical characterization of gunshots" in 2007 IEEE Workshop on Signal Processing Applications for Public Security and Forensics, (IEEE, Washington, DC, 2007), pp. 1–5.
- D. Stowell, T. Petrusková, M. Šálek, P. Linhart, Automatic acoustic identification of individual animals: Improving generalisation across species and recording conditions. arXiv, 1810.09273 (22 October 2018).
- 23. B. C. Pijanowski et al., Soundscape ecology: The science of sound in the landscape. Bioscience 61, 203–216 (2011).
- 24. J. Sueur, S. Pavoine, O. Hamerlynck, S. Duvail, Rapid acoustic survey for biodiversity appraisal. *PLoS One* **3**, e4065 (2008).
- A. Eldridge *et al.*, Sounding out ecoacoustic metrics: Avian species richness is predicted by acoustic indices in temperate but not tropical habitats. *Ecol. Indic.* 95, 939–952 (2018).
- S. Fuller, A. C. Axel, D. Tucker, S. H. Gage, Connecting soundscape to landscape: Which acoustic index best describes landscape configuration? *Ecol. Indic.* 58, 207–215 (2015).
- J. Sueur, "Indices for ecoacoustics" in Sound Analysis and Synthesis with R, J. Sueur, Ed. (Springer International Publishing, 2018), pp. 479–519.

2881/FRP/E5/Dit.KI/VII/2018. Data from the Republic of Congo were collected with permission of the Republic of Congo Ministry of Forestry. This project was supported by funding from the World Wildlife Fund (Biome Health Project, Malaysia data); the Sime Darby Foundation (Stability of Altered Forest Ecosystems Project, Malaysia data); Natural Environmental Research Council (NERC) Grant NE/K007270/1 (to N.S.J.); Engineering and Physical Sciences Research Council (EPSRC) Grant EP/N014529/1 (to N.S.J.); Fulbright Association of Southeast Asian Nations (ASEAN) Research Award for U.S. Scholars (to D.J.C.); Center for Conservation Bioacoustics (D.J.C., H.K., and P.H.W.); Project Janszoon (New Zealand data); and U.S. Fish and Wildlife Service International Conservation Fund (P.H.W.). We also thank Russ Charif, Jay McGowan, Cullen Hanks, Sarah Dzielski, Matt Young, and Randy Little for annotation of the ground truth data from Ithaca. S.S.S. is also supported by the Natural Environmental Research Council through the Science and Solutions for a Changing Planet Doctoral Training Program. This paper represents a contribution to Imperial College London's Grand Challenges in Ecosystems and the Environment initiative.

- C. Mammides, E. Goodale, S. K. Dayananda, L. Kang, J. Chen, Do acoustic indices correlate with bird diversity? Insights from two biodiverse regions in Yunnan Province, south China. *Ecol. Indic.* 82, 470–477 (2017).
- 29. D. Bohnenstiehl, R. Lyon, O. Caretti, S. Ricci, D. Eggleston, Investigating the utility of ecoacoustic metrics in marine soundscapes. J. Ecoacoustics 2, R1156L (2018).
- T. Bradfer-Lawrence et al., Guidelines for the use of acoustic indices in environmental research. Methods Ecol. Evol. 10, 1796–1807 (2019).
- R. K. Heikkinen, M. Marmion, M. Luoto, Does the interpolation accuracy of species distribution models come at the expense of transferability? *Ecography* 35, 276–288 (2012).
- M. C. Gavin, J. N. Solomon, S. G. Blank, Measuring and monitoring illegal use of natural resources. *Conserv. Biol.* 24, 89–100 (2010).
- M. Clavero, E. García-Berthou, Invasive species are a leading cause of animal extinctions. *Trends Ecol. Evol.* 20, 110 (2005).
- G.-R. Walther et al., Ecological responses to recent climate change. Nature 416, 389–395 (2002)
- 35. O. E. Sala et al., Global biodiversity scenarios for the year 2100. Science 287, 1770–1774 (2000).
- A. P. Hill *et al.*, AudioMoth: Evaluation of a smart open acoustic device for monitoring biodiversity and the environment. *Methods Ecol. Evol.* 9, 1199–1211 (2018).
- R. M. Ewers et al., A large-scale forest fragmentation experiment: The stability of altered forest ecosystems project. Philos. Trans. R. Soc. Lond. B Biol. Sci. 366, 3292–3302 (2011).
- L. McInnes, J. Healy, J. Melville, UMAP: Uniform manifold approximation and proiection for dimension reduction. arXiv. 1802.03426 (6 December 2018).
- N. Chinchor, "MUC-4 evaluation metrics" in Proceedings of the 4th Conference on Message Understanding, (Association for Computational Linguistics, 1992), pp. 22–29.
- T. G. Gunnarsson, J. A. Gill, J. Newton, P. M. Potts, W. J. Sutherland, Seasonal matching of habitat quality and fitness in a migratory bird. *Proc. Biol. Sci.* 272, 2319–2323 (2005).
- T. M. Aide, J. K. Zimmerman, L. Herrera, M. Rosario, M. Serrano, Forest recovery in abandoned tropical pastures in Puerto Rico. For. Ecol. Manage. 77, 77–86 (1995).
- J. Papán, M. Jurečka, J. Púchyová, "WSN for forest monitoring to prevent illegal logging" in 2012 Federated Conference on Computer Science and Information Systems (FedCSIS), (IEEE, 2012), pp. 809–812.
- M. Hrabina, M. Sigmund, "Acoustical detection of gunshots" in 2015 25th International Conference Radioelektronika (RADIOELEKTRONIKA), (IEEE, 2015), pp. 150–153.
- S. S. Sethi, R. M. Ewers, N. S. Jones, C. D. L. Orme, L. Picinali, Robust, real-time and autonomous monitoring of ecosystems with an open, low-cost, networked device. *Methods Ecol. Evol.* 9, 2383–2387 (2018).
- S. S. Sethi et al., SAFE Acoustics: An open-source, real-time eco-acoustic monitoring network in the tropical rainforests of Borneo. *Methods Ecol. Evol.*, 10.1111/2041-210X.13438.2041-210X.
- K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition. arXiv, 1409.1556 (10 April 2015).
- S. Abu-El-Haija et al, YouTube-8M: A large-scale video classification benchmark. arXiv, 1609.08675 (27 September 2016).
- L. J. Villanueva-Rivera, B. C. Pijanowski, J. Doucette, B. Pekin, A primer of acoustic analysis for landscape ecologists. *Landsc. Ecol.* 26, 1233–1246 (2011).
- N. Pieretti, A. Farina, D. Morri, A new methodology to infer the singing activity of an avian community: The Acoustic Complexity Index (ACI). *Ecol. Indic.* 11, 868–873 (2011).
 L. Breiman, Random forests. *Mach. Learn.* 45, 5–32 (2001).
- 51. R. Charif, A. Waack, L. Strickman, *Raven Pro 1.4 User's Manual*, (Cornell Lab of Ornithology, Ithaca, NY, 2010).
- 52. M. Pfeifer et al., Deadwood biomass: An underestimated carbon stock in degraded tropical forests? Environ. Res. Lett. 10, 44019 (2015).
- D. M. Blei, M. I. Jordan, Variational inference for Dirichlet process mixtures. *Bayesian Anal.* 1, 121–143 (2006).
- B. J. Frey, D. Dueck, Clustering by passing messages between data points. Science 315, 972–976 (2007).
- S. Sethi, sarabsethi/audioset_soundscape_feats_sethi2019: June 2020 release (Version 1.2). Zenodo. http://doi.org/10.5281/zenodo.3907296. Deposited 25 June 2020.
- S. S. Sethi, Data associated with audioset_soundscape_feats_sethi2019 (Nov 2019 release). Zenodo. http://doi.org/10.5281/zenodo.3530206. Deposited 16 October 2019.

ECOLOGY

Sethi et al.