

HumBugDB: Supplementary Materials

The supplementary materials include:

- Code and data licensing in Section [A](#).
- A code manual with additional discussions for the results of the main paper in Section [B](#).
- Section [C](#) supplies details on the database schema, and provides an explanation for every field of the metadata.
- The datasheet for HumBugDB is given in Section [D](#).

A Licenses

A.1 Code license

MIT License

Copyright (c) 2021 HumBug-Mosquito

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

A.2 Database license

CC-BY-4.0, <https://creativecommons.org/licenses/by/4.0/>

B Code use

B.1 Code access and structure

- The audio recordings and metadata csv are hosted on Zenodo <http://doi.org/10.5281/zenodo.4904800>. under a CC-BY-4.0 license.
- Code (and the metadata csv for completeness) is hosted on <https://github.com/HumBug-Mosquito/HumBugDB> under the MIT license.

The GitHub data directory structure as of commit [50656758594982480f568598874f79c222432e01](https://github.com/HumBug-Mosquito/HumBugDB/commit/50656758594982480f568598874f79c222432e01) is as follows:

```
HumBugDB
├── README.md
├── *requirements.txt
├── notebooks
│   ├── main.ipynb
│   ├── species_classification.ipynb
│   ├── supplement.ipynb
│   └── spec_audio_multispecies.ipynb
├── data
│   ├── metadata
│   │   └── *.csv
│   └── audio
│       └── *.wav
├── lib
│   ├── PyTorch
│   │   ├── __init__.py
│   │   ├── vggish
│   │   ├── ResNetDropoutSource.py
│   │   ├── ResNetSource.py
│   │   ├── runTorch.py
│   │   └── runTorchMultiClass.py
│   ├── Keras
│   │   ├── __init__.py
│   │   ├── config_keras.py
│   │   └── runKeras.py
│   ├── config.py
│   ├── feat_vggish.py
│   ├── feat_util.py
│   ├── evaluate.py
│   └── write_audio.py
├── outputs
│   ├── models
│   │   ├── keras
│   │   └── pytorch
│   ├── features
│   └── plots
```

A README and several requirements are included for installing Keras, PyTorch, and dependencies for the code. The metadata is located in `/data/metadata/` as a csv file.

Extract the audio from Zenodo to the folder `/data/audio/` and launch the Jupyter notebook `main.ipynb` to perform train-test splitting, feature extraction, model training, and evaluation for Task MED: mosquito event detection. The notebook imports from `lib` the necessary files depending on the choice of kernel and PyTorch or Keras. Task MSC: mosquito species classification is addressed in `species_classification.ipynb`. Remaining supplementary material is found in `supplement.ipynb` and `spec_audio_multispecies.ipynb`.

B.2 Code manual

Overview The following documentation has last been verified with the [commit 50656758594982480f568598874f79c222432e01](https://github.com/HumBug-Mosquito/HumBugDB/commit/50656758594982480f568598874f79c222432e01). Future code aims to maintain compatibility where possible. However, please visit the GitHub repository <https://github.com/HumBug-Mosquito/HumBugDB> for the most comprehensive instructions and updates. Latest development code can be found on the devel branch, and the stable version on master. Releases with large binaries (any pre-trained models or features) can be found on <https://github.com/HumBug-Mosquito/HumBugDB/releases>.

Top-level notebook (MED) `main.ipynb` performs data partitioning, feature extraction and segmentation in `get_train_test_from_df()`, model training in `train_model()`, and model evaluation in `get_results()`. The code is configured with `config.py`, where data directories are specified for the data, metadata and outputs, and feature transformation parameters are supplied. Model hyperparameters are given in `config_keras.py` or `config_pytorch.py`. The notebook supports both Keras [Chollet et al., 2015] and PyTorch [Paszke et al., 2019] with a common interface for convenience. In more detail, each top-level function is described as follows:

- `get_train_test_from_df(df_train, df_test_A, df_test_B)` extracts, reshapes, strides, and normalises features for use as tensors, and saves them to `config.dir_out`, if features with that particular configuration do not exist already. This function supports the creation of features for any feature-model combination. If one wishes to extract Feat. A, the function is imported from `feat_vggish.py`. For Feat. B, `get_train_test_from_df()` is imported from `feat_util.py`. The choice of import is specified in the notebook cells. Section B.3 discusses the features in more depth.
The data is split into train and test based on the matches of experiment ID to the audio tracks from the metadata given in `df_train`, `df_test_A`, `df_test_B`. It is important that no test recordings from these experiments are seen during training in advance, as otherwise model performance is overestimated.
- `train_model(X_train, y_train, X_val=None, Y_val=None, model=ResnetDropoutFull())` trains the BNNs on the data supplied (with validation data optional). The assumed input shape is that of the features produced by `get_train_test_from_df()`. The model argument is optional, and can take any model class defined in `runTorch.py`. The model architecture and training strategies may be changed further in `runKeras.py` or `runTorch.py`.
- `get_results(model, X, y, filename, n_samples=1)` evaluates the model object on test data `{X, y}` with the number of MC dropout samples as `n_samples`. If using deterministic networks, leaving the input argument blank will default to a single evaluation. For any option using Feat. A, the output is aggregated over neighbouring windows to produce predictions over the same window size as Feat. B for fair benchmarking. If you wish to use raw windows (e.g. for creating precise start/stop tags), you may modify this behaviour by removing `resize_window()` from `evaluate.py`. Specify the output plot directory in `config.plot_dir`, and the output filename in `filename`.

Species classification notebook (MSC) `species_classification.ipynb` is used to classify mosquito species on a subset of data from the database. The underlying functions for feature creation and model training are shared as much as possible from the same library sources. Some alterations have been performed for support with PyTorch multi-class classification, due to a difference in API for certain loss functions (BCE loss vs XEnt loss of binary vs multi-class problems).

- `get_feat_multispecies(df_all, train_fraction, random_seed)` extracts features. Configuration for feature extraction is given in `config.py`. The fraction of data used for training is given by `train_fraction`, for which we use 0.75 for MSC. The random seeds used are [5, 10, 21, 42, 100]. The outputs are returned in either list or tensor form, depending on whether features require aggregation (Feat A.) or not (Feat B.). Outputs shapes are designed to work without re-shaping for the evaluation function `get_results_multiclass()`.
- `train_model()` follows the same input and output structure as `main.ipynb`. Models are defined and selected from `runTorchMultiClass.py` or `runKeras.py`.

- `get_results_multiclass()` produces outputs of ROC-AUC scores per class with class averages, and confusion matrices in `.pdf` and `.txt` form. As before, the plot directory is specified in `config.py`. All outputs are given over 1.92 second windows.

Supplementary notebook `supplement.ipynb` is used to reproduce the plots of species distribution in this paper (Figure 11) and contains utilities that were used for debugging and visualising the data, should they be helpful for researchers using their own functions.

Spec audio multispecies notebook `spec_audio_multispecies.ipynb` is used to create the plots of Figure 8 and contains utilities for visualising mosquito samples of any species.

B.3 Feature parameters

We first need to define the number of feature windows that are used to represent a sample, $\mathbf{X}_i \in \mathbb{R}^{h \times w}$, where h is the height of the two-dimensional matrix, and w is the width. The longer the window, w , the better potential the network has of learning appropriate dynamics, but the smaller the resulting dataset in number of samples. It may also be more difficult to learn the salient parts of the sample that are responsible for the signal, resulting in a weak labelling problem [Kiskin et al., 2019]. Early mosquito detection efforts have used small windows due to a restriction in dataset size. For example, Fanioudakis et al. [2018] supplies a rich database of audio, however the samples are limited to just under a second. However, despite the mosquito’s simple harmonic structure, its characteristic sound also derives from the temporal variations, as is visible from spectrograms. We suspect this flight behaviour tone is better captured over longer windows, however we encourage researchers to experiment, for example by padding with noise to match the window size of this architecture, or by choosing a smaller window to extract features from.

The features used in the MED and MSC tasks are as follows:

1. **Feat. A:** Features with default configuration from the VGGish [GitHub](#) intended for use with VGGish: $\mathbf{X}_i \in \mathbb{R}^{64 \times 96}$, 64 log-mel spectrogram coefficients using 96 feature frames of 10 ms duration. To compare feature sets fairly, predictions are aggregated over neighbouring windows to create outputs over 1.92 second windows as used in Feat. B.
2. **Feat. B:** Features of previous acoustic mosquito detection work [Kiskin et al., 2021]: 128 log-mel spectrogram coefficients with a reduced time window of 30 (from 40) feature frames and a stride of 5 frames for training. Each frame spans 64 ms, forming a single training example $\mathbf{X}_i \in \mathbb{R}^{128 \times 30}$ with a temporal window of 1.92 s. To create an augmented dataset, we stride the input signal feature window with a step of 5 feature windows (a duration of 320 ms) Note that the training data is segmented by using overlapping strides specified with `config.step_size=5`, whereas the test data is created with no overlap. Samples that do not divide evenly into the window size are discarded (this is a very small number when using such a small step, and we prefer this option over padding with zeros or noise, though alternate solutions are welcome).

Detailed parameterisation is supplied in Table 5.

Table 5: Feature transformation parameters, in samples unless otherwise indicated. Audio processed with `librosa` for Feat. B and its own implementation in VGGish. The size of 1 frame in w is equal to `hop_length`. For the parameterisation of Feat. B this is 64 ms, resulting in an input feature slice of $512/8000 \times 30 = 1.92$ s duration and $h = 128$ height. For Feat. A, the example window width is $1600/16000 \times 96 = 0.96$ s with height $h = 64$ (log-mel coefficients). Feat. A parameters calculated from `mel_features.py` with the default values supplied in `vggish_params.py`.

Method	Sample rate	NFFT	win_size	hop_length	h (n_mels)	w (frames)	Stride
Feat. A	16,000	512	400	160	64	96	160
Feat. B	8,000	2,048	2,048	512	128	30	512

From the results of Section 5, it appears models are able to learn better representations with Feat. B over Feat. A for the purpose of MSC. However, Feat. A achieve better results in the MED task. We

encourage further study and window size experimentation for determining if there is an optimum window size for any given task.

B.4 Baseline models

MozzBNNv2 We give the full model structure in Figure 3. Lambda layers are dropout layers which are placed to perform MC dropout at test-time. This structure bares similarity to VGGish³, which uses 0.96 second log-mel spectrogram patches as inputs, and 11 weight layers (primarily convolutional layers and max-pool layers). Furthermore, this is an incremental improvement over the model used in Kiskin et al. [2021], <https://github.com/HumBug-Mosquito/MozzBNN>. To ensure this model does not have an advantage in benchmark tasks, each model structure was re-trained with the same data as detailed throughout the main text.

PyTorch ResNet-X We modify the final layers for compatibility with our data (in `runTorch.py`). Furthermore, we have augmented the construction blocks `BasicBlock()` and `Bottleneck()`, as well as the overall model construction, to feature dropout layers to act as an approximation for the model posterior at test-time. Dropout is implemented implicitly in `ResNetSource.py`, to not interfere with the behaviour of `model.eval()`, which by default disables dropout layers at test-time, removing the necessary stochastic component. We have pre-defined two configurations as `Resnet50DropoutFull` and `Resnet18DropoutFull`, which are both passed as objects to input arguments of model training or loading code. For further modifications see `runTorch.py` for MED and `runTorchMultiClass.py` for MSC. For ResNet-18, and ResNet-34 the final `self.fc1` layer is of size $[512, N]$, whereas for ResNet-50 the size is $[2048, N]$, where N is the number of classes for the cross-entropy loss function, or 1 if used with the binary cross-entropy loss. A quick way to check the requirement is to print `x.shape()` before the creation of the `fc1` layer.

PyTorch VGGish The source code for this model can be found at <https://github.com/harritaylor/torchvggish>, which is a PyTorch port of <https://github.com/tensorflow/models/tree/master/research/audioset/vggish>. The adaptation of this model class to transfer learning problems is open to interpretation. We opt to use the most straightforward method, which is to connect the output of the *embeddings* layer to a linear layer which then feeds into a sigmoid and binary cross-entropy loss function (in `runTorch.py`, or straight into a categorical cross-entropy loss function in `runTorchMultiClass.py`). We then re-train the network, with the weights pre-trained on AudioSet. Dropout is used in the last layer of this model only, to create a pseudo-BNN which has estimates of uncertainty when sampled at test time.

With native features (Feat. A), no further modifications were made, though we note that normalising the output of the embedding layer (division by 255) before connecting to a final linear layer provided a boost to performance that may be beneficial to the model overall. As described, this model is implemented as class `VGGishDropout(nn.Module)`.

When utilising Feat. B, the size of the output changes and thus a slight tweak to the final layers in the *embeddings* module was required: this involved removing layer (0) in *embeddings* as the dimension of `in_features` changed from 12,288 to 4,096.

We also note, finally, that this network proved troublesome with high (or sometimes default) values of learning rate for the Adam optimiser which we used for the PyTorch training loop. We lowered the learning rate from 0.0015 to 0.0003 which alleviated this problem, but it is worth keeping in mind that the learning rate should be tailored to the type of model and criterion utilised (in `config_pytorch.py`).

Model training To select the loss that is used to define the best performing model, edit `runTorch.py` to make use of `train_acc` (or any other metric as desired) by replacing. Similarly, amend the training epoch loop to change other metrics or properties during training. In `runKeras.py`, supply arguments and any other desired callbacks and model checkpointing strategies to `model.fit()`. For all models in the MED task, the validation accuracy on a random split of the training data has been used to checkpoint the best-performing model.

³<https://github.com/tensorflow/models/tree/master/research/audioset/vggish>

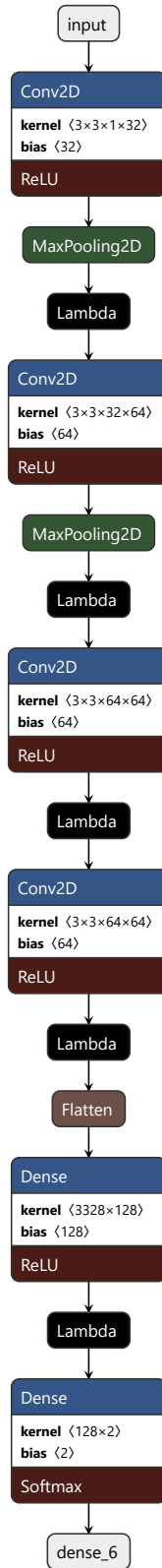


Figure 3: BCNN Keras model. Log-mel spectrograms are input with $w = 30, h = 128$, and passed through the above model. Lambda layers are dropout layers with probability 0.2. Made with <https://github.com/lutzroeder/netron>.

For MSC, no validation sets were used during the training of the models due to the way the data was partitioned and general data scarcity per class. The loss function for PyTorch models was also changed from binary cross-entropy (<https://pytorch.org/docs/stable/generated/torch.nn.BCELoss.html>) to categorical cross-entropy (<https://pytorch.org/docs/stable/generated/torch.nn.CrossEntropyLoss.html>). **Be aware** that BCELoss **does not include** a sigmoid/softmax layer before input to the loss function, which resulted in the creation of separate models for the binary and multi-class classification problems. The models are therefore stored separately in `runTorch.py` and `runTorchMultiClass.py` for the binary (MED) and multi-class (MSC) problems respectively. The Keras model remained unaffected and shares the exact same training code between the MSC and MED tasks.

Hardware The code was developed on Ubuntu 20.04 with an i7-8700K CPU, 32 GB RAM and a Titan Xp GPU with 12 GB VRAM, but models were trained and optimised with lower end hardware (Windows 10, Intel i7-4790K CPU with 16 GB RAM and a GTX970 GPU with 4 GB VRAM).

Memory optimisation Note that the default settings require at least 16 GB RAM to load into memory for ResNet-50 processing, as channels are replicated 3 times to match the pre-trained weights model. To reduce the strain on memory, increase the `step_size` parameter in `config.py` to reduce the number of windows created by feature extraction. This reduces the overlap between samples.

Alternatively, it is possible to use a non-pretrained architecture and change the tensor creation code in `build_data_loader()` from `runTorch.py` to remove `.repeat(1, 3, 1, 1)` as there will be no need to copy over identical data over three channels.

Note that once the tensors have been created, VRAM is not an issue due to the batching over the dataloaders (this code has been run on a GTX970 with 3.5 GB useable VRAM).

A further alternative is creating dataloaders which stream either the raw files and create features on the fly, or stream pre-stored features.

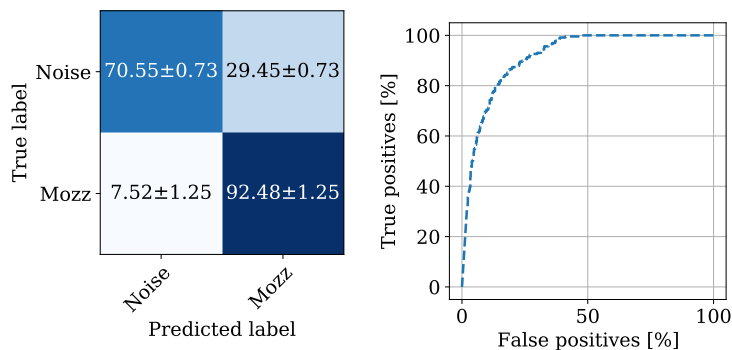
Hyperparameters Configure the hyperparameters in `config_pytorch.py` and `config_keras.py`. The number of epochs was set by observing the learning rate of the network. For MED, within a few epochs, the models began to strongly overfit, with the training accuracy failing to improve validation accuracy. For this reason, both models are set to a low epoch number, and have a fairly low `max_overrun` counter, which determines the maximum number of steps taken for which the target metric fails to improve. The dropout rate and batch size were set to 0.2 and 32, values which are generally risk-free. We note here that the point at which we stop training the model made a fairly significant difference to the balance between true positive and true negative errors (despite a similar overall ROC AUC score). In this respect, the optimisation procedure for the models could be improved with more careful thought about the metrics used for training. If error types are important, consider using loss-calibrated approaches such as that of Cobb et al. [2018].

For MSC, the learning rates are reduced and number of epochs significantly increased. The models all took many more epochs before training stalled. There is also an advantage to using a larger batch size e.g. 128 (and thus higher probability of encountering a greater number of classes per epoch) for faster convergence if desired.

B.5 Test performance

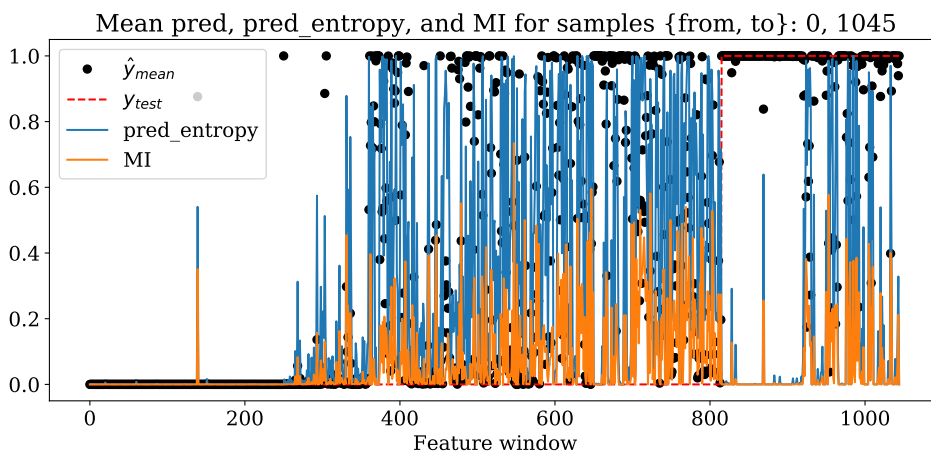
B.5.1 Verifying data integrity of Test B

To support the validity of Test B: the cage recordings conducted at Oxford Zoology, we train the Keras model on half of Test B and test on the other half (with recordings held out), with the settings: `epochs = 7`, `tau = 1.0`, `dropout = 0.2`, `validation_split = None`, `batch_size = 32`, `lengthscale = 0.01` to achieve the results of Figure 4. Figure 4c illustrates the raw output of the mean model probability across 10 MC dropout samples, alongside the predictive entropy and mutual information. The ground truth is given in red, dotted. The model produces a respectable ROC of 0.915, far outperforming the score it achieves when not trained on any part of this dataset (of 0.770). There is therefore a property of this dataset which is not captured well within the rest of the training data (perhaps the SNR or type of background noise encountered), which warrants further research.



(a) Confusion matrix

(b) ROC AUC: 0.915 ± 0.003



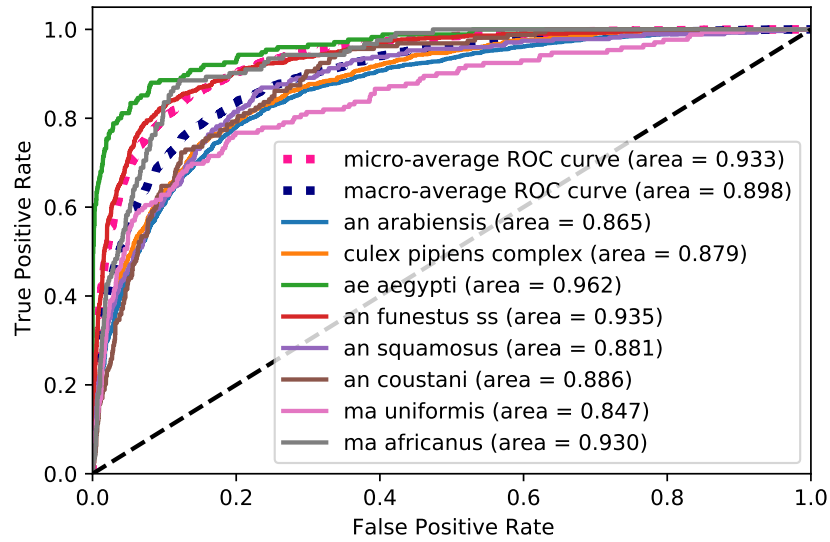
(c) Mean prediction, predictive entropy and mutual information for feature windows

Figure 4: Performance on half of a hold-out test set constructed from Test B. Confusion matrices given in normalised percentage, and ROC in the form of mean \pm standard deviation, across $N = 10$ MC dropout samples.

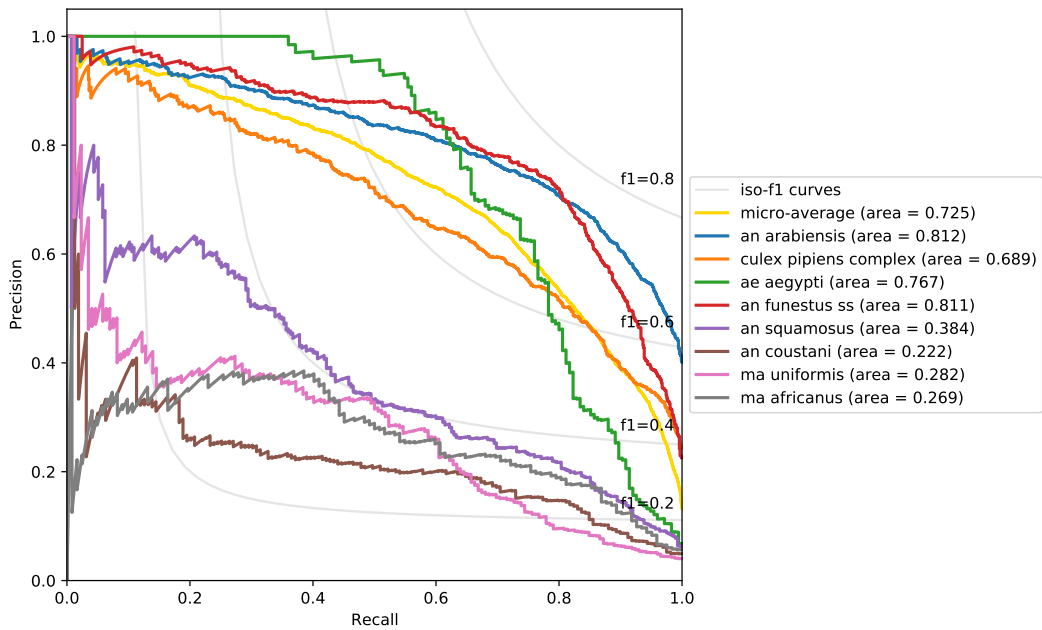
B.5.2 Mosquito Species Classification (MSC)

This section continues the discussion of 5.2. To illustrate typical model performance, we will take a single random seed from the best-performing model, MozzBNNv2, trained on Features B. We show the ROC curves in Figure 5a, and the PR curves in Figure 5b. All of the curves are plotted as a result of derived metrics from the raw softmax outputs averaged over BNN samples in Figure 7. We note overall healthy model outputs, with the probability space well occupied on a scale from 0 to 1. We can also inspect that, generally, predictions are concentrated in the correct regions: the model output per class mimics the true underlying label well, across all classes. This results in healthy ROC curves for all the classes. However, due to the class imbalance encountered, and possibly the difficulty of distinguishing certain species, PR scores leave room for improvement, especially in classes 5, 6 and 7, where the most confusion between species occurs (see confusion matrix of Figure 6). Confusion occurs between species of similar physical characteristics, resulting in similar flight tones. We illustrate a random sample of spectrograms created from audio clips for all species from the IHI Tanzanian cup data in Figure 8. We also supply a [Jupyter tool](#) to visualise and to listen to mosquito samples of any species conveniently (illustrated in Figure 9).

It should be noted that as the mosquitoes are all wild individuals, it is natural that the variation within their species produces some difficulty for the models. Nevertheless, the ROC curves demonstrate that choosing model thresholds in combination with uncertainty estimation from the BNN arm us with the ability to perform species classification.



(a) ROC curves and areas.



(b) Precision-Recall curves, alongside isometric $F1$ score curves. Note that a measure of quality of modelling is given by the PR-AUC areas supplied in the legend. A ‘good’ area is a score that is significantly above the prevalence of the class. For a list of prevalences consult Table 4.

Figure 5: Receiver Operating Characteristic and Precision Recall curves of multi-species performance on IHI Tanzanian cup data of wild mosquito data. Results generated with random seed of 42, with MozzBNNv2 Feat. B. Corresponding confusion matrices and the raw softmax outputs per class are given in Figures 6 and 7 respectively.

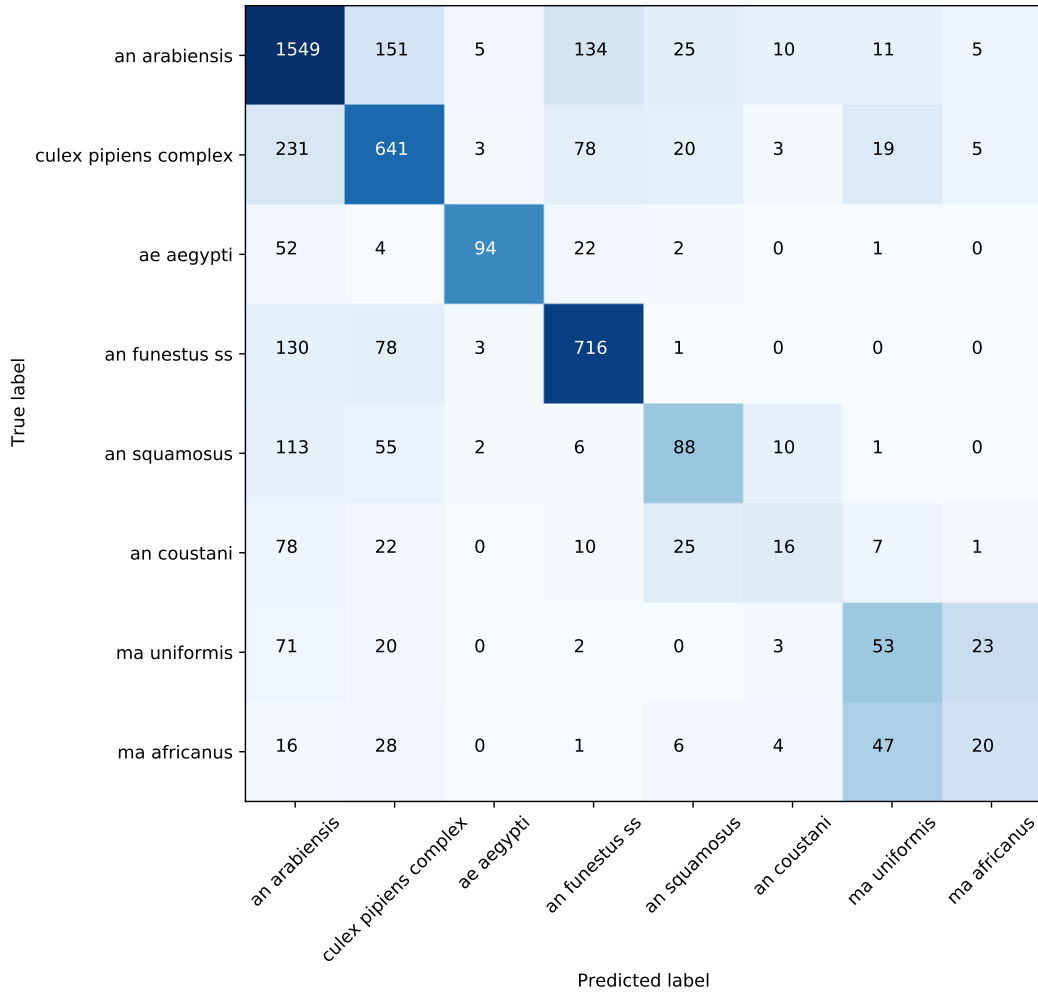


Figure 6: Confusion matrix per species of IHI Tanzanian cup data of wild mosquito data. Results generated with random seed of 42, with MozzBNNv2 Feat. B. Each entry corresponds to a 1.92 second data sample, created with no overlap during feature extraction for the test data. The species of this matrix were arranged in neighbouring classes of similarity to more clearly illustrate where the confusion occurs. Confusion occurs between species of similar physical characteristics.

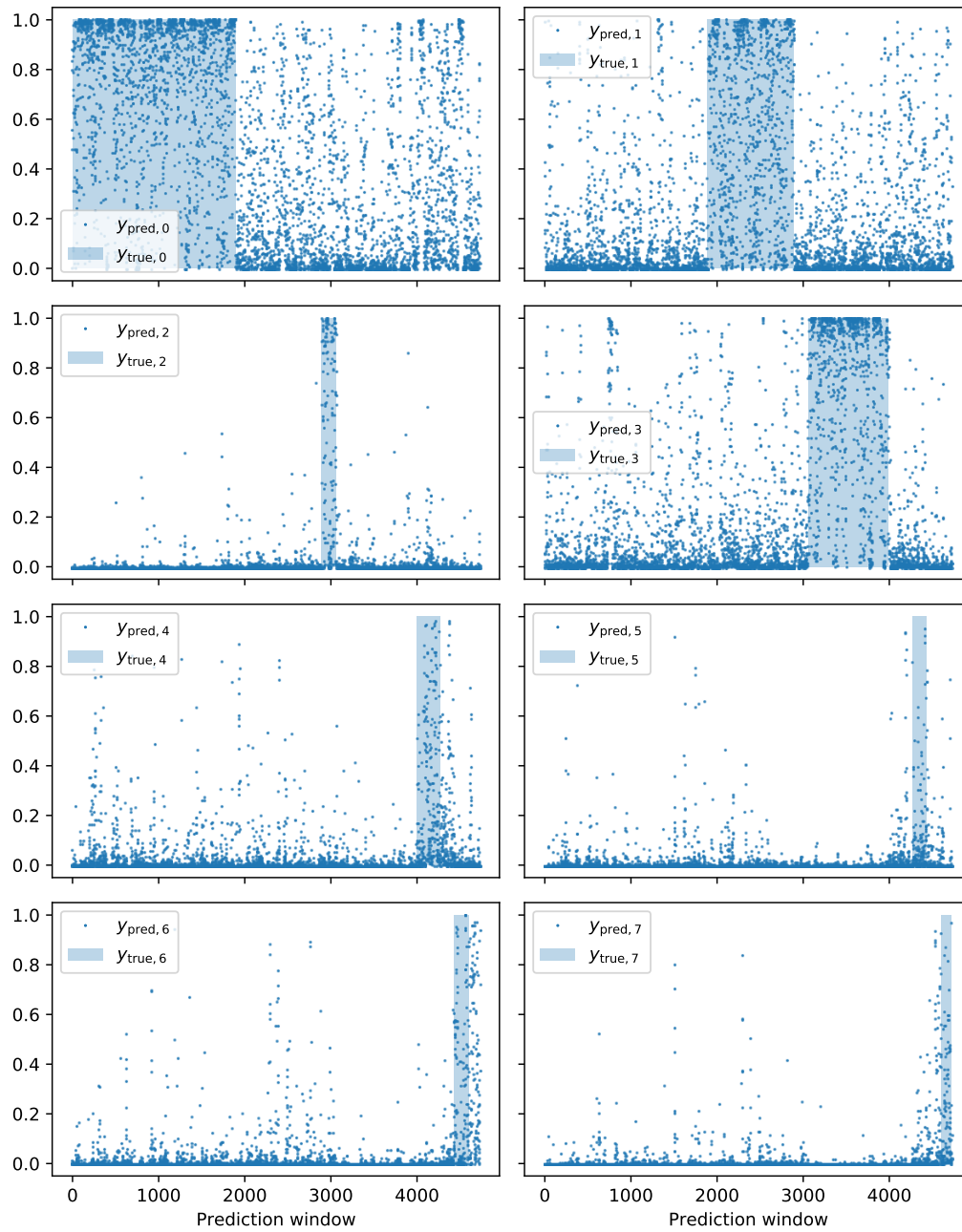


Figure 7: Raw softmax output, averaged across MC dropout samples, per species of IHI Tanzanian cup data of wild mosquito data. The shaded region represents the correct feature windows per class, and the dots are the algorithm predictions. Each plot represents a single class output of the final softmax layer. Results generated with random seed of 42, with MozzBNNv2 Feat. B. Species are numbered in order of classes of the confusion matrix.

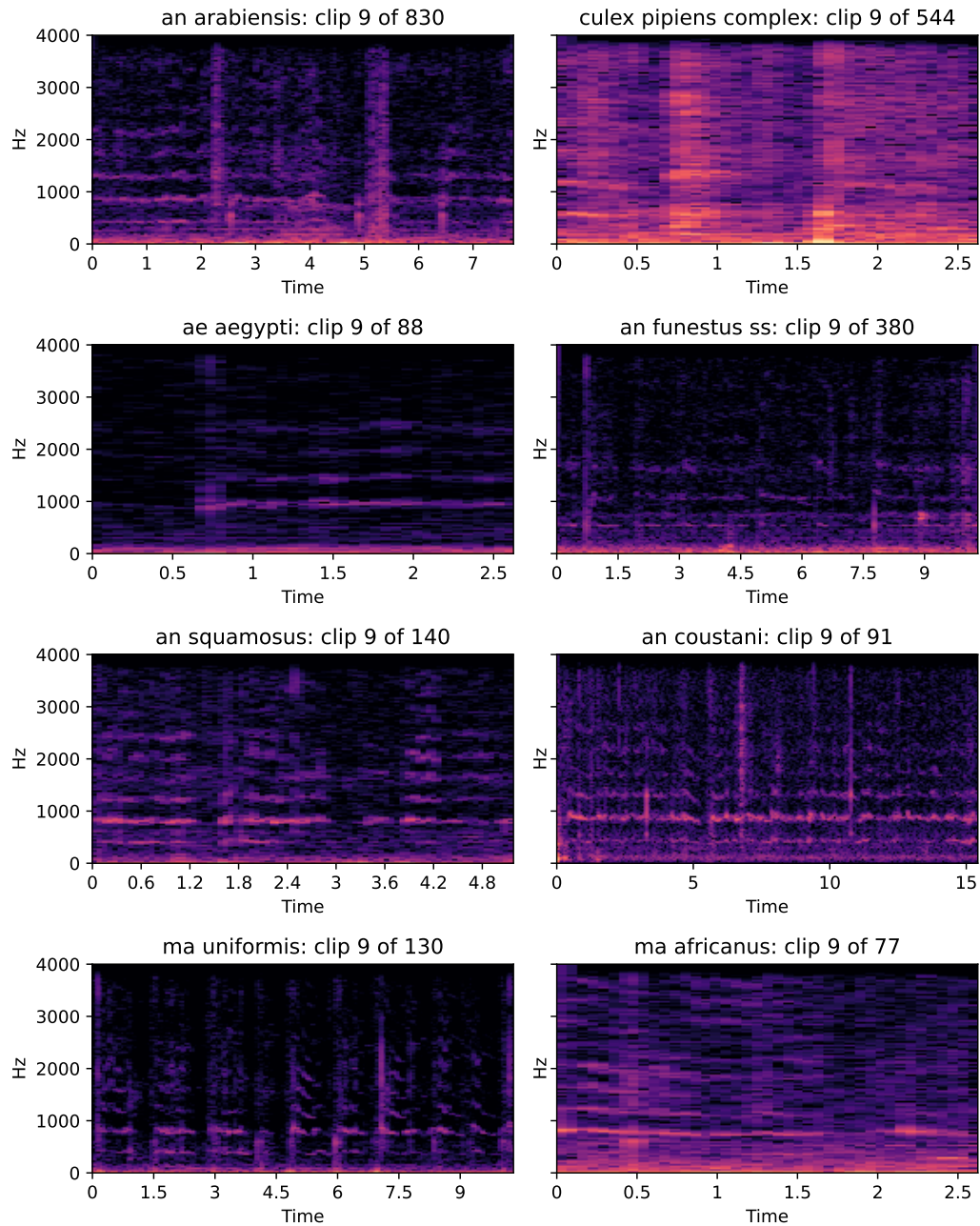


Figure 8: A comparison of randomly chosen clips of mosquito species of the IHI Tanzanian cup recordings. As seen from Figure 6, model confusion occurs commonly between *ma. uniformis* and *ma. africanus* classes. On the other hand, *ae. aegypti* are most easily identified, despite their relatively infrequent occurrence in the dataset. Spectrograms constructed by re-sampling audio to 8 kHz, matching the sample rate of Feat B. and the native sample rate of the smartphone app. In this format, it is also easier to visualise differences between species. For higher sample rate visualisations and audio, consult Figure 9.

Play audio and visualise spectrogram

Let us now visualise the audio with a spectrogram, and also play the corresponding audio clip. Recommended default to re-sample the features to `rate` of 16,000 Hz to more clearly visualise mosquito harmonics and ignore high-frequency noise. The audio will be played re-sampled to match the visual representation. The native sample rate of this audio is 44,100 Hz, retrieved when `rate = None`.

Choose any sample number within the range as illustrated in the cell above. As these data have been tagged by a BCNN with our [open-source utility](#), some sections may contain entomologist speech or other sound/silence. For more information about the model used, please consult our publication [Automatic Acoustic Mosquito Tagging with Bayesian Neural Networks](#).

```
display_audio_spec_mozz(df_all, 8, classes[0], 16000)
```

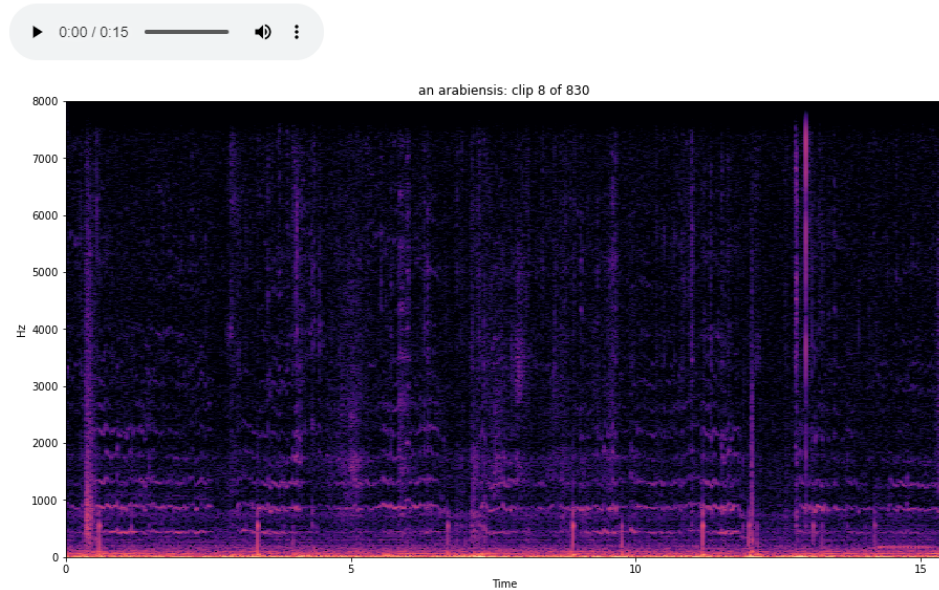


Figure 9: Interactive data visualisation tool on https://github.com/HumBug-Mosquito/HumBugDB/blob/master/notebooks/spec_audio_multispecies.ipynb. The user supplies a sample number to index an audio clip of any mosquito species from a pandas DataFrame. A spectrogram is displayed, alongside an audio playback component.

C PostgreSQL Database

C.1 Database metadata

The data presented in this paper are regularly maintained in a PostgreSQL database. For completeness, we include the full schema in Figure 10. We note that since data upload is a constant work in progress, some fields have not yet been populated sufficiently to be useful upon data extraction. We thus restrict the metadata to the fields that have been verified, and are most likely to be of greatest use. The command we use to extract all the metadata for this paper is as follows:

```
1 \copy (SELECT label.id, fine_end_time-fine_start_time, name,
2       sample_rate, record_datetime, sound_type, species, gender,
3       fed, plurality, age, method, mic_type, device_type, country
4       , district, province, place, location_type
5 FROM label
6 LEFT JOIN mosquito ON (label.mosquito_id = mosquito.id)
7 RIGHT JOIN audio ON (label.audio_id = audio.id)
8 RIGHT JOIN device ON (audio.dev_id = device.id)
9 WHERE type = 'Fine'
10 AND fine_start_time IS NOT NULL AND sound_type in
11 ('mosquito', 'background', 'audio', 'wasp', 'fly') AND
12 (path LIKE '%Kenya%'
13 OR path LIKE '%Thai%'
14 OR path LIKE '%Tanzania%'
15 OR path LIKE '%LSTMH%'
16 OR path LIKE '%CDC%'
17 OR path LIKE '%Culex%'))
18 ORDER BY path) to '/data/export/neurips_2021_zenodo_0_0_1.csv'
19 csv header;
```

We will now break down each metadata field in the data release by the table it originated from and its column heading.

Label:

- `label.id` selects the column `id` from the table `label`, which is joined to `audio` on `label.audio_id = audio.id`. This allows us to now extract a labelled section of audio as indicated by the start and end times of the label.
- `fine_start_time`, `fine_end_time` are the tags for start and end of the audio label, with reference to the original audio recording. Once audio is extracted, we assign the labelled section the filename set to the `label.id`, and define a column `length` which takes the value of `fine_start_time - fine_end_time` for each new label.

Audio:

- `name`: The original filename of the recording (including file extension).
- `sample_rate`: The sample rate of the recording.
- `record_datetime` The time of recording, as SQL `DATETIME` object (easy to parse with either `pandas` or built-in `datetime.datetime`). For newer data, this timestamp is exact, however data collected prior may only be correct to the month.

Mosquito:

- `species` is the species of the mosquito, either the species complex, or more specifically the species if available (e.g. *An. arabiensis* of the complex *An. gambiae s.s.*). If no species information is available, this field is blank (or `NaN` when imported by `pandas` with default settings). A full breakdown of the available species per experimental group is given in Figure 11 and Table 6.
- `gender`: Gender of mosquitoes (M or F) or blank if not known.

- `fed`: Whether mosquito has been fed (t or f) or blank.
- `plurality`: The quantity of mosquitoes recorded at one instant: `single`, `plural` or `blank` if unknown.
- `age`: The age of mosquito in days.
- `sound_type`: denotes whether the label corresponds to a mosquito event if `mosquito`, but can take the value of `background` for corresponding background, `audio` for sections of dense audio events not containing mosquito or wasp and fly. When parsing data, a binary distinction between `mosquito` and `NOT mosquito` can be made safely.

Device:

- `method`: The method of capture of mosquitoes, taking values `HBN`, `LT`, `ABN`, `LC`, `HLC` or `none` if not known (or applicable). Human-baited nets (HBN) are a form of mosquito intervention where humans are surrounded by a mosquito net. As part of the HumBug project, adapted bednets were used where an additional canopy to hold smartphones for recording was sewn on (from 2020 onwards) [Sinka et al., 2021].
Animal-baited nets follow the same concept but involve an animal as the main attractant for mosquitoes.
CDC light traps (LT) use several attractants to lure mosquitoes into the collection chamber. Light is the primary source, but bottled CO₂, gas or dry ice can also be used.
Larval collections, where the eggs of young mosquitoes are collected, are denoted `LC`.
Human landing catches, where mosquitoes that landed on humans are caught, are denoted `HLC`.
For mosquitoes raised from culture and not released into the wild and/or near any nets, this field is `blank`.
- `mic_type`: The microphone used. Takes values `telinga`, `phone` to denote the microphone type. Use this field to filter audio by the type of sound produced, if you wish to check for bias arising from recording device. Further refine the search with the phone model as specified in `device_type`.
- `device_type`: the device to which the microphone was connected. E.g. the field microphone (Telinga) was connected to a Tascam or Olympus recorder. If a smartphone was used, the device is the phone model (e.g. `itel A16` or `Alcatel 4015X`).

Location:

- `location_type`: The environment in which the mosquitoes were recorded in, taking values `cup` for mosquitoes recorded in sample cups, `field` for mosquitoes recorded free-flying in the field (applicable to Tanzania 2020 bednet recordings), or `culture` for mosquitoes recorded in culture cages.
- `country`, `district`, `province`, `place`: The country, district, province, and name of the recording site (e.g. `USA`, `Georgia`, `Atlanta`, `CDC insect culture`, `Atlanta`). Use these values combined with `location_type` to filter data by recording experiment.



Figure 10: Relational tables of the full PostgreSQL database which was used to generate the data for this paper. The structured nature of the database enforces a standard in label format, ensuring we can efficiently mix and match data from a wide range of experiments with differing protocols. For example, if we wish to investigate the effect of mosquito gender or microphone type on the ability to detect mosquitoes, we may sub-select data with the appropriate metadata with one query. Database schema generated with dbdiagram.io from with `pg_dump -s`.

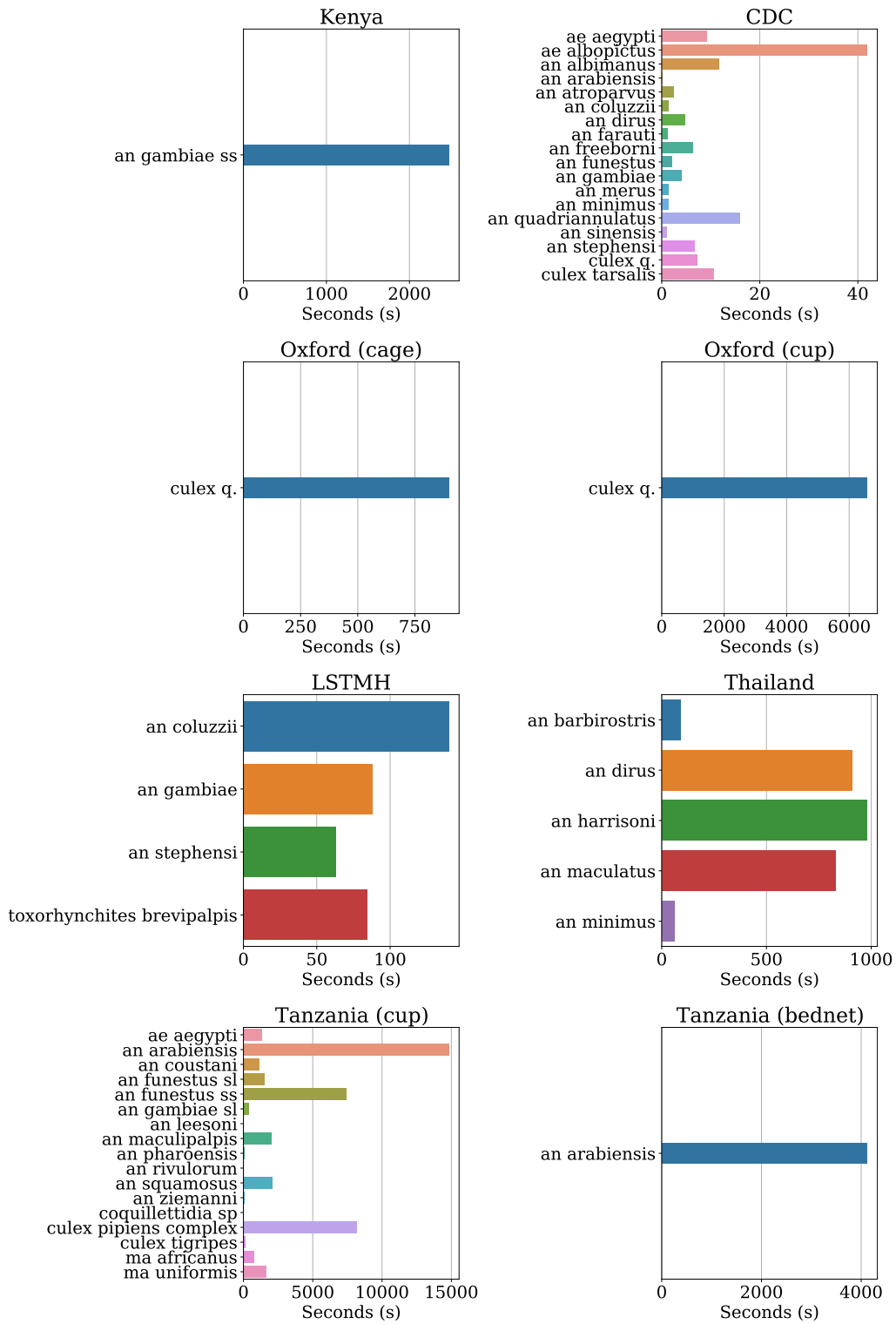


Figure 11: Species distribution per experiment corresponding to Table 7.

Table 6: Species distribution per experimental group corresponding to Table 7 and Figure 11.

Species	Species per location (seconds)								Total
	Kenya	USA (CDC)	Oxford (cup)	Oxford (cage)	LSTMH	Thailand	Tanzania (cup)	Tanzania (bednet)	
<i>Ae. aegypti</i>	0	9.1	0	0	0	0	1322.4	0	1333.6
<i>Ae. albopictus</i>	0	41.9	0	0	0	0	0	0	41.9
<i>An. albimanus</i>	0	11.6	0	0	0	0	0	0	11.6
<i>An. arabiensis</i>	0	0.1	0	0	0	0	14815.2	4118.2	18933.6
<i>An. atroparvus</i>	0	2.3	0	0	0	0	0	0	2.4
<i>An. barbirostris</i>	0	0	0	0	0	87.8	0	0	87.8
<i>An. coluzzii</i>	0	1.3	0	0	140.0	0	0	0	141.3
<i>An. coustani</i>	0	0	0	0	0	0	1140.6	0	1140.6
<i>An. dirus</i>	0	4.7	0	0	0	909.8	0	0	914.5
<i>An. farauti</i>	0	1.1	0	0	0	0	0	0	1.1
<i>An. freeborni</i>	0	6.3	0	0	0	0	0	0	6.2
<i>An. funestus</i>	0	2.1	0	0	0	0	0	0	2.1
<i>An. funestus s.l.</i>	0	0	0	0	0	0	1542.1	0	1542.1
<i>An. funestus s.s.</i>	0	0	0	0	0	0	7414.2	0	7414.2
<i>An. gambiae</i>	0	4.0	0	0	88.2	0	0	0	92.2
<i>An. gambiae s.l.</i>	0	0	0	0	0	0	406.7	0	406.7
<i>An. gambiae s.s.</i>	2474.2	0	0	0	0	0	0	0	2474.2
<i>An. harrisoni</i>	0	0	0	0	0	980.4	0	0	980.4
<i>An. lesoni</i>	0	0	0	0	0	0	43.5	0	43.5
<i>An. maculatus</i>	0	0	0	0	0	830.4	0	0	830.4
<i>An. maculipalpis</i>	0	0	0	0	0	0	2013.0	0	2013.0
<i>An. merus</i>	0	1.4	0	0	0	0	0	0	1.4
<i>An. minimus</i>	0	1.4	0	0	0	61.5	0	0	63.0
<i>An. pharoensis</i>	0	0	0	0	0	0	56.3	0	56.3
<i>An. quadriannulatus</i>	0	15.9	0	0	0	0	0	0	15.9
<i>An. rivulorum</i>	0	0	0	0	0	0	5.1	0	5.1
<i>An. sinensis</i>	0	1.0	0	0	0	0	0	0	1.0
<i>An. squamosus</i>	0	0	0	0	0	0	2091.8	0	2091.8
<i>An. stephensi</i>	0	6.7	0	0	63.1	0	0	0	69.9
<i>An. ziemanni</i>	0	0	0	0	0	0	110.0	0	110.0
<i>Coquillettidia sp.</i>	0	0	0	0	0	0	25.6	0	25.6
<i>Culex pipiens</i>	0	0	0	0	0	0	8157.8	0	8157.8
<i>Culex q.</i>	0	7.3	6573.1	898.1	0	0	0	0	7478.5
<i>Culex tarsalis</i>	0	10.5	0	0	0	0	0	0	10.5
<i>Culex tigripes</i>	0	0	0	0	0	0	158.7	0	158.7
<i>Ma. africanus</i>	0	0	0	0	0	0	785.2	0	785.2
<i>Ma. uniformis</i>	0	0	0	0	0	0	1654.5	0	1654.6
<i>Toxorhynchites br.</i>	0	0	0	0	84.6	0	0	0	84.6

C.2 Miscellaneous commands

To generate the metadata of Table 7, we include a list of commands used to generate one row for completeness.

Count total length of labelled audio for a certain path and sound type:

```
1 SELECT SUM(fine_end_time - fine_start_time)
2 FROM label
3 LEFT JOIN mosquito ON (label.mosquito_id = mosquito.id)
4 RIGHT JOIN audio ON (label.audio_id = audio.id)
5 RIGHT JOIN location ON (audio.loc_id = location.id)
6 WHERE path LIKE '%Thai%' and sound_type='mosquito';
```

Count number of audio files for a certain path and sound type:

```
1 SELECT COUNT (DISTINCT path)
2 FROM label
3 LEFT JOIN mosquito ON (label.mosquito_id = mosquito.id)
4 RIGHT JOIN audio ON (label.audio_id = audio.id)
5 RIGHT JOIN location ON (audio.loc_id = location.id)
6 WHERE path LIKE '%Thai%' and sound_type='mosquito';
```

Return location, device types, and recording methods for dataset:

```
1 SELECT DISTINCT country, location_type, method, mic_type,
2 device_type
3 FROM audio
4 RIGHT JOIN location ON (audio.loc_id = location.id)
5 RIGHT JOIN device ON (audio.dev_id = device.id)
6 WHERE path LIKE '%Tanzania%';
```

D Datasheet for HumBugDB

We follow the structure outlined in Datasheets for Datasets by Gebru et al. [2018]. D.1 gives the motivation for the data. D.2 describes the composition of the data. D.3 describes the collection process. D.4 describes the preprocessing involved in the data curation. D.5 lists past uses, and suggests a range of future uses in depth. D.6 describes potential sources of data bias and relevant mitigation strategies. Database maintenance policies are given in D.7.

D.1 Motivation

For what purpose was the dataset created? This dataset was created for academic research, and applications of machine learning for global health. One such application is the monitoring of deadly mosquito species from their acoustic signature, for which quality training data is required to capture the variation that species may exhibit.

Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)? This dataset was curated by the Machine Learning Research Group of the University of Oxford. Data was collected by the Department of Zoology, University of Oxford, the Centers for Disease Control and Prevention, Atlanta, the United States Army Medical Research Unit in Kenya (USAMRU-K), at the London School of Tropical Medicine and Hygiene, the Dept of Entomology, Kasesart University, Bangkok, and by the Ifakara Health Institute in Tanzania.

Who funded the creation of the dataset? A Google Impact Challenge Award 2014, The Bill and Melinda Gates Foundation (2019–present), available on <https://www.gatesfoundation.org/about/committed-grants/2019/07/opp1209888> (last accessed: June 2021).

D.2 Composition

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? This dataset is a collection of acoustic recordings in wav PCM format. We also supply all the metadata, generated in PostgreSQL to a csv file.

How many instances are there in total (of each type, if appropriate)? 9,295 wav audio files, 1 csv.

Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? The audio files are a sub-sample of complete audio recordings, with the recordings corresponding to one complete label defined with a label ID, extracted from the original audio with the markers `start_time`, `end_time`. We are unable to release the full unlabelled audio due to potential issues with privacy and personally identifiable information. The metadata is a curated sub-sample of all available metadata, where fields which were not sufficiently populated or unverified are excluded.

What data does each instance consist of? Each instance corresponds to a labelled section of audio with the event times originally tagged in the original recording with a `start_time`, `end_time`, either manually by human domain experts, or by machine learning models. The label type is supplied in the metadata.

Is there a label or target associated with each instance? Yes, every recording matches a label.

Is any information missing from individual instances? Though every single sample has a label, some recordings have greater availability of metadata than others; see the metadata csv for details.

Are there recommended data splits (e.g., training, development/validation, testing)? Yes, see Table 7. The splits are carried out to increase the chance of generalisation to recordings conducted in varying conditions. The validation split is part of the challenge of this benchmark, left to the discretion of the users. The test data is automatically split in the supplied code. The two tasks that the data splits encouraged are defined as follows:

- Mosquito Event Detection (MED): distinguishing mosquitoes of any species from their background surroundings, such as other insects, speech, urban, and rural noise.
- Mosquito Species Classification (MSC): the classification of detected mosquitoes into their respective species.

Table 7: Key audio metadata and division into train/test for the tasks of MED: Mosquito Event Detection, and MSC: Mosquito Species Classification. ‘Wild’ mosquitoes captured and placed into paper ‘cups’ or attracted by bait surrounded by ‘bednets’. ‘Culture’ mosquitoes bred specifically for research. Total length (in seconds) of mosquito recordings per group given, with the availability of species meta-information in parentheses. Total length of corresponding non-mosquito recordings, with matching environments, given as ‘Negative’. Full metadata documented in Appendix C.

Tasks: Train/Test	Mosquito origin	Site Country	Method (year)	Device (sample rate)	Mosquito (s) (with species)	Negative (s)
MSC: Train/Test MED: Train	Wild	IHI Tanzania	Cup (2020)	Telinga 44.1 kHz	45,998 45,998	5,600
MED: Train	Wild	Kasetsart Thailand	Cup (2018)	Telinga 44.1 kHz	9,306 2,869	7,896
MED: Train	Culture	OxZoology UK	Cup (2017)	Telinga 44.1 kHz	6,573 6,573	1,817
MED: Train	Culture	LSTMH (UK)	Cup (2018)	Telinga 44.1 kHz	376 376	147
MED: Train	Culture	CDC USA	Cage (2016)	Phone 8 kHz	133 127	1,121
MED: Train	Culture	USAMRU Kenya	Cage (2016)	Phone 8 kHz	2,475 2,475	31,930
MED: Test A	Culture	IHI Tanzania	Bednet (2020)	Phone 8 kHz	4,118 4,118	3,979
MED: Test B	Culture	OxZoology UK	Cage (2016)	Phone 8 kHz	737 737	2,307
Total					71,286 64,843	53,227

Are there any errors, sources of noise, or redundancies in the dataset? To our knowledge there are no redundancies, duplicate files, corrupt files or unintended bugs. Despite comprehensive manual checks, label errors due to human entry and ambiguity in sound type may remain.

Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)? The data is self-contained, generated from a PostgreSQL database which is hosted on University of Oxford servers. The data itself is hosted on Zenodo, and the code is accessible on GitHub.

Does the dataset contain data that might be considered confidential? No, explicit permission was obtained where speech is present.

Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? The audio recordings of mosquitoes may cause distress or discomfort to individuals with medical issues that pertain to mosquito sound.

Does the dataset identify any subpopulations (e.g., by age, gender)? The metadata identifies subpopulations of species complexes by species, and further by gender, age and plurality type (for example, if there was more than one mosquito recorded at a label). Further discriminating factors are described in Appendix C.

Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset? Yes, the speakers may announce the recording ID at the start of a recording, however explicit consent was obtained. It may be possible to trace to the person conducting the experiment indirectly.

D.3 Collection Process

How was the data associated with each instance acquired? The data was collected globally at numerous research facilities. We summarise the data collection efforts below:

- **UK, Kenya, USA:** Recordings from laboratory cultures at the London School of Tropical Medicine and Hygiene (LSTMH), the United States Army Medical Research Unit-Kenya (USAMRU-K); Center for Diseases Control and Prevention (CDC), Atlanta as well as with mosquitoes raised from eggs at the Department of Zoology, University of Oxford. Mosquitoes were recorded by placing a recording device into the culture cages where one or multiple mosquitoes were flying, or by placing individual mosquitoes into large sample cups and holding these close to the recording devices.
- **Tanzania i):** Mosquitoes recorded at Ifakara Health Institute’s semi-field facility (*‘Mosquito City’*) at Kining’ina. The facility houses six chambers containing purpose-built experimental huts, built using traditional methods and representing local housing constructions, with grass roofs, open eaves and brick walls. Four different configurations of the HumBug Net, each with a volunteer sleeping under the net, were set up in four chambers. Budget smartphones were placed in each of the four corners of the HumBug Net. Each night of the study, 200 laboratory cultured *An. arabiensis* were released into each of the four huts and the MozzWear app began recording.
- **Tanzania ii)** A collection and recording project in the Kilombero Valley, Tanzania. HBNs, larval collections and CDC-LTs were used to sample wild mosquitoes and record them in sample cups in the laboratory. *Anopheles gambiae* and *An. funestus* (another highly dangerous mosquito found across sub-Saharan Africa), are also siblings within their respective species complexes. Thus, standard PCR identification techniques [Scott et al., 1993] were used to fully identify mosquitoes from these groups. The Tanzanian sampling has collected 17 different species including: *An. arabiensis* (a member of the *gambiae* complex), *An. coluzzii*, *An. funestus*, *An. pharoensis* (see Appendix C, Figure 11 for a full breakdown).
- **Thailand:** Mosquitoes were sampled using ABNs, HBNs and larval collections over a period of two months during peak mosquito season (May to October 2018). Sampling was conducted in Pu Teuy Village (Sai Yok District, Kanchanaburi Province, Thailand) at a vector monitoring station owned by the Kasetsart University, Bangkok. The mosquito fauna at this site include a number of dominant vector species, including *An. dirus* and *An. minimus* alongside their siblings (*An. baimaii* and *An. harrisoni*) respectively (Appendix C, Figure 11 gives a species histogram for this dataset). Sampling ran from 6 pm to 6 am, as most anopheline vectors prefer to bite during the night. Mosquitoes were collected at night, carefully placed into large sample cups and recorded the following day using the high-spec Telinga field microphone and a budget smartphone.

What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)? A summary of equipment is as follows:

- Smartphone (IteI, Alcatel, and others) audio recording with the MozzWear application. Smartphone devices may have variable sample rates, as denoted by the sample rate column of the metadata. The version of MozzWear used in the curation of this dataset recorded audio in 8,000 Hz mono wave format.
- Telinga EM-23 field microphone, and Tascam, Olympus recording devices recording at 44,100 Hz. The Telinga is a very sensitive, low-noise microphone which was widely adopted in bioacoustic studies.
- Human labelling with Excel.
- Human labelling with Audacity (GNU GPLv2 license).
- Labels produced by a Bayesian convolutional neural network (our own, MIT license, included in paper).
- Voice activity detection and removal with WebRTC (BSD license).
- Python (BSD-style license), MongoDB (Server Side Public License), Django (BSD license), Apache (GPLv3 license), PostgreSQL (BSD/MIT-like license), Unix for databases, HTML dashboards, and post-processing.

Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)? Researchers

from the locations previously mentioned, paid salary from their respective institutions, through the grants disclosed previously.

Over what timeframe was the data collected? 2015 to 2020 (and ongoing).

Were any ethical review processes conducted (e.g., by an institutional review board)? We have obtained the ethics approval from the following committees:

- Oxford Tropical Research Ethics Committee (OxTREC Ref. 548-19) – University of Oxford (UK).
- Ifakara Health Institute (IHI)-IRB – Tanzania
- National Institute for Medical Research – Tanzania
- School of Public Health at the University of Kinshasa (KSPH) – DRC

D.4 Preprocessing/cleaning/labeling

Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? The data underwent rigorous curation, from manual adjustment to labels supplied in text files, to commands in the database to deal with incorrectly entered label times resulting in missing data. To encourage reproducibility and compatibility for future data release, all the label and audio quality control is performed before uploading to the database, and within the dataset itself.

Example of quality control code to check that the label end does not exceed the length (which happens frequently when labels are entered by hand into Audacity with end times longer than the recording and then exported to a text file):

```
1 SELECT path, fine_start_time, fine_end_time, sound_type, length
2 FROM label
3 LEFT JOIN mosquito ON (label.mosquito_id = mosquito.id)
4 RIGHT JOIN audio ON (label.audio_id = audio.id)
5 RIGHT JOIN location ON (audio.loc_id = location.id)
6 WHERE fine_end_time > length;
```

Sources with low estimated label quality were either removed or manually re-labelled and amended in the database.

Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? Yes, all data that may have future utility (and has not been yet used for that purpose) has been released. Unprocessed, and currently unlabelled data is also all stored on the database server, but requires further curation and data entry to the specific data tables before release. We plan to periodically update the database as more data becomes available.

Is the software used to preprocess/clean/label the instances available? The software to do so included Audacity, PostgreSQL, Python, Excel, and is available and well-maintained. We will make use of it in future for future data curation.

D.5 Uses

Has the dataset been used for any tasks already? A subset of this data (recorded in Thailand, Kenya, UK, USA) has been used to train a machine learning model to distinguish and detect a mosquito from its acoustic signature. The model was a 4-layer Bayesian convolutional neural network implemented in Keras. The predictive entropy and mutual information were used to screen predictions over thousands of hours of data. Hand labels were added to correct predictions, and the labels were fed back into the database [Kiskin et al., 2021]. Code for the training and resulting predictive pipeline is available on <https://github.com/HumBug-Mosquito/MozzBNN>.

Other past use cases and publications can be found in related works from the link of the following section. We summarise these here as:

- Bioacoustic classification with wavelet-conditioned neural networks [Kiskin et al., 2017, 2018].
- Cost-sensitive mosquito detection [Li et al., 2017a]
- A case study of species classification with field recordings [Li et al., 2018]
- A release of a subset of this database for crowdsourcing (with baseline mosquito detector model) [Kiskin et al., 2019, 2020]

Is there a repository that links to any or all papers or systems that use the dataset? Yes, the Zenodo data directory <https://zenodo.org/record/4904800> contains all the references to projects, papers and code which are associated with this dataset.

What (other) tasks could the dataset be used for? A list of use cases is not limited to, but may include:

1. **Validating species classification models from the literature.**
2. **Frequency analysis.** Identifying the fundamental and harmonic frequencies of flight tone for a particular species, to improve upon the understanding of bioacoustics literature, and entomological research.
3. **Examining inter-species (or similar) variability.** For example, the effect on the sound of flight as a result of age, gender, or any field supported in the database.

We now expand upon each point:

1. **Validating species classification models from the literature.** As a result of procuring curated data with species meta-information of both wild and lab mosquitoes, this dataset serves as an ideal test-bed to verify the effectiveness of existing species classification approaches. We encourage researchers to validate their models by making use of these data to form their own test sets without re-training their models on any parts of this dataset. Strong species discrimination performance would signify a great opportunity to utilise acoustics as a wide-scale surveillance tool.
It would also be very useful to examine transfer learning approaches, where previous models are re-trained and tested on the suggested splits of the data for either task. If you encounter any issues, or require further information do not hesitate to contact the database maintainers (Appendix D.7).
2. **Frequency analysis.** Earlier works in the literature proposed more hand-crafted approaches to building detection or classification models. These may be especially useful in very low-power embedded devices which require fast and efficient algorithms. These approaches were typically centered around specific harmonic inter-peak ratios (See Kiskin [2020, Sec. 3.2] for an overview of relevant prior work). Frequency analysis may be performed on any parts of this dataset, including on species which are under-represented. In particular, the CDC dataset contains a wide range of unique species which are sparsely labelled, however the labelled sections have very high signal-to-noise ratio. As with previous suggested use cases, we recommend trialling approaches on disjoint sets of experiments (or at the very least individual mosquito recording within an experimental set). Once again, there exists an excellent opportunity to validate models from the literature on their ability to distinguish species on this dataset.
3. **Examining the effect of species variability on their flight tone.** It is well known that mosquitoes exhibit significant variability in their physical (and therefore acoustic) properties within a species. These occur due to a multitude of factors including the age, wingspan, gender. Additionally, confounding factors such as the temperature, humidity, and potentially their fed status, can increase the difficulty in distinguishing individuals within and across species. As we maintain as much metadata as possible, this dataset provides the opportunity to examine such factors. In future releases, temperature and humidity will be added where possible, and this data is expected to be available in an update on the Tanzanian cup recordings which has already good metadata coverage including species, age, gender, fed, method. If you wish to have early access to additional metadata, please contact us and we will make the availability of such metadata a higher priority.

Is there anything about the composition of the dataset or the way it was collected and pre-processed/cleaned/labeled that might impact future uses? No, the dataset is specifically organised in PostgreSQL in a way to be consistent with future release. However, in future more metadata may become available for legacy datasets, and larger subsets may become available upon addition of labels.

Are there tasks for which the dataset should not be used? No.

D.6 Data bias

Data is collected with varying recording paradigms, and is sampling a broad (and not fully understood) population of mosquitoes. This induces inherent biases which may affect an algorithm’s performance when acting as either a mosquito detector or species discriminator. We attempted to capture biases as well as possible with comprehensive metadata coverage, which we encourage users to explore for their own use cases. We discuss potential sources of bias and suggest mitigation strategies as follows:

1. **Nature of mosquitoes: lab or wild.** These are denoted by `location_type`. The controlled conditions of laboratory cultures produce uniformly sized fully-developed adult mosquitoes which are used for a variety of purposes, including trialling new insecticides or examining the genome of these insects. Models trained on purely lab cultures run the risk of overfitting to this artificial subpopulation, encountering difficulty when generalising even to the same species. Wild mosquitoes on the other hand exhibit greater variation, at the cost of a much more laborious collection procedure. When constructing models, it is advised to train on wild data, but caution needs to be taken when testing on mosquitoes from uniform subpopulations.
2. **Location bias.** Distributions of mosquitoes may vary across regions, which would be an interesting avenue for further work. This database is expected to form training data for an upcoming challenge, where a hidden test set will consist of recordings conducted at the DRC, carried out with as similar recording methodology to Tanzanian cup data used in MSC. At that stage it will be possible to consider in more depth whether the location is a factor causing significant bias and therefore difficulty for the models.
3. **Recording device** corresponding to metadata from `device`, It is crucial that datasets are not constructed in a way where one device is used for only positive or negative instances (e.g. all noise is from one device, and all mosquito from another). If trained in such a manner, it will be easy for a high-dimensional model such as a neural network to learn the characteristics of the microphone response and use this confounding factor for classification. To mitigate this, we have included a negative control group for each experiment, and therefore also device. This issue becomes especially critical for species classification, where different species may be captured with different devices. Careful consideration and construction of data with the use of the device metadata will help avoid, or at the very least alert to possible confounding. If it is not possible to control the device, it may be desirable to use features which are (more) invariant to microphone type, e.g. MFCCs or high-level pre-trained feature representations such as VGGish embeddings [Gemmeke et al., 2017].
4. **Data imbalance.** Biased models for either species classification or mosquito detection may arise when trained naively without balancing distributions of species, or positive and negative samples. In the case of mosquito detection, a predominance of one species will likely increase the model’s ability to detect mosquitoes of that particular species, while performing worse on less well-represented groups. This is a potential source of improvement worth investigating, especially for the data split suggested in Table 7. Additionally, a closer look at species-specific performance may reveal areas for further model improvement. We recommend benchmarking against the baselines supplied to investigate areas of improvement.

For our MSC tasks we weight class samples by the inverse of their frequency. There are however well-known drawbacks to this. Bayesian models which take into account asymmetrical cost functions aim to alleviate this problem [Cobb et al., 2018]. A further option is to use different step functions in the data partitioning/augmentation pipelines. A starting point would be to modify `step_size` in `feat_util.py` in a class-specific function, to artificially balance the relative frequency of data samples of desired classes. We

have found oversampling or undersampling to produce worse per-class ROC than inverse weighting, however there is room for improvement in future work.

D.7 Database maintenance

Who is supporting/hosting/maintaining the dataset? Please contact Dr. Ivan Kiskin at ivankiskin1@gmail.com, who is maintaining the dataset. Alternative contacts include Professor Steve Roberts at sjrob@robots.ox.ac.uk at the University of Oxford Machine Learning Research Group.

How will the dataset be updated and what is the maintenance policy? The data will be updated as new data and/or metadata from new trials is obtained and curated. We expect the following updates:

1. Recordings of wild captured mosquitoes in DRC:
 - **Date:** Q2/Q3 2021
 - **Summary:** 15 species, at minimum 2000 wild captured individual mosquitoes
 - **Metadata:** species (with PCR identification where appropriate), gender, fed status, temperature, humidity, collection method, recording device information, time of collection
2. Additional metadata for IHI Tanzanian cup recordings:
 - **Date:** Q2/Q3 2021
 - **Summary:** Additional metadata
 - **Metadata:** Temperature, humidity, wing span images (and wing lengths)

To ensure the documentation is up to date, any additional metadata and code will be documented in this supplement and the datasheet for datasets. The supplement is available on GitHub alongside the baseline code at <https://github.com/HumBug-Mosquito/HumBugDB/tree/master/docs>. Database revisions will be incremented in the format of X.X.X (current version 0.0.1). The main url, <https://zenodo.org/record/4904800>, will always resolve to the latest version. If you intend to use a specific version you may select the version from the main page. Both the data and metadata are fully supported with versioning.

Updates will be communicated as follows:

- GitHub commits (mailing list), and releases.
- Posts on the HumBug official twitter account <https://twitter.com/oxhumbug>.
- Updates on our official website on <https://humbug.ox.ac.uk/news/>.
- Follow-up publications utilising additional data (arXiv and proceedings where appropriate).

If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. If you would like to contribute to this data, please contact the database host and supervising professor. We would be happy to curate data and provide requirements which would help qualify a dataset for hosting. All contributions will be credited appropriately in future work.