

---

# HumBugDB: a large-scale acoustic mosquito dataset

---

**Ivan Kiskin\***  
University of Oxford

**Marianne Sinka<sup>†</sup>**  
University of Oxford

**Adam D. Cobb<sup>||</sup>**  
SRI International

**Waqas Rafique\***  
University of Oxford

**Lawrence Wang\***  
University of Oxford

**Davide Zilli<sup>¶</sup>**  
Mind Foundry Ltd

**Ben Gutteridge<sup>§</sup>**  
University of Oxford

**Rinita Dam<sup>‡</sup>**  
University of Oxford

**Theodoros Marinos<sup>††</sup>**  
University of Surrey

**Yunpeng Li<sup>††</sup>**  
University of Surrey

**Dickson Msaky<sup>‡</sup>**  
IHI Tanzania

**Emmanuel Kaindoa<sup>‡</sup>**  
IHI Tanzania

**Gerard Killeen<sup>\*\*</sup>**  
UCC, BEES

**Kathy Willis<sup>†</sup>**  
University of Oxford

**Steve Roberts\***  
University of Oxford

\*Dept. Eng. Science: {ikiskin, waqas, sjrob}@robots.ox.ac.uk,  
lawrence.wang@eng.ox.ac.uk, <sup>†</sup>Dept. Zoology: {marianne.sinka,  
kathy.willis,rinita.dam}@zoo.ox.ac.uk, <sup>||</sup>adam.cobb@sri.com,  
<sup>††</sup>{tm00591, yunpeng.li}@surrey.ac.uk, <sup>\*\*</sup>gerard.killeen@ucc.ie,  
<sup>¶</sup>davide.zilli@mindfoundry.ai, <sup>§</sup>benjamin.gutteridge@new.ox.ac.uk  
<sup>‡</sup>Ifakara Health Institute: {dmsaky, ekaindoa}@ihi.or.tz.

## Abstract

1 This paper presents the first large-scale multi-species dataset of acoustic recordings  
2 of mosquitoes tracked continuously in free flight. Mosquitoes are well-known  
3 carriers of diseases such as malaria, dengue and yellow fever. The motivation  
4 for collecting such a large dataset comes from the need to gather information,  
5 help predict outbreaks, and inform data-driven policy. The task of detecting  
6 mosquitoes from their wingbeats is made challenging due to the difficulty in  
7 collecting recordings from realistic scenarios. To address this, as part of the  
8 HumBug project, we have conducted global experiments to record mosquitoes  
9 ranging from those bred indoors in culture cages to mosquitoes captured in the wild.  
10 As a result, the audio recordings vary widely in signal-to-noise ratio and contain  
11 a broad range of indoor and outdoor background environments from Tanzania,  
12 Thailand, Kenya, the USA and the UK. The audio recordings have been labelled  
13 by domain experts, aided by Bayesian neural networks. As a result, we present 20  
14 hours of mosquito audio recordings expertly labelled with tags precise in time, of  
15 which 18 hours are annotated from 36 different species. We provide our data from  
16 a regularly maintained database, which captures important metadata such as the  
17 capture method, age, feeding status and gender of the mosquitoes. Additionally, we  
18 provide code to extract features and train Bayesian convolutional neural networks  
19 that can distinguish mosquito sounds from their corresponding background. Our  
20 contribution is to provide a dataset that is both challenging to machine learning  
21 researchers focusing on acoustic identification, and critical to entomologists, geo-  
22 spatial modellers and other domain experts to understand mosquito behaviour,  
23 model their distribution, and manage the threat they pose to humans.

## 24 1 Introduction

25 There are over 100 genera of mosquito in the world containing over 3,500 species and they are found  
26 on every continent except Antarctica [Harbach, 2013]. Only one genus (*Anopheles*) contains species  
27 capable of transmitting the parasites responsible for human malaria. *Anopheles* contain over 475  
28 formally recognised species, of which approximately 75 are vectors of human malaria, and around 40  
29 are considered truly dangerous [Sinka et al., 2012]. These 40 species are inadvertently responsible  
30 for more human deaths than any other creature. In 2019, for example, malaria caused around 229  
31 million cases of disease across more than 100 countries resulting in an estimated 409,000 deaths  
32 [World Health Organization, 2020]. It is imperative therefore to accurately locate and identify the  
33 few dangerous mosquito species amongst the many benign ones to achieve efficient mosquito control.  
34 Mosquito surveys are used to establish vector species’ composition and abundance, human biting  
35 rates and thus the potential to transmit a pathogen. Traditional survey methods, such as human  
36 landing catches, which collect mosquitoes as they land on the exposed skin of a collector, can be  
37 time consuming, expensive, and are limited in the number of sites they can survey. They can also be  
38 subject to collector bias, either due to variability in the skill or experience of the collector, or in their  
39 inherent attractiveness to local mosquito fauna. These surveys can also expose collectors to disease.  
40 Moreover, once the mosquitoes are collected, the specimens still need to undergo post sampling  
41 processing for accurate species identification. Consequently, an affordable automated survey method  
42 that detects, identifies and counts mosquitoes could generate unprecedented levels of high-quality  
43 occurrence and abundance data over spatial and temporal scales currently difficult to achieve. It is  
44 for this reason that we utilise low-cost smartphones as acoustic mosquito sensors to solve this task.  
45 The exponential increase in smartphone ownership is a worldwide phenomenon. Governments and  
46 independent companies are continuing to extend connectivity across the African continent [Friederici  
47 et al., 2017]. More than half of sub-Saharan Africa is expected to be connected to a mobile service by  
48 2025 [GSMA, 2020]. With this expanding coverage of mobile phone networks across Africa, there is  
49 an emerging opportunity to collect huge datasets, as exemplified by the World’s Bank Listening to  
50 Africa Initiative [World Bank Organisation, 2017]. Our target application (Section 3.1) uses a free  
51 downloadable app, which means that every smartphone can be a mosquito monitor.

52 **Our contribution** In order to assist research in methods utilising the acoustic properties of  
53 mosquitoes, as part of the HumBug project (described in Section 3.1) we contribute:

- 54 • **Data:** <http://doi.org/10.5281/zenodo.4904800>: A vast database of 20 hours of  
55 finely labelled mosquito sounds, and 15 hours of associated non-mosquito control data,  
56 constructed from carefully defined recording paradigms. Data was collected over the course  
57 of five years in a global collaboration with mosquito entomologists. Recordings were  
58 captured from 36 species (or species complexes<sup>1</sup>) with a mix of low-cost smartphones  
59 and professional-grade recording devices, to capture both the most accurate noise-free  
60 representation, as well as the sound that is likely to be recorded in areas most in need. A  
61 diverse range of wild and lab culture mosquitoes is included to capture the biodiversity of  
62 naturally occurring species. Our data is stored and maintained in a PostgreSQL database,  
63 ensuring label correctness and data integrity. We export all of the audio across a vast range  
64 of experiments with a single line in Python, and the metadata we require for experiments  
65 with a single SQL query (Appendix C). This allows us to add to our database and re-release  
66 data in a reliable and efficient manner.
- 67 • **Code:** <https://github.com/HumBug-Mosquito/HumBugDB>: Detailed tutorial code for  
68 training state-of-the-art baseline Bayesian neural network models (a range of ResNet and  
69 deep CNN models) for the task of distinguishing mosquitoes of any species from their  
70 background surroundings, such as other insects, speech, urban, and rural noise. This baseline  
71 model was used to automatically tag a subset of mosquito recordings in this database with  
72 a very low false positive rate, by making use of uncertainty metrics such as the predictive  
73 entropy and mutual information [Kiskin et al., 2021].
- 74 • To ensure learnt models are tested on diverse and realistic data splits, we withheld two  
75 test sets: one which captures free-flying mosquitoes around specifically adapted bednets

---

<sup>1</sup>Species complexes are closely related sibling species that are morphologically identical but can have hugely diverse behaviours that allows one to be a prominent and dangerous vector, and another to be harmless.

(mimicking the intended target application as closely as possible), and another which contains caged mosquitoes recorded in free flight in very challenging noisy conditions.

The rest of the paper is structured as follows. Section 2 details related datasets and describes how ours contributes to the literature uniquely. Section 3 shows the primary intended use case for the data and model released in this paper for our overall aims to assist in the eradication of insect-borne diseases. Section 4 describes in detail the sources and collection methods of data present, as well as how and why we perform our train-test split. Section 5 suggests additional use cases for the data, and details the steps taken to train a benchmark model, including an overview of feature extraction, model training and evaluation code. We discuss the results that our models achieve, and the open challenges remaining that our test sets motivate. We conclude by summarising our contribution to various communities in Section 6.

We provide comprehensive instructions for using our baseline models and feature extraction code in Appendix B, and supply additional details on all the metadata in Appendix C. The datasheet (Appendix D) details the dataset’s composition (D.2), the data acquisition process (D.3), preprocessing (D.4), past and suggested use cases (D.5), sources of data bias and mitigation strategies (D.6), and database maintenance policies (D.7).

## 2 Related work

Mosquitoes have particularly short, truncated wings allowing them to flap their wings faster than any other insect of equivalent size – up to 1,000 beats per second [Simões et al., 2016, Bomphrey et al., 2017]. This produces their very distinct flight tone and has led many researchers to try and use their sound to attract, trap or kill them [Perevozkin and Bondarchuk, 2015, Johnson and Ritchie, 2016, Jakhete et al., 2017, Fanioudakis et al., 2018, Mukundarajan et al., 2017]. However, there have been very few large datasets released to the public to aid this research. We summarise key statistics of a range of datasets available publicly in Table 1, and discuss the varying sensor modalities separately due to their inherent differences in acoustic properties.

Table 1: A comparison of related mosquito acoustic and pseudo-acoustic datasets released publicly. The ‘Average mosquito length’ is the approximate length of audible mosquito recording per sample. This length can not be estimated for Mukundarajan et al. [2017], as the data is crowdsourced, unlabelled and uncured. Crowdsourced data recording or labels are marked with (\*). ‘Type’ format: majority, (minority), represents if the mosquitoes have been captured as individuals in the wild, or grown and reproduced in controlled conditions in lab colonies. Where not known, ‘Mosquito’ is estimated from the mosquito average mosquito sample duration multiplied by the number of positive samples in dataset.

Dataset	Sensor	Mosquito (Background)	Average mosquito length	Species	Type
Chen et al. [2014, UCR]	Opto-acoustic	17 min (N/A)	$\approx 0.02$ s	6	Lab
Fanioudakis et al. [2018]	Opto-acoustic	39 hr (N/A)	$\approx 0.5$ s	6	Lab
Vasconcelos et al. [2020]	Acoustic	15 min (N/A)	0.3 s	3	Lab
Mukundarajan et al. [2017] (*)	Acoustic	N/A (N/A)	N/A	20	Lab, (wild)
Kiskin et al. [2019, 2020] (*)	Acoustic	2 hr (20 hr)	1 s	N/A	Lab, (wild)
<b>HumBugDB</b>	Acoustic	20 hr (15 hr)	9.7 s	36	Wild, (lab)

**Opto-acoustic approaches** ‘Wingbeats’ [Fanioudakis et al., 2018] and ‘UCR Flying Insect Classification’ [Chen et al., 2014] are high-SNR pseudo-acoustic datasets collected via optical sensors. We note this is a different, but complementary, approach. Due to the directionality of the recording

method, typical sample durations are encountered from “only a few hundredths of a second” [Chen et al., 2014] to approximately half a second [Fanioudakis et al., 2018]. The approach therefore does not capture the acoustical properties of mosquito sound in free flight which aid mosquito detection in purely acoustic approaches [Vasconcelos et al., 2020]. Furthermore, these datasets survey lab-grown mosquito colonies which do not capture the biodiversity of mosquitoes encountered in the wild [Huh et al., 2007, Hoffmann and Ross, 2018].

**Acoustic approaches** The authors of a recent acoustic mosquito dataset [Vasconcelos et al., 2020] motivated its release by stating that none of the published datasets include environmental noise, which is essential to fully characterise mosquitoes in real-world scenarios. Their dataset consists of 300 ms snippets, amounting to a total of 15 minutes of mosquito recordings. This is an excellent first step. However, for deep learning algorithms the dataset is not readily useable due to its size. Moreover, state-of-the-art models for acoustic classification use training example sizes of at least 0.96 seconds for a variety of audio event detection tasks [Hershey et al., 2017] and often greater depending on the importance of long-range temporal context [Pons et al., 2017, Pons and Serra, 2019, Shimada et al., 2020]. Our dataset consists of mosquito samples with an average duration of 10 seconds and, additionally, we supply equal quantities of corresponding background to form a balanced class distribution of mosquito and noise (see Section 4).

Mukundarajan et al. [2017] have released an acoustic dataset recorded in free flight with smartphones. However, due to a lack of a rigorous recording protocol, the subsequent quality of the recordings is inconsistent, and there is a lack of metadata recording external factors which influence mosquito sound. There are no labels to exactly timestamp the mosquito events in files where mosquito sound is only sporadic, detracting from the overall utility of the dataset. Our database is specifically designed to eliminate these issues based on previous experience with acoustic mosquito recordings.

Kiskin et al. [2019, 2020] released extensive data spanning 22 hours of audio recordings, with crowdsourced labels covering overlapping two-second sections. However, of these, only 2 hours were labelled as containing mosquito sound. In addition, the accuracy of the labels is unknown, and the task of labelling was made difficult as clips were presented in isolation, lacking the expert knowledge and relevant background information that specialists utilised for their labels. Curated data of that release is a subset of the release of this paper, in which we improve upon the past release thanks to a dedicated joint effort between the zoological and machine learning communities.

Nevertheless, we do stress that experimentation which combines information from all of the datasets found in the literature is highly encouraged, and may help find solutions to cover multiple recording modalities, such as opto-acoustic and smartphone acoustic sensors.

## 3 Data for mosquito-borne disease prevention

### 3.1 The HumBug project

The HumBug project is a collaboration between the University of Oxford, Royal Botanic Gardens, Kew, and mosquito entomologists worldwide [HumBug, 2021]. One of the goals of the project is to develop a mosquito acoustic sensor that can be deployed into the homes of people in malaria-endemic areas to help monitor and identify the mosquito species, allowing targeted and effective vector control. Due to the rarity of mosquito events, as part of the pipeline we require a robust method for distinguishing mosquito events from background noise. This constitutes the primary use case for the baseline models of Section 5. We discuss alternate use cases further in Section 5 and Appendix D.5. In the following paragraphs we describe the role of our overall pipeline of Figure 1 by each component.

**Capturing mosquito with smartphones** We developed a power-efficient app to record mosquito flight tone using the in-built microphone on a smartphone (MozzWear [Marinos et al., 2021]). We used 16-bit mono PCM wave audio sampled at 8,000 Hz, based on prior acoustic low-cost smartphone recording solutions for mosquitoes [Li et al., 2017b, Kiskin et al., 2018].<sup>2</sup> To make mosquitoes fly close enough to a smartphone, we have developed an adapted bednet that utilises the inherent behaviour of host-seeking mosquitoes (Figure 2) [Sinka et al., 2021, Sec. 2.1.2]. The combination of

<sup>2</sup>The latest version records in 32 kbps aac in Tanzanian rural areas where bandwidth is critically limited.

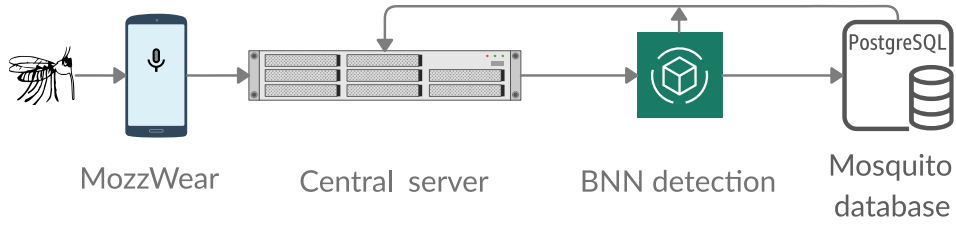


Figure 1: Schematic of project workflow. MozzWear is the mobile phone application used to capture the audio. The app synchronises to a central server, where audio enters the BNN model. Successful detections are used to update a curated database. Information feeds back to improve the model.

the bednets and smartphones constitutes the intended use case, for which we construct Test set A (see Table 2).

**Central server** Following app recording, audio is synchronised by the app, automatically or initiated by the user, to a central file server for the storage of sound recordings, and a MongoDB [MongoDB Inc, 2021] instance for the storage of metadata. The server possesses a frontend dashboard where recordings and predictions fed back from the model can be accessed. The unstructured nature of the NoSQL engine allows for additional flexibility in storing metadata, especially when new information becomes available.

**BNN detection** The classification engine deploys a Bayesian convolutional neural network (BCNN), which provides predictions with uncertainty metrics [Kiskin et al., 2021] with Monte Carlo (MC) dropout [Gal and Ghahramani, 2016]. The raw predictions of the model are fed back to the central server, and positive predictions alongside uncertainty estimates are accessible via an HTML dashboard. Positive predictions are then filtered by the probability, mutual information and predictive entropy [Houlsby et al., 2011], screened, and stored in a curated database. This drastically reduces the time spent labelling by domain experts – for our bednet data recorded in Tanzania, we estimate 1 to 2 % of 2,000 hours of recorded data contained mosquito events. Finding these events without assistance from the model was infeasible due to the vast quantity of data.

**PostgreSQL database** Due to the complex requirements of variables and data storage, we designed a relational database in PostgreSQL [PostgreSQL Global Development Group, 2021], which ensures a standardisation in the labelling and metadata process. The main concept is that all audio is stored on a data server, and each recording is uploaded with a unique ID (the full specifics are included in the database documentation provided in Appendix C). The rigorous structure of this database allows us to validate data input and ensure consistency throughout the schema. This mitigates a major cause of data quality issues and time costs in field studies. Recordings are stored in wave format at their respective sample rates, and all the metadata in csv format. For our maintenance policy, details of ethics agreements, and detailed documentation refer to the datasheet for datasets (Appendix D).

**Privacy** As a subset of data from the database may contain human speech, and other types of personal data (e.g. data recorded during trials where smartphones were actively listening continuously), we include in this paper only audio which has been assigned an explicit label of ‘mosquito’, ‘audio’, ‘background’, or otherwise full consent from members was obtained (for example where entomology experts state a recording ID, and ambient conditions etc.). Additionally, since labels have been generated both by hand and with the use of mosquito detection algorithms, to ensure no speech that has not had explicit consent for release was included in the dataset, we performed voice activity detection using Google’s WebRTC project [Ramirez et al., 2007], which is open-source, lightweight, reliable and fast [Ali, 2018, Karrer, 2020]. Sahoo [2020] tested the WebRTC VAD method over 396 hours of data, across multiple recording types. The approach was between 77 % and 99.8 % accurate. Any mosquito labels which overlapped with speech labels were removed, without truncating or re-sampling any audio to keep the format of the data in the database consistent.

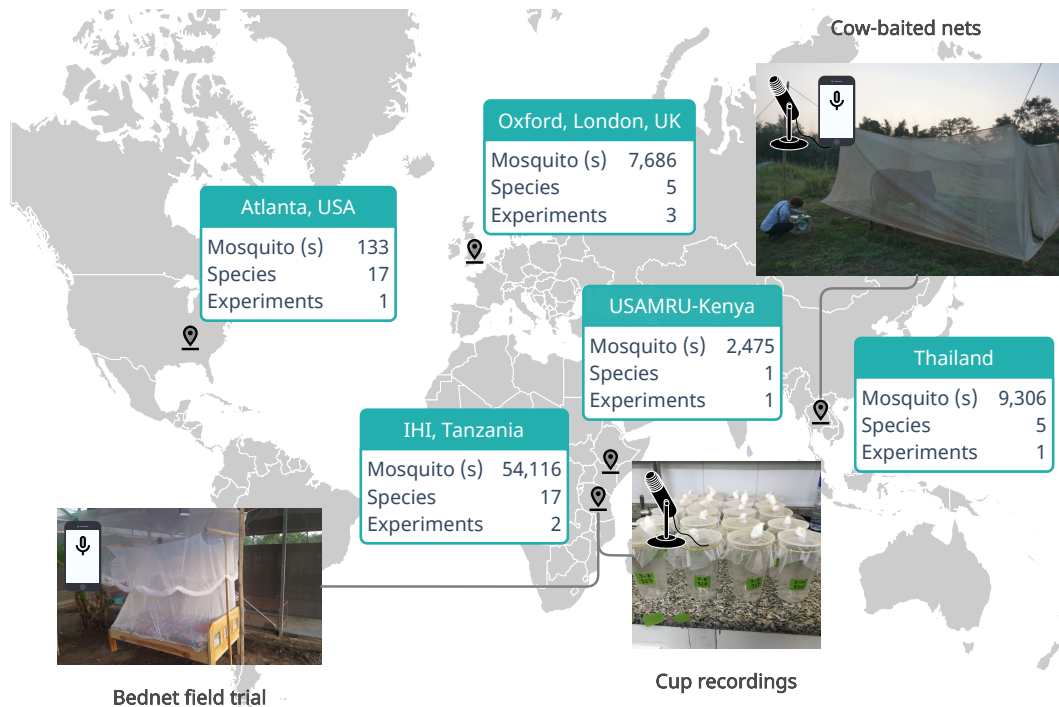


Figure 2: Map of aggregated data acquisition sites.

## 4 The HumBugDB dataset

### 4.1 Summary

Our large-scale multi-species dataset contains recordings of mosquitoes collected from multiple locations globally, as well as via different collection methods. Figure 2 shows the different locations, with the availability of labelled mosquito sound (in seconds) and number of species, and the number of experiments conducted at each location. In total, we present 71,286 seconds (20 hours) of labelled mosquito data with 53,227 seconds (15 hours) of corresponding background noise to aid with the scientific assessment process, recorded at the sites of 8 experiments. Of these, 64,843 seconds contain species metadata, consisting of 36 species (or species complexes) with the distributions illustrated in Appendix C, Figure 6 and Table 6. Table 2 gives a more detailed summary of the type of mosquitoes that were captured, and Appendix C gives a complete explanation of every field in the metadata.

In the following section we break down the data sources according to the nature of mosquitoes – bred within laboratory culture (Section 4.2.1) or wild (Section 4.2.2). We discuss the recording device and the environment the mosquitoes were recorded in – free flying in culture cages, free flying in cups or free flying in bednets (HumBug adapted bednets [Sinka et al., 2021, Sec. 2.1.2]). We also detail the methods of capture (applicable to wild mosquitoes only). These involve traditional mosquito sampling methods, including larval collection, human-baited nets (HBN), adapted Center for Disease Control Light Traps (CDC-LTs) and animal-baited nets (ABN). The method of capture is documented in more detail in Appendix C. We also make clear which dataset is used for training, and which set of experiments is used for testing the models of Section 5.

### 4.2 Data collection

#### 4.2.1 Laboratory culture mosquitoes

Many institutes that conduct research into mosquito-borne diseases hold laboratory cultures of common vector species. These include primary malaria vectors (e.g. *Anopheles gambiae*, *An. arabiensis*), arbovirus vectors including primary vectors of dengue virus (*Aedes albopictus*), yellow fever virus (*Aedes aegypti*) and west Nile virus (*Culex quinquefasciatus*). The controlled conditions



Table 2: Key audio metadata and train-test partition. ‘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 given in Appendix C.

Data (mosquitoes)	Site (country)	Recorded in	Device (sample rate)	Mosquito (s) (with species)	Negative (s)
<b>Train</b> (wild)	Kasetsart (Thailand)	cup (2018)	Telinga (44.1 kHz)	9,306 (2,869)	7,896
<b>Train</b> (wild)	IHI (Tanzania)	cup (2020)	Telinga (44.1 kHz)	45,998 (45,998)	5,600
<b>Train</b> (culture)	Zoology (Oxford, UK)	cup (2017)	Telinga (44.1 kHz)	6,573 (6,573)	1,817
<b>Train</b> (culture)	LSTMH (UK)	cup (2018)	Telinga (44.1 kHz)	376 (376)	147
<b>Train</b> (culture)	CDC (USA)	cage (2016)	phone (8 kHz)	133 (127)	1,121
<b>Train</b> (culture)	USAMRU (Kenya)	cage (2016)	phone (8 kHz)	2,475 (2,475)	31,930
<b>Test A</b> (culture)	IHI (Tanzania)	bednet (2020)	phone 8 kHz	4,118 (4,118)	3,979
<b>Test B</b> (culture)	Zoology (Oxford, UK)	cage (2016)	phone (8 kHz)	737 (737)	2,307
All	All	All	All	71,286 (64,843)	53,227

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.

**UK, Kenya, USA** Although the intrinsic variability found amongst natural populations of mosquitoes is not present in laboratory cultures, they do provide access easily to multiple species of concern. Thus we made recordings from the laboratory cultures at the London School of Tropical Medicine and Hygiene (LSTMH), the United States Army Medical Research Unit-Kenya (USAMRU-K), the Center for Diseases Control and Prevention (CDC), Atlanta, as well as with mosquitoes raised from eggs in our own laboratories at the Department of Zoology, University of Oxford. These primary recordings allowed us to quickly evaluate whether flight tone could allow us to distinguish between different species [Li et al., 2018]. 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 cups and holding these close to the recording devices.

We reserve one set of these recordings taken in culture cages by Zoology, Oxford, as one of our test datasets (denoted Test B in Table 2), as past models were able to achieve excellent mosquito detection performance when trained on data held out from the same experiment [Kiskin et al., 2018, 2017]. In this paper we treat this experiment as disparate from the remaining data, increasing the difficulty of the detection task considerably.

**Tanzania** To fulfill the aim of targeted vector control through the deployment in people’s homes, we need to be able to passively capture the mosquito’s flight tone. Therefore, in our database we include mosquitoes passively recorded in the Ifakara Health Institute’s semi-field facility (‘*Mosquito City*’) at Kining’ina, that most closely resembles the intended use of the HumBug system. It is for this reason that a labelled subset (by an expert zoologist with the help of positive BCNN predictions) of this data forms our primary test set, also marked as Test A in Table 2.

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 [Sinka et al., 2021], 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 (Figure 2). Each night of the study, 200 laboratory cultured *An. arabiensis* were released into each of the four huts and the MozzWear app began recording.

#### 4.2.2 Wild captured mosquitoes

Wild mosquitoes naturally exhibit far greater intra-specific variability. To study how this affects our ability to distinguish different species, we conducted experiments in Thailand and Tanzania.

**Thailand** Across the malaria endemic world, Asia has more dominant vector species (mosquitoes whose abundance or propensity to bite humans makes them particularly efficient vectors of disease) and species complexes anywhere else. Mosquitoes were sampled using ABNs (cow-baited nets in Figure 2), 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 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 6 and Table 6 show the exact species distribution). 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 (Appendix D.3 for device details).

**Tanzania** While Asia has the most diverse vector community, sub-Saharan Africa has the most dangerous and efficient mosquito species, namely *An. gambiae*. This is the species often referred to as the ‘most dangerous animal in the world’ and as a consequence, sub-Saharan Africa has the highest transmission of human malaria in the world, and the highest number of deaths [World Health Organization, 2020]. Using the methodology trialled in Thailand and with the help of our collaborators at the Ifakara Health Institute, we began 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. *An. gambiae* and *An. funestus* (another highly dangerous mosquito found across sub-Saharan Africa), are also siblings within their respective species complexes. Thus, standard polymerase chain reaction (PCR) identification techniques [Scott et al., 1993] were used to fully identify mosquitoes from these groups.<sup>3</sup> For all the cup recordings in Thailand and Tanzania, environmental conditions (temperature, humidity) were monitored throughout the recording process. 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 6, Table 6 for a full breakdown).

## 5 Benchmark

To showcase the utility of the data, we supply baseline models that function as acoustic mosquito event detectors. Other use cases include, but are not limited to, species classification, harmonic analysis, and the study of inter-species variability. For a more thorough consideration of these use cases refer to Appendix D.5. We discuss possible data biases arising from species imbalance, mosquito types, and multiple recording devices, and suggest mitigation strategies in Appendix D.6. For the task of mosquito event detection, we hold out Test set A of labelled field data which most closely resembles the target application. Achieving good performance on that set does not guarantee good scalability to other use cases in itself, and for this reason we use Test set B – a shorter, but very difficult low-SNR dataset as a performance marker. The prominent species in this experiment is also not as well represented, providing a further challenge. The statistics of the training and test sets are given in the rows of Table 2. In the upcoming section we will give an overview of the code we supply for our benchmarks. In Section 5.2 we describe the steps taken to train our models, and in Section 5.3 we detail how we define the performance metrics and evaluate the models supplied.

### 5.1 Code use

The top-level Jupyter notebook (Appendix B for data directory tree, code access, and layout) performs data partitioning, feature extraction and segmentation in `get_train_test_from_df()`, model

<sup>3</sup>The database gives the PCR identification within the `species` column, or the genus/complex if not available.



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 `librosa` features for use as tensors, and saves them to `config.dir_out`, if features with that particular configuration do not exist already. 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. Appendix B.3, Table 5 shows the result of feature extraction with baseline feature parameters.
- `train_model(X_train, y_train, X_val=None, Y_val=None)` 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 architecture and training strategies may be changed in `runKeras.py` or `runTorch.py`.
- `get_results(model, X, y, 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.

## 5.2 Model architecture and training

We extract 128 log-mel spectrogram features with a time window of 30 feature frames and a stride of 5 frames for training. Each frame spans 64 ms, forming a single training example  $X_i \in \mathbb{R}^{128 \times 30}$  with a temporal window of 1.92 s. Test data is strided with the stride length equal to the window size. We list all our parameters affecting the feature transformation in Appendix B.3, Table 4, and include a discussion with general recommendations for feature parameterisation. We supply two benchmark BNN model classes for this dataset:

- **Keras BNN:** A CNN with four convolutional, two max-pooling, and one fully connected layer augmented with dropout layers (shown in Appendix B.4, Figure 3). Its structure is based on prior models that have been successful in assisting domain experts in curating parts of this dataset by thresholding with uncertainty metrics [Kiskin et al., 2021].
- **PyTorch ResNet BNN:** ResNet has achieved state-of-the-art performance in audio tasks [Palanisamy et al., 2020] motivating its use as a baseline model in this paper. We augment the model with dropout layers in the appropriate building blocks to approximate a BNN. We opt to use the pre-trained model for a warm start to the weight approximations. We describe our modifications to the model class in Appendix B.4.

For both models the validation accuracy on a random split of the training data has been used to checkpoint the best-performing model. 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). We give the number of epochs, the learning rate, dropout rate, the batch size, and discuss ways to further optimise the memory usage in Appendix B.4.

## 5.3 Test results

As a benchmark, we define the test performance with three metrics: the receiver operating characteristic area-under-curve score (ROC AUC), the true positive rate (TPR), also known as the recall, and the true negative rate (TNR), to account for class imbalances in the test sets. These are evaluated over 1.92 second audio chunks. The number of audio samples in each test set following test feature extraction is given in column one of Table 3. Test features are strided by the length of the window to evaluate non-overlapping sections. To simplify the problem, edge cases where the data cannot be partitioned into full 1.92 second sections are removed from the test set. On feature extraction, all

Table 3: Test performance of the four-conv-layer Keras CNN, and two ResNet configurations over the two test sets. The number of 1.92 second samples over which the scores are evaluated is given for mosquitoes by  $N_{\text{mozz}}$  and for noise as  $N_{\text{noise}}$  respectively. Scores are reported as the mean  $\pm$  standard deviation over 10 MC dropout samples.

Data	Metric	BNN-Keras-4conv	BNN-ResNet-50	BNN-ResNet-18
Test A	ROC AUC	<b><math>0.960 \pm 0.003</math></b>	$0.959 \pm 0.001$	$0.918 \pm 0.001$
$N_{\text{mozz}} = 1,714$	TPR (%)	$71.0 \pm 0.71$	<b><math>95.6 \pm 0.24</math></b>	$72.64 \pm 0.41$
$N_{\text{noise}} = 2,068$	TNR (%)	<b><math>98.0 \pm 0.25</math></b>	$73.4 \pm 0.43$	$90.86 \pm 0.22$
Test B	ROC AUC	$0.349 \pm 0.055$	$0.545 \pm 0.004$	<b><math>0.670 \pm 0.006</math></b>
$N_{\text{mozz}} = 430$	TPR (%)	$2.16 \pm 0.48$	<b><math>2.70 \pm 0.50</math></b>	$1.42 \pm 0.22$
$N_{\text{noise}} = 1,015$	TNR (%)	<b><math>99.8 \pm 0.07</math></b>	$99.4 \pm 0.25$	$99.71 \pm 0.03$

labels shorter than that window duration are not included in the test set, though this is an area that is left for future work. When comparing performance, we suggest using a test set which has the window size as currently implemented in the code (within `get_feat()` in `feat_util.py`).

Table 3 shows the results that our baselines models were able to achieve. For the intended use case of Test A, all of the models were able to achieve ROC AUC above 0.91. The choice of model to deploy would depend on the preference over error types. For example, ResNet-50 performs better at recalling mosquito events, at the expense of a 26 % false positive rate. On the other hand, the Keras model achieves a false positive rate of only 2 %, but at the expense of missing 29 % of mosquito events. However, performance on Test B is unacceptable by all models, with all of the models categorising nearly all the audio as noise. To verify that the issue does not lie in the test set, after manually verifying each label resulting from feature extraction, we trained the models on half of Test B’s recordings, and predicted on the second half, to achieve an ROC AUC of 0.915 (Appendix B.5, Figure 4). Furthermore, prior work was able to achieve ROC AUCs of 0.871 to 0.952 with smaller neural networks which were optimised for use with scarce data [Kiskin et al., 2017]. The task presented in this paper, however, is to be able to achieve good performance over Test B, in addition to Test A, without the model having access to any data (or covariates) from both Test A and Test B.

## 6 Conclusion

In this paper we present a vast database of 20 hours of finely labelled mosquito sounds, and 15 hours of associated non-mosquito control data, constructed from carefully defined recording paradigms. Our recordings capture a diverse mixture of 36 species of mosquitoes from controlled conditions in laboratory cultures, as well as mosquitoes captured in the wild. The dataset is a result of a global co-ordination as part of the HumBug project. The HumBug project is ongoing and the robust recording pipeline described in this paper means that the database will continue to grow in the coming years. A major contribution of this paper has therefore been to link together all the moving parts, from the smartphone sensors and in-house apps, to the curation of a PostgreSQL database with the help of Bayesian neural networks.

Despite decades of work, mosquito-borne diseases are still dangerous and prevalent, with malaria alone contributing to hundreds of thousands of death each year. Therefore a further contribution of this work is to make available mosquito data that is still a scarce commodity. In addition, we have highlighted that our dataset contains real field data collected from smartphones, as well as varying background environments and different experimental settings. As a result, this multi-species data set will continue to help domain-experts in the bio-sciences study the spread of mosquito-carrying diseases, as well as the myriad of factors that affect acoustic flight tone.

Finally, our dataset will be of interest to machine learning researchers working with acoustic data, both in the availability of a real-world acoustic dataset, as well as in the way that we use Bayesian neural networks in the labelling pipeline. We provide simple functions for data manipulation and baseline models in both Keras and PyTorch, alongside extensive documentation. As a result, we make it easy for researchers to start building their own models. It is our aim, by releasing this dataset, to encourage further work in the detection of mosquitoes leading to improved models and better mosquito detection algorithms in the future.

## Checklist

### 1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** Claim: first large-scale multi-species dataset, supported with evidence in Section 2. Claim: BNNs for labelling, supported with evidence in Section 5, with code instructions. Further detail is given in Appendix B.
- (b) Did you describe the limitations of your work? **[Yes]** We describe the limitations of the baseline models in Section 5.3. We also describe how we had to withhold certain data due to potential privacy issues in Section 3.
- (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** We discuss how we mitigated potential negative impacts by incorporating a paragraph on privacy (Section 3.1). We mitigate the risk of people misusing models from a misunderstanding of performance generalisation (e.g. by making claims they have may have solved the task of mosquito detection and seek to deploy in countries without a fail-safe) by ensuring a robust train-test split of data. An assertion check in the code is performed ensure no audio recordings feature in both train and test sets, and we explain in detail how performance figures can be misinterpreted on the test sets in question.
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]** This project involves the study of potentially lethal mosquitoes, and therefore, explicit permission was obtained from the relevant Ethics committees for research. These are listed in the datasheet for datasets in Appendix D. Any personally identifiable information was removed, and explicit consent was obtained from all individuals that may feature in audio recordings throughout (see section on Privacy 3).

### 2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? **[N/A]** Results are experimental and empirical.
- (b) Did you include complete proofs of all theoretical results? **[N/A]**

### 3. If you ran experiments (e.g. for benchmarks)...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** The links to code, data, and instructions are given in Section 1. Additionally, we supply extra meta analysis to assist with code useage in Appendix B. We also describe the reasoning for our metadata format by explaining the underlying database schema and commands used to generate the metadata in Appendix C.
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** The data splits are a key factor of performance and are clearly described in Section 4 and Section 5.3. Our reasoning and the selection of hyperparameters is given in Appendix B.4.
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[Yes]** The randomness resulting from stochastic predictions with BNNs is described with a mean and standard deviation in Section 5.3. Due to the nature of random initialisation of weights during model training, we also include the trained models used to generate the predictions, and all random seeds used for data manipulation in the codebase.
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** We describe the computational resources for development and testing in Section 5.

### 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? **[Yes]** All software packages were credited to the developers (e.g. Keras, PyTorch, Audacity)
- (b) Did you mention the license of the assets? **[Yes]** The licenses of any software used are given in the datasheet for datasets in Appendix D.3.
- (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** Yes, in both Section 1 and Appendix B.1.

- 438 (d) Did you discuss whether and how consent was obtained from people whose data you're  
439 using/curating? [N/A] The dataset is original, and consent was obtained from the  
440 relevant ethics reviews and members of the teams (see datasheet for datasets, Appendix  
441 D.3).
- 442 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
443 information or offensive content? [Yes] Discussed in the Privacy paragraph of Section  
444 3, as well as in the datasheet for datasets, Appendix D.
- 445 5. If you used crowdsourcing or conducted research with human subjects...
- 446 (a) Did you include the full text of instructions given to participants and screenshots, if  
447 applicable? [N/A] For the data collection of this paper, our collaborators were working  
448 closely with us, the research was done by humans and not on human subjects.
- 449 (b) Did you describe any potential participant risks, with links to Institutional Review  
450 Board (IRB) approvals, if applicable? [N/A]
- 451 (c) Did you include the estimated hourly wage paid to participants and the total amount  
452 spent on participant compensation? [N/A]

## 453 **A Licenses**

### 454 **A.1 Code license**

455 MIT License

456 Copyright (c) 2021 HumBug-Mosquito

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

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

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

### 471 **A.2 Database license**

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

## 473 B Code use

### 474 B.1 Code access and structure

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

479 The GitHub data directory structure is as follows:

```
480 HumBugDB
├── README.md
├── *requirements.txt
├── notebooks
│   ├── main.ipynb
│   └── supplement.ipynb
├── data
│   ├── metadata
│   │   └── *.csv
│   └── audio
│       └── *.wav
├── lib
│   ├── config.py
│   ├── config_keras.py
│   ├── config_pytorch.py
│   ├── feat_util.py
│   ├── runKeras.py
│   ├── runTorch.py
│   ├── ResNetDropoutSource.py
│   ├── evaluate.py
│   └── write_audio.py
└── outputs
    ├── models
    │   ├── keras
    │   └── pytorch
    ├── features
    └── plots
```

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

483 Extract the audio from Zenodo to the folder /data/audio/ and launch the Jupyter notebook  
484 main.ipynb to perform train-test splitting, feature extraction, model training, and evaluation. The  
485 notebook imports from lib the necessary files depending on the choice of kernel and PyTorch or  
486 Keras.

### 487 B.2 Code manual

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

- 495 • `get_train_test_from_df(df_train, df_test_A, df_test_B)` extracts, reshapes,  
496 strides, and normalises librosa features for use as tensors, and saves them to  
497 `config.dir_out`, if features with that particular configuration do not exist already. The



Table 4: Feature transformation parameters, in samples. Audio processed with librosa at 8,000 Hz. The size of 1 frame in  $w$  is equal to `hop_length`. For our parameterisation this is 64 ms, resulting in an input feature slice of  $w = 1.92$  s duration and  $h = 128$  height.

NFFT	win_size	hop_length	$h$ (n_mels)	$w$	Stride
2048	2048	512	128	30	512

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. Appendix B.3, Table 5 shows the result of feature extraction with baseline feature parameters.

- `train_model(X_train, y_train, X_val=None, Y_val=None)` 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 architecture and training strategies may be changed in `runKeras.py` or `runTorch.py`.
- `get_results(model, X, y, 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.

**Supplementary notebook** `supplement.ipynb` is used to reproduce the plots of species distribution in this paper (Figure 6) and contains utilities that were used for debugging and visualising the data, should they be helpful for researchers using their own functions.

### 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, since we achieved more robust results with  $w = 30, h = 128$ , corresponding to 30 frames per window, each of 64 ms duration for a total audio slice of 1.92 seconds per sample. Nevertheless, we encourage researchers to make use of any data they wish to augment their model, for example by padding with noise to match the window size of this architecture, or by choosing a smaller window to extract features from.

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).

### B.4 Baseline models

**Keras** 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 VGGish<sup>4</sup>, which uses 0.96 second log-mel spectrogram patches as inputs, and 11 weight layers (primarily convolutional layers and max-pool layers).

<sup>4</sup><https://github.com/tensorflow/models/tree/master/research/audioset/vggish>

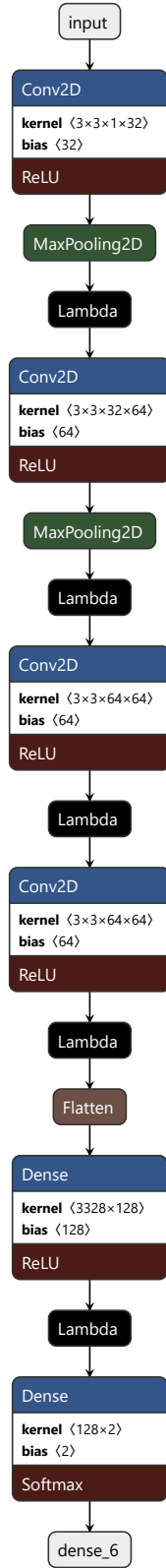


Figure 3: BCNN Keras model. Log-mel spectrograms are input with  $w = 40, h = 128$ , and passed through the above model. Lambda layers are dropout layers with probability 0.2. Made with <https://github.com/lutzroeder/netron>.

Table 5: Statistics of samples passed to models. Shown in number of data samples  $X$  resulting from log-mel feature transformation with a window size of 1.92 s, a step size of 0.32 s over the data of Table 7. All the variables affecting the size of the dataset created are found in `config.py`. and documented in `main.ipynb`.

Data	Mosquito	Negative	Total
Train	164,677	138,784	303,461
Test A: Tanzania (bednet)	1,714	2,068	3,782
Test B: Oxford Zoology (caged)	430	1,015	1,445

**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. Select which version of ResNet to use by modifying the class `ResnetDropoutFull` in `runTorch.py`. For ResNet-18, 34 the final `self.fc1` layer is of size `[512, 1]`, whereas ResNet-50 the size is `[2048, 1]`. A quick way to check this is to print `x.shape()` before the creation of the `fc1` layer.

**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 line 126. Similarly, amend the training epoch loop starting at line 85 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()` in line 105.

**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.

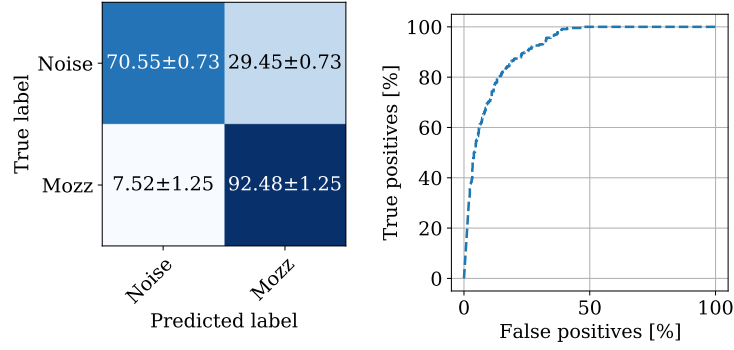
Alternatively, it is possible to use a non-pretrained architecture and change the tensor creation code in `build_data_loader()` from `runTorch.py` to remove `.repeat(1, 3, 1, 1)` as there will be no need to copy over identical data over three channels.

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).

**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. 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_ouerrun` 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].

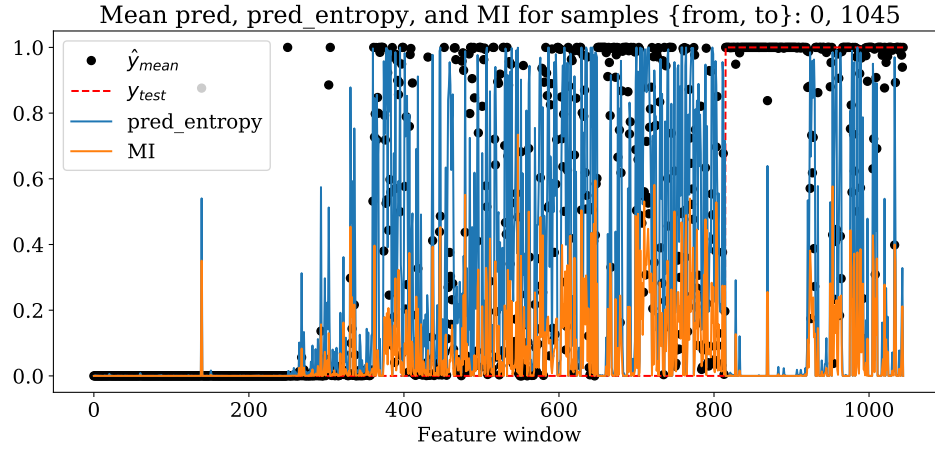
## B.5 Test performance

**Verifying data integrity of Test B** To support the validity of this data, 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.



(a) Confusion matrix

(b) ROC AUC:  $0.915 \pm 0.003$



(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.

## 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 5. 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:

```
copy (SELECT label.id, fine_end_time-fine_start_time, name,
sample_rate, record_datetime, sound_type, species, gender, fed, plurality,
age, method, mic_type, device_type, country, district, province, place,
location_type
FROM label
LEFT JOIN mosquito ON (label.mosquito_id = mosquito.id)
RIGHT JOIN audio ON (label.audio_id = audio.id)
RIGHT JOIN device ON (audio.dev_id = device.id)
WHERE type = 'Fine'
AND fine_start_time IS NOT NULL AND sound_type in
('mosquito', 'background', 'audio', 'wasp', 'fly') AND
(path LIKE '%Kenya%'
OR path LIKE '%Thai%'
OR path LIKE '%Tanzania%'
OR path LIKE '%LSTMH%'
OR path LIKE '%CDC%'
OR path LIKE '%Culex%'))
ORDER BY path) to '/data/export/neurips_2021_zenodo_0_0_1.csv' 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 6 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.

#### 616 **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  $\text{CO}_2$ , 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`).

#### 638 **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.



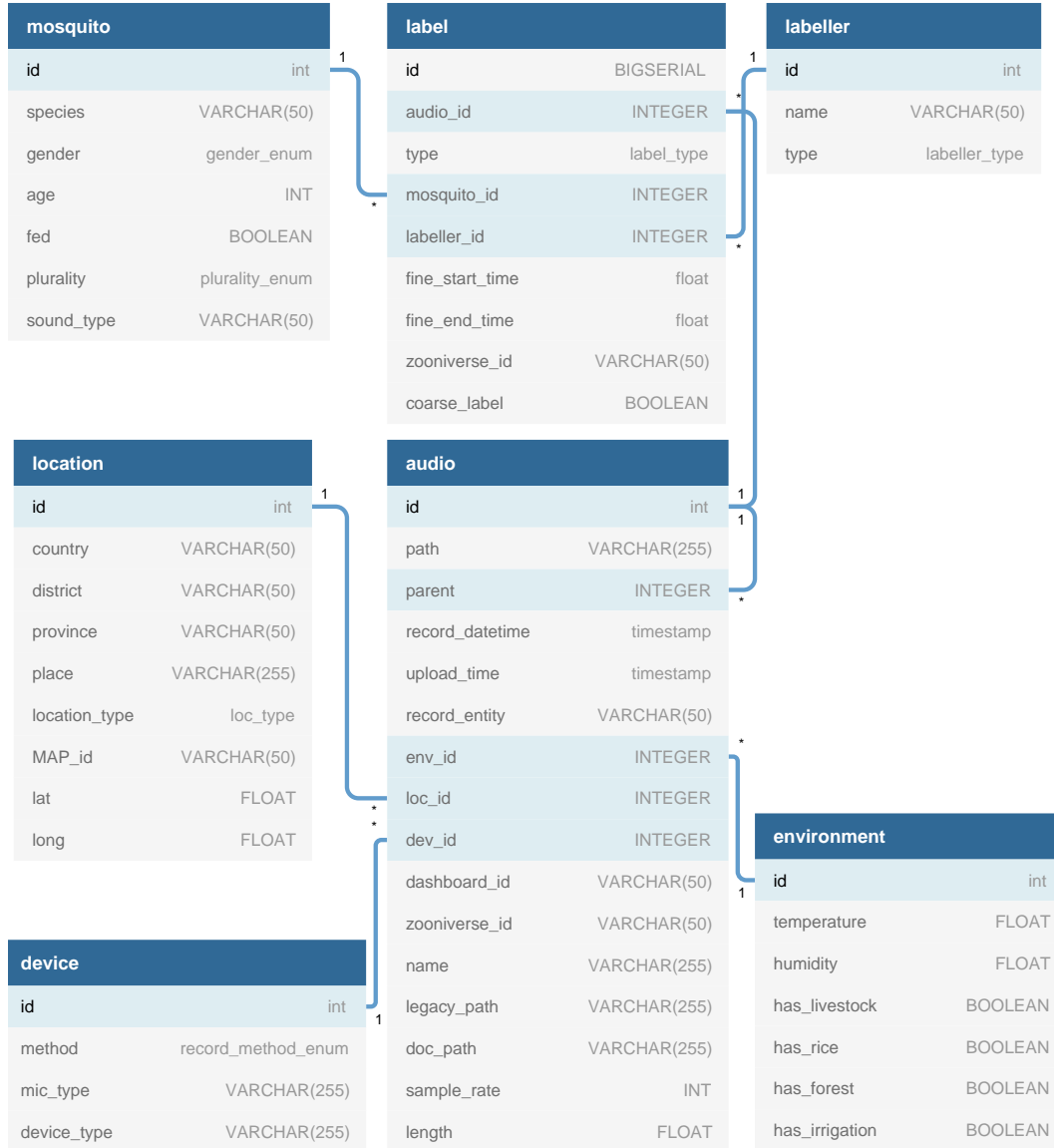


Figure 5: Relational tables of the full PostgreSQL database which was used to generate the data for this paper. The structured nature of the database enforces a standard in label format, ensuring we can efficiently mix and match data from a wide range of experiments with differing protocols. For example, 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](https://dbdiagram.io) from with `pg_dump -s`.

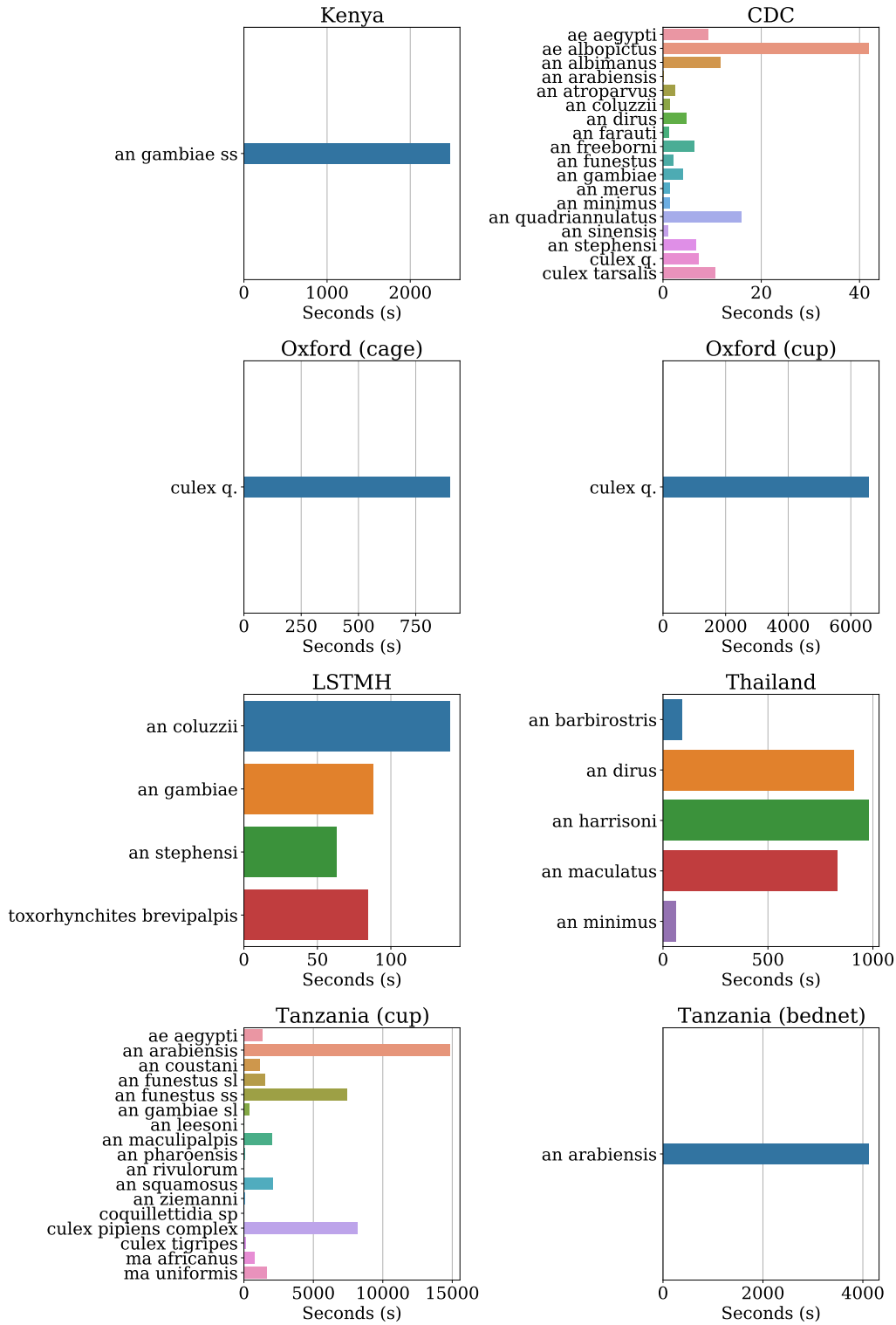


Figure 6: Species distribution per experiment corresponding to Table 7.

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

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

646 **C.2 Miscellaneous commands**

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

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

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

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

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

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

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

## D Datasheet for dataset

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, Kasetsart 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.

Table 7: Key audio metadata and train-test partition. ‘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 is given in Appendix C.

Data (mosquitoes)	Site (country)	Recorded in	Device (sample rate)	Mosquito (s) (with species)	Negative (s)
<b>Train</b> (wild)	Kasetsart (Thailand)	cup (2018)	Telinga (44.1 kHz)	9,306 (2,869)	7,896
<b>Train</b> (wild)	IHI (Tanzania)	cup (2020)	Telinga (44.1 kHz)	45,998 (45,998)	5,600
<b>Train</b> (culture)	Zoology (Oxford, UK)	cup (2017)	Telinga (44.1 kHz)	6,573 (6,573)	1,817
<b>Train</b> (culture)	LSTMH (UK)	cup (2018)	Telinga (44.1 kHz)	376 (376)	147
<b>Train</b> (culture)	CDC (USA)	cage (2016)	phone (8 kHz)	133 (127)	1,121
<b>Train</b> (culture)	USAMRU (Kenya)	cage (2016)	phone (8 kHz)	2,475 (2,475)	31,930
<b>Test A</b> (culture)	IHI (Tanzania)	bednet (2020)	phone 8 kHz	4,118 (4,118)	3,979
<b>Test B</b> (culture)	Zoology (Oxford, UK)	cage (2016)	phone (8 kHz)	737 (737)	2,307
All	All	All	All	71,286 (64,843)	53,227

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

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

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

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

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

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



### 716 D.3 Collection Process

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

- 719 • **UK, Kenya, USA:** Recordings from laboratory cultures at the London School of Tropical  
 720 Medicine and Hygiene (LSTMH), the United States Army Medical Research Unit-Kenya  
 721 (USAMRU-K); Center for Diseases Control and Prevention (CDC), Atlanta as well as  
 722 with mosquitoes raised from eggs at the Department of Zoology, University of Oxford.  
 723 Mosquitoes were recorded by placing a recording device into the culture cages where one  
 724 or multiple mosquitoes were flying, or by placing individual mosquitoes into large sample  
 725 cups and holding these close to the recording devices.
- 726 • **Tanzania i):** Mosquitoes recorded at Ifakara Health Institute’s semi-field facility (*‘Mosquito*  
 727 *City’*) at Kining’ina. The facility houses six chambers containing purpose-built experimental  
 728 huts, built using traditional methods and representing local housing constructions, with  
 729 grass roofs, open eaves and brick walls. Four different configurations of the HumBug  
 730 Net, each with a volunteer sleeping under the net, were set up in four chambers. Budget  
 731 smartphones were placed in each of the four corners of the HumBug Net. Each night of the  
 732 study, 200 laboratory cultured *An. arabiensis* were released into each of the four huts and  
 733 the MozzWear app began recording.
- 734 • **Tanzania ii)** A collection and recording project in the Kilombero Valley, Tanzania. HBNs,  
 735 larval collections and CDC-LTs were used to sample wild mosquitoes and record them in  
 736 sample cups in the laboratory. *Anopheles gambiae* and *An. funestus* (another highly dan-  
 737 gerous mosquito found across sub-Saharan Africa), are also siblings within their respective  
 738 species complexes. Thus, standard PCR identification techniques [Scott et al., 1993] were  
 739 used to fully identify mosquitoes from these groups. The Tanzanian sampling has collected  
 740 17 different species including: *An. arabiensis* (a member of the *gambiae* complex), *An.*  
 741 *coluzzii*, *An. funestus*, *An. pharoensis* (see Appendix C, Figure 6 for a full breakdown).
- 742 • **Thailand:** Mosquitoes were sampled using ABNs, HBNs and larval collections over a  
 743 period of two months during peak mosquito season (May to October 2018). Sampling  
 744 was conducted in Pu Teuy Village (Sai Yok District, Kanchanaburi Province, Thailand) at  
 745 a vector monitoring station owned by the Kasetsart University, Bangkok. The mosquito  
 746 fauna at this site include a number of dominant vector species, including *An. dirus* and *An.*  
 747 *minimus* alongside their siblings (*An. baimaii* and *An. harrisoni*) respectively (Appendix C,  
 748 Figure 6 gives a species histogram for this dataset). Sampling ran from 6 pm to 6 am, as  
 749 most anopheline vectors prefer to bite during the night. Mosquitoes were collected at night,  
 750 carefully placed into large sample cups and recorded the following day using the high-spec  
 751 Telinga field microphone and a budget smartphone.

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

- 755 • Smartphone (Iitel, Alcatel, and others) audio recording with the MozzWear application.  
 756 Smartphone devices may have variable sample rates, as denoted by the sample rate column  
 757 of the metadata. The version of MozzWear used in the curation of this dataset recorded  
 758 audio in 8,000 Hz mono wave format.
- 759 • Telinga EM-23 field microphone, and Tascam, Olympus recording devices recording at  
 760 44,100 Hz. The Telinga is a very sensitive, low-noise microphone which was widely adopted  
 761 in bioacoustic studies.
- 762 • Human labelling with Excel.
- 763 • Human labelling with Audacity (GNU GPLv2 license).
- 764 • Labels produced by a Bayesian convolutional neural network (our own, MIT license, in-  
 765 cluded in paper).
- 766 • Voice activity detection and removal with WebRTC (BSD license).

767 • Python (BSD-style license), MongoDB (Server Side Public License), Django (BSD license),  
768 Apache (GPLv3 license), PostgreSQL (BSD/MIT-like license), Unix for databases, HTML  
769 dashboards, and post-processing.

770 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors)**  
771 **and how were they compensated (e.g., how much were crowdworkers paid)?** Researchers  
772 from the locations previously mentioned, paid salary from their respective institutions, through the  
773 grants disclosed previously.

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

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

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

#### 782 **D.4 Preprocessing/cleaning/labeling**

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

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

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

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

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

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

#### 803 **D.5 Uses**

804 **Has the dataset been used for any tasks already?** A subset of this data (recorded in Thailand,  
805 Kenya, UK, USA) has been used to train a machine learning model to distinguish and detect a

806 mosquito from its acoustic signature. The model was a 4-layer Bayesian convolutional neural network  
807 implemented in Keras. The predictive entropy and mutual information were used to screen predictions  
808 over thousands of hours of data. Hand labels were added to correct predictions, and the labels were  
809 fed back into the database [Kiskin et al., 2021]. Code for the training and resulting predictive pipeline  
810 is available on <https://github.com/HumBug-Mosquito/MozzBNN>.

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

- 813 • Bioacoustic classification with wavelet-conditioned neural networks [Kiskin et al., 2017,  
814 2018].
- 815 • Cost-sensitive mosquito detection [Li et al., 2017a]
- 816 • A case study of species classification with field recordings [Li et al., 2018]
- 817 • A release of a subset of this database for crowdsourcing (with baseline mosquito detector  
818 model) [Kiskin et al., 2019, 2020]

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

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

- 824 1. **Species classification with HumBugDB.** Training a machine learning model to distinguish  
825 between various species.
- 826 2. **Validating species classification models from the literature.**
- 827 3. **Frequency analysis.** Identifying the fundamental and harmonic frequencies of flight tone  
828 for a particular species, to improve upon the understanding of bioacoustics literature, and  
829 entomological research.
- 830 4. **Examining inter-species (or similar) variability.** For example, the effect on the sound of  
831 flight as a result of age, gender, or any field supported in the database.

832 We now expand upon each point:

- 833 1. **Species classification with HumBugDB.** There a multitude of mosquito species, however  
834 only certain species have potential to transmit certain pathogens. It is imperative therefore to  
835 accurately locate and identify the few dangerous mosquito species amongst the many benign  
836 ones to achieve efficient mosquito control. Dangerous species include primary malaria  
837 vectors (e.g. *Anopheles gambiae*, *An. arabiensis*), arbovirus vectors including primary  
838 vectors of dengue virus (*Aedes albopictus*), yellow fever virus (*Aedes aegypti*) and west nile  
839 virus (*Culex quinquefasciatus*). The species available and their associated occurrence per  
840 experimental group is given in Figures 6 and Table 6 in Appendix C.

841 When designing experiments, it would be useful to take into account the method of capture  
842 and whether mosquitoes are from lab cultures or captured as individuals in the wild. The  
843 experimental groups are given in Table 7, and further details can be found in the database  
844 metadata as detailed in Appendix C. As a starting point, we recommend using the Tanza-  
845 nian cup recordings of wild individuals, denoted by the metadata `country = Tanzania`,  
846 `location_type = field`. Each recording (denoted by the name metadata field) is of an  
847 individual mosquito. You may partition the train/validation/test splits by recording name, to  
848 ensure appropriately disjoint subsets (avoiding repetition of individuals when validating or  
849 testing models).

850 A very useful property of this dataset is that there is overlap in species recorded in entirely  
851 different experimental conditions (Table 6). This gives the opportunity to test the robustness  
852 of the model when generalising across datasets.

- 853 2. **Validating species classification models from the literature.** As a result of procuring  
854 curated data with species meta-information of both wild and lab mosquitoes, this dataset  
855 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. If you encounter any issues, or require further information do not hesitate to contact the database maintainers (Appendix D.7).

3. **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.

4. **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 preprocessed/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. **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].

3. **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.

If using multi-class classifiers, it is possible to begin by weighting 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.

## 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.

**Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?** The data will be updated as new data from new trials is obtained and curated. We expect updates to occur every few months in 2021. Updates will be communicated by commits and pushes to the GitHub repository. Please also follow the link on Zenodo for versioning details, where older versions will continue to be hosted.

**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.

## References

- H. Ali. Real-time Communication Using WebRTC. Technical report, Georgia Institute of Technology, 2018.
- R. J. Bomphrey, T. Nakata, N. Phillips, and S. M. Walker. Smart wing rotation and trailing-edge vortices enable high frequency mosquito flight. *Nature*, 544(7648):92–95, 2017.
- Y. Chen, A. Why, G. Batista, A. Mafra-Neto, and E. Keogh. Flying insect classification with inexpensive sensors. *Journal of Insect Behavior*, 27(5):657–677, 2014.
- F. Chollet et al. Keras, 2015. URL <https://keras.io>. Accessed: 2018-06-07.
- A. D. Cobb, S. J. Roberts, and Y. Gal. Loss-calibrated approximate inference in Bayesian neural networks. *arXiv preprint arXiv:1805.03901*, 2018.
- E. Fanioudakis, M. Geismar, and I. Potamitis. Mosquito wingbeat analysis and classification using deep learning. In *2018 26th European Signal Processing Conference (EUSIPCO)*, pages 2410–2414, 2018.
- N. Friederici, S. Ojanperä, and M. Graham. The impact of connectivity in Africa: Grand visions and the mirage of inclusive digital development. *The Electronic Journal of Information Systems in Developing Countries*, 79(1):1–20, 2017.
- Y. Gal and Z. Ghahramani. Dropout as a Bayesian approximation: representing model uncertainty in deep learning. In *International Conference on Machine Learning*, pages 1050–1059, 2016.
- T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. Daumé III, and K. Crawford. Datasheets for datasets. *arXiv preprint arXiv:1803.09010*, 2018.
- J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter. Audio set: an ontology and human-labeled dataset for audio events. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 776–780. IEEE, 2017.
- GSMA. The mobile economy-sub-saharan africa, 2020. URL <https://www.gsma.com/mobileeconomy/sub-saharan-africa/>. Last accessed: 2021-07-08.
- R. Harbach. Mosquito taxonomic inventory, 2013. URL <http://mosquito-taxonomic-inventory.info/>. Last accessed: 2021-06-07.
- S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold, et al. CNN architectures for large-scale audio classification. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 131–135. IEEE, 2017.
- A. A. Hoffmann and P. A. Ross. Rates and Patterns of Laboratory Adaptation in (Mostly) Insects. *Journal of Economic Entomology*, 111(2):501–509, 03 2018. ISSN 0022-0493. doi: 10.1093/jee/toy024. URL <https://doi.org/10.1093/jee/toy024>.
- N. Houlsby, F. Huszár, Z. Ghahramani, and M. Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
- B. Huho, K. Ng’habi, G. Killeen, G. Nkwengulila, B. Knols, and H. M. Ferguson. Nature beats nurture: a case study of the physiological fitness of free-living and laboratory-reared male *Anopheles gambiae* sl. *Journal of Experimental Biology*, 210(16):2939–2947, 2007.
- HumBug. The HumBug Project, 2021. URL <https://humbug.ox.ac.uk/>. Accessed: 2021-06-21.
- S. Jakhete, S. Allan, and R. Mankin. Wingbeat frequency-sweep and visual stimuli for trapping male *Aedes aegypti* (Diptera: Culicidae). *Journal of medical entomology*, 54(5):1415–1419, 2017.
- B. J. Johnson and S. A. Ritchie. The siren’s song: exploitation of female flight tones to passively capture male *Aedes aegypti* (Diptera: Culicidae). *Journal of medical entomology*, 53(1):245–248, 2016.



994 R. Karrer. Google WebRTC Voice Activity Detection module, 2020. URL [https://github.com/](https://github.com/rafaelkarrer/mex-webrtcvad/releases/tag/v0.1)  
995 [rafaelkarrer/mex-webrtcvad/releases/tag/v0.1](https://github.com/rafaelkarrer/mex-webrtcvad/releases/tag/v0.1). Accessed: 2021-06-05.

996 I. Kiskin. *Machine learning for acoustic mosquito detection*. PhD thesis, University of Oxford, 2020.

997 I. Kiskin, B. P. Orozco, T. Windebank, D. Zilli, M. Sinka, K. Willis, and S. Roberts. Mosquito  
998 detection with neural networks: the buzz of deep learning. *arXiv preprint arXiv:1705.05180*, 2017.

999 I. Kiskin, D. Zilli, Y. Li, M. Sinka, K. Willis, and S. Roberts. Bioacoustic detection with wavelet-  
1000 conditioned convolutional neural networks. *Neural Computing and Applications: Special Issue on*  
1001 *Deep Learning for Music and Audio*, Aug 2018. ISSN 1433-3058.

1002 I. Kiskin, U. Meepegama, and S. Roberts. Super-resolution of time-series labels for bootstrapped  
1003 event detection. *Time-series Workshop at the International Conference on Machine Learning*,  
1004 2019.

1005 I. Kiskin, L. Wang, A. Cobb, et al. Humbug Zooniverse: a crowd-sourced acoustic mosquito dataset.  
1006 *International Conference on Acoustics, Speech, and Signal Processing 2020, NeurIPS Machine*  
1007 *Learning for the Developing World Workshop 2019*, 2019, 2020.

1008 I. Kiskin, A. D. Cobb, M. Sinka, and S. J. Roberts. Automatic acoustic mosquito tagging with  
1009 Bayesian neural networks. *The European Conference on Machine Learning and Principles and*  
1010 *Practice of Knowledge Discovery in Databases*, 2021.

1011 Y. Li, I. Kiskin, D. Zilli, M. Sinka, H. Chan, K. Willis, and S. Roberts. Cost-sensitive detection  
1012 with variational autoencoders for environmental acoustic sensing. *NeurIPS Workshop on Machine*  
1013 *Learning for Audio Signal Processing*, 2017a.

1014 Y. Li, D. Zilli, H. Chan, I. Kiskin, M. Sinka, S. Roberts, and K. Willis. Mosquito detection with  
1015 low-cost smartphones: data acquisition for malaria research. *NeurIPS Workshop on Machine*  
1016 *Learning for the Developing World*, 2017b.

1017 Y. Li, I. Kiskin, M. Sinka, D. Zilli, H. Chan, E. Herreros-Moya, T. Chareonviriyaphap, R. Tisgratog,  
1018 K. Willis, and S. Roberts. Fast mosquito acoustic detection with field cup recordings: an initial  
1019 investigation. *Detection and Classification of Acoustic Scenes and Events*, 2018.

1020 T. Marinos, S. Lin, D. Zilli, and H. Chan. MozzWear, 2021. URL [https://github.com/](https://github.com/HumBug-Mosquito/MozzWear)  
1021 [HumBug-Mosquito/MozzWear](https://github.com/HumBug-Mosquito/MozzWear). Pending update on Google Play store, GitHub private, accessed:  
1022 2021-06-05.

1023 MongoDB Inc. Mongoddb, 2021. URL <https://www.mongodb.com/>. Accessed: 2021-06-05.

1024 H. Mukundarajan, F. J. H. Hol, E. A. Castillo, C. Newby, and M. Prakash. Using mobile phones  
1025 as acoustic sensors for high-throughput mosquito surveillance. *eLife*, 6:e27854, Oct 2017. ISSN  
1026 2050-084X.

1027 K. Palanisamy, D. Singhania, and A. Yao. Rethinking CNN models for audio classification. *arXiv*  
1028 *preprint arXiv:2007.11154*, 2020.

1029 A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein,  
1030 L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy,  
1031 B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance  
1032 deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox,  
1033 and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages  
1034 8024–8035. Curran Associates, Inc., 2019. URL [http://papers.neurips.cc/paper/](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)  
1035 [9015-pytorch-an-imperative-style-high-performance-deep-learning-library.](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)  
1036 [pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).

1037 V. P. Perevozkin and S. S. Bondarchuk. Species specificity of acoustic signals of malarial mosquitoes  
1038 of anopheles maculipennis complex. *International Journal of Mosquito Research*, 2(3):150–155,  
1039 2015.

1040 J. Pons and X. Serra. musicnn: Pre-trained convolutional neural networks for music audio tagging.  
1041 *arXiv preprint arXiv:1909.06654*, 2019.

1042 J. Pons, O. Nieto, M. Prockup, E. Schmidt, A. Ehmann, and X. Serra. End-to-end learning for music  
1043 audio tagging at scale. *arXiv preprint arXiv:1711.02520*, 2017.

1044 PostgreSQL Global Development Group. PostgreSQL, 2021. URL [https://www.postgresql.](https://www.postgresql.org/docs/9.3/app-psql.html)  
1045 [org/docs/9.3/app-psql.html](https://www.postgresql.org/docs/9.3/app-psql.html). Accessed: 2021-06-05.

1046 J. Ramirez, J. M. Górriz, and J. C. Segura. Voice activity detection. fundamentals and speech  
1047 recognition system robustness. *Robust speech recognition and understanding*, 6(9):1–22, 2007.

1048 A. Sahoo. Voice activity detection for low-resource settings. *Department of Electrical Engineering,*  
1049 *Stanford University*, 2020.

1050 J. A. Scott, W. G. Brogdon, and F. H. Collins. Identification of single specimens of the anopheles  
1051 gambiae complex by the polymerase chain reaction. *The American journal of tropical medicine*  
1052 *and hygiene*, 49(4):520–529, 1993.

1053 K. Shimada, N. Takahashi, S. Takahashi, and Y. Mitsufuji. Sound event localization and detection  
1054 using activity-coupled cartesian doa vector and rd3net. Technical report, DCASE2020 Challenge,  
1055 July 2020.

1056 P. M. Simões, R. A. Ingham, G. Gibson, and I. J. Russell. A role for acoustic distortion in novel rapid  
1057 frequency modulation behaviour in free-flying male mosquitoes. *Journal of Experimental Biology*,  
1058 219(13):2039–2047, 2016.

1059 M. Sinka, D. Zilli, I. Kiskin, Y. Li, D. Kirkham, W. Rafique, H. Chan, B. Gutteridge, E. Herreros-  
1060 Moya, H. Portwood, S. J. Roberts, and K. J. Willis. HumBug – An Acoustic Mosquito Monitoring  
1061 Tool for use on budget smartphones. *Methods in Ecology and Evolution*, 2021. doi: 10.1111/  
1062 2041-210X.13663.

1063 M. E. Sinka, M. J. Bangs, S. Manguin, Y. Rubio-Palis, T. Chareonviriyaphap, M. Coetzee, C. M.  
1064 Mbogo, J. Hemingway, A. P. Patil, W. H. Temperley, et al. A global map of dominant malaria  
1065 vectors. *Parasites & vectors*, 5(1):1–11, 2012.

1066 D. Vasconcelos, N. J. Nunes, and J. Gomes. An annotated dataset of bioacoustic sensing and features  
1067 of mosquitoes. *Scientific Data*, 7(1):1–8, 2020.

1068 World Bank Organisation. Listening to Africa, 2017. URL [https://www.worldbank.org/en/](https://www.worldbank.org/en/programs/listening-to-africa)  
1069 [programs/listening-to-africa](https://www.worldbank.org/en/programs/listening-to-africa). Last accessed: 2021-07-08.

1070 World Health Organization. Fact Sheet, 2020. URL [https://www.who.int/news-room/](https://www.who.int/news-room/fact-sheets/detail/vector-borne-diseases)  
1071 [fact-sheets/detail/vector-borne-diseases](https://www.who.int/news-room/fact-sheets/detail/vector-borne-diseases). Accessed: 2020-01-26.