

Comparison of VGGish embedding and MFCC feature 1 in bee colony sound classification (#73841)

First revision

Guidance from your Editor

Please submit by **2 Dec 2022** for the benefit of the authors .



Structure and Criteria

Please read the 'Structure and Criteria' page for general guidance.



Raw data check

Review the raw data.



Image check

Check that figures and images have not been inappropriately manipulated.

Privacy reminder: If uploading an annotated PDF, remove identifiable information to remain anonymous.

Files

Download and review all files from [materials page](#).

the

- 1 Tracked changes manuscript(s)
- 1 Rebuttal letter(s)
- 11 Figure file(s)
- 5 Table file(s)
- 1 Other file(s)

For assistance email peer.review@peerj.com

Structure and Criteria

2



Structure your review

The review form is divided into 5 sections. Please consider these when composing your review:

1. BASIC REPORTING
2. EXPERIMENTAL DESIGN
3. VALIDITY OF THE FINDINGS
4. General comments
5. Confidential notes to the editor

You can also annotate this PDF and upload it as part of your review When ready

[submit online](#).

Editorial Criteria

Use these criteria points to structure your review. The full detailed editorial criteria is on your [guidance page](#).

BASIC REPORTING

✓ Clear, unambiguous, professional ~~English~~~~Original~~~~English~~ language used through ~~out the journal~~~~out~~. The

✓ Intro & background to show context. Research question & relevant. & meaningful. It is stated how the Structure

✓ fills an identified knowledge gap. discipline norm or improved for clarity. performed to a Figures are relevant, high quality, wellhigh technical & ethical standard.

✓ Raw data supplied (see [PeerJ policy](#)).

EXPERIMENTAL DESIGN

✓ ~~Original~~ primary research within [Scope of](#)

✓ well defined, relevant Literature well referenced conforms to [PeerJ standards](#), research

✓ Rigorous investigation

✓ information to replicate.

✓ labelled & described. Methods described with sufficient detail &

VALIDITY OF THE FINDINGS

i Impact and novelty not assessed. Conclusions are replication encouraged where original research literature is clearly supporting results. stated.

✓ well stated, linked to *Meaningful* question & limited to rationale & benefit to

✓ All underlying data have been provided; they are robust, statistically sound, & controlled.

Standout reviewing tips

3



The best reviewers use these techniques

Tip

Support criticisms with evidence from the text or from other sources

Example

Smith et al (J of Methodology, 2005, V3, pp 123) have shown that the analysis you use in Lines 241-250 is not the most appropriate for this situation. Please explain why you used this method.

Give specific suggestions on how to improve the manuscript

Your introduction needs more detail. I suggest that you improve the description at lines 57- 86 to provide more justification for your study (specifically, you should expand upon the knowledge gap being filled).

Comment on language and grammar issues

The English language should be improved to ensure that an international audience can clearly understand your text. Some examples where the language could be improved include lines 23, 77, 121, 128 – the current phrasing makes comprehension difficult. I suggest you have a colleague who is proficient in English and familiar with the subject matter review your manuscript or contact a professional editing service.

Organize by importance of the issues, and number your points

1. *Your most important issue*
2. *The next most important item*
3. *...*
4. *The least important points*

Please provide constructive criticism, and avoid personal opinions

I thank you for providing the raw data, however your supplemental files need more descriptive metadata identifiers to be useful to future readers. Although your results are compelling, the data analysis should be

Comment on strengths (as well as weaknesses) of the manuscript

improved in the following ways: AA, BB, CC

I commend the authors for their extensive data set, compiled over many years of detailed fieldwork. In addition, the manuscript is clearly written in professional, unambiguous language. If there is a weakness, it is in the statistical analysis (as I have noted above) which should be improved upon before Acceptance.

Comparison of VGGish embedding and MFCC feature in bee colony sound classification

Nayan Di ^{Corresp., 1, 2}, Muhammad Zahid Sharif ^{1, 2}, Zongwen Hu ^{3, 4}, Renjie Xue ^{1, 2}, Baizhong Yu ^{1, 2}

¹ Anhui Institute of Optics and Fine Mechanics, Hefei Institute of Physical Science, Chinese Academy of Sciences, Hefei, China

² University of Science and Technology of China, Hefei, China

³ Eastern Bee Research Institute, College of Animal Science and Technology, Yunnan Agricultural University, Kunming, China

⁴ The Sericultural and Apicultural Research Institute, Yunnan Academy of Agricultural Sciences, Mengzi, China

Corresponding Author: Nayan Di
Email address: dny_yan@126.com

Background. Bee colony sound is a continuous, low-frequency buzzing sound that varies with the environment or the colony's behavior and is considered meaningful. Bees use sounds to communicate within the hive, and bee colony sounds investigation can reveal helpful information about the circumstances in the colony. **Therefore, one crucial step in analyzing bee colony sounds is to extract appropriate acoustic features.**

Methods. This paper uses VGGish embedding and MFCC feature generated from three bee colony sound datasets to train several machine learning algorithms to determine which acoustic feature performs best in bee colony sound recognition.

Results. The results showed that VGGish embedding performs better than or on par with MFCC feature in all three datasets.

1 Comparison of VGGish embedding and MFCC feature in 2 bee colony sound classification

3 Nayan Di^{1,2,*}, Muhammad Zahid Sharif^{1,2}, Zongwen Hu^{3,4}, Renjie Xue^{1,2}, Baizhong Yu^{1,2}

4 ¹ Anhui Institute of Optics and Fine Mechanics, Hefei Institute of Physical Science, Chinese

5 Academy of Sciences, Hefei, 230031, China

6 ² University of Science and Technology of China, Hefei, 230026, China

7 ³ Eastern Bee Research Institute, College of Animal Science and Technology, Yunnan

8 Agricultural University, Kunming, China

9 ⁴ The Sericultural and Apicultural Research Institute, Yunnan Academy of Agricultural Sciences,

10 Mengzi Yunnan, China

Commented [JB1]: Define the aconyms VGGish and MFCC here

11

12 **Corresponding Author:*

13 *Nayan Di¹*

14 *350 Shushanhu Road, Hefei, Anhui, 230031, P. R. China*

15 *Email address: dny_yan@126.com*

16 Abstract

17 **Background.** Bee colony sound is a continuous, low-frequency buzzing sound that varies with the
18 environment or the ~~colony~~'s ~~colony's~~ behavior and is considered meaningful. Bees use sounds
to
19 communicate within the hive, and bee colony sounds investigation can reveal helpful 20
information about the circumstances in the colony. One crucial step in analyzing bee colony 21
sounds is to extract appropriate acoustic feature.

22 **Methods.** This paper uses VGGish embedding and ~~MFCC~~~~coefficients-feature~~, ~~extracted~~
~~generated~~ from three bee 23 colony sound datasets, to train several machine learning algorithms to
determine which acoustic 24 feature performs best in bee colony sound recognition.

25 **Results.** The results showed that VGGish embedding performs better than or on par with ~~the~~ MFCC
26 feature in all three datasets.

27 **Keywords:** Acoustic feature; Bee colony sound; VGGish embedding; MFCC feature; *Apis*
28 *cerana* *Apis mellifera*? If *Apis ceranae*, this adds application of acoustic screening to another
species of honey bee. If *Apis mellifera*, which is what was examined in the first draft of this
study, the *A. ceranae* should be replaced with *A. mellifera* throughout the manuscript.

29

30 1. Introduction

31 Honey bees play an essential role in agriculture production and are almost responsible for
32 90% of global commercial pollination service pollination(Klein, VaissiLre et al. 2007). As a
vital 33 node of the agriculture section, it is essential to ensure that the bee colonies can
provide service.

Commented [JB2]: Spell out the words from which each of these acronyms is derived.

34 To save human resources and reduce disturbance to bee colonies, a non-invasive method that can
 35 detect the intra-colonial condition of the hive without disturbing the colony is a consensus among
 36 researchers and practitioners (Meikle and Holst 2015). The internal environment of a honeybee
 37 hive includes sound, temperature, and humidity, which are relatively stable under certain
 38 conditions (Murphy, Magno et al. 2015). By monitoring these indicators in the hive and
 39 establishing the association between these indicators, we could learn a lot about the status of the
 40 colony (Ferrari, Silva et al. 2008, Braga, Gomes et al. 2020). Among these indicators, beehive
 41 sound is critical. Bee buzzing carries information on colony behavior and phenology. Honey
 42 bees emit specific sounds when exposed to stressors such as pest infection (Qandour, Ahmad et
 43 al. 2014), airborne toxicants (Zhao, Deng et al. 2021), and failing queens (Cejrowski,
 Szymański
 44 et al. 2018). Using both statistical and A.I. analysis of colony sounds, ~~Jerry Bromenshenk et al.,~~
in their patents (2009) and in their review paper (2015) showed that their Artificial Intelligence
(A.I.) could detect a diverse
~~45 developed a smartphone app (Bee Health Guru KS) which could detect~~ a diverse variety of
 46 chemicals and eight colony health variables, inside beehives by simply ~~putting~~ putting a
microphone into the bottom of a beehive and recording bee colony sounds for 30 or 60
seconds. In 2019, they released a cellphone app (Bee Health Guru) that can run the
diagnostic programs, record and analyze the results, and upload the data, visual inspections,
and app analyses to a cloud-based site, which automatically generates a report with the GPS
location shown on a map. the cellphone near the beehive and recording the bee ~~47 colony~~
sound for a few seconds. Currently, the app is being calibrated for a variety of phone
operating systems for bee sounds from around the world (www.beehealth.guru).
 48 One of the critical phases in analyzing the bee colony sound would be extracting
 49 appropriate feature from the bee colony sound for machine learning or deep learning
 50 algorithms.
 51 Traditionally we use frequency domain or time domain feature of sound, such as
 soundcape
 indices and low-frequency signal features (Sharif, Wario et al. 2020). Mel Frequency
 Cepstrum

Commented [JB3]: The 2015 review is available on line and should be cited. The acoustic discrimination includes both chemical and biological endpoints.

52 Coefficient (MFCC) is one of the most commonly used features. It is characterized by using
a set

53 of critical coefficients to create Mel cepstrum, which makes its cepstrum more similar to the

54 nonlinear human auditory system (Muda, Begam et al. 2010). Due to the nonlinear

55 correspondence between Mel frequency and Hz frequency, the calculation accuracy of
MFCC

56 decreases with the increase of frequency. This characteristic makes MFCC more suitable for
bee 57 colony sound than other feature extraction methods in the past because the sound
signal in the 58 colony is concentrated in the low-frequency part (Dietlein 1985).

59 Thanks to the rapid development of artificial intelligence, Convolutional Neural Net (CNN)

60 and Recurrent Neural Networks (RNN) have been widely applied in audio recognition
(Kumar

61 and Raj 2017). Experimental results showed that the recognition method based on CNN is
prior

62 to the method based on machine learning models (Kulyukin, Mukherjee et al. 2018). Visual

63 Geometry Group (VGG) is one of the most popular CNN models. ~~It was proposed by~~
~~Simonyan~~ ~~Simonyan and Zisserman proposed it~~

64 ~~and Zisserman~~ in 2014 and is named after the Visual Geometry Group (Simonyan and
Zisserman

65 2014). VGGish is a TensorFlow definition of a VGG-like audio classification model. ~~The~~
VGGish

66 model is a derivative network of the VGG network trained on a large YouTube dataset

67 (Gemmeke, Ellis et al. 2017). Its structure is consistent with VGG11, including eight

68 convolutional layers, five pooling layers, and three fully connected layers. Each
convolutional

69 layer uses a 3x3 convolution kernel. VGGish converts audio input feature into a
semantically

70 meaningful, high-level 128-dimensional embedding, which can be fed as input to a
downstream

71 classification model. Due to the scale and diversity of the YouTube dataset, the resulting
acoustic

Commented [JB4]: This definition needs to appear the first time it is mentioned in the paper, then use MFCC. MFCC is a set of calculations that generate a coefficient. Drop the redundant and confusing phrase MFCC coefficient feature. The word feature is confusing - MFCC is a coefficient, not a feature.

Commented [JB5]: The original article used YouTube derived data, this revision suggests that it wasn't used. If used, the YouTube data is likely to be for *Apis mellifera*, not *Apis ceranae*.

7271 features are both very general and of high resolution, placing each audio sample in a high73
dimensional feature space that is unlikely to show ecosystem-specific bias. This 128-
dimensional

74 embedding characteristic is helpful in various identification contexts, including monitoring 75
anomalous events in an ecosystem (Sethi, Jones et al. 2020) and sound-based disease detection 76
(Shi, Du et al. 2019).

77 In this article, we contribute to the body of research on audio beehive monitoring by
78 comparing VGGish embedding and standard MFCC statistics (?) feature in classifying
audio samples from
79 microphones deployed inside beehives. We tested the VGGish embedding and MFCC
feature on 80 three different classification tasks and compared these two-feature using four
machine-learning 81 algorithms.

82 In particular, Section two will describe the hardware and software configuration to obtain
83 bee colony sound and report the detail of the three bee colony datasets we used in this
paper.

84 Section three will give the performance of VGGish embedding and MFCC feature
coefficients in bee colony

8485 sound classification, as well as the effects of different dimensional reduction algorithms.
Section 86 four will report conclusions and a future perspective.

Commented [JB6]: Overall, still a very small data set. Most investigators use recordings from 100 or more colonies, splitting the data into a training and a testing group. One can not train on all samples, then test on the same set of samples.,

87 2. Materials and Methods

88 2.1 Hardware

89 The hardware and software systems for obtaining bee colony sound are as follows: A
 90 microphone inside the beehive (PCK200, TAKSTAR) was placed about 15cm from the
 bottom.
 91 The microphone has a frequency range of 30-20 kHz ~~30 Hz to 20 kHz (30 kHz to 20 kHz~~
~~makes no sense)~~ and a sensitivity of -35 dB. A digital sound
 92 card (UM2, BEHRINGER) was used to convert the analog signal into a digital signal. The
 digital
 93 signal was transmitted to a personal computer (HP 2170p, Windows 7), The software
 Audacity 94 was used to record the sound, and the sound sampling rate was set to 44.1kHz,
 mono. Sound files
 95 were saved on the hard disk in .wav format. The hardware structure is illustrated in
 Figure1.

96 2.2 Audio data

97 The experiment was carried out at the Sericulture and Apiculture Research Institute of
 98 Yunnan Academy of Agricultural Sciences (23.5144N,103.4043E) from November 2020 to
 June
 99 2021. The institute is located in Caoba Town, Mengzi City, Yunnan Province, China. We
 100 collected three collections of honey bees (*Apis* ~~cerana~~ *Serena*) colony sounds and named
 them dataset
 101 one, two and three, respectively. A detailed description of these datasets is given below.
 Every
 102 bee colony lived in standard wood beehive with a queen of 10 months old. All the bee
 colonies
 103 were ~~are~~ healthy without any sign of attack by pests, emerging diseases, and viruses.

Commented [JB7]: Cerana or mellifera?

2.2.1 Dataset one

Dataset one contains the colony sound of three experimental groups. Each group was treated with unique odorous compounds.

Honeybees were trained with syrup solutions containing different volatile compounds to visit artificial feeding sites approximately 200 meters ~~away~~ from the hive. A feeder containing 50% sucrose solution was placed approximately five meters from the hive, and the marked foragers were caught in a glass tube at the door of the hive. The foragers were gently let out to feed on the feeder. When the foragers had eaten enough syrup, they returned to the hive after hovering over the feeder a few times. This was repeated several times, and when visited by a larger number of foragers, the feeder was slowly placed approximately 10 meters from the hive, and so on, gradually moving the feeder to the target position. When a large number of marked bees were feeding at the target distance, the sound inside the colony was recorded for 10min. Before changing the compounds added to the sucrose solution, we stopped feeding for two days, waiting for the colony to be depleted of food and odors before starting another treatment. The sound files were collected from three different colonies, each colony with two frames. The number of recordings and duration were shown in Table 1. In this dataset, all the colony sound files were collected during winter from November 2020 to January 2021, and very few food sources were available outside. In this way, the artificial food source we provide may be the only food ~~sources~~ for honeybees.

This dataset contains the colony sound of three experimental groups, which were treated with unique odorous compounds at a mass ratio of 0.1% in 50% (w/w) sucrose solution, sucrose solution with 50% concentration was used as blank control. The compound used were ethyl acetate and acetone. The colony sound was labeled ~~"_blank,"~~ ~~"_acetone"~~ ~~"_blank,"~~ ~~"_acetone,"~~ and ~~"_ethyl,"~~ respectively.

127 2.2.2 Dataset two

128 Dataset two collects bee colony sounds concerning the ~~queen's~~ status. The object is to
129 use the colony sound to detect whether there is a queen pupa and whether the pupa has hatched.

This 130 dataset includes honey bee sounds under three scenarios.

131 This work was carried out in June 2021, alternating between spring and summer. It
132 simulated the occurrence of a new queen cell in the colony before swarming. We selected
133 two groups of healthy and strong colonies of *Apis cerana*, each with six frames of honeybees
134 and a normal breeding queen. In the first scenario, we caged the queen and collected colony
135 sounds. In the second stage, we introduced a mature queen pupa into this colony. The original queen
136 was still in the cage and, therefore, would not attack the new queen pupa. Collecting sound data
137 began after a day. In the third stage, we opened the hive every night, checked the pupa condition,
138 and recorded the next day after the new queen emerged. All recordings started around 11:00 am.
139 In this way, we obtained colony sounds in three different queen states. They were labeled as
140 blank, ~~queen~~, ~~queen~~ pupa, ~~queen~~ or ~~queen~~ new queen.

141 2.2.3 Dataset three

142 This dataset contains sounds from bee colonies of different colony sizes. We investigated
143 six bee colonies, including two colonies with two frames, two with four, and two with six.
144 The bee colony sound was recorded at 9:00 am for about three to ten minutes in each of the
145 colonies, and the recorded sound files were labeled as C2-, C4-, and C6, respectively. We estimate the
146 number of bees by weighing the colony. The weight of an empty hive is measured first,
then the

~~146~~ ~~147~~ whole swarm of bees is shaken off into the empty hive, and the mass is measured again. The

~~147~~ 148 mass difference obtained is the total weight of the swarm. We estimated the number of bees per

~~148~~ 149 colony based on the average honeybee weight, which was 94.9mg (Table 2).

2.3 Data processing

The data processing was based on python 3.5.1 and Scikit-learn 1.0.2 (Pedregosa et al., 2012).

2.3.1 Feature extraction

VGGish Embedding. The audio sample was first split into segments of 0.96s. Each 0.96s segment was first resampled to ~~16 kHz~~ using a Kaiser window, and a log-scaled Mel-frequency spectrogram was generated (96 temporal frames, 64 frequency bands). Each audio sample was then passed through CNN from Google's AudioSet project (Gemmeke, Ellis et al. 2017, Hershey, Chaudhuri et al. 2017) to generate a 128-dimensional embedding of the audio. Figure 2 shows the structure of the VGGish network and the work-flow of extracting VGGish embedding.

Mel-frequency Cepstral Coefficient (MFCC). MFCCs are based on the known variation of the human ear's critical bandwidths with frequency. The MFCC technique uses two types of filters: linearly spaced and log arithmetically spaced. The signal is expressed in the Mel frequency scale to capture the phonetically important characteristics of speech. This scale has a linear frequency spacing below 1000 Hz and a logarithmic spacing above 1000 Hz. The MFCC feature extraction procedures are as follows: windowing the sound signal, applying the FFT (Fast Fourier Transform), taking the log, and then warping the frequencies on a Mel scale, followed by applying the inverse DCT (Discrete Cosine Transform). The 13-dimensional MFCC will be combined with the first-order difference coefficients and second-order coefficients difference to

~~166~~ ~~168~~ get the 39-dimensional MFCC feature.

Commented [JB8]: Why resample and drop the frequency? Colony sounds extent into the ultrasonic range.

Commented [JB9]: There is no reason to suspect that bee hearing approximates that of the human ear. Is the MFCC simply correcting the signal produced by the sound card and software like Audacity which is maximized to pick out sounds the are discernible to the human ear?

169 A block diagram of the structure of an MFCC processor is given in Figure 3.

170 2.3.2 Dimension reduction

171 Since the features extracted from the raw data are high-dimensional, it is not conducive to
 172 visualization. It is necessary to use the technique for dimensionality reduction to get 2D
 points

173 from a high-dimensional input vector.

Commented [JB10]: Why not use an X,Y,Z, 3-D visualization?

174 To estimate the impact of dimension reduction, we experimented with the following
 175 dimensionality reduction algorithms: (R1) uniform manifold approximation and projection
 176 (UMAP). UMAP works by learning approximate manifolds from higher dimensional
 Spaces and
 177 mapping them into lower dimensional Spaces (McInnes, Healy et al. 2018); (R2) t-
 distributed
 178 stochastic neighbor embedding(t-SNE) (Van der Maaten and Hinton 2008). This technique
 is a
 179 variation of Stochastic Neighbor Embedding (Becht, McInnes et al. 2019, Diaz-Papkovich,
 180 Anderson-TrocmØ et al. 2019).

181 The multidimensional colony sound feature were narrowed down to two by
 the two

182 algorithms. Machine learning algorithms then classify the reduced feature set.

183 2.3.3 Training classifiers

184 In this paper, we trained four well-known machine learning algorithms, namely decision
 185 tree (DT), K-nearest neighbors (KNN), support vector machine (SVM) classification, and
 186 random forests (RF). DT is a tree-structured classifier. The internal nodes represent the
 features.

187 The branches represent the rules, and each leaf node represents the outcome. KNN (Altman
 188 1992) is a supervised learning model. A majority vote classifies its neighbors in vector
 space,

189 and the data is assigned to the class with the nearest neighbors. SVM classification (Hong
and
190 Cho 2008) aims to create the best decision boundary(which is called a hyperplane) that can
191 segregate n-dimensional space into classes so that the new data point can be put into the
correct
192 category. RF(Breiman 2001) is a classifier that contains a bunch of decision trees. It takes
the
193 prediction from each tree and predicts the final output based on the majority votes of
predictions
194 from those decision trees.

195 We trained all four models on the same feature vectors automatically extracted from the
196 raw audio files in three bee colony datasets. The following feature: (F1) VGGish
embedding;(F2)
197 Mel frequency cepstral coefficients (MFCC) (Davis et al., 1980) are used in training all
four
198 models. We used the mean of the test accuracy as a summary of the ~~model-model's~~
performance. Then
199 the paired Student s t-test was used to check if the difference in the mean accuracy
between the 200 two models is statistically significant.

201 The labeled feature were split into a training set (70%) and a testing test (30%) with the
202 training_test_split procedure from the Python sklearn.model_selection library (Pedregosa,
203 Varoquaux et al. 2011). All these classification models were trained with the training set on
an
204 Intel Xeon E5-2676V3@2.40 GHz x 12 processor with 64 GiB of RAM and 64-bit
Windows 10.

205 2.3.4 Model evaluation

206 Classification accuracy and F1 score were used to evaluate the performance of the ML
207 models. The classification accuracy is the percentage of correct predictions. The F1 score

Commented [JB11]: Is a student T-test the proper statistical test for data such as this?

Commented [JB12]: 79% training compared to 20 testing heavily weights the analysis model towards success with a small set of recordings.

208 integrates information regarding both precision and recall(Chinchor and Sundheim 1993).
 The
 209 balanced accuracy of the classifier on the test set was reported as an average F1 score for
 each 210 class to account for sample-size imbalances among classes.
 211 The data processing work-flow is presented in Figure 4.

3. Results

3.1 The performance of models on dataset one

3.1.1 Audio signal

Two different compounds were added separately into the sucrose solution. Figure 5 presents the log spectrogram of the bee colony sound. We can see that: 1) after being treated with a compounds-sucrose solution, the low-frequency sound in the bee colony increased; 2) the bee

colony sound increased more significantly when feeding with the acetone-sucrose solution than

219 when feeding with ethyl-sucrose solution, and there was a significant increase in bee colony 220 sound around 130hz.

3.1.2 Dimensional reduction of audio feature

Figure 6 shows the output of VGGish embedding and MFCC feature dimensionality reduction in dataset one. In the two-dimensional diagram, it is evident that the MFCC feature 224 overlaps after dimensionality reduction, while the VGG embedding can better distinguish the 225 sound in these three situations.

3.1.3 Model evaluation

Table 3 and Table 4 summarize the results of four machine learning methods. VGGish embedding performs significantly better than the MFCC feature ($P < 0.005$) and shows an advantage of about 30% over MFCC feature in all four machine learning methods, among which KNN 230 performs best, achieves achieving an accuracy of 94.79%.

3.2 The performance of models on dataset two

3.2.1 Audio signal

From the log spectrogram of the bee colony sound (Figure 7), the colony with a queen

Commented [JB13]: 1 drop of toluene in a multi-story, A. mellifera colony, produces and immediate, and easy to hear roaring sound. Your dosing levels are very high for an insect that can detect many odors at the parts per trillion range.

pupae seemed more active than the colonies in the other two conditions. The signal around 250hz 235 and 500hz are stronger in the sound collection Queen pupa and New queen than in the sound 236 collection Blank.

237 3.2.2 Dimensional reduction of audio feature

238 Compared with the MFCC dimensionality reduction diagram (Figure 8), the scatter plot of 239 VGG embedding after dimensionality reduction has less overlap.

240 3.2.3 Model evaluation

241 The MFCC feature performs slightly better than VGGish embedding and shows an
242 advantage of about 4 percent in all four machine learning methods (Table 3, Table 4), ~~but~~
243 ~~Still~~, the
244 difference was not statistically significant ($P > 0.05$). Moreover, KNN performed best, and
244 achieved an accuracy of 90%.

245 3.3 The performance of models on dataset three (Identifying colony size)

246 3.3.1 Audio signal in dataset three

247 This dataset includes bee colony sounds from 3 different colony sizes: C2) bee colony size
248 of
248 about 7500 work bees; C4) bee colony size of about 11000 work bees; C6) bee colony size
249 of ~~of~~ about 17000 work bees. Figure 9 presents the log spectrogram of the bee colony
250 sound signals of 250 ~~in~~ this dataset.

251 3.3.2 Dimensional reduction of audio feature

252 The output of UMAP (Figure 10) exhibits the VGGish embedding and MFCC ~~feature~~ of
253 colony sound in dataset three.

254 3.3.3 Model evaluation

255 The accuracy of four machine learning models using different colony sound features on

dataset three is shown in Table 3. VGGish embedding has an advantage over ~~the~~ MFCC feature of

about 20 percent in all four machine learning methods, and the difference was statistically significant ($P < 0.05$). Moreover, KNN performed best and achieved an accuracy of 91%.

3.4 The influence of different dimensionality reduction methods

~~In order to~~ To test the effects of different dimensionality reduction algorithms on the accuracy of the models, ~~we~~ We have chosen two dimensionality reduction algorithms, namely UMAP and t-SNE.

Figure 11 exhibits the results of dimensionality reduction of dataset one using the t-SNE algorithm, compared with the output of the UMAP algorithm (Figure 6), UMAP performs better than t-SNE feature in separating bee colony sounds. Table 4 shows the accuracy of four machine learning methods trained by ~~two-dimension~~ two-dimension factors obtained by UMAP and t-SNE. The original sound feature used by those dimensional reduction algorithms were the MFCC feature. The results show that UMAP performs better than t-SNE in almost all datasets and all machine learning methods.

270 4. Discussion

271 Hive monitoring based on colony sound has made a lot of research achievements in recent
 272 years (Terenzi, Cecchi et al. 2020) and has become increasingly popular with ~~in~~ many
 273 international companies such as Arnia, Bee Hero, Nectar, and Broodminder 274
 (https://www.umn.edu/bee/monitoringconference_2020/).

275 In this paper, We compared the performance of VGGish embedding and MFCC feature of
 276 bee colony sound in four classification algorithms. The result in Table 3 indicated that all
 four
 277 classification algorithms could generate prediction accuracy percentages that are better than
 278 ~~'chance'~~ 'chance' based percentages. In all classification methods, the VGGish feature can
 guarantee more
 279 than 80% testing accuracy, among which KNN has the best performance of 94%. The
 testing
 280 accuracy of the MFCC feature varies a lot between different datasets. In ~~dataset~~ one and
 three,
 281 the MFCC could only achieve an accuracy of about 69%, while in dataset two, it achieved
 an
 282 accuracy of 90%. Results (Table 3, Table 5) show that the difference between the two
 features in 283 datasets one and three is statistically significant ($P < 0.005$). At the same time,
 in dataset two, there 284 is ~~not any~~ significant difference between the two models ($P > 0.005$).

285 We confirm that the VGGish embedding applies to bee colony sound classification and
 286 performs more stability than the MFCC feature among different datasets. This may be due
 to the
 287 MFCC ~~being is~~ highly dependent on data and feature, which causes weak generalization
 ability due
 288 to insufficient bee colony data and the similarity of bee colony sound. The VGGish network
 is
 289 trained on a more extensive and general Audio set, which means a better generalization
 ability.

290 Our results suggest that different compounds do lead to different responses in the bee

colony (Figure 6, Table 3, Table 4), ~~it~~which further confirms the results of previous studies (Bromenshenk, Henderson et al. 2009, Sharif, Wario et al. 2020, Zhao, Deng et al. 2021, Yu, Huang et al. 2022), and ~~moreover~~, verifies the applicability of VGGish embedding for the classification of bee colony sounds. As seen from the log spectrum of bee colony sounds (Figure 295 5), the acetone-sucrose solution and acetone ethyl-sucrose solution would agitate the colony compared to the sucrose solution. The ~~low-low~~-frequency amplitude was much larger when treated with acetone than when treated with sucrose solution. This may be due to the fact that acetone stimulates bee colonies more strongly than ethyl acetate at the same concentration, and low concentrations of ethyl acetate were mildly attractive to bees (Schmidt and Hanna 2006).

The MFCC feature performs better in dataset two (Table 3, Table 4). This may be because of the fact that the sound changes fundamentally during bee swarming (Michelsen, Kirchner et al. 1986). Thus, it is easier for the standard MFCC ~~feature~~ to capture the character in colony sounds. Dataset three is ~~relatively small~~small. The total duration of sound in dataset three is less than one hour, and the machine learning models trained by the VGGish embedding could still achieve an accuracy of around 90%, which may be because the VGGish could better capture the distinctions among the ~~dataset~~datasets. We have compared two different dimensionality reduction algorithms (Figure 11, Table 5), and UMAP performs better than the t-SNE in ~~almost every~~every situation. The secret of UMAP lies in its ability to infer local and global structures while maintaining relative global distances in low-dimensional space. The result also shows that UMAP performed better in separating different colony sounds.

310 ~~311~~—In summary, the results of this paper indicate that the combination of VGGish embedding 311 ~~312~~ and the KNN method has achieved the highest accuracy on the testing set of all three datasets 313 (Table 3, Table 4, Table 5).

314 Several ways in which this research can be improved are given below:

315 1) Beehive sound samples are ~~relatively few~~few, and only one type of microphone is used
for

316 collecting the sound, which causes a lack of data diversity and affects the ~~model'~~'s-model's
317 generalizability. A more comprehensive data set must be attained in future work to train the

~~317~~318 318 system and improve the ~~model'~~'s-model's generalizability.

~~319~~319 2) Expand the application of the model: in this study, we applied VGGish
embedding in the

~~320~~ 320 ~~classification~~classifications of three datasets. Beehive sound can be influenced by
many other factors, such as

321 the invasion of natural enemies and parasites. Subsequent studies can check how VGGish 322
embedding performs in these scenarios.

323 Acknowledgments

324 The authors would like to thank the Sericultural & Apicultural Research Institute for
 325 permission to conduct this study and ~~for~~ help during data collection. We also thank Xuewen
 326 Zhang, Chuntao Zhou, Chunhui Miao, and Xinqiu Huang, who have generously shared
 their time
 325326 327 and expertise.

328

329 References

330 Altman, N. S. (1992). "An introduction to kernel and nearest-neighbor nonparametric regression." The American 331
Statistician **46**(3): 175-185.

332

333 Becht, E., et al. (2019). "Dimensionality reduction for visualizing single-cell data using UMAP." Nature biotechnology
334 **37**(1): 38-44.

335

336 Braga, A. R., et al. (2020). "A method for mining combined data from in-hive sensors, weather and apiary inspections
337 to forecast the health status of honey bee colonies." Computers and Electronics in Agriculture **169**: 105161.

338

339 Breiman, L. (2001). "Random forests." Machine learning **45**(1): 5-32.

340

341 Bromenshenk, J. J., et al. (2009). Honey bee acoustic recording and analysis system for monitoring hive health, 342
Google Patents.

343

344 Cejrowski, T., et al. (2018). Detection of the bee queen presence using sound analysis. Asian Conference on 345
Intelligent Information and Database Systems, Springer.

346

347 Chinchor, N. and B. M. Sundheim (1993). MUC-5 evaluation metrics. Fifth Message Understanding Conference
(MUC348 5): Proceedings of a Conference Held in Baltimore, Maryland, August 25-27, 1993.

349

350 Diaz-Papkovich, A., et al. (2019). "UMAP reveals cryptic population structure and phenotype heterogeneity in large
351 genomic cohorts." PLoS genetics **15**(11): e1008432.

352

353 Dietlein, D. G. (1985). "A method for remote monitoring of activity of honeybee colonies by sound analysis." Journal
354 of Apicultural Research **24**(3): 176-183.

355

356 Ferrari, S., et al. (2008). "Monitoring of swarming sounds in bee hives for early detection of the swarming period." 357
Computers and electronics in agriculture **64**(1): 72-77.

358

359 Gemmeke, J. F., et al. (2017). Audio set: An ontology and human-labeled dataset for audio events. 2017 IEEE 360
 international conference on acoustics, speech and signal processing (ICASSP), IEEE.

361

362 Hershey, S., et al. (2017). CNN architectures for large-scale audio classification. 2017 IEEE 363
 international conference on acoustics, speech and signal processing (ICASSP), IEEE.

364

365 Hong, J.-H. and S.-B. Cho (2008). "A probabilistic multi-class strategy of one-vs.-rest support vector machines for 366
 cancer classification." Neurocomputing **71**(16-18): 3275-3281.

367

368 Klein, A.-M., et al. (2007). "Importance of pollinators in changing landscapes for world crops." Proceedings of the 369
royal society B: biological sciences **274**(1608): 303-313.

370

371 Kulyukin, V., et al. (2018). "Toward audio beehive monitoring: Deep learning vs. standard machine learning in 372
 classifying beehive audio samples." Applied Sciences **8**(9): 1573.

373

374 Kumar, A. and B. Raj (2017). "Deep cnn framework for audio event recognition using weakly labeled web data." arXiv
 375 preprint arXiv:1707.02530.

376

377 McInnes, L., et al. (2018). "Umap: Uniform manifold approximation and projection for dimension reduction." arXiv
 378 preprint arXiv:1802.03426.

379

380 Meikle, W. and N. Holst (2015). "Application of continuous monitoring of honeybee colonies." Apidologie **46**(1): 10-
 381 22.

382

383 Michelsen, A., et al. (1986). "The tooting and quacking vibration signals of honeybee queens: a quantitative analysis."
384 Journal of Comparative Physiology A **158**(5): 605-611.

385

386 Muda, L., et al. (2010). "Voice recognition algorithms using mel frequency cepstral coefficient (MFCC) and dynamic
387 time warping (DTW) techniques." arXiv preprint arXiv:1003.4083.

388

389 Murphy, F. E., et al. (2015). b+ WSN: Smart beehive for agriculture, environmental, and honey bee health 390 monitoring
Preliminary results and analysis. 2015 IEEE Sensors Applications Symposium (SAS), IEEE.

391

392 Pedregosa, F., et al. (2011). "Scikit-learn: Machine learning in Python." the Journal of machine Learning research
12:
393 2825-2830.

394

395 Qandour, A., et al. (2014). "Remote beehive monitoring using acoustic signals."

396

397 Schmidt, J. O. and A. Hanna (2006). "Chemical nature of phagostimulants in pollen attractive to honeybees." Journal
398 of Insect Behavior **19**(4): 521-532.

399

400 Sethi, S. S., et al. (2020). "Characterizing soundscapes across diverse ecosystems using a universal acoustic feature 401
set." Proceedings of the National Academy of Sciences **117**(29): 17049-17055.

402

403 Sharif, M. Z., et al. (2020). "Soundscape indices: new features for classifying beehive audio samples." Sociobiology 404
67(4): 566-571. 405

406 Shi, L., et al. (2019). "Lung sound recognition algorithm based on vggish-bigru." IEEE Access **7**: 139438-139449.

407

408 Simonyan, K. and A. Zisserman (2014). "Very deep convolutional networks for large-scale image recognition." arXiv
 409 preprint arXiv:1409.1556. 410

411 Terenzi, A., et al. (2020). "On the importance of the sound emitted by honey bee hives." Veterinary Sciences 7(4):
 412 168.

413

414 Van der Maaten, L. and G. Hinton (2008). "Visualizing data using t-SNE." Journal of machine learning research 9(11).
 415

416 Yu, B., et al. (2022). "A matter of the beehive sound: Can honey bees alert the pollution out of their hives?" 417
Environmental Science and Pollution Research: 1-11.

418

419 Zhao, Y., et al. (2021). "Based investigate of beehive sound to detect air pollutants by machine learning." Ecological
 420 Informatics 61: 101246.

421

422

Figure 1

The hardware system used to obtain bee colony sound.

The microphone is placed inside the beehive, ~~then t.~~ The sound signal captured by the microphone is converted to a digital signal by the digital sound card, then transmitted to the PC and saved on a hard disk for further analysis.

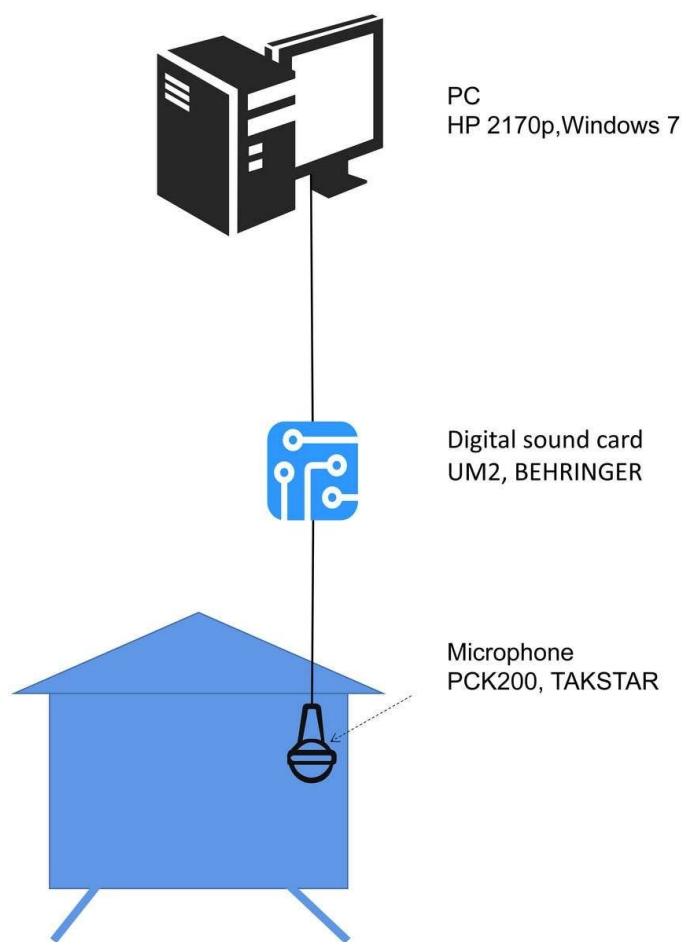


Figure 2

An overview of the structure of the VGGish network.

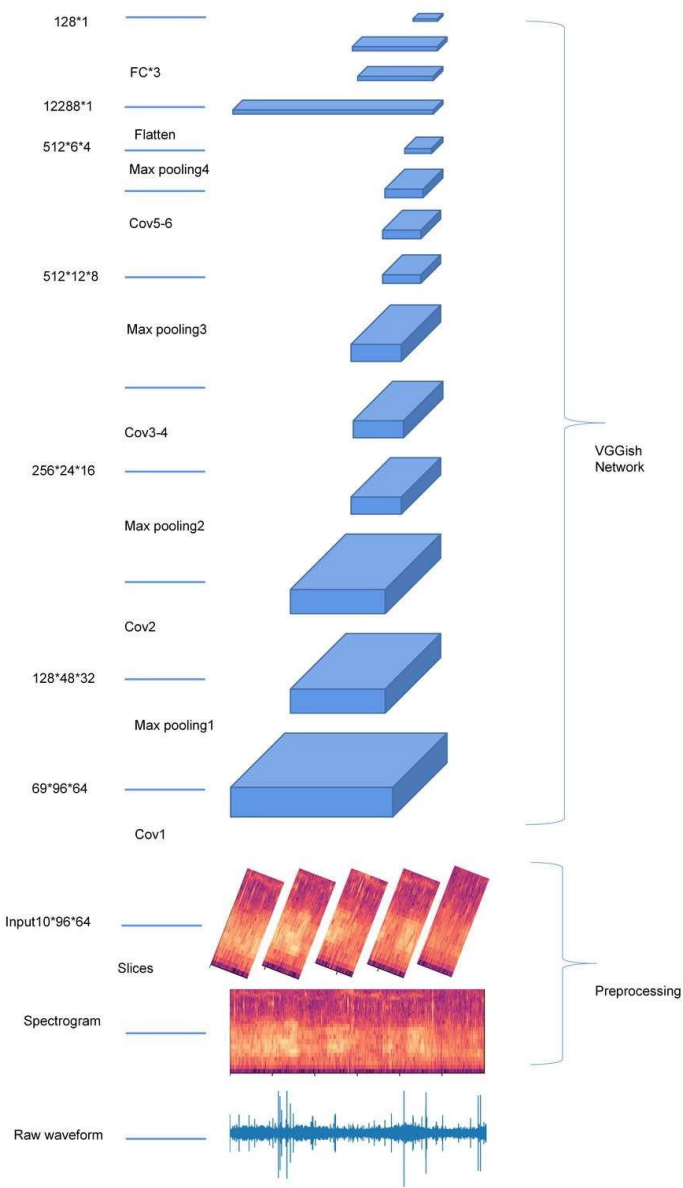


Figure 3

Block diagram of the MFCC processor.

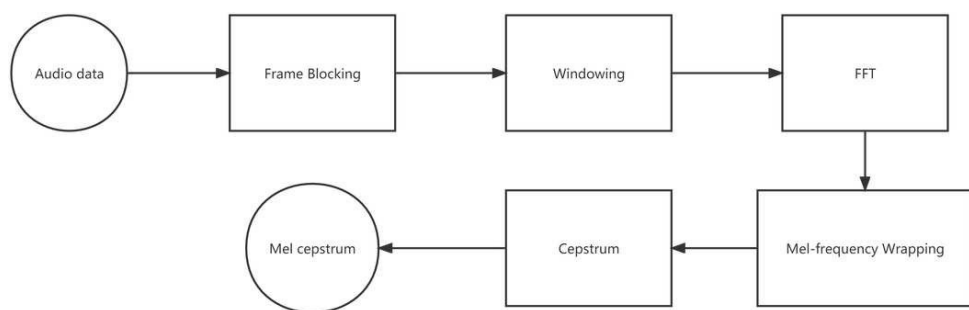


Figure 4

Overview of the approach adopted for the acoustic classification of beehive sounds workflow [\(work-flow\)](#).

The original audio files [\(files\)](#) (.wav format) containing recordings of beehive sounds were manually classified [\(classified\)](#) into corresponding scenarios. Then, the MFCC and VGGish embedding were used to extract the audio features, respectively. Dimensionality reduction was performed using the UMAP method for the two sets of feature data. After that, the resulting data set was split into 70% for the training/development set and 30% for the testing data set. ~~The Finally, the~~ test data set was used to evaluate the performance of the classifiers [\(classifiers\)](#) in correctly assigning the beehive sound to the respective scenario.

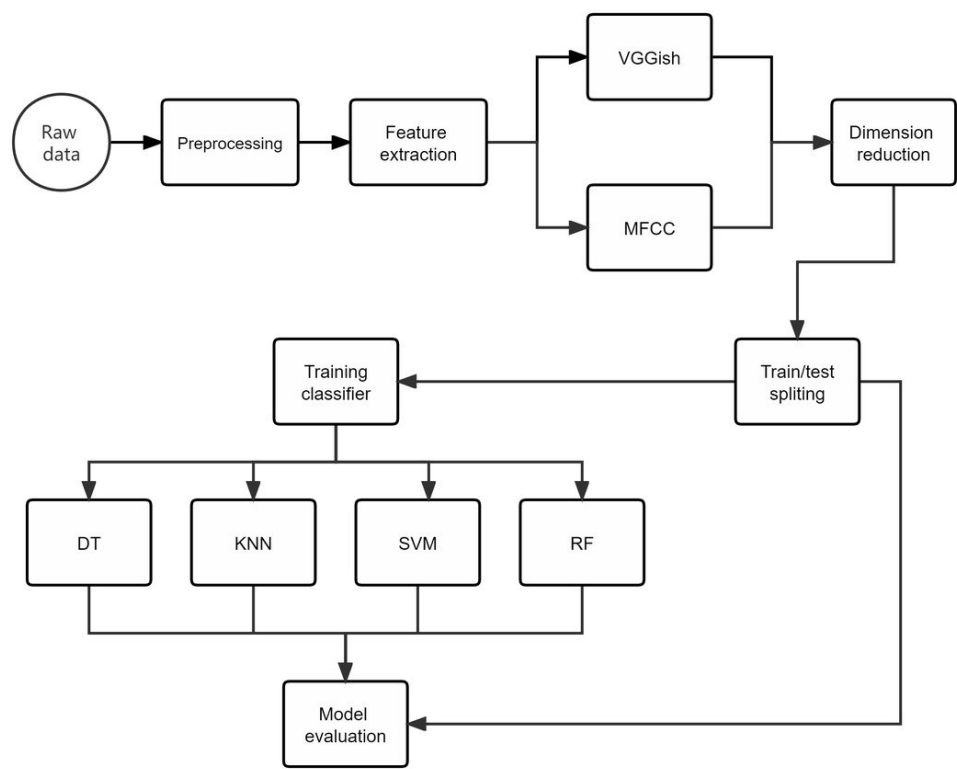


Figure 5

Log spectrum of bee colony sounds from dataset one.

Left: Acetone(treated with acetone-sucrose solution); Middle: Ethyl(treated with ethyl acetate-sucrose solution); Left: Blank(treated with sucrose solution).

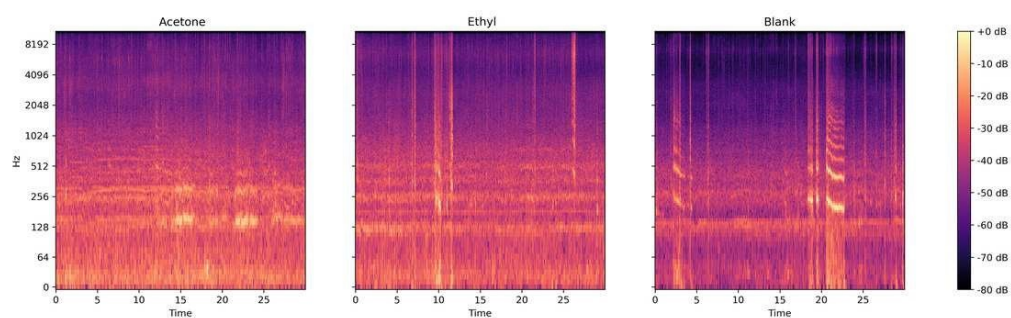
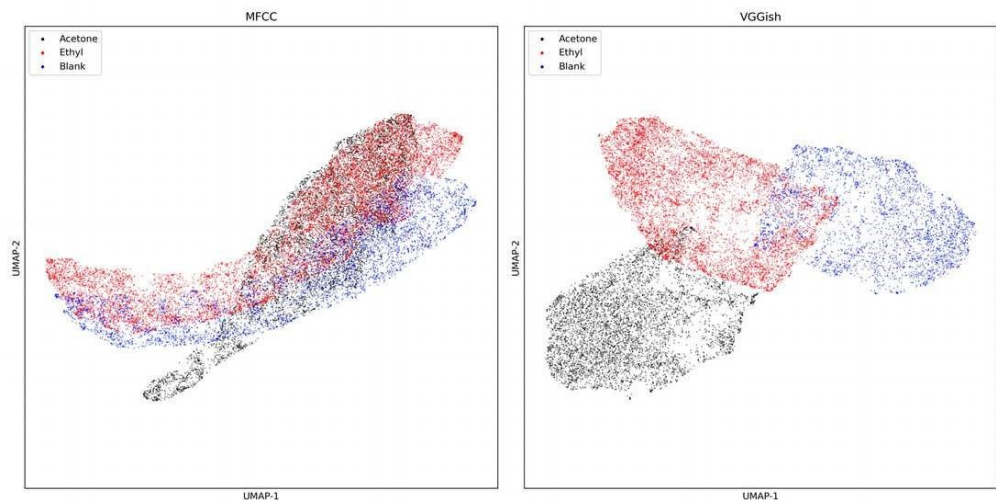


Figure 6

UMAP dimension reduction of sound features from dataset one.



Log spectrum of bee colony sounds of dataset two.

Left: Normal situation; Middle: Queen pupa inside colony; Left: New queen emerged(two queens in the colony).

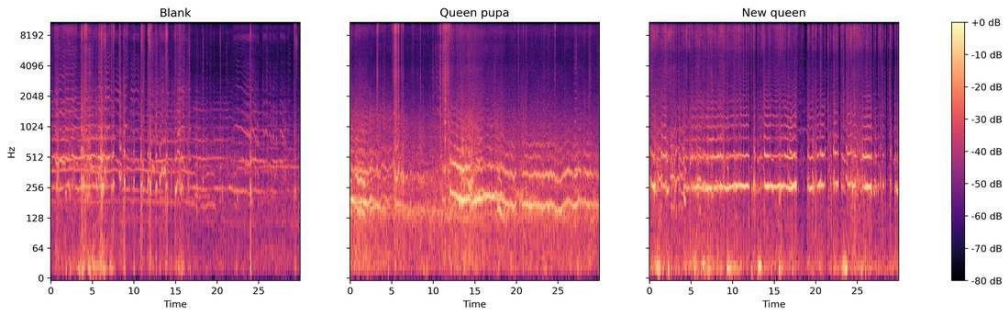
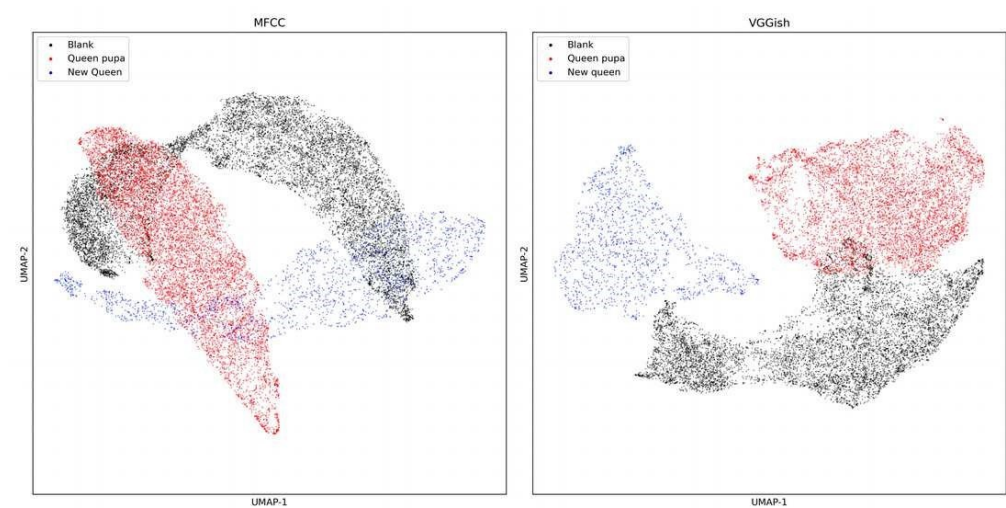


Figure 7

UMAP dimension reduction of sound features from dataset two.



Log spectrum of bee colony sounds for dataset three.

Left: Colony size of around 7500 bees(C2); Middle: Colony size of around 11000 bees(C4); Right: Colony size of around 17000 bees(C6).

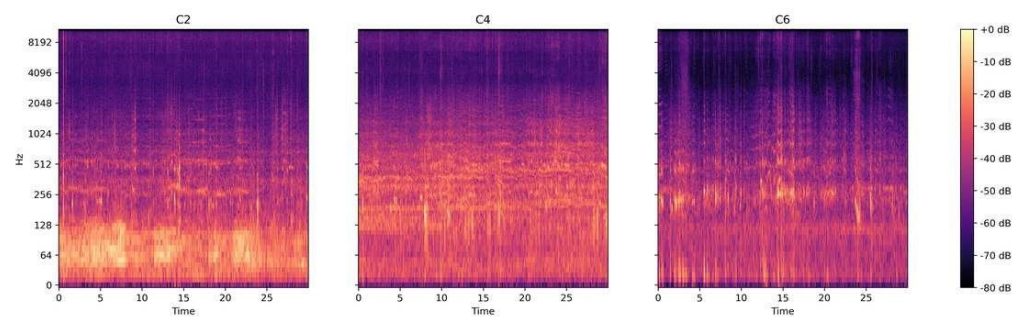
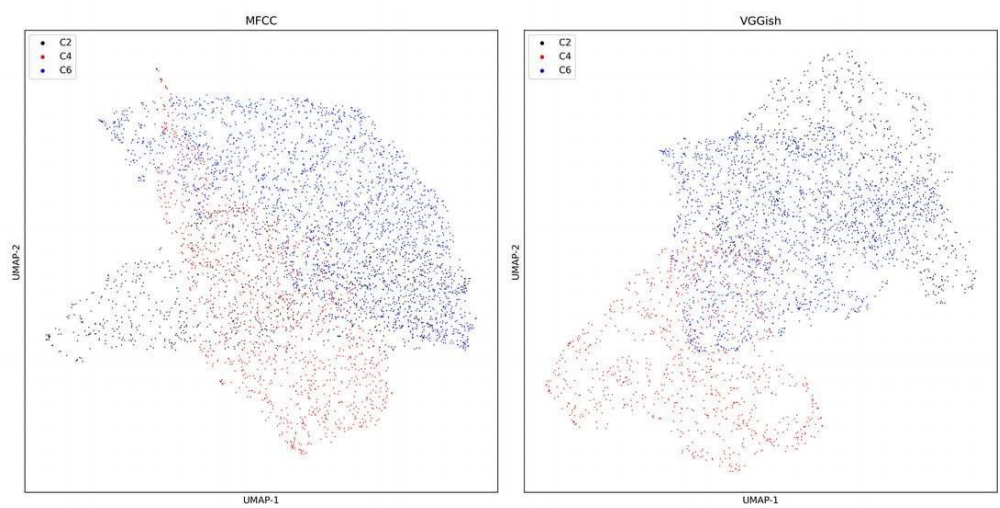


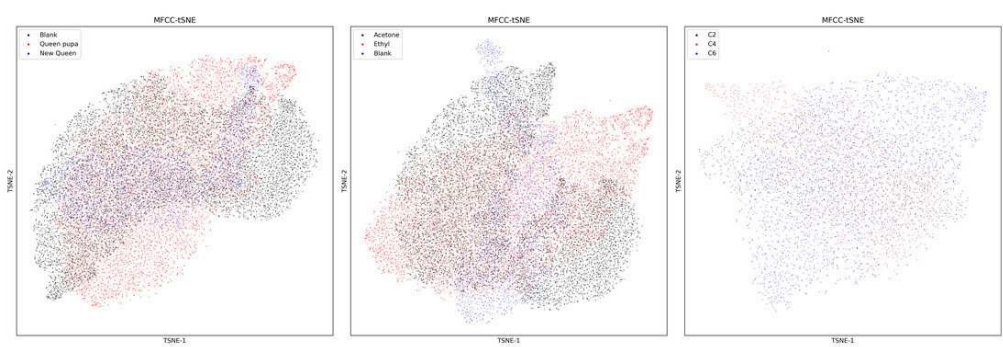
Figure 8

UMAP dimension reduction of sound features for dataset three.



MFCC features of three datasets after t-SNE dimensionality reduction.

Left: MFCC feature using t-SNE dimensionality reduction on dataset two; Middle: MFCC feature using t-SNE dimensionality reduction on dataset one; Right: MFCC feature using t-SNE dimensionality reduction on dataset three.



1

Table 1

The size of each colony used in dataset three.

"N frames" denotes the number of frames in the colony; "Total bee weight(Kg)" represents the total weight of each colony; <N worker bees= denotes the approximate number of worker bees in each colony.

(on next page)

Colony	N frames	Total bee weight(Kg)	N worker bees
1#	2	0.723	7619
2#	2	0.685	7218
3#	4	1.010	10643
4#	4	1.095	11538
5#	6	1.650	17387
6#	6	1.580	16649

2
3
4

1

Table 2

An overview of the datasets collected in order to identify compounds in nectar and queen's presence.

<Scenario= "N recordings" denotes the number of individuals with buzzing sounds recorded; "Total duration" represents the total recording time in each case; N colonies denotes the number of colonies in which we recorded sounds; N frames represent the colony size.

(on next page)

Datasets	Scenario	N colonies	N frames	N Recordings	Total Duration
Dataset one Identify compounds		3	2	6	50min
	Acetone	3	2	9	90min
	Ethyl acetate	3	2	11	111min
Dataset two Identify queen state	Blank	2	6	12	131min
	New queen pupa	2	6	9	101min
	New queen	2	6	3	23min
Dataset three Identify colony size	C2	2	2	2	12min
	C4	2	4	2	15min
	C6	2	6	2	29min

1

Table 3

Accuracy of machine learning models using different (different) colony sound features on three |
colony sound datasets

(on next page)

Datasets	Dataset 1				Dataset 2				Dataset 3			
Algorithm	KNN	DT	RF	SVM	KNN	DT	RF	SVM	KNN	DT	RF	SVM
VGGish	94.79%	93.45%	94.43%	91.56%	86.58%	85.14%	85.94%	81.46%	91.08%	88.81%	89.23%	89.15%
MFCC	69.09%	66.28%	69.17%	68.29%	90.48%	88.45%	89.95%	87.25%	66.04%	65.78%	65.13%	68.05%

¹
Table 4

F1-score of machine learning models using di erent (different) colony sound features on three colony sound datasets

(on next page)

Datasets	Dataset 1				Dataset 2				Dataset 3			
Algorithm	KNN	DT	RF	SVM	KNN	DT	RF	SVM	KNN	DT	RF	SVM
VGGish	94.79%	93.45%	94.42%	91.55%	86.58%	85.17%	85.93%	81.49%	91.06%	88.82%	89.21%	89.03%
MFCC	68.24%	66.32%	68.49%	65.26%	90.13%	88.44%	89.63%	85.41%	65.73%	65.74%	64.80%	65.09%

¹
Table 5

Comparison of different (different) dimensionality reduction methods

|

(on next page)

Datasets	Dataset 1				Dataset 2				Dataset 3			
Algorithm	KNN	DT	RF	SVM	KNN	DT	RF	SVM	KNN	DT	RF	SVM
t	69.09%	66.28%	69.17%	68.29%	90.48%	88.45%	89.95%	87.25%	66.04%	65.78%	65.13%	68.05%
	51.62%	54.85%	55.07%	56.63%	62.64%	65.38%	66.83%	66.42%	52.24%	55.04%	57.45%	60.18%