# ICSD: An Open-source Dataset for Infant Cry and Snoring Detection

Qingyu Liu[a], Longfei Song[a], Dongxing Xu[b], Yanhua Long[a,*]

[a]*Shanghai Engineering Research Center of Intelligent Education and Bigdata, Shanghai Normal University, Shanghai, China*
[b]*Unisound AI Technology Co., Ltd., Beijing, China*

## Abstract

The detection and analysis of infant cry and snoring events are crucial tasks within the field of audio signal processing. While existing datasets for general sound event detection are plentiful, they often fall short in providing sufficient, strongly labeled data specific to infant cries and snoring. To provide a benchmark dataset and thus foster the research of infant cry and snoring detection, this paper introduces the **I**nfant **C**ry and **S**noring **D**etection (ICSD) dataset, a novel, publicly available dataset specially designed for ICSD tasks. The ICSD comprises three types of subsets: a real strongly labeled subset with event-based labels annotated manually, a weakly labeled subset with only clip-level event annotations, and a synthetic subset generated and labeled with strong annotations. This paper provides a detailed description of the ICSD creation process, including the challenges encountered and the solutions adopted. We offer a comprehensive characterization of the dataset, discussing its limitations and key factors for ICSD usage. Additionally, we conduct extensive experiments on the ICSD dataset to establish baseline systems and offer insights into the main factors when using this dataset for ICSD research. Our goal is to develop a dataset that will be widely adopted by the community as a new open benchmark for future ICSD research.

*Keywords:* Infant Cry, Snoring, Sound event detection, Dataset, Data collection

## 1. Introduction

In the field of machine perception, datasets are of critical importance. Particularly in the domain of audio, each research direction has its specialized dataset, such as the LibriSpeech dataset [1] for Automatic Speech Recognition (ASR), the VoxBlink dataset [2] for speaker verification, the Voice-Bank+DEMAND dataset [3] for speech enhancement, and the LJSpeech dataset [4] for speech synthesis, etc. In the realm of sound event detection, the Detection and Classification of Acoustic Scenes and Events (DCASE) has initiated and provided specialized datasets, such as the DESED dataset for DCASE task4 [5]. This dataset is targeted for applications in domestic environments, focusing on 10 classes of sound events such as doorbell rings, dog barks, speech, and so on. However, for application scenarios involving the detection of infant crying and snoring sounds, there are currently no publicly available datasets in the literature.

The detection of infant crying and snoring has significant implications for both health and well-being. Infant crying is a critical form of communication for infant, indicating various needs such as hunger, discomfort, or distress. Continuous monitoring and prompt response to infant crying are essential, particularly during the night when parents are sleeping. On the other hand, snoring is a common symptom of sleep disorders such as sleep apnea, which can lead to serious health complications if left untreated. Regular detection and analysis of snoring patterns can provide valuable insights for medical diagnosis and treatment. In a domestic environment, the development of effective detection systems for these sounds can greatly enhance the quality of life for families, improve infant care, and contribute to sleep disorder diagnosis and treatment.

Numerous existing datasets have collected audio recordings of infant cry and snoring for potential acoustic event detection detection or classification studies. For instance, Google Research extracted 120,459 clips from the original Audioset[6] in 2020 for strong label annotation, including 556 infant cry audio files and 465 snoring audio files. However, the quantity of these strongly labeled data is not sufficient for sound event detection specific to infant crying and snoring. ESC-50 [33] comprises 40 five-second audio clips for each of 50 different classes, including categories for infant crying and snoring. However, rather than for detection tasks, this dataset has been specifically designed for multi-label sound event classification tasks. Baby Chillanto database [15] includes 2267 one-second audio clips of infant crying. However, this dataset is primarily designed for direct medical pathologies classification and identification rather than routine monitoring. Two snoring datasets from Kaggle[1], one for general snoring [31] and another for female and male snoring [32], provide 1500 one-second snoring audio clips totally. However, the majority of the samples in these two datasets are set against quiet backgrounds, which does not reflect the complexity of real-world scenarios.

The existing datasets are either too small or not readily adaptable for research purposes in the joint sound event detection of

---

*The first two authors (Qingyu Liu, Longfei Song) contributed equally to this work. Yanhua Long is the corresponding author. e-mail: yanhua@shnu.edu.cn

[1]https://www.kaggle.com/

Table 1: Basic Description of Source Datasets Used for ICSD Dataset Creation

| Database | Creator | Data | Papers |
|---|---|---|---|
| Audioset | Google Research | **Total 3911** (Infant cry 1815, Snoring 2096) | [7, 8, 9, 10, 11, 12, 13, 14] |
| BCD | National Institute of Astrophysics and Optical Electronics, CONACYT Mexico | **Total 2268** (Infant cry) | [27, 16, 17, 18, 19, 20, 21, 22, 23, 24] |
| Donate A Cry | `https://github.com/gveres/donateacry-corpus` | **Total 457** (Infant cry) | [26, 27, 28, 29, 30] |
| Self-collected database | Collected by Tareq Khan | **Total 500** (Snoring) | / |
| FM Snoring Dataset | Self-recorded | **Total 1000** (Snoring) | / |
| ESC-50 | Warsaw University of Technology, Institute of Electronic Systems | **Total 80** (Infant cry 40, Snoring 40) | [34, 11, 12] |
| SINS | DCASE 2018 Challenge Task5 | **Total 72984** (Absence 18860, Cooking 5124, Dish washing 1424, Eating 2308, Other 2060, Social activity 4944, Vacuum cleaning 972, Watching TV 18648, Working 18644) | [5] |
| MUSAN | Center for Language and Speech Processing, The Johns Hopkins University | **Total 10831** (Music 660, Noise 930, Speech 426) | [37] |

infant crying and snoring sounds. These limitations of current available resources highlight the need for a large and easily accessible dataset to facilitate more advanced research in the field of infant cry and snoring detection. In this study, we aim to create a new open-source dataset that can foster significant advancements in the field of infant cry and snoring detection. The main contributions are as follows:

1) We present the open-source Infant Cry and Snoring Detection (ICSD) dataset, a unique resource that includes weakly labeled, synthetic, and real strongly labeled data annotations.

2) We establish three baseline systems on the ICSD dataset, two of which are derived from the baselines of DCASE Task 4 Challenge [5], while the third is a novel approach proposed by our team. This setup provides a comprehensive reference point for future ICSD studies.

3) We provide a detailed analysis and discussion of the results and challenges in infant cry and snoring detection. By presenting the ICSD dataset and our preliminary findings, we hope to promote the future development in the detection of infant cry and snoring.

The ICSD dataset and the baseline are publicly available at `https://github.com/QingyuLiu0521/ICSD/`.

## 2. Source Datasets

Table 1 introduces the source datasets utilized in this research. These datasets covers a variety of audio samples, including the background sounds from various environments, the foreground snoring and infant cries. Each dataset is unique in its composition, providing a comprehensive range of audio samples that significantly aid the data wide-domain coverage of our research. The table provides details about the creator of each dataset, the total number of data clips/recordings it contains, and references to papers that have previously utilized the dataset. This diversity of data have allowed us to develop and validate our models effectively. A brief description of each source dataset is presented below.

**Audioset**: AudioSet [6] is a comprehensive audio event dataset comprising over 2 million human-annotated, 10-second clips from YouTube. The dataset uses a hierarchical ontology of 632 event classes for annotation, allowing for complex labelling of sounds. The aim of AudioSet is to aid in developing high-performance audio event recognizers, similar to ImageNet's impact in the image domain. In 2020, additional strong labelling was performed on select clips, with annotators marking every distinct sound event and indicating their start and end times. AudioSet's strongly labeled data includes 456 classes. Given our research focus on Infant Cry and Snoring Detection, we only extracted the Infant Cry and Snoring categories from AudioSet, including 1713 Snoring and 1391 Infant Cry clips of weakly labeled data, and 383 Snoring and 424 Infant Cry clips of event-based strongly labelled data.

**Baby Chillanto Database (BCD)**: The Baby Chillanto Database [15] is a publicly available resource compiled by the National Institute of Astrophysics and Optical Electronics, CONACYT, Mexico. It is specifically designed for infant cry pathology classification tasks, featuring a diverse array of infant crying sounds. Each infant cry audio was segmented into one-second duration, and is grouped into five categories, including asphyxia, deaf, hunger, normal and pain. In total, the BCD comprises 2,268 infant cry samples, making it a comprehensive resource for researchers studying infant cries.

**Donate A Cry**: Donate A Cry [25] is a database designed to

Table 2: Overall Statistics of ICSD dataset (#clips, 10s/clip).

| Set | Train | | Validation | | Test | |
|---|---|---|---|---|---|---|
| | InfantCry | Snoring | InfantCry | Snoring | InfantCry | Snoring |
| Weakly labeled | 1699 | 1577 | 189 | 176 | None | None |
| Real strongly labeled | 338 | 305 | 43 | 39 | 43 | 39 |
| Synthetic strongly labeled | 4000 | 4000 | 500 | 500 | 500 | 500 |

aid in identifying the needs of infants through their crying patterns. The cries in the dataset were recorded from babies aged between 0 and 2 years old. Following data cleaning, the dataset was streamlined to include five primary categories: hunger, burping, belly pain, discomfort, and tiredness. The dataset comprises a total of 457 files, with each audio clip having a duration of 7 seconds.

**Self-collected database**: This self-recorded snoring dataset [31] is an organized collection of audio samples, all of which are one-second clips, divided into two primary classes: one for snoring sounds and the other for non-snoring sounds. Each class has 500 samples and self-collected from various online sources. Among the 500 snoring samples, 363 samples consist of snoring sounds of children, adult men and adult women without any background noise. The remaining 137 samples include snoring superimposed on non-snoring background sounds.

**FM Snoring Dataset**: The FM Snoring Dataset [32] is Female and Male Snoring Dataset is a well-structured collection of audio samples, each of which is a one-second clip with a sampling rate of 44,100 Hz. This dataset is categorized into two distinct classes, each containing 500 audio clips, representing snoring sounds from females or males.

**ESC-50**: ESC-50 [33] is a collection of 2000 environmental recordings evenly distributed across 50 classes of various common sound events. These are grouped into 5 categories: animal sounds, natural soundscapes, human sounds, interior/domestic sounds, and exterior/urban noises, with each category containing 10 classes. Each class contains 40 clips with 5 seconds long for each. In our research, we only extracted the portion of Infant Cry and Snoring data from the ESC-50 dataset.

**SINS**: SINS [35] is a comprehensive collection of continuous recordings from a person spending a week in a holiday home. This dataset is meticulously organized with each file segmented into 60-second intervals, and the events categorized into specific activities such as absence, cooking, dish washing, eating, social activities, vacuum cleaning, watching TV, and working. For our research, we utilized a derivative of the SINS dataset, specifically used in the DCASE 2018 Challenge Task 5 [49]. This subset provided us with 72,984 clips of 10-second duration each, yielding approximately 200 hours of data.

**MUSAN**: MUSAN [36] is a dataset that holds approximately 109 hours of audio comprising music, speech, and noise. It embraces diversity with music from various genres, speech from twelve different languages, and a broad range of noises.

## 3. ICSD Dataset

Table 2 provides an overview of our proposed ICSD dataset. The entire dataset is divided into three subsets: train, validation, and test sets for building the ICSD system. Three types of data clips are collected to construct these subsets: weakly labeled clips with clip-level event annotations, real strongly labeled clips with manually annotated event-based time-stamps, and synthetic clips with synthesized event-based time-stamps. Since we focus solely on the ICSD event detection problem rather than an audio tagging task, we did not include the weakly labeled data in the test set.

As shown in Table 2, we have achieved a relatively balanced sample distribution for the infant cry and snoring detection system in each subset. For instance, we have 1,699 and 1,577 weakly labeled clips of 10-second duration for the infant cry and snoring training sets, respectively. For the real strongly labeled data, there are 338 clips for infant cry model training and 305 clips for snoring model training. In this version of ICSD, we have created 8,000 clips for training and 1,000 clips for model validation and testing.

All the data clips in Table 2 are derived from the eight source datasets as listed in Table 1. Due to the variations in audio formats, sampling rates, recording channels, and audio quality among different sources, these datasets were not directly used for ICSD data generation. Instead, a data format unifying process was conducted to clean up all these source samples. The following sections provide detailed descriptions of the data cleaning process and the characteristics of the weakly labeled, real strongly labeled, and synthetic strongly labeled data.

### 3.1. Data Format Unifying

For all the eight source datasets in Table 1, we use ffmpeg [38] to convert those multi-channel audio files into single-channel ones, and down-sampling them into 16 KHz sampling rate.

### 3.2. Weakly Labeled Data Creation

Based on the cleaned-up source data clips presented in Section 3.1, we first created the weakly labeled data (WLD) portion of the ICSD dataset. The detailed statistics are summarized in Table 3. The Infant Cry category of our WLD is derived from three sources: Audioset, Donate A Cry, and ESC-50. Meanwhile, the Snoring category is derived from two sources: Audioset and ESC-50.This results in a total of 1888 clips for Infant Cry and 1753 clips for Snoring. We divided the samples within this part of the dataset into a 9:1 ratio for training (3276

samples) and validation (365 samples). This diverse collection of WLD covers a wide range of sound variations, potentially enhancing the ICSD system's generalization ability and effectiveness on unseen data.

Table 3: Detail Statistics of Weakly Labeled Data (#clip, 10s/clip).

| Category | Data Source | #Clips |
|---|---|---|
| Infant Cry | Audioset | 1391 |
| | Donate A Cry | 457 |
| | ESC-50 | 40 |
| Snoring | Audioset | 1713 |
| | ESC-50 | 40 |

### 3.3. *Real Strongly Labeled Data Creation*

The real strongly labeled data in our ICSD dataset is directly sourced from Audioset, all these data clips are with human labeled target event-based time-stamps. The detail statistics are shown in Table 4. We divide the clips within this part of dataset into a ratio of 8:1:1 for training (647 clips), validation (80 clips) and testing (80 clips).

Table 4: Statistics of Real Strongly Labeled Data (#clip, 10s/clip).

| Category | Data Source | Train | Validation | Test |
|---|---|---|---|---|
| Infant Cry | Audioset (424) | 340 | 42 | 42 |
| Snoring | Audioset (383) | 307 | 38 | 38 |

In addition, in Fig.1, we plot the duration histogram of both infant cry and snoring events as described in Table 4, along with their time intervals between events. Fig.1(a) and (b) show the duration distribution of the target events, while Fig.1(c) and (d) illustrate the distribution of time intervals between two events in each recording clip. These figures aid in understanding the temporal characteristics of the sound events, which are fundamental for effective feature extraction and model training.
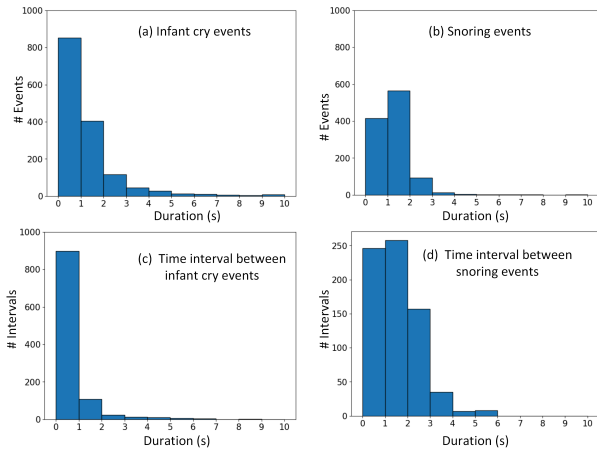


Figure 1: Duration histogram of events and time intervals.

As illustrated in Fig.1(a), the infant cry events mostly have durations within the first two seconds. There is a notable peak

at the outset, indicating a high occurrence of very brief cries. The sharp decline in event frequency as duration increases highlights the transient nature of infant cries. These brief events imply the necessity of feature extraction techniques tailored to capturing these rapid, intense bursts of sound effectively. In Fig.1(b), we see the snoring events tend to have slightly longer durations, with the most common durations ranging from one to four seconds. This suggests that snoring, while still brief, includes somewhat longer sounds compared to infant cries, necessitating a distinct approach to feature analysis capable of accommodating these long acoustic characteristics.

Examining the duration distribution of intervals between events, Fig.1(c) highlights that intervals of infant cry events are primarily less than one second, indicating a high frequency of closely spaced cry events. This pattern suggests the need for algorithms capable of detecting rapid sequences of sound, which could signify specific behaviors or needs. In contrast, the time interval distribution for Snoring events shown in Fig.1(d) demonstrates more extended periods between occurrences, with intervals more evenly distributed up to three seconds. This regularity in event spacing provides a clear rhythmic pattern that feature extraction or acoustic modeling techniques must capture to effectively differentiate snoring from other sounds.

### 3.4. *Synthetic Strongly Labeled Data Creation*

The synthetic strongly labeled (SSL) data in our ICSD dataset is specifically generated to augment the existing data and provide a more robust training set for the ICSD system training.

#### 3.4.1. *Source Material*

The source material used to synthesize the SSL dataset can be divided into two groups: the foreground events and background sounds. Details of these sources are presented as bellow.

**Foreground Events**: In our study, foreground sounds are individual sound events including both events segments of infant cry and snoring (they are also called target events in ICSD system). As shown in Table 5, there target events are all extracted from the source datasets as listed in Table 1, specifically, the infant cry events are extracted from BCD and Audioset, while the snoring events are extracted from both the self-collected database and FM Snoring dataset that directly downloaded from the Kaggle website [2]. To ensure the quality of synthetic data, we discarded foreground events with poor quality including the asphyxia and deaf segments of BCD as well as any foreground sounds with duration less than 250ms. As in Table 5, all these cleaned-up foreground events are then divided into a ratio of 9:1 for the synthesize of training and validation sets, and the synthesize of the test set.

**Background Sounds**: The background sounds used in our study consist of individual ambient sound clips, each segmented into 10-second intervals. To avoid the dominance of speech and music backgrounds, we did not use all the cleaned-up clips from
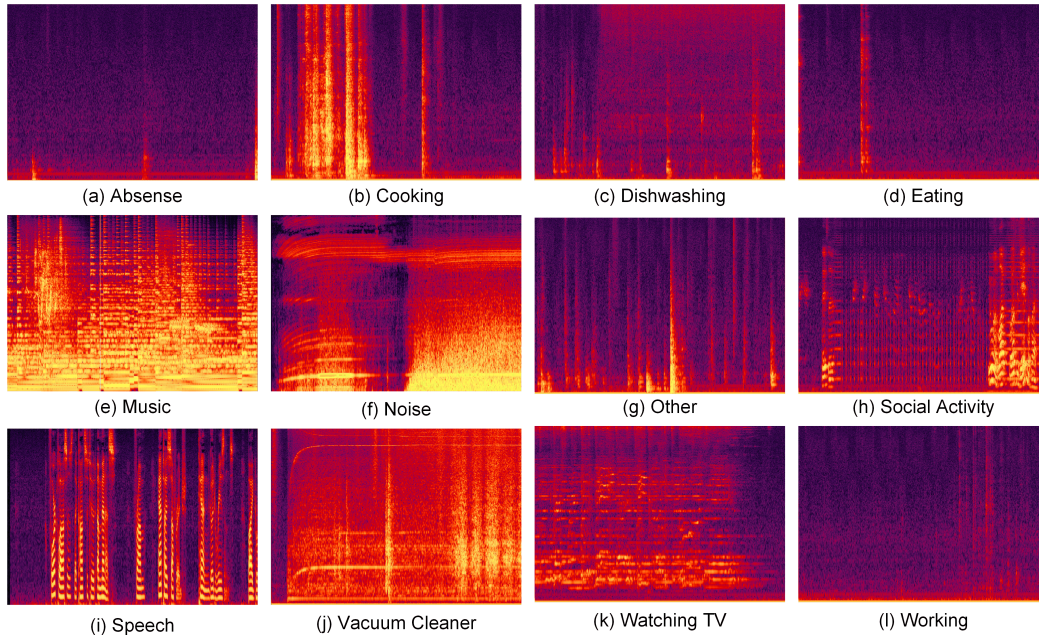
---

[2]https://www.kaggle.com/

Figure 2: Spectrogram samples of 12 types of background sounds.

Table 5: Statistics of Foreground Events.

| Category | Source(#Events) | Train & Validation | Test |
|---|---|---|---|
| Infant Cry | BCD (1048) | 943 | 105 |
| | Audioset (1354) | 1211 | 143 |
| Snoring | Kaggle (1500) | 1350 | 150 |
| | Audioset (1058) | 946 | 112 |

Table 6: Statistics of Background Sound Segments.

| Category | # Segments (10s) | # hrs |
|---|---|---|
| Absence | 18860 | 52.4 |
| Cooking | 5124 | 14.2 |
| Dish washing | 1424 | 4.0 |
| Eating | 2308 | 6.4 |
| Social activity | 4944 | 13.7 |
| Vacuum cleaning | 972 | 2.7 |
| Watching TV | 18648 | 51.8 |
| Working | 18644 | 51.8 |
| Speech | 5400 | 15 |
| Music | 3600 | 10 |
| Noise | 1831 | 5.1 |
| Other | 2060 | 5.7 |
| Total | 83815 | 232.8 |

the SINS and MUSAN datasets (in Table 2) for our ICSD backgrounds in SSL data synthesis. Instead, we randomly extracted 15 hours of speech data and 10 hours of music data from the MUSAN dataset and combined them with the SINS dataset as background sounds. Detailed information is provided in Table 6, which includes 11 different types of background environment sounds and an additional category named 'Other' for miscellaneous background sounds. During SSL synthesizing, all the backgrounds in Table 6 are divided into a ratio of 9:1 for the synthesis of training plus validation sets, and the synthesis of the test set. Fig. 2 illustrates a randomly picked sample background spectrogram for each category of background.

### 3.4.2. SSL Data Synthesizing Process

Given the above provided foreground events and background sound clips, we synthesized the SSL subset of the ICSD dataset using the Scaper toolkit, as illustrated in Fig. 3. Scaper [42] is a library specifically designed for the synthesis of strongly labeled data and has been widely utilized in the DCASE Challenges [5]. Scaper can create a rich, customizable dataset from collections of foreground and background sounds. It employs probabilistic definitions to generate diverse audio scenes, allowing for the manipulation of individual sounds with transformations such as pitch shifting and time stretching. This toolkit

is particularly valuable for sound event detection in diverse environments, addressing the scarcity of strongly labeled audio data necessary for training and evaluating machine listening systems.

The parameters for Scaper to generate the synthetic data SSL were carefully chosen to control the characteristics of the soundscapes. These include the duration of each synthetic clip, the minimum and maximum number of events in each clip, the level of reverberation, the foreground-background signal-to-noise ratio (FB-SNR), and the range of pitch shifts applied to the events. The specific values of these parameters are detailed in Table 7.

Based on the above Scaper settings, we finally generated a total of 10,000 synthetic clips for the current SSL dataset, as detailed in Table 2. This SSL dataset includes 4,000 training
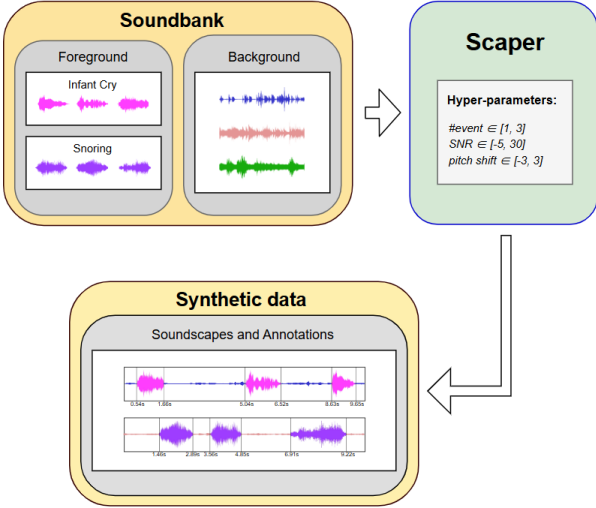
5

Figure 3: Illustration of SSL data synthesizing.

Table 7: Parameters in Scaper for SSL Data Generation.

| Parameter | Value |
|-----------|-------|
| Duration | 10s |
| Min event | 1 |
| Max event | 3 |
| Reverb | 0.1 |
| FB-SNR | (uniform, -5, 30) |
| Pitch shifts | (uniform, -3.0, 3.0) |

clips each for infant cry and snoring, along with 500 clips each for validation and testing in both categories.

### 3.5. Ontology Data Release

Fig.4 illustrates the organized structure where audio files are stored in the audio folder and event time-stamp annotations in the metadata folder, each further categorized into train, validation, and test subfolders. Furthermore, we have also released the source materials for generating synthetic strongly labeled data, including foreground events and background events.

## 4. ICSD Systems

We build three main ICSD systems on our ICSD dataset to serve as initial references for further research on infant cry and snoring detection. These systems include an MT-CRNN system, a CRNN+BEATs system, and an enhanced CRNN+BEATs system. The first two architectures are adapted from the DCASE 2023 Challenge Task 4 baseline systems [41], while the competitive system is our improved version of CRNN+BEATs. Details of these three systems are presented in the following sections.

### 4.1. MT-CRNN

The MT-CRNN ICSD system is based on the Mean-Teacher Convolutional Recurrent Neural Network [39], it consists of a
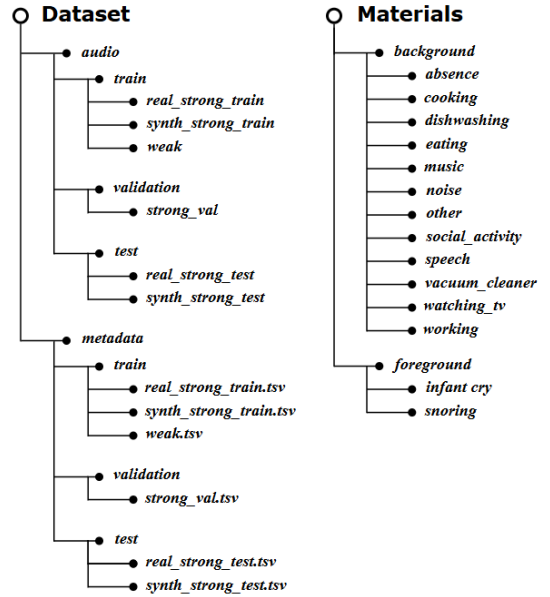


Figure 4: The structure of the ICSD dataset ontology.

student model and a teacher model, both sharing same architecture but serving different purposes. The model architecture is a fusion of a Convolutional Neural Network (CNN) and a Recurrent Neural Network (RNN), topped with an attention layer. The RNN output offers strong frame-level predictions, while the attention layer output provides segment/clip-level predictions (weak predictions).

The student model is trained on a mix of strongly and weakly labeled data. Binary cross entropy (BCE) is used to compute the supervised learning loss at the frame level for strongly labeled data and at the clip level for weakly labeled data. However, the teacher model is not directly trained. Instead, the teacher model's parameters are updated using an Exponential Moving Average (EMA) approach[40], where the weights are a moving average of the student model's weights. This procedure is defined as:

$$\theta'_t = \alpha \cdot \theta'_{t-1} + (1 - \alpha) \cdot \theta_t \qquad (1)$$

where $\alpha$ is the decay hyperparameter used to control how far the EMA reaches into the previous training history, $\theta'_t$ and $\theta_t$ are the current model weights of the teacher and student model, respectively.

In this MT-CRNN, the teacher model contributes to the training of the student model via a consistency loss (mean-squared error) for both strong and weak predictions. The whole MT-CRNN model training is performed using four distinct losses combination as defined in Eq.(2),

$$\mathcal{L} = \mathcal{L}_{BCE}^{weak} + \mathcal{L}_{BCE}^{strong} + \gamma \cdot \left( \mathcal{L}_{MSE}^{weak} + \mathcal{L}_{MSE}^{strong} \right) \qquad (2)$$

where $\gamma$ denotes the loss weight, $\mathcal{L}_{BCE}^*$ and $\mathcal{L}_{MSE}^*$ are the supervised BCE loss and teacher-student MSE consistency loss, respectively. All the losses with 'weak' subscript are computed using the weakly labeled training data, while the losses with 'strong' subscripts are computed using the strongly labeled training data. More details can be found in [5].
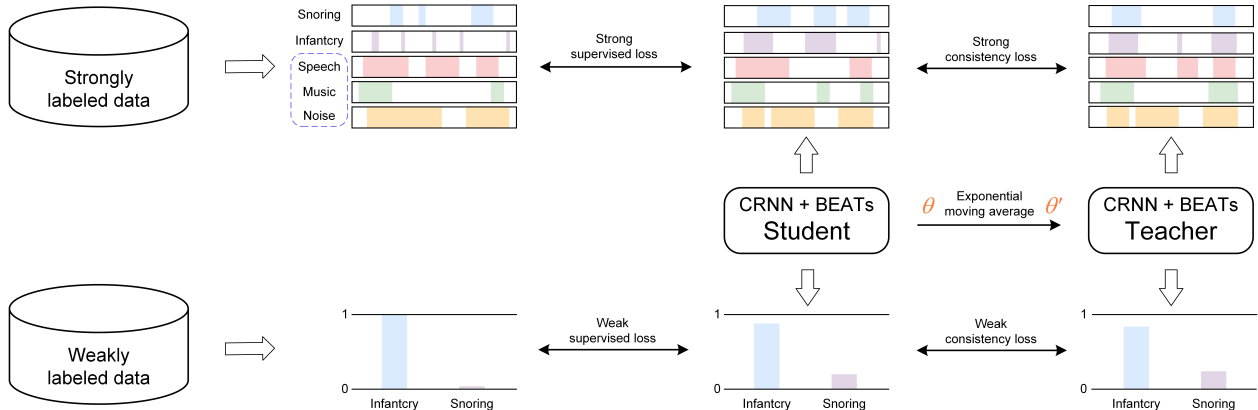
Figure 5: Training procedure illustration of the Competitive CRNN+BEATs model.

## 4.2. *CRNN+BEATs*

The CRNN+BEATs model builds upon the MT-CRNN model by incorporating the BEATs features. BEATs [11], is an iterative audio pre-training framework to learn Bidirectional Encoder representation from Audio Transformers. The BEATs model is trained to extract high-level non-speech audio semantics information using iterative self-supervised learning. In our study, the CRNN+BEATs system architecture is the same as the one used in DCASE 2023 Task 4 Challenge baselines [41], but we train it on our ICSD dataset for infant cry and snoring detection.

In the CRNN+BEATs model, the integration of BEATs feature representation and the CNN output is performed by a linear transformation followed by a layer normalization. The transformed feature map and CNN output are then concatenated and passed through another linear transformation to match the dimension of the GRU input of the CRNN. Our CRNN+BEATs model is trained in the same way as the MT-CRNN model. The inclusion of BEATs feature representations provide additional complementary information to the CRNN latent representations, which can potentially improve the performance of the model in sound event detection tasks.

## 4.3. *Competitive CRNN+BEATs*

The Competitive CRNN+BEATs model is an improved version of the CRNN+BEATs model. It incorporates three additional categories, speech, music and noise, into the real strongly labeled data. This means we train a five-class SED system to perform the ICSD event detection, rather than focusing solely on infant cry and snoring detection. For each of speech and music categories, we downloaded 400 clips of strongly labeled data from Audioset, dividing them into 350 clips for training and 50 clips for validation. For the noise category, we directly pick 140 samples from the background clips for model training. The training process for the Competitive CRNN+BEATs model follows a similar procedure to that of the CRNN+BEATs model, as detailed in Fig.5. In this case, the $\mathcal{L}_*^{strong}$ in Eq.(2) are with all the five types of strongly labeled data (for the noise category, we labeled the whole clip as noise strong label in this case).

By incorporating additional categories for speech, music and noise, the model is better equipped to handle acoustic confusion between speech, music, noise and target sound events, potentially improving its performance in ICSD tasks. The motivation behind the proposed Competitive CRNN+BEATs model stems from our extensive preliminary experiments, which showed that infant cry sounds have the most strong acoustic confusion with music, speech and noise, leading to high false positive rates during non-target environmental background sound testing.

## 5. Experiments and Results

### 5.1. *Configurations*

All our models were trained using the Adam optimizer. The batch size of the MT-CRNN system was set to 36, including 24 for strongly labeled data and 12 for weakly labeled data. Meanwhile, the batch size of the CRNN+BEATs system was set to 48, with 24 strongly labeled data and 24 weakly labeled data. The learning rate was warmed up for 50 epochs, and then decayed to 0.001. The models were trained over a total of 200 epochs. The training loss weights $\gamma$ in Eq.(2) are defined in Eq.(3):

$$\gamma = C \times e^{\beta \times \left(1 - \frac{s}{N \times L}\right)^2} \tag{3}$$

where $C = 2$, represents the constant maximum value used for the self-supervised loss weight, $\beta$ is the exponent parameter set to -5.0, $s$ is the current training step, $N$ is the number of warm-up epochs, set to 50. $L$ is the length of each epoch in terms of training steps.

### 5.2. *Evaluation Metrics*

To validate the performance of our models on the proposed ICSD, we use both intersection-based F1 [44, 48] and segment-based F1 (Seg-F1) scores [46, 47] as primary evaluation metrics. The intersection-based F1 (Inter-F1) is calculated using a detection tolerance criterion (DTC) and ground truth intersection criterion (GTC) set to 0.5, and a cross-trigger tolerance criterion (CTTC) set at 0.3. The Inter-F1 is employed to provide a robust measure of the overlap between the predicted and

actual labels of target sound events, while the Seg-F1 evaluates the model's performance in accurately identifying sound events within fixed-length (1s in this study) segments of an audio stream. These metrics have been widely used for evaluating sound event detection systems as in [43, 44, 45].

In addition to these F1 measures, we also use a false positive event (FPE) count per hour, similar to the false alarm (FA) count per hour used in keyword spotting (KWS) tasks, to evaluate the error trigger rate of ICSD systems in long-time background environments without any target events. In practical applications, a very low FPE count is crucial to ensure the reliability and efficiency of infant cry and snoring detection systems. This minimizes the chances of false alarms, thereby improving overall system performance.

## 5.3. Results

### 5.3.1. Overall Results

Table 8: Performance Comparison of Different Models on the Real Test Set.

| System | Infant Cry | | Snoring | |
|---|---|---|---|---|
| | Inter-F1 | Seg-F1 | Inter-F1 | Seg-F1 |
| MT-CRNN | 0.8203 | 0.8497 | 0.8089 | 0.8529 |
| CRNN+BEATs | 0.8475 | 0.8990 | **0.8522** | 0.8763 |
| Competitive CRNN+BEATs | **0.8805** | **0.9174** | 0.8254 | 0.8697 |
| CRNN+BEATs-14 Class | 0.8534 | 0.8919 | 0.8448 | **0.8765** |
| CRNN+BEATs+others | 0.8299 | 0.8732 | 0.8341 | 0.8636 |

Tables 8 and 9 demonstrate the ICSD results on five different systems for the real and synthetic test sets, respectively. By comparing the results of the MT-CRNN and CRNN+BEATs systems in Table 8, we observe an absolute performance improvement of around 2%-5% in both infant cry and snoring detection, as measured by both intersection-based F1 and segment-based F1 metrics. This indicates that introducing additional feature maps derived from the pre-trained BEATs model provides significant complementary information to the basic mean-teacher CRNN model.

Furthermore, when comparing the proposed competitive CRNN+BEATs results with those from the CRNN+BEATs system, we achieve an absolute 3.3% improvement in intersection-based F1 for infant cry detection. This suggests that modeling music, speech and noise together with the target infant cry and snoring as separate categories can effectively reduce acoustic confusion between infant cry and background music, speech and noise. However, we observe a slight performance degradation in snoring detection. This may be due to the model balancing acoustic discrimination and confusion between noise and snoring. As a result, some snoring events are misclassified as noise during the inference on the real snoring test set.

In Table 9, all the results on the synthetic test tests are presented. Compared to Table 8, all five models perform much better on the synthetic test sets than on the real ICSD test sets. For example, both the intersection-based F1 and segment-based

F1 metrics of five different models are all above 95% in both infant cry and snoring detection. This indicates that the quantity of synthetic strongly labeled data in the train set far exceeds that of real strongly labeled data, which could lead to potential better models to these synthetic data, hence resulting in the better performance. Moreover, real strongly labeled data often contain more complex factors, such as variations in sound quality and event intervals, which may affect model performance. In contrast, synthetic datasets are usually simpler and clearer, enabling the models to learn and predict better.

Table 9: Performance Comparison on the Synthetic Test Set.

| System | Infant Cry | | Snoring | |
|---|---|---|---|---|
| | Inter-F1 | Seg-F1 | Inter-F1 | Seg-F1 |
| MT-CRNN | 0.9519 | 0.9527 | 0.9744 | 0.9566 |
| CRNN+BEATs | **0.9878** | 0.9769 | 0.9821 | **0.9722** |
| Competitive CRNN+BEATs | 0.9845 | 0.9755 | 0.9815 | 0.9672 |
| CRNN+BEATs-14 Class | 0.9832 | 0.9760 | 0.9850 | 0.9690 |
| CRNN+BEATs+others | **0.9878** | **0.9786** | **0.9901** | 0.9702 |

In addition, when comparing the results of the competitive CRNN+BEATs system with the CRNN+BEATs system, there is almost no performance difference between them on synthetic test sets. However, on the real infant cry test set, the competitive CRNN+BEATs significantly outperforms CRNN+BEATs. This indicates that including more categories (speech, music, noise) helps the model better generalize to real-world scenarios where such sounds are present. For the real snoring test set, the competitive CRNN+BEATs shows slightly lower performance than CRNN+BEATs. This could be because the inclusion of more categories may cause some snoring events to be misclassified as noise.

Furthermore, motivated by the competitive CRNN+BEATs performance, we also performed the 'CRNN+BEATs-14 Class' and 'CRNN+BEATs+Others' systems for a fair system comparison in both Table 8 and 9. 'CRNN+BEATs-14 Class' means train a 14-class CRNN+BEATs model using balanced training data of 14 classes (two target classes: infant cry and snoring, 12 background classes as shown in Table 6.). The 'CRNN+BEATs+Others' model, on the other hand, was trained as a 3-class system that combines the 12 background classes into a single universal background class labeled 'Others'. Detailed data organization for these two systems can be found on our ICSD open-source website.

By comparing the results of last three lines in Table 8 and 9, we see that on the real infant cry test set, the competitive CRNN+BEATs model shows the best performance with the highest Inter-F1 (0.8805) and Segment-F1 (0.9174) scores, suggesting that modeling five categories enhances the model's ability to differentiate between these sounds and the infant cry in real-scenarios. Treating all background sounds as a single category in 'CRNN+BEATs+Others' may not be as effective in distinguishing infant cries from various background noises. And introducing more background sounds categories as modeling targets in 'CRNN+BEATs-14 Class' may bring more mis-

classified possibility of target infant cry detection. On the real snoring test set, these findings differ slightly from those observed for the infant cry test set. Here, the CRNN+BEATs model achieves the highest F1 scores, suggesting that it handles snoring detection more effectively when focusing solely on the two target categories. This may be attributed to the significant acoustic differences between the two target sounds, infant cries, and snoring.

### 5.3.2. *Ablation Study of Different Training Data Sets*

Considering the significant role that training data plays in shaping the performance of a model, Table 10 provides an ablation study using different combinations of training data to observe their effects on system performance. All the experiments are performed on the CRNN+BEATs system.

In Table 10, we observe a significant performance gap between the synthetic and real ICSD test sets. The F1 measures on the synthetic set are much better than those on the real set, indicating that the synthetic data is much easier and cleaner compared to the real test set. By comparing A0 with A1, we see a large acoustic mismatch between the synthetic and real ICSD data. In A1, training the model using only the synthetic training set (shown in Table 2) results in high F1 measures on the synthetic test set, while the F1 measures on the real test set are significantly reduced. In A2, adding the synthetic training data to the real training set improves performance on both the synthetic and real test sets. This suggests that, despite being simpler, synthetic data can still provide complementary information to the real training data. In A3, adding weakly labeled data to the strongly labeled set yields the best results on the real test set, while maintaining almost the same good results on the synthetic test set as in A2. This indicates that both weakly labeled data and synthetic training data are important for augmenting the limited real training data set. Therefore, in all of Table 8, Table 9, and Fig.6, the models are trained using the entire training set, as shown in Table 2, to include all strongly labeled and weakly labeled data.

### 5.4. *Results on Non-target Test Set*

As the importance of false alarm in keyword spotting tasks, the false positive rate on non-target test set is also extremely important for the ICSD tasks in real-world applications. Because a high false positive rate for target events will greatly reduce the user's satisfaction in actual product applications. For example, if an alert is sent to parents every time an infant's cry is detected, a high false positive event (FPE) rate will end up irritating the users. Fig.6 presents the FPEs per hour on a 'Non-target Test Set' with 12 background classes, with 1 hour data for each class randomly collected from a total of 232.8 hrs of background sounds as listed in Table 6. Fig.6 (a) shows the FPEs for infant cry while (b) shows the results for snoring FPEs.

For the CRNN+BEATs model, the most common false positive detection were music, noise, and speech, indicating that the model often misclassify these non-target sounds as infant cry or snoring. In contrast, the Competitive CRNN+BEATs model

Figure 6: False Positive Events per Hour on Non-target Test Set for (a) Infant Cry and (b) Snoring.
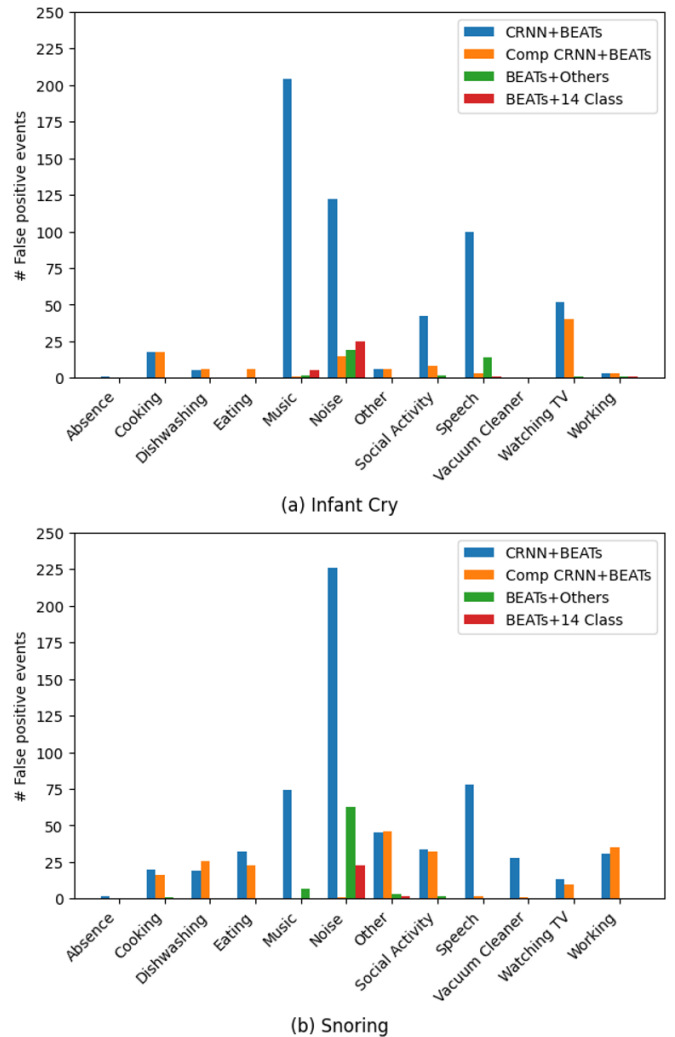


(a) Infant Cry



(b) Snoring

Table 10: Training Data ablation experiments of CRNN+BEATs on the Synthetic and Real Test Sets (including both infant cry and snoring).

| ID | Training data | Synthetic Test Set | | Real Test Set | |
|----|---------------|---------|---------|---------|---------|
| | | Inter-F1 | Seg-F1 | Inter-F1 | Seg-F1 |
| A0 | Real strongly labeled only | 0.8819 | 0.8835 | 0.8104 | 0.8521 |
| A1 | Synthetic strongly labeled only | 0.9858 | 0.9773 | 0.7459 | 0.8413 |
| A2 | Synthetic + Real strongly labeled | **0.9896** | **0.9774** | 0.8122 | 0.8635 |
| A3 | All strongly + weakly labeled | 0.9827 | 0.9724 | **0.8425** | **0.8815** |

shows a significant reduction in FPEs for these classes, indicating the competitive model shows better discrimination ability between target and non-target sounds, thereby significantly minimizing false positive detection counts in music, noise, and speech conditions. However, when comparing the FPEs between 'Competitive CRNN+BEATs' and the 'CRNN+BEATs-14 Class', 'CRNN+BEATs+Others' systems, we find that the 'CRNN+BEATs-14 Class' achieves the lowest overall FPEs of background sounds for both infanct cry and snoring sounds. However, from Table 8 and 9, we see that the 'CRNN+BEATs-14 Class' may bring more target sounds misclassifying to other non-target background sounds at some extent, therefore, there is a performance trade-off between F1 measures and FPEs in different real application scenarios. We hope these preliminary results can provide a basic reference for other researchers in the ICSD field.

## 6. Conclusion

In this study, we introduced the ICSD dataset, a new publicly available resource designed to bridge the gap in the availability of high-quality, strongly and weakly labeled data for the detection of infant cries and snoring sounds. Besides the open-source dataset, we also provided several baseline ICSD systems to show extensive preliminary ICSD results on this dataset as reference points for future ICSD studies, including a mean-teacher based CRNN, a CRNN+BEATs, and competitive CRNN+BEATs systems. Based on these baseline systems, we performed detailed analysis on both the real and synthetic test sets, as well as on a purely non-target background test set to evaluate the degree of false positive rates in different types of background acoustic environments. We hope that our ICSD dataset will be beneficial for the advancement of sound event detection in home environments, and all these preliminary experimental findings and observations on this released dataset can serve as a solid foundation for future innovations in this area. In addition, it is worth noting that the current usage of our released synthetic and background datasets provided in the ICSD is just a basic example. Researchers can re-organize and use our data simulation script to synthesize any size/any types of synthetic dataset for their ICSD system building. Our future work will focus on exploring new methods to improve ICSD performance and continuously enhancing this dataset to achieve larger and wider domain coverage.

## References

[1] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an ASR corpus based on public domain audio books," in *Proc. ICASSP*, 2015, pp. 5206–5210.

[2] Y. Lin, X. Qin, G. Zhao, M. Cheng, N. Jiang, H. Wu, and M. Li, "Voxblink: A large scale speaker verification dataset on camera," in *Proc. ICASSP*, 2024, pp. 10 271–10 275.

[3] J. Thiemann, N. Ito, and E. Vincent, "The diverse environments multi-channel acoustic noise database (demand): A database of multichannel environmental noise recordings," in *Proc. Acoustics*, vol. 19, no. 1, 2013.

[4] K. Ito and L. Johnson, "The LJ speech dataset," Available at: https://keithito.com/LJ-Speech-Dataset/, 2017.

[5] N. Turpault, R. Serizel, A. P. Shah, and J. Salamon, "Sound event detection in domestic environments with weakly labeled data and soundscape synthesis," in *DCASE Workshop*, 2019.

[6] 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 *Proc. ICASSP*, 2017, pp. 776–780.

[7] J. Shor, A. Jansen, R. Maor, O. Lang, O. Tuval, F. d. C. Quitry, M. Tagliasacchi, I. Shavitt, D. Emanuel, and Y. Haviv, "Towards learning a universal non-semantic representation of speech," in *Proc. INTERSPEECH*, 2020.

[8] A. Saeed, D. Grangier, and N. Zeghidour, "Contrastive learning of general-purpose audio representations," in *Proc. ICASSP*, 2021, pp. 3875–3879.

[9] A. Agostinelli, T. I. Denk, Z. Borsos, J. Engel, M. Verzetti, A. Caillon, Q. Huang, A. Jansen, A. Roberts, M. Tagliasacchi *et al.*, "MUSICLM: Generating music from text," *arXiv preprint arXiv:2301.11325*, 2023.

[10] F. Kreuk, G. Synnaeve, A. Polyak, U. Singer, A. Défossez, J. Copet, D. Parikh, Y. Taigman, and Y. Adi, "AudioGen: Textually guided audio generation," in *Proc. ICLR*, 2023.

[11] S. Chen, Y. Wu, C. Wang, S. Liu, D. Tompkins, Z. Chen, and F. Wei, "Beats: Audio pre-training with acoustic tokenizers," in *Proc. ICML*, 2023, pp. 5178–5193.

[12] Q. Kong, Y. Cao, T. Iqbal, Y. Wang, W. Wang, and M. D. Plumbley, "Panns: Large-scale pretrained audio neural networks for audio pattern recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing,*, vol. 28, pp. 2880–2894, 2020.

[13] Q. Kong, Y. Xu, W. Wang, and M. D. Plumbley, "Audio set classification with attention model: A probabilistic perspective," in *Proc. ICASSP*, 2018, pp. 316–320.

[14] Q. Kong, C. Yu, Y. Xu, T. Iqbal, W. Wang, and M. D. Plumbley, "Weakly labelled audioset tagging with attention neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing,*, vol. 27, no. 11, pp. 1791–1802, 2019.

[15] O. F. Reyes-Galaviz, S. D. Cano-Ortiz, and C. A. Reyes-García, "Evolutionary-neural system to classify infant cry units for pathologies identification in recently born babies," in *Proc. MICAI*, 2008, pp. 330–335.

[16] O. F. Reyes-Galaviz, E. A. Tirado, and C. A. Reyes-Garcia, "Classification of infant crying to identify pathologies in recently born babies with ANFIS," in *Proc. ICCHP*, 2004, pp. 408–415.

[17] C. Ji, X. Xiao, S. Basodi, and Y. Pan, "Deep learning for asphyxiated infant cry classification based on acoustic features and weighted prosodic features," in *2019 International conference on Internet of Things (iThings) and IEEE green computing and communications (GreenCom) and IEEE cyber, physical and social computing (CPSCom) and IEEE smart data (SmartData)*, 2019, pp. 1233–1240.

[18] M. Moharir, M. Sachin, R. Nagaraj, M. Samiksha, and S. Rao, "Identification of asphyxia in newborns using gpu for deep learning," in *Proc. I2CT*, 2017, pp. 236–239.

[19] C. C. Onu, I. Udeogu, E. Ndiomu, U. Kengni, D. Precup, G. M. Sant'Anna, E. Alikor, and P. Opara, "Ubenwa: Cry-based diagnosis of birth asphyxia," *arXiv preprint arXiv:1711.06405*, 2017.

[20] M. Sachin, R. Nagaraj, M. Samiksha, S. Rao, and M. Moharir, "GPU based deep learning to detect asphyxia in neonates," *Indian J. Sci. Technol,*, vol. 10, no. 3, 2017.

[21] O. M. Badreldine, N. A. Elbeheiry, A. N. M. Haroon, S. ElShehaby, and E. M. Marzook, "Automatic diagnosis of asphyxia infant cry signals using wavelet based mel frequency cepstrum features," in *Proc. ICENCO*, 2018, pp. 96–100.

[22] H. B. Sailor and H. A. Patil, "Auditory filterbank learning using ConvRBM for infant cry classification," in *Proc. INTERSPEECH*, 2018, pp. 706–710.

[23] J. Chunyan, M. Chen, L. Bin, and Y. Pan, "Infant cry classification with graph convolutional networks," in *Proc. ICCCS*, 2021, pp. 322–327.

[24] J. Saraswathy, M. Hariharan, V. Vijean, S. Yaacob, and W. Khairunizam, "Performance comparison of daubechies wavelet family in infant cry classification," in *Proc. CSPA*, 2012, pp. 451–455.

[25] G. Veres, "Donateacry-corpus," https://github.com/gveres/donateacry-corpus.

[26] P. Kulkarni, S. Umarani, V. Diwan, V. Korde, and P. P. Rege, "Child cry classification-an analysis of features and models," in *Proc. I2CT*, 2021, pp. 1–7.

[27] T. Ozseven, "Infant cry classification by using different deep neural network models and hand-crafted features," *Biomedical Signal Processing and Control,*, vol. 83, p. 104648, 2023.

[28] K. Sharma, C. Gupta, and S. Gupta, "Infant weeping calls decoder using statistical feature extraction and gaussian mixture models," in *Proc. ICCCNT*, 2019, pp. 1–6.

[29] E. Sutanto, F. Fahmi, W. Shalannanda, and A. Aridarma, "Cry recognition for infant incubator monitoring system based on internet of things using machine learning," *International Journal of Intelligent Engineering & Systems,*, vol. 14, no. 1, 2021.

[30] L. Jiang, Y. Yi, D. Chen, P. Tan, and X. Liu, "A novel infant cry recognition system using auditory model-based robust feature and gmm-ubm," *Concurrency and Computation: Practice and Experience,*, vol. 33, no. 11, p. e5405, 2021.

[31] T. Khan, "A deep learning model for snoring detection and vibration notification using a smart wearable gadget," *Electronics,*, vol. 8, no. 9, p. 987, 2019.

[32] O. NAHUM, "Female and Male Snoring Dataset," Available at: https://www.kaggle.com/datasets/orannahum/female-and-male-snoring, 2022.

[33] K. J. Piczak, "ESC: Dataset for environmental sound classification," in *Proc. ACM international conference on Multimedia*, 2015, pp. 1015–1018.

[34] H. Akbari, L. Yuan, R. Qian, W.-H. Chuang, S.-F. Chang, Y. Cui, and B. Gong, "Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text," *Advances in Neural Information Processing Systems,*, vol. 34, pp. 24 206–24 221, 2021.

[35] G. Dekkers, S. Lauwereins, B. Thoen, M. W. Adhana, H. Brouckxon, B. Van den Bergh, T. Van Waterschoot, B. Vanrumste, M. Verhelst, and P. Karsmakers, "The SINS database for detection of daily activities in a home environment using an acoustic sensor network," in *DCASE Workshop*, 2017, pp. 32–36.

[36] D. Snyder, G. Chen, and D. Povey, "MUSAN: A music, speech, and noise corpus," *arXiv preprint arXiv:1510.08484*, 2015.

[37] Y. Wang, H. Lv, D. Povey, L. Xie, and S. Khudanpur, "Wake word detection with alignment-free lattice-free mmi," in *Proc. INTERSPEECH*, 2020.

[38] J. Newmarch and J. Newmarch, "FFmpeg/Libav," *Linux sound programming,*, pp. 227–234, 2017.

[39] U. Singh, D. D. Dash, M. Sharma, S. Mishra, S. Malarvizhi, S. Tiwari, and R. T. Shankarappa, "Polyphonic sound event detection and classification using convolutional recurrent neural network with mean teacher," in *Proc. ICCCNT*, 2021, pp. 1–4.

[40] A. Lawrance and P. Lewis, "An exponential moving-average sequence and point process (EMA1)," *Journal of Applied Probability,*, vol. 14, no. 1, pp. 98–113, 1977.

[41] DCASE Community, "DCASE 2023 Task: Sound Event Detection with Weak Labels and Synthetic Soundscapes," Available at: https://dcase.community/challenge2023/task-sound-event-detection-with-weak-labels-and-\synthetic-soundscapes#baseline-system, 2023, [Online; accessed 10-October-2023].

[42] J. Salamon, D. MacConnell, M. Cartwright, P. Li, and J. P. Bello, "Scaper: A library for soundscape synthesis and augmentation," in *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2017, pp. 344–348.

[43] R. Serizel and N. Turpault, "Sound event detection from partially annotated data: Trends and challenges," in *IcETRAN conference*, 2019.

[44] G. Ferroni, N. Turpault, J. Azcarreta, F. Tuveri, R. Serizel, Ç. Bilen, and S. Krstulović, "Improving sound event detection metrics: insights from dcase 2020," in *Proc. ICASSP*, 2021, pp. 631–635.

[45] T. Khandelwal, R. K. Das, and E. S. Chng, "Is your baby fine at home? Baby cry sound detection in domestic environments," in *Proc. APSIPA ASC*, 2022, pp. 275–280.

[46] A. Mesaros, T. Heittola, and T. Virtanen, "Metrics for polyphonic sound event detection," *Applied Sciences,*, vol. 6, no. 6, p. 162, 2016.

[47] A. Mesaros, T. Heittola, and D. Ellis, "Datasets and evaluation," *Computational Analysis of Sound Scenes and Events,*, pp. 147–179, 2018.

[48] Ç. Bilen, G. Ferroni, F. Tuveri, J. Azcarreta, and S. Krstulović, "A framework for the robust evaluation of sound event detection," in *Proc. ICASSP*, 2020, pp. 61–65.

[49] J. Chew, Y. Sun, L. Jayasinghe, and C. Yuen, "Dcase 2018 challenge: Solution for task 5," *arXiv preprint arXiv:1812.04618*, 2018.