

# MODULATION SPECTRAL SIGNAL REPRESENTATION AND I-VECTORS FOR ANOMALOUS SOUND DETECTION

## Technical Report

Parth Tiwari<sup>1, 3</sup>, Yash Jain<sup>2</sup>, Anderson Avila<sup>3</sup>, João Monteiro<sup>3</sup>,  
Shruti Kshirsagar<sup>3</sup>, Amr Gaballah<sup>3</sup>, Tiago H. Falk<sup>3</sup>

<sup>1</sup> Department of Industrial and Systems Engineering, IIT Kharagpur, India

<sup>2</sup> Department of Mathematics, IIT Kharagpur, India

<sup>3</sup> MuSAE Lab, Institut National de la Recherche Scientifique - Centre EMT, Montreal, Canada

### ABSTRACT

This report summarizes our submission for Task-2 of the DCASE 2020 Challenge. We propose two different anomalous sound detection systems, one based on features extracted from a modulation spectral signal representation and the other based on i-vectors extracted from mel-band features. The first system uses a nearest neighbour graph to construct clusters which capture local variations in the training data. Anomalies are then identified based on their distance from the cluster centroids. The second system uses i-vectors extracted from mel-band spectra for training a Gaussian Mixture Model. Anomalies are then identified using their negative log likelihood. Both these methods show significant improvement over the DCASE Challenge baseline AUC scores, with an average improvement of 6% across all machines. An ensemble of the two systems is shown to further improve the average performance by 11% over the baseline.

**Index Terms**— i-Vectors, Amplitude-Modulation Spectrums, Graph, Clustering, Gaussian Mixture Models

### 1. INTRODUCTION

Monitoring industrial machinery can prevent the production of faulty products and decrease the chances of machine breakdown. Anomalous sounds can indicate symptoms of unwanted activity, hence, Anomalous Sound Detection (ASD) systems can potentially be used for real time monitoring of machines. Unsupervised ASD systems rely on only “normal” sounds for identifying anomalies. Developing ASD systems in an unsupervised manner is essential, as: (i) the nature of anomalies may not be known beforehand, and (ii) deliberately destroying expensive devices is impractical from a development cost perspective. In addition, the frequency at which anomalies occur is low and the variability in the type of anomaly is high, therefore, developing balanced datasets for supervised learning is difficult.

In our proposed systems, we focus on using features which are able to capture anomalous behaviour. Simple machine learning algorithms when used on top of these features are able to beat the baseline performance[1]. In our first system, we propose an outlier detection method which is similar to a nearest neighbour search. In this method, clusters of normal sounds are formed using a nearest neighbour graph to capture variations in the normal working sounds of a machine. Anomalies are then identified based on their distance from these clusters. Modulation spectrum features are used for this

system. The features are extracted from the so-called modulation spectrum (MS) signal representation, which quantifies the rate of change of the signal spectral components over time. These features have previously been useful for stress detection in speech [2], for speech enhancement [3], and room acoustic characterization [4], to name a few applications.

In our second system, in turn, we use i-vectors and Gaussian Mixture Models (GMM) for anomaly detection. i-vectors have been widely used for speech applications, including speech, speaker, language, and accent recognition. We extract i-vectors from MFCC features and use them to train GMMs for anomaly detection. The negative log likelihood for a sample is used as its anomaly scores. Lastly, an ensemble of these two systems is also experimented with.

### 2. SYSTEM DESCRIPTION

#### 2.1. System 1 - Graph Clustering using Modulation Spectrograms

##### 2.1.1. Pre-processing and Feature Extraction

Modulation spectrum corresponds to an auditory spectro-temporal representation that captures long-term dynamics of an audio signal. The pipeline proposed in [5] is used to extract modulation spectrograms.

Prior to feature extraction, noise reduction is performed on the signal through a spectral gating method (using *noisereduce*<sup>1</sup> in Python), described as follows: 100 normal training sound clips for a machine-id are averaged and used as a noise clip for that machine-id. An FFT is calculated over this noise clip and statistics including the mean power, are tabulated for each frequency band. A threshold for each frequency band is calculated based upon the statistics. An FFT is calculated over the signal. A mask is determined by comparing the signal FFT to the threshold. The mask is smoothed with a filter over frequency and time. The mask is applied to the FFT of the signal, and is inverted.

The speech activity level is normalized to -26 dBov (dB overload), after noise removal thus eliminating unwanted energy variations caused by different loudness levels in the speech signal. Next, the pre-processed speech signal  $\hat{x}(n)$  is filtered by a 60-channel gammatone filterbank, simulating cochlear processing [6]. The first filter of the filterbank is centered at 125 Hz and the last one at at just below half of the sampling rate. Each filter bandwidth follows the

<sup>1</sup><https://pypi.org/project/noisereduce/>

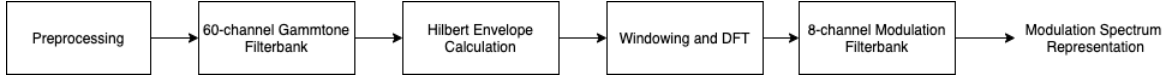


Figure 1: Block diagram describing steps for computing the modulation spectral representation

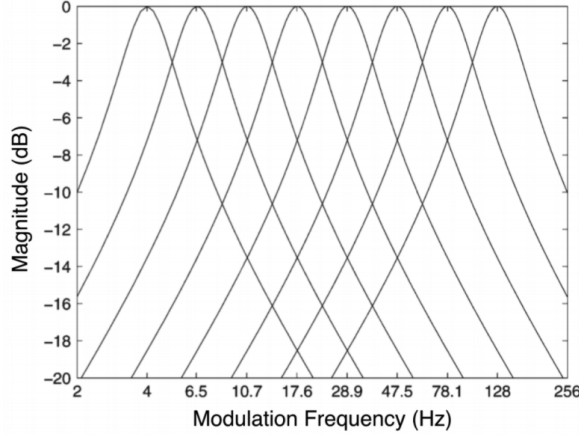


Figure 2: Frequency responses of the 8-channel modulation filterbank

equivalent rectangular bandwidth (ERB), which is an approximation of the bandwidths of the filters in human hearing, as described below:

$$ERB_j = \frac{f_j}{Q_{\text{ear}}} + B_{\text{min}}, \quad (1)$$

where  $f_j$  represents the center frequency of the  $j$ -th filter.  $Q_{\text{ear}}$  represents the asymptotic filter quality at high frequencies and  $B_{\text{min}}$  is the minimum bandwidth for low frequencies. They are set, respectively, to 9.265 and 24.7.

The temporal envelope  $e_j(n)$  is then computed from  $\hat{x}_j(n)$ , the output of the  $j$ -th acoustic filter, via the Hilbert transform:

$$e_j(n) = \sqrt{\hat{x}_j(n)^2 + \mathcal{H}\{\hat{x}_j(n)\}^2} \quad (2)$$

where  $\mathcal{H}\{\cdot\}$  denotes the Hilbert Transform. Temporal envelopes  $e_j(n)$ ,  $j = 1, \dots, 60$  are then windowed with a 256-ms Hamming window and shifts of 40 ms. The discrete Fourier transform  $\mathcal{F}\{\cdot\}$  of the temporal envelope  $e_j(m; n)$  ( $m$  indexes the frame) is then computed in order to obtain the modulation spectrum  $E_j(m, f_m)$ , i.e.,

$$E_j(m; f_m) = \|\mathcal{F}\{e_j(m; n)\}\| \quad (3)$$

where  $m$  represents the  $m$ -th frame obtained after every Hamming window multiplication and  $f_m$  designates modulation frequency. The time variable  $n$  is dropped for convenience. Lastly, following recent physiological evidence of a modulation filterbank structure in the human auditory system [6], an auditory-inspired modulation filterbank is further used to group modulation frequencies into eight bands. These are denoted as  $\mathcal{E}_{(j,k)}(m)$ ,  $k = 1, \dots, 8$ , where  $j$  indexes the gammatone filter and  $k$  the modulation filter. Figure 3 depicts the frequency response for the 8-channel modulation filterbank used in our system. Note that the filter center frequencies are equally spaced in the logarithmic scale from 4 to 128 Hz. The modulation spectral representation frames obtained over time are averaged for

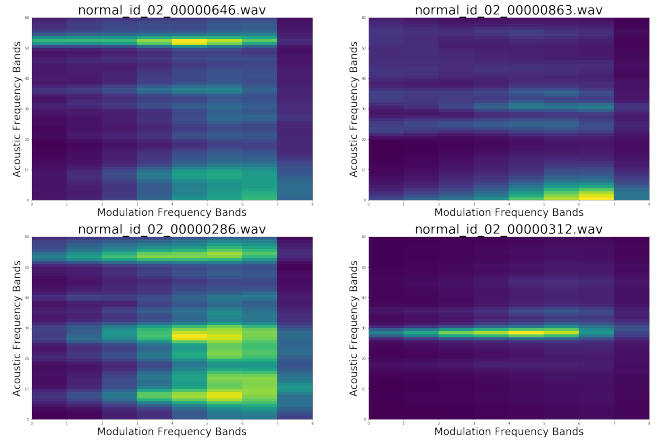


Figure 3: Modulation spectrograms for four normal training samples from Pump machine-id 2

all our experiments. This results in a  $60 \times 8$  modulation spectral representation (i.e., modulation spectrogram) for each sound clip.

### 2.1.2. Anomaly Detection

Figure 3 shows four modulation spectrograms for the normal training samples of Pump machine-id 2. It can be seen that a significant amount of variability exists within the same machine-id. We capture this variability by a graph-based clustering approach using modulation spectra as features.

Consider a graph  $G = (V, E)$  where  $V$  is the set of nodes comprising of the normal training sound clips.  $E$  is the set of edges connecting the nodes. Two nodes  $v_p, v_q \in V$  share an edge  $e_{pq} \in E$  such that  $q = \text{argmin}(\{D(v_p, v_r) \mid r = 1, \dots, p-1, p+1, \dots, |V|\})$ . Here  $D(v_p, v_r)$  is the L1 distance between  $v_p, v_r$ . The graph  $G$ , when constructed in this manner, consists of several disjoint subgraphs i.e.  $G = g_1 \cup g_2 \cup g_3 \dots g_n$  such that  $g_l \cap g_m = \phi \forall l, m = 1, 2, \dots, n \mid l \neq m$ . Each of these subgraphs is treated as a separate cluster. A centroid  $\mu_l$  and standard deviation  $\sigma_l$  are calculated corresponding to each cluster by taking the mean and standard deviation of all frequency bins.  $\mu_l, \sigma_l$  both have dimensions  $60 \times 8$ . This graph  $G$  and the corresponding cluster centroids, standard deviations are computed separately for each machine-id.

Anomaly score for a sound clip of a machine-id is calculated using the standard-deviation normalized distance from each cluster centroid corresponding to that machine-id. For a given sample in the test dataset  $v_t \notin V$ , its anomaly score  $A_t$  is given by  $A_t = \min(\{z_{t,l} \mid \forall l\})$

$$z_{t,l} = \sum_{j=1}^{60} \sum_{k=1}^8 \frac{|z_{t,j,k} - \mu_{l,j,k}|}{\sigma_{l,j,k}}. \quad (4)$$

Here,  $z_{t,j,k}$  is the energy value at the  $j^{\text{th}}$  gammatone filterbank

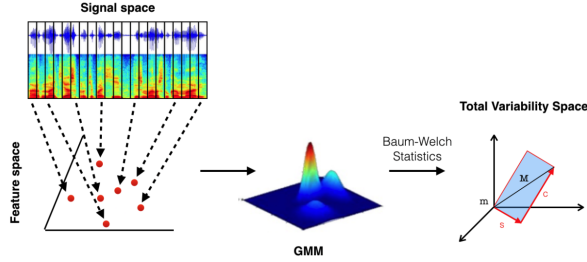


Figure 4: Diagram describing the steps for i-vector extraction

and  $k^{\text{th}}$  modulation filterbank. The intuition behind this strategy for finding anomaly scores is that normal samples will lie closer to the one of the cluster centroids in comparison to anomalous samples. This method does not require any training and can be seen as a KNN based ASD system where instead of computing distances from each training sample, we only find the distances from the cluster centroid.

## 2.2. System 2 - GMMs using MFCC i-Vectors

### 2.2.1. Feature Extraction

The i-vector framework maps a list of feature vectors,  $O = \{o_t\}_{t=1}^N$ , where  $o_t \in \mathbb{R}^F$ , and  $N$  is the frame index. Typically Mel-frequency cepstral coefficients (MFCC's) extracted from an utterance, into a fixed-length vector,  $n \in \mathbb{R}^D$ . In order to achieve that, a Gaussian mixture model (GMM),  $\lambda = (\{w_k\}, \{m_k\}, \{\sigma_k\})$ , is used. The GMM, trained on multiple utterances, is referred to as the universal background model (UBM), and is used to collect Baum-Welch statistics from each utterance [7]. Such statistics are computed for each mixture component  $k$ , resulting in the so-called supervector  $M \in \mathbb{R}^{FK}$ , where  $F$  represents the feature dimension and  $K$  is the number of Gaussian components. As in the Joint Factor Analysis (JFA) [8], the i-vector framework also considers that speaker and channel variability lies in a lower subspace of the GMM supervectors [9]. The main difference between the two approaches is that the i-vector projects both speaker and channel variability into the same subspace, namely total variability space, represented as follows:

$$M = m + Tw, \quad (5)$$

where  $M$  is the dependent supervector (extracted from a specific utterance) and  $m$  is the independent supervector (extracted from the UBM),  $T$  corresponds to a rectangular low-rank total variability matrix and  $w$  is a random vector with a normal distribution, the so-called i-vector. In our experiments, a 100-dimensional i-vector was adopted extracted on top of MFCC features.

*Mel frequency spectrum coefficients* - Prior to their extraction, the input signals (sampled at 16 kHz) are normalized to -26 dBOV. The signals also undergo a pre-emphasis filter of coefficient 0.95, which is meant to balance low and high frequency magnitudes. A 30-ms Hamming window with 50% overlap is applied before extracting the MFCCs. The Hamming window is used to remove edge effects [10]. The cepstral feature vector can then be extracted from each frame according to:

$$c_n = \sum_{m=1}^M [Y_m] \cos \left[ \frac{\pi n}{M} \left( m - \frac{1}{2} \right) \right], n = 1, 2, 3, \dots, N, \quad (6)$$

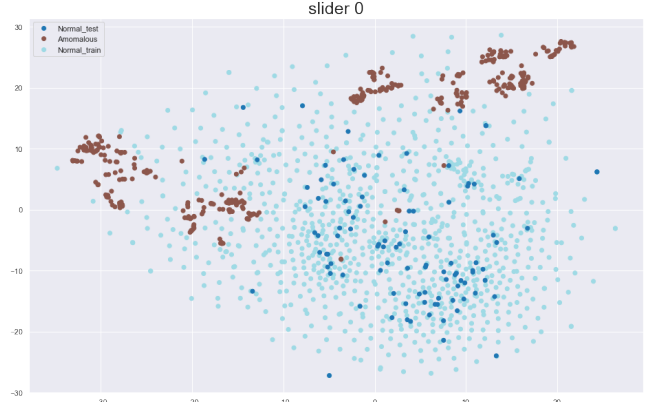


Figure 5: 2D t-SNE projections of i-Vectors corresponding to machine-slider, machine-id 0

where  $c_n$  is the  $n^{\text{th}}$  mel-cepstral coefficient and  $Y_m$  refers to the log-energy of the  $m^{\text{th}}$  filter. In this work, a set of 13 coefficients together with log energy, delta and delta-delta coefficients form the feature vector from each frame.

### 2.2.2. Anomaly Detection

The 100 dimensional i-vectors extracted from the normal training sounds are used to train a Gaussian Mixture Model using scikit-learn [11]. Ten mixture components are used for all machines with each component having its own general full covariance matrix. The ability of i-vectors to capture anomalous behaviour is depicted in Figure 5. Here, 2D t-SNE embeddings show that i-vectors corresponding to anomalous sounds are well separated from the normal sounds. i-vectors seem to be spread according to a Gaussian density, which is the key motivation behind using GMMs for anomaly detection.

The negative log-likelihood of sample  $x$  is given by:

$$-\log P(x|\pi, \mu, \sigma) = -\log \left\{ \sum_{k=1}^K \pi_k \mathbf{N}(x|\mu_k, \sigma_k) \right\}, \quad (7)$$

where  $\pi_k$  is the mixing coefficient for the  $k^{\text{th}}$  component of the GMM and  $\mu_k, \sigma_k$  are the corresponding mean and co-variance matrices. These values are used as the anomaly score.

## 2.3. System 3 - Graph-i-vector Ensemble

Lastly, an ensemble of the two proposed system shows performance improvement in several cases. Anomaly scores obtained from the two systems are first normalized using their minimum and maximum values. The ensemble anomaly scores are then computed by taking the geometric mean of the normalized values.

## 3. RESULTS

The results are shown in Table 1. When trained on development data[12, 13], the graph clustering system outperforms the baseline[1] by an average of 6% and 32% AUC for machines slider and valve respectively. The i-vector GMM system outperforms baseline and graph clustering system for some of the machine IDs

Table 1: AUC and pAUC scores for all machines in the development dataset. Modspec Graph, iVGmm, Ensemble correspond to Systems 1-3, respectively. The best scores for each case have been shown in bold.

Machine	Mid	Baseline AUC	Modspec Graph AUC	iVGmm AUC	Ensemble AUC	Baseline pAUC	Modspec Graph pAUC	iVGmm pAUC	Ensemble pAUC
ToyCar	1	81.36%	78.24%	75.04%	<b>81.64%</b>	<b>68.40%</b>	64.69%	57.54%	66.75%
	2	<b>85.97%</b>	89.06%	83.30%	<b>91.72%</b>	77.72%	76.14%	67.00%	<b>79.78%</b>
	3	63.30%	67.16%	79.47%	<b>78.21%</b>	55.21%	52.58%	59.52%	<b>56.37%</b>
	4	84.45%	89.40%	94.84%	<b>96.44%</b>	68.97%	63.54%	82.94%	<b>84.80%</b>
	Avg	78.77%	80.96%	83.16%	<b>87.00%</b>	67.58%	64.24%	66.75%	<b>71.92%</b>
ToyConveyor	1	<b>78.07%</b>	62.56%	55.51%	64.62%	<b>64.25%</b>	51.59%	52.82%	52.24%
	2	<b>64.16%</b>	54.03%	53.80%	56.65%	<b>56.01%</b>	49.99%	50.95%	50.23%
	3	<b>75.35%</b>	59.10%	59.09%	64.06%	<b>61.03%</b>	50.31%	52.82%	52.25%
	Avg	<b>72.53%</b>	58.57%	56.13%	61.78%	<b>60.43%</b>	50.63%	52.20%	51.58%
fan	0	54.41%	63.37%	<b>67.85%</b>	67.12%	49.37%	49.73%	<b>57.38%</b>	52.92%
	2	73.40%	79.32%	70.39%	<b>80.48%</b>	54.81%	57.16%	<b>61.93%</b>	59.21%
	4	61.61%	71.76%	73.52%	<b>78.07%</b>	53.26%	50.68%	<b>57.53%</b>	53.99%
	6	73.92%	74.00%	81.15%	<b>81.90%</b>	52.35%	49.38%	<b>56.31%</b>	49.23%
	Avg	65.83%	72.11%	73.23%	<b>76.89%</b>	52.45%	51.74%	<b>58.29%</b>	53.84%
pump	0	67.15%	86.66%	74.99%	<b>86.95%</b>	56.74%	<b>82.52%</b>	67.10%	78.32%
	2	61.53%	62.44%	<b>74.91%</b>	70.06%	58.10%	64.77%	60.08%	<b>65.72%</b>
	4	88.33%	84.16%	<b>92.02%</b>	90.73%	67.10%	59.95%	73.74%	<b>68.00%</b>
	6	74.55%	81.64%	71.10%	<b>82.65%</b>	58.02%	66.20%	51.70%	<b>66.56%</b>
	Avg	72.89%	78.72%	78.26%	<b>82.60%</b>	59.99%	68.36%	63.15%	<b>69.65%</b>
slider	0	96.19%	<b>99.91%</b>	83.92%	98.72%	81.44%	<b>99.53%</b>	50.04%	93.44%
	2	78.97%	<b>84.36%</b>	56.93%	77.93%	63.68%	<b>73.86%</b>	47.84%	52.89%
	4	94.30%	<b>97.83%</b>	87.84%	95.93%	71.98%	<b>88.59%</b>	62.71%	79.69%
	6	69.59%	<b>79.03%</b>	59.04%	71.40%	49.02%	<b>55.47%</b>	49.91%	52.28%
	Avg	84.76%	<b>90.28%</b>	71.93%	86.00%	66.53%	<b>79.36%</b>	52.63%	69.57%
valve	0	68.76%	<b>100.00%</b>	79.33%	98.80%	51.70%	<b>100.00%</b>	52.94%	95.62%
	2	68.18%	<b>99.88%</b>	85.35%	98.69%	51.83%	<b>99.34%</b>	56.27%	93.29%
	4	74.30%	<b>98.26%</b>	84.10%	95.88%	51.97%	<b>91.32%</b>	56.32%	80.26%
	6	53.90%	<b>89.22%</b>	69.84%	85.01%	48.43%	<b>72.59%</b>	49.91%	59.65%
	Avg	66.28%	<b>96.84%</b>	79.65%	94.59%	50.98%	<b>90.81%</b>	53.86%	82.21%

in pump and fan. The ensemble of graph clustering system and i-vector GMM system outperforms the baseline by an average AUC score of 8%, 11% and 10% for machines ToyCar, fan, and pump, respectively. Interestingly, the performance for the ToyConveyor case was lower than that achieved by the benchmark system for all three proposed systems. This may be due to anomalies which occur for a very small time interval and are not being captured by the proposed longer-term features. We provide our implementation here<sup>2</sup>.

#### 4. REFERENCES

- [1] Y. Koizumi, Y. Kawaguchi, K. Imoto, T. Nakamura, Y. Nikaido, R. Tanabe, H. Purohit, K. Suefusa, T. Endo, M. Yasuda, and N. Harada, "Description and discussion on DCASE2020 challenge task2: Unsupervised anomalous sound detection for machine condition monitoring," in *arXiv e-prints: 2006.05822*, June 2020, pp. 1–4. [Online]. Available: <https://arxiv.org/abs/2006.05822>
- [2] A. R. Avila, S. R. Kshirsagar, A. Tiwari, D. Lafond, D. O'Shaughnessy, and T. H. Falk, "Speech-based stress classification based on modulation spectral features and convolutional neural networks," in *2019 27th European Signal Processing Conference (EUSIPCO)*. IEEE, 2019, pp. 1–5.
- [3] S. Karimian-Azari and T. H. Falk, "Modulation spectrum based beamforming for speech enhancement," in *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2017, pp. 91–95.
- [4] T. H. Falk and W.-Y. Chan, "Temporal dynamics for blind measurement of room acoustical parameters," *IEEE Transactions on Instrumentation and Measurement*, vol. 59, no. 4, pp. 978–989, 2010.
- [5] T. H. Falk and W. Chan, "Modulation spectral features for robust far-field speaker identification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 1, pp. 90–100, 2010.
- [6] M. Slaney, "An efficient implementation of the patterson-holdsworth auditory filter bank," 1997.
- [7] D. Garcia-Romero and C. Y. Espy-Wilson, "Analysis of i-vector length normalization in speaker recognition systems," in *Twelfth Annual Conference of the International Speech Communication Association*, 2011.

<sup>2</sup><https://github.com/parth2170/DCASE2020-Task2>

- [8] P. Kenny, G. Boulianne, P. Ouellet, and P. Dumouchel, "Joint factor analysis versus eigenchannels in speaker recognition," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 4, pp. 1435–1447, 2007.
- [9] J. H. L. Hansen and T. Hasan, "Speaker recognition by machines and humans: A tutorial review," *IEEE Signal processing magazine*, vol. 32, no. 6, pp. 74–99, 2015.
- [10] B. Logan and et al., "Mel frequency cepstral coefficients for music modeling," in *Ismir*, vol. 270, 2000, pp. 1–11.
- [11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [12] Y. Koizumi, S. Saito, H. Uematsu, N. Harada, and K. Imoto, "ToyADMOS: A dataset of miniature-machine operating sounds for anomalous sound detection," in *Proceedings of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, November 2019, pp. 308–312. [Online]. Available: <https://ieeexplore.ieee.org/document/8937164>
- [13] H. Purohit, R. Tanabe, T. Ichige, T. Endo, Y. Nikaido, K. Suefusa, and Y. Kawaguchi, "MIMII Dataset: Sound dataset for malfunctioning industrial machine investigation and inspection," in *Proceedings of the Detection and Classification of Acoustic Scenes and Events 2019 Workshop (DCASE2019)*, November 2019, pp. 209–213. [Online]. Available: [http://dcase.community/documents/workshop2019/proceedings/DCASE2019Workshop\\\_Purohit\\\_21.pdf](http://dcase.community/documents/workshop2019/proceedings/DCASE2019Workshop\_Purohit\_21.pdf)