

NATIONAL RESEARCH UNIVERSITY
HIGHER SCHOOL OF ECONOMICS

Faculty of Computer Science
Bachelor's Programme "Data Science and Business Analytics"

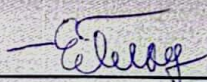
Research Project Report on the Topic:

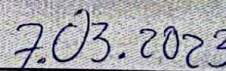
"Structural breaks identification in sound signals"

Title of the Project

Fulfilled by:

Student of the Group БПАД201
Gladilin Egor Vyacheslavovich


(signature)


(date)

Checked by the Project Supervisor:

Lukyanchenko Petr Pavlovich
Senior Lecturer
Faculty of Computer Science, HSE University



(signature)

7.03.2023

(date)

Moscow 2023

Contents

1	Introduction	2
2	Dataset	3
3	Models	4
3.1	Group-MADE + MobileNetV2 ensemble	4
3.1.1	Data preprocessing	4
3.1.2	How Group-MADE works	4
3.1.3	Classification	5
3.2	Ensemble of Auto-Encoder Based and WaveNet like Systems	5
3.2.1	Heteroskedastic Variational Auto-Encoder (HVAE) . .	6
3.2.2	ID Conditioned Auto-Encoder (IDCAE)	6
3.2.3	FREAK	7
3.3	Outlier-exposed classification	8
3.3.1	Data preprocessing	8
3.3.2	ResNet for anomaly detection	8
3.4	Ensemble of statistical methods	9
4	Models' benchmarking	10

1 Introduction

Structural breaks identification is not a new challenge, which relates to many spheres: finance, econometrics, physics, medicine, climate forecasting. In statistics and econometrics structural breaks are known as unexpected changes over time in the parameters of regression models. There were invented many methods to identify them like Chow test or CUSUM test.

Within the epoch of AI development this challenge went further and now it is addressed to finding anomalies in videos, photos, text information, time series using neural networks and has other applications. With the usage of deep learning methods there were introduced multiple interesting solutions, which helps to overcome: unknown new data (as anomalies are distinct from the majority of data, model has to generalise known data efficiently and identify whether the difference exists), imbalance of classes (as anomalies occurs rarely), heterogeneity of different classes of objects (depending on the class the difference between anomaly and normal object can distinguish). Solving these issues, programmers faces low accuracy scores classifying data as abnormal/normal, lack of marked-up data, dealing with high-dimensional data, anomaly detection of a particular abnormal group, difficulty of results interpretation.

Sound processing is considered as a modern field in deep learning, however many approaches from CV and NLP are well adopted and widely used there. Hence, they are used to detect anomalies in weather changes (rain, wind, storm, hail), malfunctions in the operation of mechanisms, unusual road traffic, out-of-tune musical instrument, echolocation of ships and submarines, anomalies of air flows in air conditioning systems.

In this study anomaly detection in sound signals is explored. However, solutions are highly dependent on the given dataset, so it was decided to choose a particular dataset and explore models that were constructed specially for it. In the article I discuss the dataset, study the models and give the comparative analysis based on the metric's scores.

2 Dataset

First of all, before going to the models, one has to choose and describe the properties of the dataset. For this study, it was decided to unite 2 datasets: ToyADMOS [7] and MIMII [11]. Both dataset contains audio recordings of mechanisms' performance stored in "wav" files. There are 6 types of mechanisms: Toy-car (ToyADMOS), Toy-conveyor (ToyADMOS), Valve (MIMII), Pump (MIMII), Fan (MIMII), Slider (MIMII). For each type of mechanisms there were chosen several distinct machines, which id's are labeled on Figure 1. The united dataset consists of development dataset, additional training dataset and evaluation dataset.

Development dataset contains nearly 1000 normal sound recordings of each machine for training and around 200 recordings of each normal and abnormal to test.

Additional training dataset contains around 1000 normal recordings for each machine form evaluation dataset.

In evaluation dataset Machine IDs coincide with machine IDs from additional training dataset. This one is used only for estimating models and contains around 400 unlabeled recording for each machine.

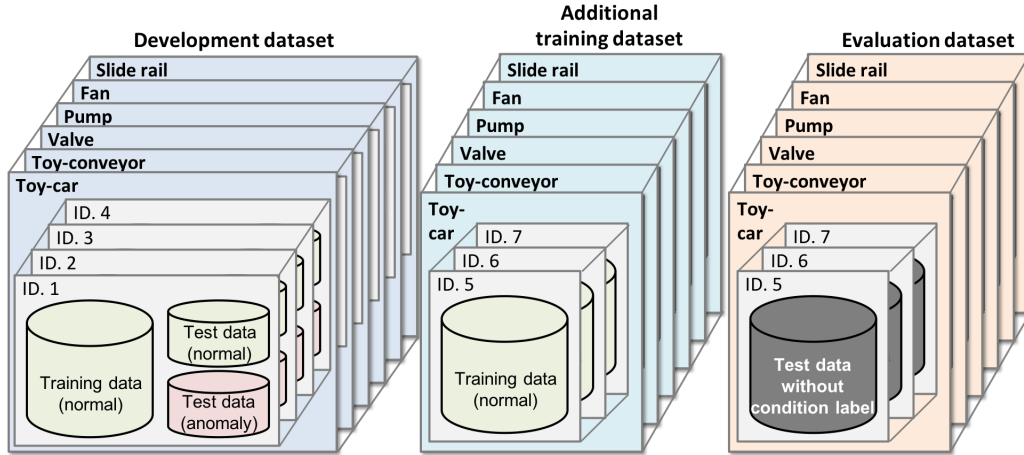


Figure 1: The dataset visualisation

3 Models

There already exists models that well-enough deal with anomaly detection for the dataset. In this section properties of the models and their metric's scores are discussed. The metrics for evaluating models' performance are Area Under the ROC Curve (AUC) and partial AUC (pAUC).

3.1 Group-MADE + MobileNetV2 ensemble

The model architecture and usage called Group-MADE was described in [4] and [3].

3.1.1 Data preprocessing

The authors propose to use log-Mel spectrograms as inputs constructed in the following way: each recording is split in 64ms frames having 32ms hop length between frames, then for each the log-Mel spectrogram is computed with 1024-FFT and 128 Mel bins and finally each 5 frames are concatenated resulting in $5 \times 128 = 640$ dimensional input.

3.1.2 How Group-MADE works

The core idea of MADE is to mask the connections between the layers of a vanilla autoencoder. Thus, the probability of a subsequent prediction depends on the previous ones and the result of the autoencoder is considered autoregressive for a specific sequence.

Thus, having autoregressive autoencoder, the probability density, \mathbf{x} is calculated in the following way:

$$p(\mathbf{x}) = \prod_{d=1}^D p(x_d | x_{<d}) \quad (1)$$

Where $x_{<d} = [x_1, \dots, x_{d-1}]^T$.

Defining $p(x_d = 1 | x_{<d}) = \hat{x}_d$ and $p(x_d = 0 | x_{<d}) = 1 - \hat{x}_d$, the loss-function becomes a negative log-likelihood:

$$-\log p(\mathbf{x}) = \sum_{d=1}^D -\log p(x_d | x_{<d}) \quad (2a)$$

$$= \sum_{d=1}^D -x_d \log p(x_d = 1 | x_{<d}) - p(x_d = 0 | x_{<d}) \quad (2b)$$

$$= \ell(x) \quad (2c)$$

Thus \hat{x}_d is a function that takes $x_{<d}$ as an input and outputs $p(x_d)$. This allows us to refer to the autoregression, as each next probability depends only on previous object.

The authors constructed Group-MADE in such way that each further prediction depends on previous objects in a frame, which are inputs, instead of the whole probability distribution for each dimension.

Having an input $t = [t_{i+1}, t_{i+2}, t_{i+3}, t_{i+4}, t_{i+5}]$, where i^{th} frame is t_i , the probability of t_i is computed as following:

$$p(\mathbf{t}) = \prod_{i=1}^5 p(t_i | t_{<i}) = \prod_{i=1}^5 \prod_{j=128}^5 p(t_{ij} | t_{<i}) \quad (3)$$

So, each Mel bin is dependent on the previous ones, but is not dependent on the other ones.

The autoencoder is constructed with fully connected layers with the following features: [128, 128, 128, 128, 32, 128, 128, 128, 128]

Negative log-likelihood was used as loss function and for estimating the normality for test samples. There is used Adam optimizer with learning rate set at 0.001 to train. As for training data, there were used all mechanism types without regard the machine IDs.

3.1.3 Classification

For classification task the authors decided to use self-supervised classification strategy. For each mechanism type, it was chosen to train MobileNetV2 [12] on normal data do not considering machine IDs to identify machine's ID and to distinguish real sample from synthetically generated versions.

Finally, there was added a softmax layer and negative softmax score on the test sample is considered to be the anomaly score.

The final solution represents the ensemble of Group-MADE and 2 MobileNetV2 with appended softmax layer.

3.2 Ensemble of Auto-Encoder Based and WaveNet like Systems

The solution proposed by the authors is the ensemble consisting of several methods, whose weighted anomaly score is aimed to maximize AUC and pAUC metrics. These combined 3 methods are separately explained below.

3.2.1 Heteroskedastic Variational Auto-Encoder (HVAE)

For each mechanism type was constructed its own HVAE [1], so HVAE is either trained on OpenL3 embeddings or on log-Mel spectrograms. The key modification proposed by the authors is changing the distribution. In vanilla VAE [6] the loss function is

$$D_{KL}(\mathcal{N}(\mu, \sigma^2) || \mathcal{N}(0, 1)) + MSE(X, \hat{X}) \quad (4)$$

Where D_{KL} stands for kl-divergence, $(\mu, \log(\sigma)) = E(X)$, $E(X)$ - encoder's output, \hat{X} - decoder's output (reconstruction). So the probability distribution $p(x)$ of the feature space is approximated by maximizing:

$$\log(p(x)) \geq E_{q(Z|X)} + D_{KL}(q(Z|X) || p(Z)) \quad (5)$$

$q(Z|X)$ is normal distribution with parameters set by encoder's output, $p(Z)$ is $\mathcal{N}(0, 1)$ and $p(X|Z)$ is normal distribution with $(\mu, \sigma) = (\hat{X}, 1)$. So the modification considers changing $p(X|Z)$, instead of setting 1 for σ , authors propose $\hat{\sigma} = D(Z)$. And with this change, not only reconstruction is estimated, but also its precision. Thus the new loss function:

$$D_{KL}(\mathcal{N}(\mu, \sigma^2) || \mathcal{N}(0, 1)) + \beta(wMSE(X, \hat{X}, \hat{\sigma}) - \sum_{i=1}^n \log(\hat{\sigma}_i)) \quad (6)$$

So the authors constructed several HVAE for each mechanism type, each HVAE has the same number of hidden layers. Estimation of $\hat{\sigma}$ is realised either by doubling an output layer, like in encoder, or by building separate decoder for variance, this approach is called Big Precision (BP). The precise architectures are described in the table in [1].

3.2.2 ID Conditioned Auto-Encoder (IDCAE)

The model uses basic principles of encoding and decoding previously described in article devoted to C2AE [9], however the input for decoder is changed. The authors firstly propose to encode machine IDs using one-hot approach, after that use some functions $H_\gamma, H_\beta : \mathbb{Y} \rightarrow \mathbb{Z}$, which take one-hot encoded label ell and map it to latent vector Z from \mathbb{Z} . Thus the decoder looks so: $D(H_\gamma(\ell) \times Z(E(X)) + H_\beta(\ell))$.

In IDCAE architecture DenseBlock implies combination of Dense, batch-normalization and relu. Having sound signals, they are constructed to Short Time Fourier Transform (STFT) with 1024-window and 512 hop length, then STFT is transformed to power mel-spectrogram with M mels and take log

Encoder(E)	H_γ, H_β	Decoder(E)
Input(F,M)	Input one-hot IDs	Input 16
Flatten	Dense 16	DenseBlock 128
DenseBlock 128	Sigmoid	DenseBlock 128
DenseBlock 64	Dense 16	DenseBlock 128
DenseBlock 32		DenseBlock 128
DenseBlock 16		Dense F · M
		Reshape(F, M)

Table 1: IDCAE architecture

with base 10 and multiply by 10. The input shape is (F, M) where, F is frame size and M is number of mels. For training Adam optimizer with basic, parameters, exponential learning rate and 100 epochs are used.

3.2.3 FREAK

The Freak model contains 3 ResidualBlocks that were taken from WaveNet [13] and One CausalConv1D. One ResidualBlock contains one causal convolution, 2 independent convolutions, then results of these convolutions are multiplied and go to one more convolution layer. After each convolution there is normalization and activation function, except for the last layer. Causal convolution is implemented by doing element-wise multiplication with the convolution kernel. For 1D outputs are realised just by shifting forward convolutions by a few elements. The input is represented by computed STFT with 2048-window and 512 hop length, then it is transformed to mel-spectrograms with either 64 or 128 bins. Afterward logarithm is applied to each mel-spectrogram. Finally, mean and variance based on the data are computed and standardization is applied. Adam is set to be optimizer with $\alpha = 0.001$, $\beta_1 = 0.85$, $\beta_2 = 0.999$, the batch size is equal to 32.

layer	channels	dilation	kernel	group
ResidualBlock	m*bins	1	3	4
ResidualBlock	m*bins	2	3	4
ResidualBlock	m*bins	4	3	4
CausalConv1D	bins	8	3	4

Table 2: FREAK architecture

3.3 Outlier-exposed classification

3.3.1 Data preprocessing

This method is based on the exposure of abnormal class. As training data contains only normal data, for classification we are lack of abnormal. The solution is to take into consideration that are neither belong to normal nor to abnormal classes - outliers. This allows to treat the problem as classification task and label outliers as abnormal objects.

After that the data transformation takes place. Firstly, audios are normalized to zero mean and unit variance. Then we calculate mel-spectrograms with sampling rate = 16000Hz, 1024 window, 512 hop length, 128 filters. Finally, by taking logarithm log-mel spectrograms are obtained.

3.3.2 ResNet for anomaly detection

The chosen model for classification is ResNet [5] which is widely used for classification tasks. Loss-function is Binary Cross Entropy (BCE), optimizer - Adam, having $\beta_1 = 0.9$, $\beta_2 = 0.999$, the batch size is equal to 64.

Layers	channels	Kernel size	ResidualBlock (RB)	Kernel size
Conv2d	5	5	Conv2d	ks1
BN			BN	
RB	5	(ks1=3, ks2=3)	Conv2d	ks2
MaxPool	5	2	BN	
RB	5	(ks1=3, ks2=3)	result + input	
MaxPool	5	2		
RB	5	(ks1=3, ks2=3)		
RB	5	(ks1=3, ks2=3)		
MaxPool	5	2		
RB	5	(ks1=1, ks2=1)		
RB	5	(ks1=1, ks2=1)		
RB	5	(ks1=1, ks2=1)		
Conv2d	1	1		
BN				
AvgPool				

Table 3: ResNet architecture

3.4 Ensemble of statistical methods

The authors [14] suggest to apply model based on geometric transformations for anomaly detection, however facing obstacles they add 2 statistical methods: scaling method [8] and probability aggregation method [2]. Authors assumes that the anomaly score of sub-model i denoted $s_i(x)$ is distributed by Gamma-distribution $s_i(x) \sim \Gamma(\alpha_i, \beta_i)$. For obtaining scaled scores one can use CDF $F_i(x; \alpha_i, \beta_i)$ over the training dataset distribution.

$$\hat{s}_i(x) = F_i(x; \alpha_i, \beta_i) \quad (7a)$$

$$= \frac{\gamma_i(\alpha_i, \beta_i, s_i(x))}{\Gamma(\alpha_i)} \quad (7b)$$

Ultimately, the anomaly score is obtained by calculating:

$$\hat{s}_e(x) = 1 - \prod_{i=1}^n (1 - \hat{s}_i(x)) \quad (8)$$

Where $\hat{s}_e(x)$ is the anomaly score and $\hat{s}_i(x)$ are scaled scores calculated in (7a), (7b).

4 Models' benchmarking

Here we can see the results that are demonstrated by the models. Overall AUC scores are high, from 92.7% up to 94,6%. The results can be explained by several reasons.

First of all, let's refer to data preparation and augmentation. Due to complex data transformation the models are able to generalise data in more efficient, this increase robustness of the models when the inputs are normal and vise versa in case id data is abnormal the outputs significantly differ from the normal and it is easier to detect them. In particular, [1] used the denoising properties of UNet to clear recordings from unwanted external sounds. This was used under assumption that abnormal sounds are distinct from normal and to clarify this all sounds except for machine sounds has to be removed. Another complex data augmentation was done by [10]. The authors did not the complex, just basic ResNet, and training this model on the dataset, which contains just normal, is likely to fail on test dataset (low metric scores). However taking into account the outliers and considering them as abnormal data makes the model more robust to test samples. As outliers are more similar to abnormal data ResNet shows low loss values during training, consequently the model is even more confident during test as abnormal data differs more from normal than outliers.

On the other side authors tried to improve models and maximize score due to complex models and ensembles. Thus, we see that [1], besides augmentation, uses an ensemble of several auto-encoder based models and diverse classification models, [4] used an ensemble of auto-encoder, that takes into account sequences and IDs, and classification model. [14] used a statistical methods to solve the tasks, nevertheless authors applied an ensemble and succeeded as well.

All the graphs below are results of benchmarking models described above.

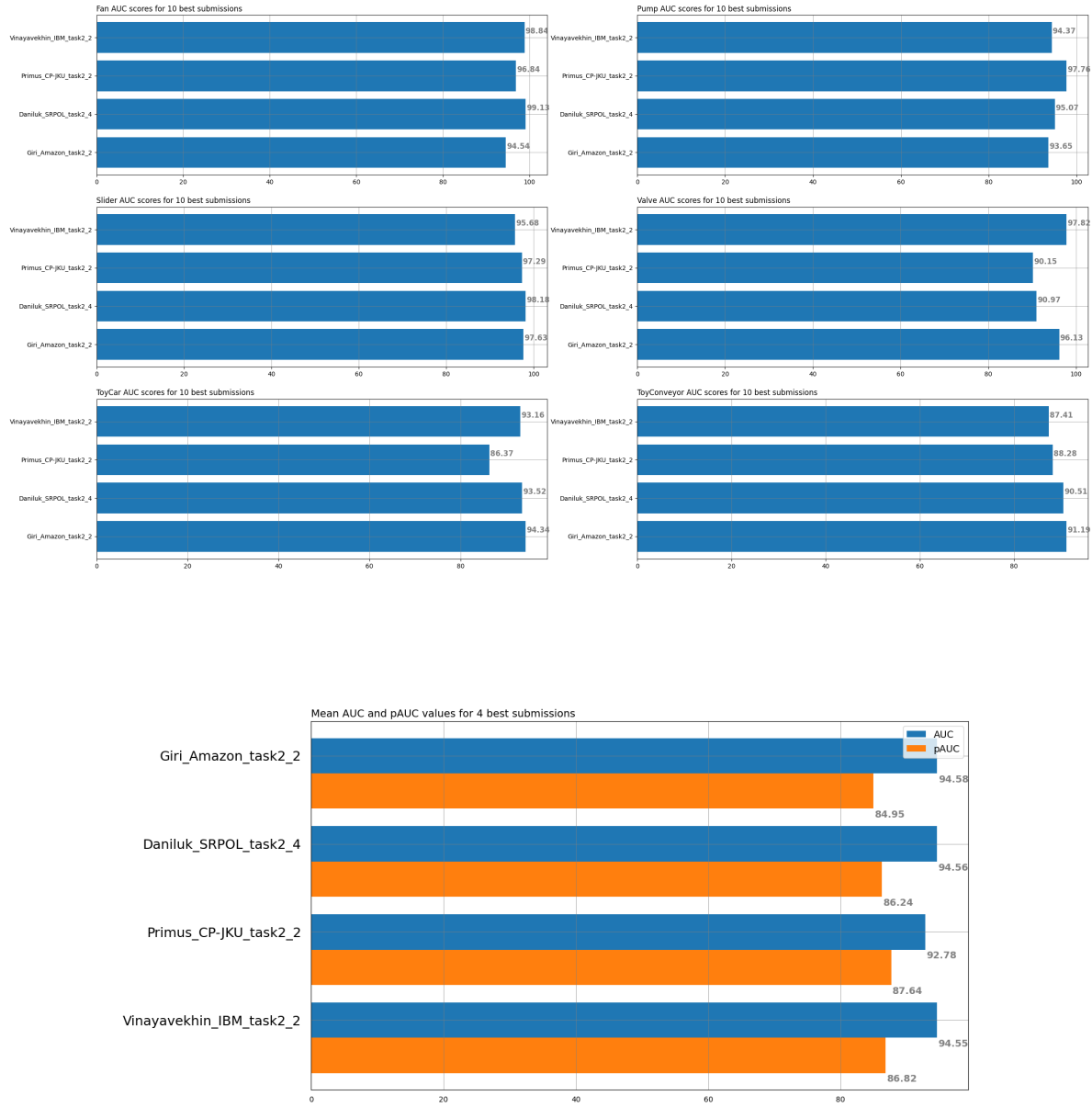


Figure 2: Graphical representation of the results

References

- [1] Pawel Daniluk, Marcin Gozdziwski, Slawomir Kapka, and Michal Kosmider. Ensemble of auto-encoder based systems for anomaly detection. Technical report, DCASE2020 Challenge, July 2020.
- [2] Jing Gao and Pang-ning Tan. Converting output scores from outlier detection algorithms into probability estimates. In *Sixth International Conference on Data Mining (ICDM'06)*, pages 212–221, 2006.
- [3] Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. MADE: Masked Autoencoder for Distribution Estimation. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *JMLR Proceedings*, pages 881–889. JMLR.org, 2015.
- [4] Ritwik Giri, Srikanth V. Teneti, Karim Helwani, Fangzhou Cheng, Umut Isik, and Arvinth Krishnaswamy. Unsupervised anomalous sound detection using self-supervised classification and group masked autoencoder for density estimation. Technical report, DCASE2020 Challenge, July 2020.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [6] Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.
- [7] Yuma Koizumi, Shoichiro Saito, Hisashi Uematsu, Noboru Harada, and Keisuke Imoto. Toyadmos: A dataset of miniature-machine operating sounds for anomalous sound detection. In *2019 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 313–317, 2019.
- [8] Hans-Peter Kriegel, Peer Kroger, Erich Schubert, and Arthur Zimek. *Interpreting and Unifying Outlier Scores*, pages 13–24.
- [9] Poojan Oza and Vishal M. Patel. C2ae: Class conditioned auto-encoder for open-set recognition. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2302–2311, 2019.
- [10] Paul Primus. Reframing unsupervised machine condition monitoring as a supervised classification task with outlier-exposed classifiers. Technical report, DCASE2020 Challenge, July 2020.

- [11] Harsh Purohit, Ryo Tanabe, Kenji Ichige, Takashi Endo, Yuki Nikaido, Kaori Suefusa, and Yohei Kawaguchi. MIMII dataset: Sound dataset for malfunctioning industrial machine investigation and inspection. *CoRR*, abs/1909.09347, 2019.
- [12] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks, 2019.
- [13] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *CoRR*, abs/1609.03499, 2016.
- [14] Phongtharin Vinayavekhin, Tadanobu Inoue, Shu Morikuni, Shiqiang Wang, Tuan Hoang Trong, David Wood, Michiaki Tatsubori, and Ryuki Tachibana. Detection of anomalous sounds for machine condition monitoring using classification confidence. Technical report, DCASE2020 Challenge, July 2020.