

CnnSound: Convolutional Neural Networks for the Classification of Environmental Sounds

OZKAN INIK *

Department of Computer Engineering, Tokat Gaziosmanpasa University, Tokat, Turkey

HUSEYIN SEKER

School of Computing and Digital Technologies, Staffordshire University, Stoke-on-Trent, UK

The classification of environmental sounds (ESC) has been increasingly studied in recent years. The main reason is that environmental sounds are part of our daily life, and associating them with our environment that we live in is important in several aspects as ESC is used in areas such as managing smart cities, determining location from environmental sounds, surveillance systems, machine hearing, environment monitoring. The ESC is however more difficult than other sounds because there are too many parameters that generate background noise in the ESC, which makes the sound more difficult to model and classify. The main aim of this study is therefore to develop more robust convolution neural networks architecture (CNN). For this purpose, 150 different CNN-based models were designed by changing the number of layers and values of their tuning parameters used in the layers. In order to test the accuracy of the models, the Urbansound8k environmental sound database was used. The sounds in this data set were first converted into an image format of 32x32x3. The proposed CNN model has yielded an accuracy of as much as 82.5% being higher than its classical counterpart. As there was not that much fine-tuning, the obtained accuracy has been found to be better and satisfactory compared to other studies on the Urbansound8k when both accuracy and computational complexity are considered. The results also suggest further improvement possible due to low complexity of the proposed CNN architecture and its applicability in real-world settings.

CCS CONCEPTS • Computing methodologies~Artificial intelligence~Computer vision•Computing methodologies~Machine learning~Machine learning approaches •Computing methodologies~Machine learning~Learning paradigms~Supervised learning •Applied computing~Computer forensics

Additional Keywords and Phrases: Environmental sound classification (ESC), Deep Learning, Convolutional Neural Networks (CNN), Urbansound8k

1 INTRODUCTION

Sound data contains more semantic information than visual data [1]. In particular, sound data becomes more important to obtain information about an environment. In order to realize some applications in daily life, it is necessary to use environmental sounds unlike speech and music sounds. For this reason, studies on the classification of urban sounds have intensified in recent years. Environmental Sounds Classification (ESC) is regarded as one of the most important issues of the non-speech voice classification task [2]. ESC is of critical importance in many problems such as; noise pollution analysis [3, 4], surveillance systems [5-7], context-aware applications [1, 8-13], machine hearing [14-17], environment monitoring [18], crime alert systems [19], soundscape assessment [20, 21], and smart city [22, 23]. Different data sets have been created for ESC task. ESC-10, ESC-50 [24] and Urbansound8k (US8K) [25] datasets are used extensively. Different statistical and machine learning methods have been used for ESC in the literature [1, 26-33].

* Corresponding author: ozkan.inik@gop.edu.tr.

The success rates of these methods are relatively low compared to deep learning-based studies in recent years. Deep learning [34] achieved a high success rate in the imageNet [35] competition in 2012. Due to this success, deep learning models have been used in different areas [36-42]. Also, deep learning based methods have been used for ESC task in recent years [43-56]. In general, it has been observed that the success rates obtained with deep learning models have better results than other artificial intelligence methods. The main reason for this can be summarized as an automatic feature discovery in deep learning models. Recently, it is seen that CNN models [2, 36-38, 41-46, 48-50] are used for ESC task. There are a lot of parameters that need to be adjusted in the design of CNNs. Therefore, the best CNN model can be found in different layer depths and different parameters. In this study, the task is to find the most suitable CNN model for ESC with an appropriate number of layers and values of their tuning parameters. The designed CNN model has been found to perform well in the ESC task compared to earlier works.

This paper is organized as follows. In Section 2, information about the features of the Urbansound8k ESC data set is given. In Section 3, the proposed CNN model is explained. The details of the experimental studies are provided in Section 4. Finally, the paper is concluded in Section 5.

2 DATA SET

In this study, Urbansound8k [25] data set is used for ESC task. Urbansound8k data set was obtained from real environment according to 4 seconds recording time. Environmental noise is present in the records obtained. The data set consists of 10 classes. They are air conditioner, car horn, children playing, dog bark, drilling, engine idling, gun shot, jackhammer, siren, and street music. These sounds are transformed into images by adapting the method of scalogram. The scalogram is the absolute value of the continuous wavelet transform (CWT) of a signal plotted as a function of time and frequency. Wavelet Toolbox of Matlab R2020b software was used for the conversion process. The transformed form of each class in the data set into sound signals and images is given in Figure 1. There are a total of 8732 records in the data set. The image resolution for training the CNN model is set to 32x32x3. For the statistical validation, 80% of the dataset was used for training, and 10% for validation and the remaining 10% for testing. The total number of images for each class for training, validation and testing are given in Table 1. Further information on the data set can be found in [57].

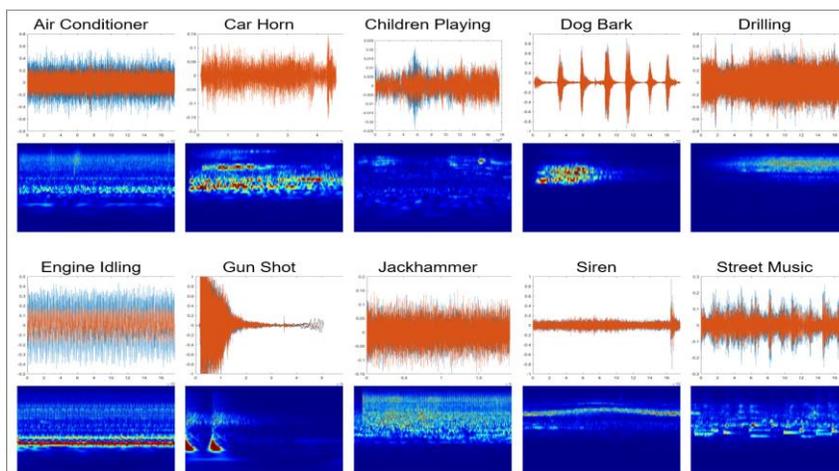


Figure 1. Sound(up) and translated image according to scalogram (down) of classes in Urbansound8k data set

Table 1. Number of records used for training, validation and testing in the Urbansound8k data set

Class	Number of images	Train	Validation	Test
Air Conditioner	1000	600	200	200
Car Horn	429	257	86	86
Children Playing	1000	600	200	200
Dog Bark	1000	600	200	200
Drilling	1000	600	200	200
Engine Idling	1000	600	200	200
Gun Shot	374	224	75	75
Jackhammer	1000	600	200	200
Siren	929	557	186	186
Street Music	1000	600	200	200

3 PROPOSED METHOD

In this study, the most suitable CNN model for the ESC task is obtained through a grid search. CNN models were designed and trained according to the layer depth and the parameter values used in the layers. The layer structure of the model that gives the best result among these models and the parameter values used in the layers are given in Figure 2. The proposed CNN model consists of 3 convolution layers, 1 pooling layer and 2 fully connected layers. There are 79 filters in the first convolution layer of the proposed model. After the training, the feature maps created by these filters and the effect of the filters on the input image is given in Figure 3.

Name	Type	Activations	Learnables	Total Learnables
input 32x32x3 images with 'zero-center' normalization	Image Input	32x32x3	-	0
conv_1 79 11x11x3 convolutions with stride [1 1] and padding 'same'	Convolution	32x32x79	Weights 11x11x3x79 Bias 1x1x79	28756
relu_1 ReLU	ReLU	32x32x79	-	0
conv_2 256 11x11x79 convolutions with stride [1 1] and padding 'same'	Convolution	32x32x256	Weigh... 11x11x79x2... Bias 1x1x256	2447360
relu_2 ReLU	ReLU	32x32x256	-	0
AveragePooling 2x2 average pooling with stride [6 6] and padding [0 0 0 0]	Average Pooling	6x6x256	-	0
conv_3 187 11x11x256 convolutions with stride [1 1] and padding 'same'	Convolution	6x6x187	Weigh... 11x11x256x1... Bias 1x1x187	5792699
relu_3 ReLU	ReLU	6x6x187	-	0
FC 797 fully connected layer	Fully Connected	1x1x797	Weights 797x6732 Bias 797x1	5366201
relu_4 ReLU	ReLU	1x1x797	-	0
droupOut 50% dropout	Dropout	1x1x797	-	0
FC_2 10 fully connected layer	Fully Connected	1x1x10	Weights 10x797 Bias 10x1	7980
softmax softmax	Softmax	1x1x10	-	0
classoutput crossentropyex with 'Air Conditioner' and 9 other classes	Classification Output	-	-	0

Figure 2. Architecture of the proposed CNN model and parameter information used in each layer

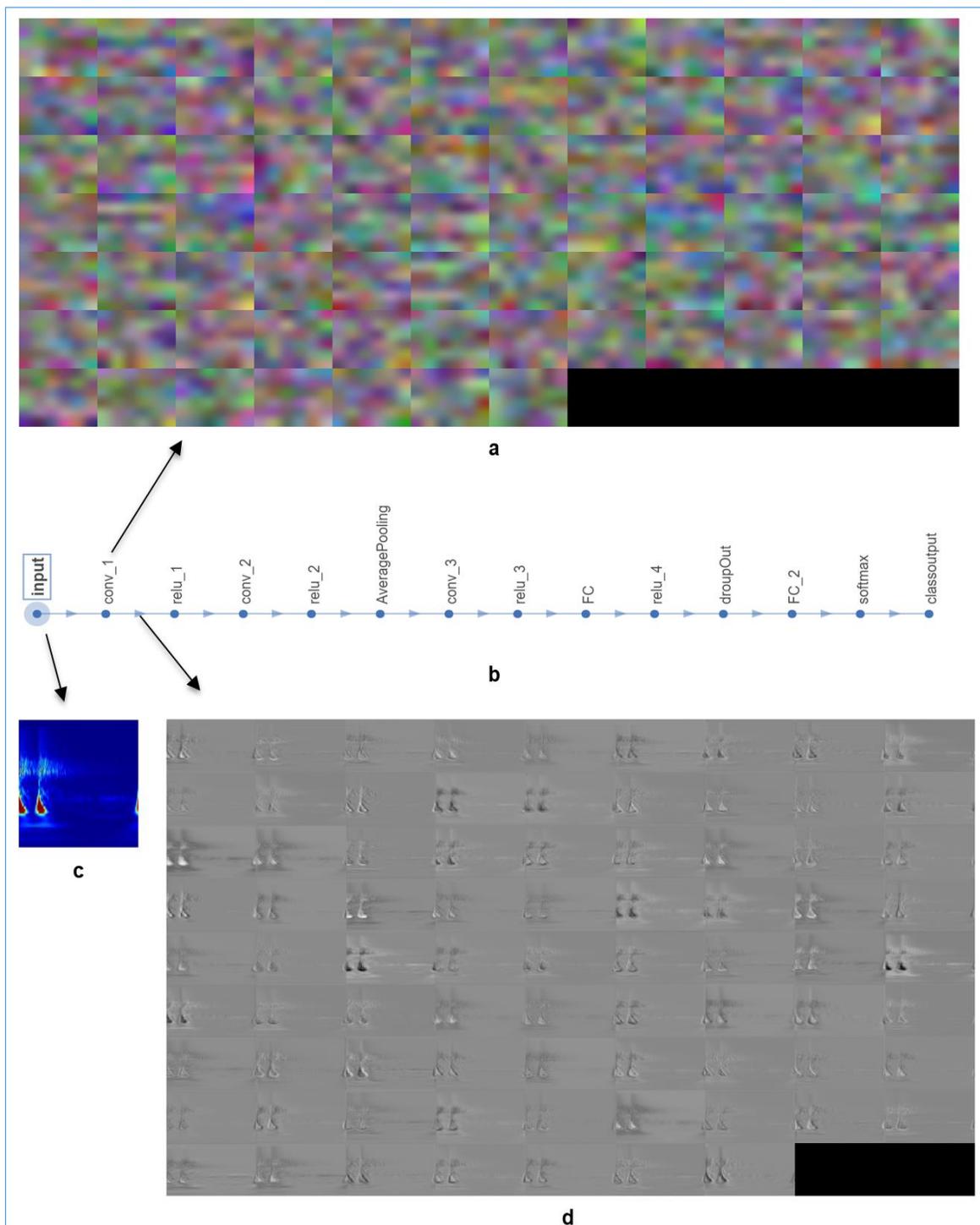


Figure 3. Layer architecture of the proposed CNN model (b), trained filters in the first convolution layer(a), the test image of the Gun shot class (c), the effect of each filter in the first convolution layer on the input image(d).

4 EXPERIMENTAL STUDIES

Experimental studies have been carried out on a computer with Intel® Core™ i9-7900X 3.30GHz×20 processor, 64 GB Ram and 2 x GeForce RTX2080Ti graphic card. Matlab R2020a 64bit (win64) has used as the software platform. The parameters used for the training of CNN model are given in Table 2.

Table 2. CNN training parameters

Parameters	Value
Optimizer	SGDM
Epochs	50
Learning rate drop factor	0.1
Learning rate drop period	10
Dropout rate	0.5
Mini Batch Size	256
Initial learning rate	0.001
Validation Frequency	50

In the studies, the accuracy of the most suitable CNN model was 82.26%. The confusion matrix obtained by this model is given in Figure 4. When confusion matrix is examined, it is seen that the most confused class with each other are Children Playing and Street Music. While the highest classification performance was achieved in the Car Horn class, the lowest classification success was achieved in the Engine Idling class. The graph of accuracy and validation values according to the epoch in the training phase of the CNN model is given in Figure 5, and the accuracy and validation loss graph is given in Figure 6. When Figures 5 and 6 are examined, it is seen that the model at the training stage reach the optimum performance approximately after the 15th epoch.

Output Class	Air Conditioner	Car Horn	Children Playing	Dog Bark	Drilling	Engine Idling	Gun Shot	Jackhammer	Siren	Street Music	Accuracy
Air Conditioner	178 10.2%	3 0.2%	7 0.4%	4 0.2%	0 0.0%	3 0.2%	0 0.0%	0 0.0%	2 0.1%	5 0.3%	88.1%
Car Horn	0 0.0%	48 2.7%	1 0.1%	1 0.1%	3 0.2%	0 0.0%	0 0.1%	1 0.1%	1 0.2%	4 0.2%	81.4%
Children Playing	8 0.5%	5 0.3%	142 8.1%	17 1.0%	2 0.1%	3 0.2%	2 0.1%	0 0.0%	10 0.6%	16 0.9%	69.3%
Dog Bark	2 0.1%	2 0.1%	11 0.6%	154 8.8%	1 0.1%	1 0.1%	1 0.1%	1 0.2%	3 0.5%	8 0.5%	83.7%
Drilling	2 0.1%	7 0.4%	6 0.3%	5 0.3%	167 9.6%	0 0.0%	0 0.0%	3 0.2%	3 0.2%	10 0.6%	82.3%
Engine Idling	2 0.1%	0 0.0%	1 0.1%	1 0.1%	1 0.1%	187 10.7%	1 0.1%	1 0.1%	0 0.0%	5 0.3%	94.0%
Gun Shot	0 0.0%	1 0.1%	1 0.1%	2 0.1%	0 0.0%	0 0.0%	69 3.9%	0 0.0%	0 0.0%	0 0.0%	94.5%
Jackhammer	3 0.2%	3 0.2%	10 0.6%	2 0.1%	13 0.7%	2 0.1%	2 0.1%	189 10.8%	0 0.0%	9 0.5%	81.1%
Siren	2 0.1%	6 0.3%	4 0.2%	4 0.2%	3 0.2%	1 0.1%	0 0.0%	0 0.0%	163 9.3%	3 0.2%	87.6%
Street Music	3 0.2%	11 0.6%	17 1.0%	10 0.6%	10 0.6%	3 0.2%	0 0.0%	5 0.3%	4 0.2%	140 8.0%	69.0%
	89.0%	55.8%	71.0%	77.0%	83.5%	93.5%	92.0%	94.5%	87.6%	70.0%	82.3%
	11.0%	44.2%	29.0%	23.0%	16.5%	6.5%	8.0%	5.5%	12.4%	30.0%	17.7%
	Air Conditioner	Car Horn	Children Playing	Dog Bark	Drilling	Engine Idling	Gun Shot	Jackhammer	Siren	Street Music	
											Target Class

Figure 4. Confusion matrix of proposed CNN model

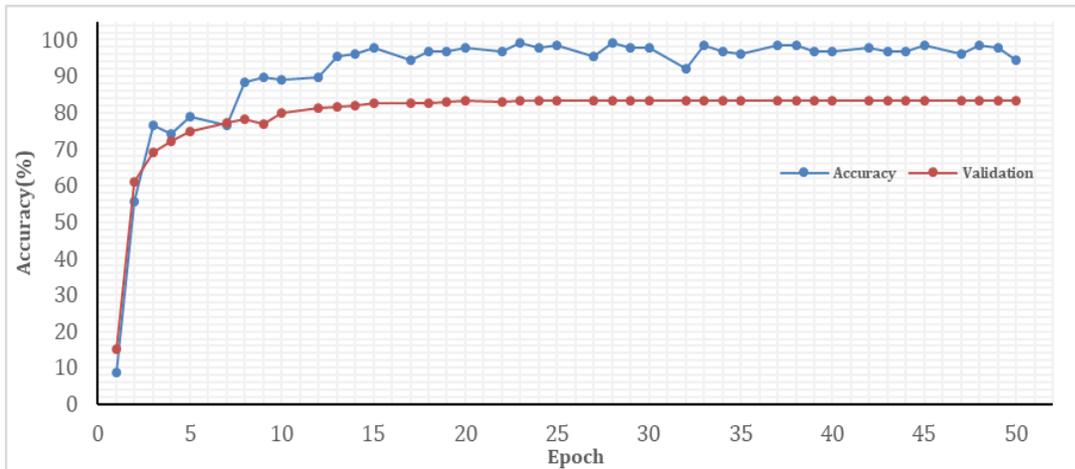


Figure 5. Accuracy and validation graph of the proposed CNN model during training.

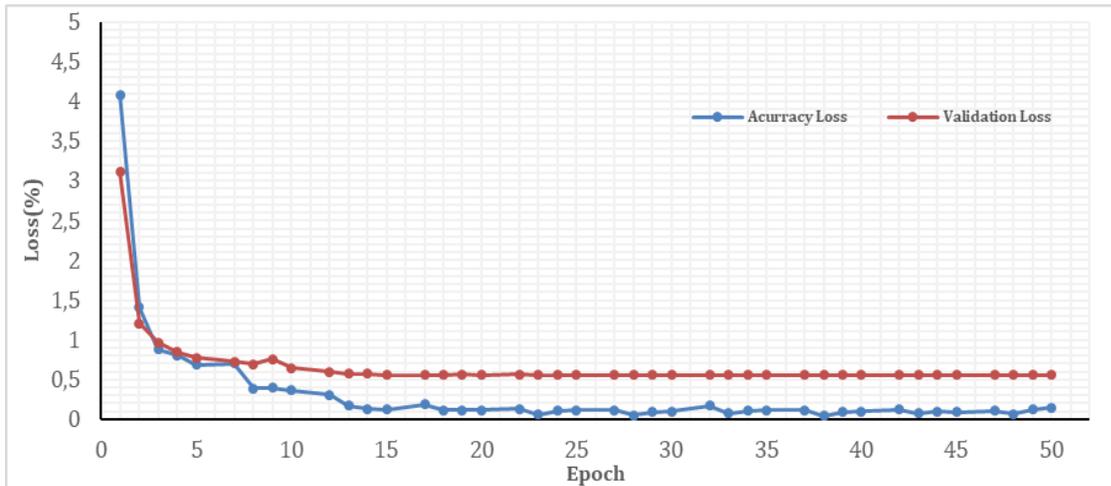


Figure 6. Accuracy loss and validation loss graph of the proposed CNN model during training.

4.1 Comparison with other studies

Different studies based on the deep learning have been conducted on the ESC data set. For the sake of comparison, accuracy values obtained in earlier studies are given in Table 3. It can be observed that the proposed CNN method achieves a very good performance. The method only performed lower than GoogLeNet and AlexNet. The reason for this is related to the image size obtained during the transformation of the data set. GoogLeNet input image size is 224x224x3 and AlexNet input image size is 227x227x3, being much higher than that of our current study in the proposed CNN models (the input image size of 32x32x32). The large input image size enables the model to capture more information and discover more features. Thus, it enables the model to be more successful.

Table 3. Comparison of the accuracy value obtained by the proposed method with other methods

Method	Accuracy(%)
GooLeNet and AlexNet [47]	93
Proposed method (mean)	82.45
D-CNN(Activation functions=LeakyReLU) [55]	81.9
CNN [21]	81.5
D-CNN(Activation functions= PReLU) [55]	81.4
D-CNN(Activation functions= ReLU) [55]	81.2
DNN [20]	79.23
SoundNet [58]	79
DCNN + augmentation SB-CNN (DA) [44]	79
D-CNN(Activation functions= ELU) [55]	78.9
EnvNet-v2 + augmentation[46]	78.3
Pyramid-Combined CNN[2]	78.1
EnvNet-v2 (Tokozume et al., 2017)[46]	78
Dilated CNN [52]	78
DCNN [59]	77.36
Unsupervised feature learning SKM (DA)[30]	76
Convolutional layers with max-pooling[43]	74
SKM[30]	74
Deep CNN[44]	74
D-CNN(Activation functions= Softplus) [55]	73.7
CNN (Baseline model) [43]	73.7
Unsupervised feature learning SKM [30]	73.6
M18 CNN (Dai et al., 2017)[60]	72
VGG (Pons & Serra, 2018)[61]	70
SVM [25]	71
Very Deep CNN[60]	69.38
Baseline system[25]	68
SVM[62]	62.4
ANN, KNN + features cascading + optimization[63]	56.4

5 CONCLUSIONS

Sound is a part of daily life and associated with the environment that we live in. Therefore, it is regarded as important as images and videos, and its analysis further leads the discovery of insights of our environment. In this study, the most suitable CNN model was obtained through a grid search of its parameters for the classification of environmental sounds. For this purpose, 150 CNN different models were designed, trained and tested over the Urbansound8k environmental sounds that were represented in an image format of 32x32x3. Among these methods developed, the best performing CNN model, which we call “CnnSound”, has achieved 82.45% predictive accuracy. When compared with similar studies in the literature, it has been observed that the CnnSound model has a satisfactory performance. As there is still a room for improvement, further research will be geared towards further improvement through pre-processing methods of sounds (e.g., noise elimination), sound representation in an image format by means of different methods, optimization methods and further fine-tuning of CNN models. This is further expected to be studied along with other sound libraries to demonstrate robustness of the sound representation of the methods and deep learning-based frameworks being developed and adapted into sound modelling and classification.

REFERENCES

- [1] S. Chu, S. Narayanan, and C.-C. J. Kuo, "Environmental sound recognition with time–frequency audio features," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, pp. 1142-1158, 2009.
- [2] F. Demir, M. Turkoglu, M. Aslan, and A. Sengur, "A new pyramidal concatenated CNN approach for environmental sound classification," *Applied Acoustics*, vol. 170, p. 107520, 2020.
- [3] P. Aumond, C. Lavandier, C. Ribeiro, E. G. Boix, K. Kamboja, E. D'Hondt, *et al.*, "A study of the accuracy of mobile technology for measuring urban noise pollution in large scale participatory sensing campaigns," *Applied Acoustics*, vol. 117, pp. 219-226, 2017.
- [4] J. Cao, M. Cao, J. Wang, C. Yin, D. Wang, and P.-P. Vidal, "Urban noise recognition with convolutional neural network," *Multimedia Tools and Applications*, vol. 78, pp. 29021-29041, 2019.
- [5] R. Radhakrishnan, A. Divakaran, and A. Smaragdis, "Audio analysis for surveillance applications," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, 2005.*, 2005, pp. 158-161.
- [6] M. Crocco, M. Cristani, A. Trucco, and V. Murino, "Audio surveillance: A systematic review," *ACM Computing Surveys (CSUR)*, vol. 48, pp. 1-46, 2016.
- [7] P. Laffitte, Y. Wang, D. Sodoyer, and L. Girin, "Assessing the performances of different neural network architectures for the detection of screams and shouts in public transportation," *Expert systems with applications*, vol. 117, pp. 29-41, 2019.
- [8] T. Heittola, A. Mesaros, A. Eronen, and T. Virtanen, "Audio context recognition using audio event histograms," in *2010 18th European Signal Processing Conference*, 2010, pp. 1272-1276.
- [9] M. Xu, C. Xu, L. Duan, J. S. Jin, and S. Luo, "Audio keywords generation for sports video analysis," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 4, pp. 1-23, 2008.
- [10] A. Waibel, H. Steusloff, and R. Stiefelhagen, "CHIL-Computers in the human interaction loop. 5th Intern," in *Workshop on Image Analysis for Multimedia Interactive Services*, 2004.
- [11] D. P. Ellis and K. Lee, "Minimal-impact audio-based personal archives," in *Proceedings of the the 1st ACM workshop on Continuous archival and retrieval of personal experiences*, 2004, pp. 39-47.
- [12] A. J. Eronen, V. T. Peltonen, J. T. Tuomi, A. P. Klapuri, S. Fagerlund, T. Sorsa, *et al.*, "Audio-based context recognition," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, pp. 321-329, 2005.
- [13] D. Barchiesi, D. Giannoulis, D. Stowell, and M. D. Plumbley, "Acoustic scene classification: Classifying environments from the sounds they produce," *IEEE Signal Processing Magazine*, vol. 32, pp. 16-34, 2015.
- [14] H. Li, S. Ishikawa, Q. Zhao, M. Ebana, H. Yamamoto, and J. Huang, "Robot navigation and sound based position identification," in *2007 IEEE International Conference on Systems, Man and Cybernetics*, 2007, pp. 2449-2454.
- [15] R. F. Lyon, "Machine hearing: An emerging field [exploratory dsp]," *IEEE signal processing magazine*, vol. 27, pp. 131-139, 2010.
- [16] S. Chu, S. Narayanan, C.-C. J. Kuo, and M. J. Mataric, "Where am I? Scene recognition for mobile robots using audio features," in *2006 IEEE International conference on multimedia and expo*, 2006, pp. 885-888.
- [17] J. Huang, "Spatial auditory processing for a hearing robot," in *Proceedings. IEEE International Conference on Multimedia and Expo*, 2002, pp. 253-256.
- [18] M. Green and D. Murphy, "Environmental sound monitoring using machine learning on mobile devices," *Applied Acoustics*, vol. 159, p. 107041, 2020.
- [19] P. Intani and T. Orachon, "Crime warning system using image and sound processing," in *2013 13th International Conference on Control, Automation and Systems (ICCAS 2013)*, 2013, pp. 1751-1753.
- [20] A. J. Torija, D. P. Ruiz, and Á. F. Ramos-Ridao, "A tool for urban soundscape evaluation applying support vector machines for developing a soundscape classification model," *Science of the Total Environment*, vol. 482, pp. 440-451, 2014.
- [21] V. P. Romero, L. Maffei, G. Brambilla, and G. Ciaburro, "Modelling the soundscape quality of urban waterfronts by artificial neural networks," *Applied Acoustics*, vol. 111, pp. 121-128, 2016.
- [22] A. Agha, R. Ranjan, and W.-S. Gan, "Noisy vehicle surveillance camera: A system to deter noisy vehicle in smart city," *Applied Acoustics*, vol. 117, pp. 236-245, 2017.
- [23] S. Ntalampiras, "Universal background modeling for acoustic surveillance of urban traffic," *Digital Signal Processing*, vol. 31, pp. 69-78, 2014.
- [24] K. J. Piczak, "ESC: Dataset for environmental sound classification," in *Proceedings of the 23rd ACM international conference on Multimedia*, 2015, pp. 1015-1018.
- [25] J. Salamon, C. Jacoby, and J. P. Bello, "A dataset and taxonomy for urban sound research," in *Proceedings of the 22nd ACM international conference on Multimedia*, 2014, pp. 1041-1044.

- [26] V. Bisot, R. Serizel, S. Essid, and G. Richard, "Feature learning with matrix factorization applied to acoustic scene classification," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, pp. 1216-1229, 2017.
- [27] D. Stowell, D. Giannoulis, E. Benetos, M. Lagrange, and M. D. Plumbley, "Detection and classification of acoustic scenes and events," *IEEE Transactions on Multimedia*, vol. 17, pp. 1733-1746, 2015.
- [28] P. Dhanalakshmi, S. Palanivel, and V. Ramalingam, "Classification of audio signals using AANN and GMM," *Applied soft computing*, vol. 11, pp. 716-723, 2011.
- [29] J. Ludena-Choez and A. Gallardo-Antolin, "Acoustic Event Classification using spectral band selection and Non-Negative Matrix Factorization-based features," *Expert Systems with Applications*, vol. 46, pp. 77-86, 2016.
- [30] J. Salamon and J. P. Bello, "Unsupervised feature learning for urban sound classification," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 171-175.
- [31] J. T. Geiger and K. Helwani, "Improving event detection for audio surveillance using gabor filterbank features," in *2015 23rd European Signal Processing Conference (EUSIPCO)*, 2015, pp. 714-718.
- [32] M. Mulimani and S. G. Koolagudi, "Segmentation and characterization of acoustic event spectrograms using singular value decomposition," *Expert Systems with Applications*, vol. 120, pp. 413-425, 2019.
- [33] J. Xie and M. Zhu, "Investigation of acoustic and visual features for acoustic scene classification," *Expert Systems with Applications*, vol. 126, pp. 20-29, 2019.
- [34] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097-1105.
- [35] J. Deng, A. Berg, S. Satheesh, H. Su, A. Khosla, and L. Fei-Fei, "Imagenet large scale visual recognition competition 2012 (ILSVRC2012)," *See net. org/challenges/LSVRC*, p. 41, 2012.
- [36] M. P. Pour and H. Seker, "Transform domain representation-driven convolutional neural networks for skin lesion segmentation," *Expert Systems with Applications*, vol. 144, p. 113129, 2020.
- [37] M. Buyukyilmaz, A. O. Cibikdiken, M. A. Abdalla, and H. Seker, "Identification of chicken Eimeria species from microscopic images by using MLP deep learning algorithm," in *Proceedings of the International Conference on Video and Image Processing*, 2017, pp. 84-88.
- [38] Ö. İnik, A. Ceyhan, E. Balcioglu, and E. Ülker, "A new method for automatic counting of ovarian follicles on whole slide histological images based on convolutional neural network," *Computers in biology and medicine*, vol. 112, p. 103350, 2019.
- [39] İ. Özkan and E. Ülker, "Derin Öğrenme ve Görüntü Analizinde Kullanılan Derin Öğrenme Modelleri," *Gaziosmanpaşa Bilimsel Araştırma Dergisi*, vol. 6, pp. 85-104, 2017.
- [40] İ. Özkan and B. TURAN, "Classification of Animals with Different Deep Learning Models," *Journal of New Results in Science*, vol. 7, pp. 9-16, 2018.
- [41] Ö. İnik, K. Uyar, and E. Ülker, "Gender Classification with A Novel Convolutional Neural Network (CNN) Model and Comparison with other Machine Learning and Deep Learning CNN Models," *Journal Of Industrial Engineering Research*, vol. 4, pp. 57-63, 2018.
- [42] Ö. İnik, E. Balcioglu, A. Ceyhan, and E. Ülker, "Using Convolution Neural Network for Classification of Different Tissue Images in Histological Sections," *Annals of the Faculty of Engineering Hunedoara*, vol. 17, pp. 101-104, 2019.
- [43] K. J. Piczak, "Environmental sound classification with convolutional neural networks," in *2015 IEEE 25th International Workshop on Machine Learning for Signal Processing (MLSP)*, 2015, pp. 1-6.
- [44] J. Salamon and J. P. Bello, "Deep convolutional neural networks and data augmentation for environmental sound classification," *IEEE Signal Processing Letters*, vol. 24, pp. 279-283, 2017.
- [45] N. Takahashi, M. Gygli, B. Pfister, and L. Van Gool, "Deep convolutional neural networks and data augmentation for acoustic event detection," *arXiv preprint arXiv:1604.07160*, 2016.
- [46] Y. Tokozume, Y. Ushiku, and T. Harada, "Learning from between-class examples for deep sound recognition," *arXiv preprint arXiv:1711.10282*, 2017.
- [47] V. Boddapati, A. Petef, J. Rasmusson, and L. Lundberg, "Classifying environmental sounds using image recognition networks," *Procedia computer science*, vol. 112, pp. 2048-2056, 2017.
- [48] S. Li, Y. Yao, J. Hu, G. Liu, X. Yao, and J. Hu, "An ensemble stacked convolutional neural network model for environmental event sound recognition," *Applied Sciences*, vol. 8, p. 1152, 2018.
- [49] Y. Su, K. Zhang, J. Wang, and K. Madani, "Environment sound classification using a two-stream CNN based on decision-level fusion," *Sensors*, vol. 19, p. 1733, 2019.
- [50] Z. Mushtaq and S.-F. Su, "Environmental sound classification using a regularized deep convolutional neural network with data augmentation," *Applied Acoustics*, vol. 167, p. 107389, 2020.

- [51] Z. Mushtaq, S.-F. Su, and Q.-V. Tran, "Spectral images based environmental sound classification using CNN with meaningful data augmentation," *Applied Acoustics*, vol. 172, p. 107581, 2021.
- [52] Y. Chen, Q. Guo, X. Liang, J. Wang, and Y. Qian, "Environmental sound classification with dilated convolutions," *Applied Acoustics*, vol. 148, pp. 123-132, 2019.
- [53] S. Abdoli, P. Cardinal, and A. L. Koerich, "End-to-end environmental sound classification using a 1D convolutional neural network," *Expert Systems with Applications*, vol. 136, pp. 252-263, 2019.
- [54] F. Medhat, D. Chesmore, and J. Robinson, "Masked Conditional Neural Networks for sound classification," *Applied Soft Computing*, vol. 90, p. 106073, 2020.
- [55] X. Zhang, Y. Zou, and W. Shi, "Dilated convolution neural network with LeakyReLU for environmental sound classification," in *2017 22nd International Conference on Digital Signal Processing (DSP)*, 2017, pp. 1-5.
- [56] M. Lim, D. Lee, H. Park, Y. Kang, J. Oh, J.-S. Park, *et al.*, "Convolutional Neural Network based Audio Event Classification," *KSII Transactions on Internet & Information Systems*, vol. 12, 2018.
- [57] K. Chng, "Classify Urban Sound using Machine Learning & Deep Learning (https://github.com/KevinChngJY/classifyurbansound_matlab), GitHub. Retrieved September 28, 2020.," 2020.
- [58] Y. Aytar, C. Vondrick, and A. Torralba, "Soundnet: Learning sound representations from unlabeled video," in *Advances in neural information processing systems*, 2016, pp. 892-900.
- [59] J. Ye, T. Kobayashi, and M. Murakawa, "Urban sound event classification based on local and global features aggregation," *Applied Acoustics*, vol. 117, pp. 246-256, 2017.
- [60] W. Dai, C. Dai, S. Qu, J. Li, and S. Das, "Very deep convolutional neural networks for raw waveforms," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 421-425.
- [61] J. Pons and X. Serra, "Randomly weighted CNNs for (music) audio classification," in *ICASSP 2019-2019 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 2019, pp. 336-340.
- [62] A. Kumar and B. Raj, "Features and kernels for audio event recognition," *arXiv preprint arXiv:1607.05765*, 2016.
- [63] B. da Silva, A. W Happi, A. Braeken, and A. Touhafi, "Evaluation of Classical Machine Learning Techniques towards Urban Sound Recognition on Embedded Systems," *Applied Sciences*, vol. 9, p. 3885, 2019.