

IDMT-Traffic: An Open Benchmark Dataset for Acoustic Traffic Monitoring Research

Jakob Abeßer
Semantic Music Technologies
Fraunhofer IDMT
 Ilmenau, Germany
 jakob.abesser@idmt.fraunhofer.de

Saichand Gourishetti
Industrial Media Applications
Fraunhofer IDMT
 Ilmenau, Germany

András Káta
Industrial Media Applications
Fraunhofer IDMT
 Ilmenau, Germany

Tobias Clauß
Industrial Media Applications
Fraunhofer IDMT
 Ilmenau, Germany

Prachi Sharma
Semantic Music Technologies
Fraunhofer IDMT
 Ilmenau, Germany

Judith Liebetrau
Industrial Media Applications
Fraunhofer IDMT
 Ilmenau, Germany

Abstract—In many urban areas, traffic load and noise pollution are constantly increasing. Automated systems for traffic monitoring are promising countermeasures, which allow to systematically quantify and predict local traffic flow in order to support municipal traffic planning decisions. In this paper, we present a novel open benchmark dataset, containing 2.5 hours of stereo audio recordings of 4718 vehicle passing events captured with both high-quality sE8 and medium-quality MEMS microphones. This dataset is well suited to evaluate the use-case of deploying audio classification algorithms to embedded sensor devices with restricted microphone quality and hardware processing power. In addition, this paper provides a detailed review of recent acoustic traffic monitoring (ATM) algorithms as well as the results of two benchmark experiments on vehicle type classification and direction of movement estimation using four state-of-the-art convolutional neural network architectures.

I. INTRODUCTION

A world-wide rise in population and a steady urbanization trend causes people to move from rural areas to bigger cities. With more and more active vehicles, travelling times increase and so do noise and air pollution levels. Intelligent transportation systems (ITS) are effective countermeasures to reduce and optimize traffic flow by adapting to local traffic situations. In the past decade, several automatic methods for traffic monitoring were developed for application scenarios such as controlling traffic light cycles, traffic accident detection, logistics monitoring, and other smart city application.

Traffic monitoring systems use various sensor modalities to measure traffic flow ranging from camera sensors for visual object detection and tracking, magnetic loop sensors for counting passing vehicles, to measurement systems based on radio waves (Radar) and light waves (Lidar). While such systems can be installed as distributed sensor networks to cover large areas, installation and maintenance costs are often high. Acoustic traffic monitoring (ATM) provides a cheaper alternative for non-intrusive traffic measurements and is the sole focus of this paper.

This paper has three main contributions. First, we present a compact state-of-the-art review of recent ATM systems.

As a second contribution, we introduce the IDMT-Traffic dataset, a novel dataset for traffic monitoring that includes around 2.5 hours of multi-microphone audio recordings with 4718 annotated passing vehicles. The dataset is intended as public benchmark to further stimulate research on traffic monitoring. Finally, we present the results of two benchmark experiments for vehicle type classification and direction of movement estimation using four different convolutional neural networks (CNNs) architectures.

This paper is structured as follows: We first review recent ATM algorithms in Section II before Section III describes the IDMT-Traffic dataset in details. Then, Section IV discusses the experimental procedure and the results of the two benchmark experiments. Finally, Section V concludes this work.

II. RELATED WORK

The audible sound on a road is a combination of several sound sources such as the engines, the exhausts, the wheels and air turbulence, which occur when vehicles pass by [1]. As a reasonable way to break down this complex audio analysis task, researches approached traffic monitoring from different perspectives. In this section, we categorize existing ATM algorithms based on the approaches for audio data acquisition and statistical modeling.

Moving sound sources such as vehicles can be detected based on their emitted sound if they are recorded with at least two microphones. Therefore, most ATM methods analyze either *stereo audio signals* [1], [2] [3] [4] [5] [6] or *multi-channel audio recordings* [7] [8] [9] [10], which are recorded with microphone arrays [7]–[9]. Microphones are commonly integrated into sensor units, which are either placed at the roadside [3] [4] [8] [2] or mounted on light poles at a height between 0.5 and 3 meters [1] [7] [6].

Traffic density is commonly measured on a two-stage scale (congested/non-congested) [5] or on a three-stage scale as either low/free (equivalent to vehicle speeds larger equal to

40 km/h), medium (20-40 km/h), or heavy/jammed (below 20 km/h) [1]–[4]. Traffic density can also be measured by *detecting and counting passing vehicles*. A common approach is to investigate run-time differences between stereo audio signals [7], [9], [10]. Moving sources exhibit a sweep-like peak contour in the temporal development of the cross-correlation function between both signals. While the contour’s diagonal alignment indicates the direction of movement, the contour angle indicates the speed of a vehicle. Ishida et al. match pre-defined templates with the cross-correlation function to detect for left-right and right-left movements [10]. In contrast, our experiments shown in Section IV indicate that a convolutional neural network can learn such characteristic patterns by itself.

Heavy traffic can lead to *traffic accidents*, which can be detected by the two sounds tire skidding and car crash [11], [12]. Another approach for traffic monitoring is to distinguish between vehicles in good and bad *mechanical condition* based on emitted sounds [2].

Different audio signal representations are used for solving ATM tasks. While most often raw spectrogram representations are used, some authors compute more advanced audio features such as the Mel-frequency Cepstral Coefficients (MFCC) prior to the modeling and classification steps [2], [3]. In order to train more robust algorithms, Gatto et al. apply data augmentation and mix recorded audio signals with additional noise [5].

ATM algorithms apply various mostly traditional classification algorithms such as Nearest Neighbor classifier [3], Bayer’s classifier [4], Random Forest classifier [5], Support Vector Machines (SVM) [2], Artificial Neural Networks (ANN) [2], [3], as well as hybrid approaches such as the Neuro-Fuzzy Classifier [1]. While most above-mentioned tasks require classification algorithms, Djukanović et al. use Support Vector Regression (SVR) [6] to predict the vehicle-to-microphone distances.

Most publications for ATM systems rely on proprietary datasets. However, some publicly available datasets are applicable for traffic monitoring. The MIVIA road audio events data set includes audio recordings of 400 sound events of the two classes tire skidding and car crashes [11]. The MAVD dataset [13] was published for sound event detection (SED) of particular traffic sounds which were recorded from the vehicle classes car, truck, bus, motorcycle in different states such as idling, accelerating, or braking. In the research field of SED, many datasets such as the FSK50k [14] or the AudioSet [15] include general sound classes such as car, truck, or train, whose recognition could be applied in ATM systems. Similarly, acoustic scene classification (ASC) datasets such as the TUT Urban Acoustic Scenes 2018 dataset [16] allow to train algorithms to detect amongst others traffic-related sound scenes such as “Street, traffic”, “Bus”, “Metro”, and “Tram”.

III. IDMT-TRAFFIC DATASET

In this section, we introduce a novel dataset for acoustic traffic monitoring (IDMT-Traffic)¹. It is intended as a public evaluation benchmark for the detection and classification of passing vehicles on inner-city and overland roads. The dataset includes time-synchronized stereo audio recordings made with both high-quality sE8 microphones² as well as lower-budget microelectro-mechanical systems (MEMS) microphones³ from four different recording locations including three city traffic locations and one country road location in and around Ilmenau, Germany. The recording scenarios include different speed limits (30, 50, and 70 km/h) as well as wet and dry road conditions.

Figure 1a and Figure 1b illustrate the recording setup, which was placed with a distance of 0.5 meters to the adjacent street. Both pairs of sE8 and MEMS microphones are fixed at a distance of 18.5 cm. Video recordings were solely made for the purpose of annotating the type and movement direction of passing vehicles. For reasons of data protection, we strictly avoided filming faces and licence plates by aiming the camera at the lower part of the vehicles as shown in Figure 1c and Figure 1d. For each of both microphone types, around 2.5 hours of audio recordings exist with a total of 4718 annotated passing vehicles. The dataset includes four classes: cars (3903 events), trucks (511 events), busses (53 events), and motorcycles (251 events). This distribution reflects the natural imbalance of vehicle types in common traffic scenarios.

IV. BENCHMARK EXPERIMENTS

Using the IDMT-Traffic dataset introduced in Section III, we conducted two benchmark experiments for different ATM tasks. Here, we only used the audio recordings recorded with the high-quality sE8 microphones.

A. Audio Representation & Pre-processing

As a first processing step, the left and right channels of the stereo audio files were averaged to a mono channel at the original sample rate of 48 kHz. In this work, we address the tasks of detecting vehicles and classifying their type and direction of movement. Therefore, we extract two types of features to be processed by the convolutional neural networks introduced in Section IV-B.

As first feature type, we extract mel-spectrograms using the *librosa* library [17]. We average the left and right audio channels and down-sample the audio signal to a sample rate of 22.05 kHz. Mel-spectrograms are computed using an FFT size, a window size, and a hop size of 2048, 1024, and 512 samples, respectively. In our experiments, we investigate the effect of changing the number of mel-bands $N_B \in \{16, 32, 64, 128\}$ on the recognition performance. Log-magnitude scaling is applied in order to compensate for the natural dynamic range of

¹The dataset will be made publicly available soon at <https://www.idmt.fraunhofer.de/en/publications/datasets.html>

²<https://www.seeelectronics.com/se8-mic>

³InvenSense ICS-43434

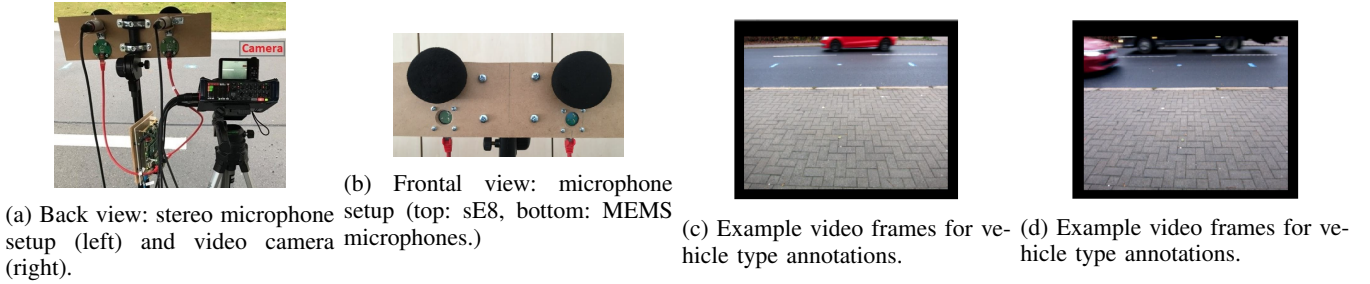


Fig. 1: Audio-video recording setup for the dataset creation.

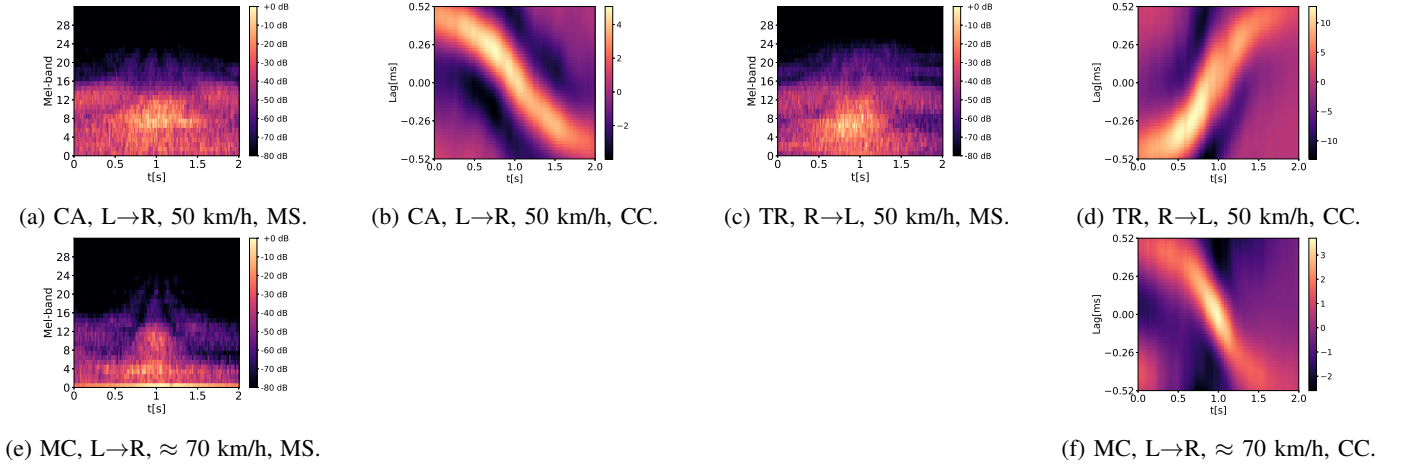


Fig. 2: Examples of two-second long patches taken from the IDMT-Traffic dataset for the vehicle type classes car (CA), truck (TR), and motorcycle (MC) for both mel-spectrogram (MS) and cross-correlation (CC) features. Direction of movement is either left-to-right (L→R) or right-to-left (R→L).

the traffic recordings. Finally, two-second long sub-sequences (patches) are extracted.

As a second feature type, we compute the local cross-correlation between the left and right audio channel at the original sample rate of 48 kHz. We extract blocks of 200 ms length with a hopsize of 25 ms from the audio signals. From the cross-correlation function between the left and right channel of the b -th block, we keep the center part $c_b \in \mathbb{R}^{51}$ with a margin of 25 lags around zero-lag index. We derive a two-dimensional feature representation by stacking such cross-correlation blocks for two-second long patches as for the mel-spectrograms.

Features are standardized (zero mean and unit variance) per bin, i. e. per frequency bin for the mel-spectrogram patches and per time lag for the cross-correlation patches, over all patches of a given dataset. This normalization procedure is performed independently for the training set, validation set, and test set. Figure 2 shows both the mel-spectrogram as well as the cross-correlation features for three examples, including passings of a car, a truck, and a motorcycle.

B. Neural Network Architectures

In our benchmark experiments, we test three different convolutional neural network architectures, which will be detailed

in the following sections. Table I summarizes the number of parameters per model.

1) *VGGNet*: The VGGNet model proposed by Takahashi et al. in [18] uses four pairs of 3x3 convolutional layers with intermediate pooling only between the layer pairs. This way, the spatial resolution is decreased while the number of filters is increased from 32 to 256. Several regularization strategies such as batch normalization, dropout as well as L2 regularization in the penultimate dense layers are applied to improve the model’s generalization towards new data.

2) *ResNet*: The ResNet is the “RN1” as proposed by Koutini et al. in [19]. It includes five residual blocks with two convolutional layers each. The network was designed to have a reduced receptive field and has been evaluated the task of acoustic scene classification.

3) *SqueezeNet*: The SqueezeNet architecture was introduced in [20] and implements several model compression strategies. As a first strategy, 3x3 filters are replaced by 1x1 filters in the convolutional layers. As a second strategy, the network includes a number fire modules, which use a squeeze-and-expand approach to reduce the depth of feature maps while maintaining their size.

4) *MobileNetMini*: The MobileNetMini model is a miniaturized version of the MobileNet architecture proposed

Model	# Parameters	Classification Task
VGGNet	1,442,788	Vehicle type
ResNet	3,259,012	Vehicle type
SqueezeNet	1,171,652	Vehicle type
MobileNetMini	15,363	Direction of movement

TABLE I: Summary of the compared neural network architectures, their number of parameters, as well as the classification tasks, they have been evaluated for.

Dataset	Car	Truck	Motorcycle	No vehicle
Training Set	2471	290	132	2393
Validation Set	275	32	15	266
Test Set	1157	189	99	1412

TABLE II: Number of patches per class in the training set, validation set, and test set for vehicle type classification.

in [21]. It includes one convolutional layer and one depth-wise convolutional layer with batch normalization and ReLU activation functions each followed by a global max pooling operations and a final softmax dense layer.

C. Experimental Procedure

From the IDMT-Traffic dataset, we first select audio files recorded at two locations with speed limits at 30 and 50 km/h. Both sets are combined, then shuffled and split into training set (90 %) and validation set (10 %). Recordings from the third location (speed limit at 70 km/h) was used as test set. Using this data partition, we aim to test the robustness of the ATM algorithms against different vehicle speeds and the corresponding changes in the vehicle sound characteristics. We trained all neural networks using the Adam optimizer [22] for 250 epochs with a learning rate of 10^{-5} . Early stopping with a patience of 50 epochs is used on the validation loss to monitor the training process.

D. Experiment 1 - Vehicle Type Classification

We consider a four-class classification scenario where we include the three vehicle types cars, trucks, and motorcycles as well a no-vehicle class, which includes spectrogram patches without any passing vehicles. The patches for the first three classes are centered around the annotated passing times. Non-vehicle patches were randomly sampled in between annotated vehicle passings in the audio recordings of the IDMT-Traffic dataset. The number of patches per class as well as their partition to training set, validation set, and test set is given in Table II. It can be observed that the classes no-vehicle and car have most patches followed by truck and motorcycle.

Table III summarizes the class-wise f-scores for the three investigated neural network architectures VGGNet, ResNet, and SqueezeNet. We observe that all models perfectly recognize the no-vehicle patches therefore allow for a robust vehicle detection (binary classification task) based on the high-quality sE8 audio recordings. Concerning the model perfor-

Model	Car	Truck	Motorcycle	No vehicle
$N_B = 16$				
VGGNet	0.94	0.5	0.96	1.0
ResNet	0.94	0.49	0.96	1.0
SqueezeNet	0.92	0.48	0.9	1.0
$N_B = 32$				
VGGNet	0.94	0.46	0.96	1.0
ResNet	0.94	0.44	0.97	1.0
SqueezeNet	0.94	0.42	0.9	1.0
$N_B = 64$				
VGGNet	0.94	0.49	0.97	1.0
ResNet	0.94	0.49	0.97	1.0
SqueezeNet	0.94	0.5	0.95	1.0
$N_B = 128$				
VGGNet	0.94	0.44	0.96	1.0
ResNet	0.94	0.45	0.95	1.0
SqueezeNet	0.91	0.53	0.97	1.0

TABLE III: Class-wise f-score results for vehicle type classification using the three neural network models for mel-spectrograms with 16 bins.

	Car	Truck	Motorcycle	No vehicle
Car	97.29	2.62	0.02	0.09
Truck	60.21	38.84	0.63	0.32
Motorcycle	3.23	1.21	95.35	0.2
No vehicle	0.23	0.01	0.11	99.65

TABLE IV: Confusion matrix for vehicle type classification using VGGNet with $N_B = 16$ (all values in percent).

mance, VGGNet and ResNet perform comparably well and slightly outperform the SqueezeNet model. Interestingly, the results show that a frequency resolution of only 16 mel-bands ($N_B = 16$) is sufficient to classify vehicle types.

Table IV illustrates as an example the confusion matrix for the VGGNet with $N_B = 16$. It becomes apparent that the truck-to-car confusion is the most prominent misclassification. We assume that since both vehicles have only small differences in their geometric size, they cause a similar acoustic footprint which complicates their distinction.

E. Experiment 2 - Direction of Movement Estimation

In this experiment, we evaluate the performance of the MobileNetMini architecture on the cross-correlation features for detecting the direction of movement. Here, we include patches across different vehicle types for the classes left-to-right and right-to-left and add no-vehicle patches as third class to simulate the detection task. The number of patches per class as well as their distribution among training, validation, and test sets is given in the Table V. As can be seen in the confusion matrix in Table VI, the direction of movement can be easily determined using a very small MobileNet architecture and the cross-correlation features. This result confirms the findings from the scientific literature [7], [9], [10]. As

Dataset	L→R	R→L	No vehicle
Training Set	1445	1448	2393
Validation Set	161	161	266
Test Set	678	767	1412

TABLE V: Number of patches per class in the training set, validation set, and test set for direction of movement estimation.

	L→R	R→L	No vehicle
L→R	96.61	0.83	2.57
R→L	0.26	98.64	1.1
No vehicle	0	0.21	99.79

TABLE VI: Confusion matrix for direction of movement estimation using the MobileNetMini (all values in percent).

the only distinction, our method relies on automatic feature learning as part of the CNN model.

V. CONCLUSIONS

In this paper we show that acoustic traffic monitoring provides a low-cost and non-invasive alternative to traffic monitoring approaches based on other sensor modalities such as vision or radar. After providing a thorough review of scientific publications on acoustic traffic monitoring, we present the novel IDMT-Traffic dataset, which is a freely-accessible benchmark dataset intended to stimulate further research in acoustic traffic monitoring.

In our baseline experiments, which use solely the high-quality audio recordings in the dataset, we show that state-of-the-art convolutional neural networks already achieve high performance scores for vehicle type classification and direction of movement estimation. Furthermore, the results show that vehicle detection can be implemented easily with either the mel-spectrogram or the cross-correlation features.

Having the goal of a real-world deployment of an ATM system in mind, several challenges need to be addressed. The first challenge arises from the microphone mismatched between high-quality and low-quality microphones used in mobile sensor devices. By including audio recordings from both high-quality sE8 microphones as well as lower-quality MEMS microphones, the IDMT-Traffic dataset provides a suitable test-bed to develop new algorithmic strategies for domain adaptation. A second challenge comes from computational performance constraints of mobile sensor devices, which might require to compress the neural network models. In addition to these challenges, future research directions include a precise speed estimation of vehicles as well as an improved classification of passing trucks.

ACKNOWLEDGEMENTS

This research was supported by the Fraunhofer Innovation Program “Trusted Resource Aware ICT” (TRAICT). The authors would like to thank Ke Shen, Yishi Sun, Jiaming Tan,

Julia Aileen, Katharina Duong, and Thanh Phong Duong for their assistance in the dataset recording process.

REFERENCES

- [1] P. Borkar and L. G. Malik, “Acoustic Signal based Traffic Density State Estimation using Adaptive Neuro-Fuzzy classifier,” in *Proceedings of the IEEE International Conference on Fuzzy Systems (FUZZ-IEEE)*, Hyderabad, India, 2013, pp. 1–8.
- [2] M. Bhandarkar and T. Waykole, “Vehicular Mechanical Condition Determination and On Road Traffic Density Estimation Using Audio Signals,” in *Proceedings of the 6th International Conference on Computational Intelligence and Communication Networks*, Bhopal, India, 2014, pp. 395–401.
- [3] V. P. Warghade and M. S. Deshpande, “Road Traffic Condition Estimation Based on Road Acoustics,” in *Proceedings of the International Conference on Computing, Communication, Control and Automation (ICCUBEA)*, Pune, India, 2017, pp. 1–5.
- [4] V. Tyagi, S. Kalyanaraman, and R. Krishnapuram, “Vehicular Traffic Density State Estimation Based on Cumulative Road Acoustics,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 13, no. 3, pp. 1156–1166, 2012.
- [5] R. C. Gatto and C. H. Q. Forster, “Audio-Based Machine Learning Model for Traffic Congestion Detection,” *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–8, 2020.
- [6] S. Djukanovic, J. Matas, and T. Virtanen, “Robust Audio-Based Vehicle Counting in Low-to-Moderate Traffic Flow,” in *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, Las Vegas, NV, USA, 2020, pp. 1608–1614.
- [7] Shipping Chen, Ziping Sun, and B. Bridge, “Traffic Monitoring Using Digital Sound Field Mapping,” *IEEE Transactions on Vehicular Technology*, vol. 50, no. 6, pp. 1582–1589, 2001.
- [8] Y. Na, Y. Guo, Q. Fu, and Y. Yan, “An Acoustic Traffic Monitoring System: Design and Implementation,” in *Proceedings of the IEEE UIC-ATC-ScalCom-CBDCCom-IoP*, Guangzhou, China, 2015, pp. 119–126.
- [9] B. Barbagli, G. Manes, and R. Facchini, “Acoustic Sensor Network for Vehicle Traffic Monitoring,” in *Proceedings of the International Conference on Advances in Vehicular Systems, Technologies and Applications (VEHICULAR)*, Venice, Italy, 2012, pp. 1–6.
- [10] S. Ishida, S. Liu, K. Mimura, S. Tagashira, and A. Fukuda, “Design of acoustic vehicle count system using DTW,” in *Proceedings of the ITS World Congress*, Melbourne, Australia, 2016, pp. 1–10.
- [11] P. Foggia, N. Petkov, A. Saggese, N. Strisciuglio, and M. Vento, “Audio Surveillance of Roads: A System for Detecting Anomalous Sounds,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 17, no. 1, pp. 279–288, 2016.
- [12] Y. Li, X. Li, Y. Zhang, M. Liu, and W. Wang, “Anomalous Sound Detection Using Deep Audio Representation and a BLSTM Network for Audio Surveillance of Roads,” *IEEE Access*, vol. 6, pp. 58 043–58 055, 2018.
- [13] P. Zinemanas, P. Cancela, and M. Rocamora, “MAVD: A Dataset for Sound Event Detection in Urban Environments,” in *Proceedings of the Workshop on Detection and Classification of Acoustic Scenes and Events (DCASE)*, New York, NY, USA, 2019, pp. 263–267.
- [14] E. Fonseca, S. Member, X. Favory, J. Pons, F. Font, and X. Serra, “FSD50K: An Open Dataset of Human-labeled Sound Events,” *arXiv preprint arXiv:2010.00475*, vol. 14, no. 8, pp. 1–21, 2020.
- [15] J. F. Gemmeke, D. P. W. 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 *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, New Orleans, LA, USA, 2017, pp. 776–780.
- [16] A. Mesaros, T. Heittola, and Tuomas Virtanen, “A Multi-Device Dataset for Urban Acoustic Scene Classification,” in *Proceedings of the Detection and Classification of Acoustic Scenes and Events (DCASE)*, Surrey, UK, 2018, pp. 9–13.
- [17] B. McFee, C. Raffel, D. Liang, D. P. W. Ellis, M. McVicar, E. Battenberg, and O. Nieto, “librosa: Audio and Music Signal Analysis in Python,” in *Proceedings of the Scientific Computing with Python conference (Scipy)*, Austin, Texas, 2015, pp. 18–24.
- [18] G. Takahashi, T. Yamada, N. Ono, and S. Makino, “Performance Evaluation of Acoustic Scene Classification using DNN-GMM and Frame-Concatenated Acoustic Features,” in *Proceedings of the 9th Asia-Pacific Signal and Information Processing Association Annual Summit*

and Conference (APSIPA), Honolulu, Hawaii, USA, 2018, pp. 1739–1743.

- [19] K. Koutini, H. Eghbal-Zadeh, M. Dorfer, and G. Widmer, “The receptive field as a regularizer in deep convolutional neural networks for acoustic scene classification,” *European Signal Processing Conference*, vol. 2019-Septe, 2019.
- [20] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, “SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, Toulon, France, 2016, pp. 1–13. [Online]. Available: <http://arxiv.org/abs/1602.07360>
- [21] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. C. Chen, “MobileNetV2: Inverted Residuals and Linear Bottlenecks,” in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, Salt Lake City, UT, USA, 2018, pp. 4510–4520.
- [22] D. P. Kingma and J. L. Ba, “Adam: A method for stochastic optimization,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, San Diego, CA, USA, 2015, pp. 1–15.