

Article

# A Review of Approaches and Challenges for Acoustic Scene Classification

Jakob Abeßer <sup>1</sup> 

<sup>1</sup> Semantic Music Technologies, Fraunhofer IDMT, Ehrenbergstraße 31, 98693 Ilmenau, Germany; jakob.abesser@idmt.fraunhofer.de

Version February 18, 2020 submitted to Appl. Sci.

**Abstract:** The number of publications on acoustic scene classification (ASC) in environmental audio recordings constantly increased over the last years. This was mainly stimulated by the annual Detection and Classification of Acoustic Scenes and Events (DCASE) competition with its first edition in 2013. All competitions so far involved one or multiple ASC tasks. With a focus on deep learning-based ASC algorithms, this article summarizes and groups existing approaches for data preparation, i. e., feature representations, feature pre-processing, and data augmentation, and for data modelling, i. e., neural network architectures and learning paradigms. Finally, the paper discusses current algorithmic limitations and open challenges in order to preview possible future developments towards the real-life application of ASC systems.

**Keywords:** acoustic scene classification, machine listening, deep neural networks

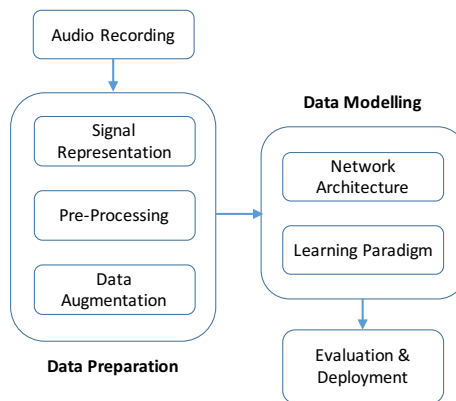
## 1. Introduction

Recognizing different indoor and outdoor acoustic environments from recorded acoustic signals is an active research field that has received a lot of attention in the last years. The task is an essential part of auditory scene analysis and involves summarizing an entire recorded acoustic signal using a pre-defined semantic description like “office room” or “public place”. Those semantic entities are denoted as *acoustic scenes*, and the task of recognizing them as acoustic scene classification (ASC) [1].

A particularly challenging task related to ASC is the detection of audio events which are temporarily present in an acoustic scene. Examples of such audio events include vehicles, car horns, and footsteps among others. This task is referred to as *acoustic event detection* (AED) and it substantially differs from ASC as it focuses on the precise temporal detection of particular sound events.

State-of-the-art ASC systems have been shown to outperform humans on this task [2]. Therefore, they are applied in numerous *application scenarios* such as context-aware wearables and hearables, hearing aids, health care, security surveillance, wild-life monitoring in nature habitats, smart cities, IoT, and autonomous navigation.

As current ASC methods are consistently based on deep neural networks [1], this article summarizes the recent methodological advances in the field of deep learning based ASC. Existing methods are summarized and categorized based on the typical processing steps illustrated in Figure 1. Section 2, Section 3, and Section 4 discuss techniques to represent, pre-process, and augment audio signals for ASC. Commonly used neural network architectures and learning paradigms are detailed in Section 5 and Section 6. Finally, Section 7 discusses open challenges and limitations of current ASC algorithms before Section 8 concludes this article. Methodologies and common datasets for evaluating



**Figure 1.** Flowchart summarizes the article structure and lists the typical processing flow of an ASC algorithm.

ASC algorithms are not further addressed in this article. The interesting reader is referred to [1,3] and the DCASE community website<sup>1</sup>.

## 2. Signal Representations

Datasets for the tasks of ASC or AED commonly contain digitized audio recordings. The resulting acoustic signals are commonly represented as waveforms that denote the amplitude of the recorded signal over discrete time samples. In most cases, ASC or AED systems perform the tasks of interest on derived signal representations, which will be introduced in the following section.

### 2.1. Monaural vs. Multi-Channel Signals

ASC algorithms commonly process monaural audio signals. Sound sources in acoustic scenes are spatially distributed by nature. If multi-channel audio recordings are available, the inherent spatial information can be exploited to better localize sound sources. The joint localization and detection of sound events has been first addressed in task 3 of the DCASE 2019 challenge.<sup>2</sup>

In addition to the left/right channels, a mid/side channel representation can be used as additional signal representation [4,5]. As an example for using a larger number of audio channels, Green and Murphey [6] classify acoustic scene recordings of 4th-order Ambisonics by combining spatial features describing the direction of sound arrival with band-wise spectral diffuseness measures. Similarly, Zirliński and Lee combine spatial features from binaural recordings with spectro-temporal features to characterize the foreground/background sound distribution in acoustic scenes [7].

### 2.2. Fixed Signal Transformations

Most neural network architectures applied for ASC require multi-dimensional input data (compare Section 5). The most commonly used time-frequency transformations are the short-time Fourier transform (STFT), the mel spectrogram, and the wavelet spectrogram. The mel spectrogram is based on a non-linear frequency scale motivated by human auditory perception and provides a more compact spectral representation of sounds compared to the STFT. ASC algorithms process only the magnitude of the Fourier transform while the phase is discarded.

Wavelets can be computed in a one-step [8,9] or cascaded fashion [10] to decompose time-domain signals into a set of basis function coefficients. The deep scattering spectrum [10] decomposes a signal using a sequential cascade of wavelet decompositions and modulation operations. The scalogram

<sup>1</sup> <http://dcase.community/>

<sup>2</sup> <http://dcase.community/challenge2019/task-sound-event-localization-and-detection>

[11,12] uses multiple parallel wavelet filters with logarithmically spaced support and bandwidth to provide invariance to both time-warping and local signal translations. Time-averaged statistics based on computer vision algorithms like Local Binary Patterns (LPB) or Histogram of Oriented Gradients (HOG) can be used to summarize such two-dimensional signal transformations [13].

In addition to basic time-frequency transformations, *perceptually-motivated signal representations* are used as input to deep neural networks. Such representations for instance characterize the distribution (e. g. Mel-Frequency Cepstral Coefficients (MFCC) [14], sub-band power distribution [15], and Gammatone Frequency Cepstral Coefficients [16]) and modulation of the spectral energy (e. g. amplitude modulation bank features [17] and temporal energy variation [18]). Feature learning techniques based on hand-crafted audio features and traditional classification algorithms such as Support Vector Machines (SVM) have been shown to underperform deep learning based ASC algorithms [19,20].

High-dimensional feature representations are often redundant and can lead to model overfitting. Therefore, before being processed by the neural network, the feature space dimensionality can be further reduced: One approach is to *aggregate subbands* of spectrograms using local binary pattern (LBP) histograms [21] or Subband Power Distribution (SBD) features [15]. A second approach is to map features to a randomized low-dimensional feature space as proposed by Jimenez et al. [22].

### 2.3. Learnable Signal Transformations

Three different approaches have been used to the best of our knowledge in ASC systems to avoid fixed pre-defined signal transformations. The first approach is to apply *end-to-end learning* where neural networks directly process raw audio samples. Examples of such network architectures are AclNet and AclSincNet [23] as well as SoundNet [24]. As a potential advantage against spectrogram-based methods, the signal phase is not discarded.

The second approach is to interpret the *signal transformation step as learnable function*, commonly denoted as “front-end”, which can be jointly trained with the classification back-end [25]. The third approach is to use *unsupervised learning* to derive semantically meaningful signal representations. Amiriparian et al. combine representations learnt using a deep convolutional generative adversarial network (DCGAN) and using a recurrent sequence to sequence autoencoder (S2SAE) [26]. Similarly, environmental audio recordings can be decomposed into suitable basis functions using well-established matrix factorization techniques such Non-Negative Matrix Factorization (NMF) [27] and Shift-Invariant Probabilistic Latent Component Analysis (SIPLCA) [28].

## 3. Pre-Processing

Feature standardization is commonly used to speed up the convergence of gradient descent based algorithms [5]. This process changes the feature distribution to have zero mean and unit variance. In order to compensate for the large dynamic range in environmental sound recordings, logarithmic scaling can be applied to spectrogram based features. Other low-level audio signal pre-processing methods include dereverberation and low-pass filtering [29].

Both ASC and AED face the challenge that foreground sound events in acoustic scenes are often overshadowed by background noises. Lostanlen et al. use Per-Channel Energy Normalization (PCEN) [30] to reduce stationary noise and to enhance transient sound events in environmental audio recordings [31]. This algorithm performs an adaptive, band-wise normalization and decorrelates the frequency bands. Wu et al. enhance edge-like structures in mel spectrograms using two edge detection methods from image processing based on Difference of Gaussians (DoG) and Sobel filtering [32]. The background drift of the mel spectrogram is removed using median filtering. Similarly, Han et al. use *background subtraction* and apply median filtering over time [4] to remove irrelevant noise components from the acoustic scene background and the recording device.

Several filtering approaches are used as pre-processing for ASC algorithms. For example, Nguyen et al. apply a Nearest Neighbor Filter based on the Repeating Pattern Extraction Technique

(REPET) algorithm [33] and replace the most similar spectrogram frames by their median prior to the classification [34]. This allows to emphasize *repetitive sound* events in acoustic scenes such as from sirens or horns. As another commonly used filtering approach, Harmonic-Percussive Source Separation (HPSS) decomposes the spectrogram into *horizontal and vertical components* and provides additional feature representations for ASC [4,29,35].

#### 4. Data Augmentation Techniques

Training deep learning models usually requires large amounts of training data to fully capture the natural variability in the data to be modeled. The size of machine listening datasets increased over the last years but lag behind computer vision datasets such as the ImageNet dataset with over 14 million images and over 21 thousand object classes [36]. The only exception to this day is the AudioSet dataset [37] with currently over 2.1 million audio excerpts and 527 sound event classes. This section summarizes techniques for *data augmentation* to address this lack of data.

The first group of data augmentation algorithms generates new training data instances from existing ones by applying various signal transformations. Basic audio signal transformation include *time stretching*, *pitch shifting*, *dynamic range compression*, as well as *adding random noise* [38–40]. Koutini et al. apply *spectral rolling* by randomly shifting spectrogram excerpts over time [41].

Several data augmentation methods allow to simulate overlap between multiple sound events and the resulting occlusion effects in the spectrogram. *Mixup* data augmentation allows to create new training instances by mixing pairs of features and their corresponding targets based on a given mixing ratio [42]. Another approach adapted from the computer vision field is *SpecAugment*, where features are temporally warped and blocks of the features are randomly masked [43]. Similarly, *random erasing* involves replacing random boxes in feature representations by random numbers [44]. In the related research task of bird audio detection, Lasseck combines several data augmentation techniques in the time domain (e. g. mosaicing random segments, time stretching, time interval dropout) and time-frequency domain (e. g. piece-wise time/frequency stretching and shifting) [45].

A second group of data augmentation synthesizes novel data instances from scratch. The most common synthesis approaches are based on Generative Adversarial Networks (GAN) [46], where class-conditioned synthesis models are trained using an adversarial training strategy by imitating existing data samples. While data synthesis is usually performed in the audio signal domain [12,47], Mun et al. instead synthesize intermediate embedding vectors [48]. Kong et al. generate acoustic scenes using the SampleRNN model architecture [49].

#### 5. Network Architectures

ASC algorithms mostly use CNN-based network architectures since they usually provide a summarizing classification of longer acoustic scene excerpts. In contrast, AED algorithms commonly use convolutional recurrent neural networks (CRNN) as they focus on a precise detection of sound events [50]. This architecture combines convolutional neural networks (CNN) as front-end for representation learning and a recurrent layer for temporal modeling. Hence, the main focus is on CNN-based ASC methods in Section 5.1. Other methods using feedforward neural networks (FNN) and CRNN are briefly discussed in Section 5.2 and Section 5.3, respectively. Network architectures and the corresponding hyper-parameters are usually optimized manually. As an exception, Roletscheck et al. automate this process and compare various architectures, which are automatically generated using a genetic algorithm [51].

##### 5.1. Convolutional Neural Networks

Traditional CNN architectures use multiple blocks of successive convolution and pooling operations for feature learning and down-sampling along the time and feature dimensions, respectively. As an alternative, Ren et al. use *atrous CNNs* which are based on dilated convolutional kernels [52]. Such kernels allow to achieve a comparable *receptive field* size without intermediate pooling operation.

Koutine et al. show that ASC systems can be improved if the receptive field is regularized by restricting its size [53].

In most CNN-based architectures, only the activations of the last convolutional layer are connected to the final classification layers. As an alternative, Yang et al. follow a *multi-scale feature* approach and further process the activations from all intermediate feature maps [54]. Additionally, the authors use the Xception network architecture, where the convolution operation is split into a depthwise (spatial) convolution and a pointwise (channel) convolution to reduce the number of trainable parameters. A related approach is to factorize two-dimensional convolutions into two one-dimensional kernels to separately model transient and long-term characteristics of sounds [16,55]. The influence of different symmetric and asymmetric kernel shapes are systematically evaluated by Wang et al. [56].

Several extensions to the common CNN architecture were proposed to improve the feature learning. Basbug and Sert adapted the *spatial pyramid pooling* strategy from computer vision, where feature maps are pooled and combined on different spatial resolutions [57]. In order to learn frequency-aware filters in the convolutional layers, Koutini et al. propose to encode the frequency position of each input feature bin within an additional channel dimension (*frequency-aware CNNs*) [41]. Similarly, Marchi et al. add the first and second order time derivative of spectrogram-based features as *additional input channels* in order to facilitate detecting transient short-term events which have a rapid increase in magnitude [58].

## 5.2. Feedforward Neural Networks

*Feedforward neural networks* (FNN) are used in several ASC algorithms. Bisot et al. use an FNN architecture to concatenate features from an NMF decomposition and a Constant-Q transform of the audio signal [59]. Takahashi et al. combine an FNN with multiple Gaussian Mixture Model (GMM) classifiers to model the individual acoustic scenes [60].

## 5.3. Convolutional Recurrent Neural Networks

The third category of ASC algorithms are based on convolutional recurrent neural networks (CRNN). Li et al. combine in two separate input branches CNN-based front-ends for feature learning with bidirectional gated recurrent units (BiGRU) for temporal feature modeling [10]. In contrast to a sequential ordering of convolutional and recurrent layers, parallel processing pipelines using long short-term memory (LSTM) layers were used in [47] and [61]. Two recurrent network types used in ASC systems require fewer parameter and less training data compared to LSTM layers—gated recurrent neural networks (GRNN) [8,9,62] and time-delay neural networks (TDNN) [17,63].

# 6. Learning Paradigms

Building up on the basic neural network architectures introduced in Section 5, approaches to further improve ASC systems are summarized in this section. After discussing methods for closed/open set classification in Section 6.1, extensions to neural networks such as multiple input networks (Section 6.2) and attention mechanisms (Section 6.3) are presented. Finally, both multitask learning (Section 6.4) and transfer learning (Section 6.5) will be discussed as two promising training strategies to improve ASC systems.

## 6.1. Closed/Open Set Classification

Most ASC algorithms assume a *closed-set classification* scenario with a fixed predefined set of acoustic scenes to classify. In real-world applications however, the underlying data distributions of acoustic scenes is often unknown and can furthermore change over time with new classes becoming relevant. This motivates the use of *open-set classification* approaches, where an algorithm can also classify a given audio recording as “unknown” class. This scenario was first addressed as part of the DCASE 2019 challenge in the task 1C “Open set Acoustic Scene Classification” [64].



Saki et al. propose the *Multi-Class Open-set Evolving Recognition* (MCOSR) algorithm to tackle open-set ASC [65]. Unknown samples are first rejected by a recognition model before the algorithm tries to identify underlying (hidden) classes in these samples in an unsupervised manner. Finally, the recognition model can be updated using the novel classes. Wilkinghoff and Kurth combine a closed-set classification algorithm and an outlier detection algorithm based on deep convolutional autoencoders (DCAE) to recognize unknown samples in an open-set ASC scenario [66]. Lehner et al. evaluated the model's classification confidence to identify unknown samples [67]. Therefore, a threshold is applied on the highest logit value at the input of the final neural network layer.

## 6.2. Multiple Input Networks

As discussed before, most ASC algorithms use a convolutional front-end to learn characteristic patterns in multi-dimensional feature representations. As a general difference to image processing, the time and frequency axes in spectrogram-based feature representations do not carry the same semantic meaning. In order to train networks to detect spectral patterns, which are characteristic for certain frequency regions, several authors split a spectrogram into two [68] or multiple [69] sub-bands and use networks with multiple input branches. Using the same idea of distributed feature learning, other ASC systems individually process the left/right or mid/side channels [5] or filtered signal variants, which are obtained using harmonic/percussive separation (HPSS) [35] or nearest neighbor filtering (NNF) [34]. Instead of feeding multiple signal representations to the network as individual input branches, Dang et al. propose to concatenate both MFCC and log mel spectrogram features along the frequency axis as input features [70].

## 6.3. Attention

The temporal segments of an environmental audio recording contribute differently to the classification of its acoustic scene. Neural *attention mechanisms* allow neural networks to focus on a specific subset of its input features. Attention mechanisms can be incorporated at different positions within neural network based ASC algorithms. Li et al. incorporate Gated Linear Units (GLU) in several steps of the feature learning part of the network ("multi-level attention") [10]. GLUs implement pairs of mutually gating convolutional layers to control the information flow in the network. Attention mechanisms can also be applied in the pooling of feature maps [71]. Wang et al. use self-determination CNNs (SD-CNNs) to identify frames with higher uncertainty due to overlapping sound events. A neural network can learn to focus on local patches within the receptive field if a *network-in-network* architecture is used [72]. Here, individual convolutional layers are extended by micro neural networks, which allows for more powerful approximations by additional non-linearities.

## 6.4. Multitask-Learning

*Multitask learning* involves learning to solve multiple related classification tasks jointly with one network [73]. By learning shared feature representations, the performance on the individual tasks can be improved and a better generalization can be achieved.

A natural approach is to train one model to perform ASC and AED in a joint manner [74] as acoustic events are the building blocks of acoustic scenes. Sound events and acoustic scenes naturally follow a hierarchical relationship. While most publications perform a "flat" classification, Xu et al. exploit a hierarchical acoustic scene taxonomy and group acoustic scenes to the three high-level scene classes "vehicle", "indoor", and "outdoor" [75]. The authors use a hierarchical pre-training approach, where the network learns to predict the high-level scene class as main task and the low-level scene class such as car or tram as auxiliary task. Using a similar scene grouping approach, Nwe et al. train a CNN with several shared convolutional layers and three three branches of task-specific convolutional layers to predict the most likely acoustic scene within each scene group [76].

## 6.5. Transfer Learning & Result Fusion

Many ASC algorithms rely on well proven neural network architectures from the computer vision domain such as AlexNet [71,77], VGG16 [9], Xception [54], DenseNet [41], GoogLeNet [77], and Resnet [35,67]. *Transfer learning* allows the finetuning of models, which were pretrained on related audio classification tasks. For instance, Huang et al. used the AudioSet dataset to pretrain four different neural network architectures and finetune them using a task-specific development set [23]. Similarly, Singh et al. take a pretrained SoundNet [78] network as basis for their experiments [24,79]. Ren et al. use the VGG16 model as seed model, which was pre-trained for object recognition in images [9]. Kumar et al. pre-train a CNN in a supervised fashion using weak label annotation of the AudioSet dataset. The authors compare three transfer learning strategies to adapt the model to novel AED and ASC target tasks [80].

Many ASC algorithms include *result fusion* steps where intermediate results from different time frames or classifiers are merged. Similar to computer visions, features learnt in different layers of the network capture different levels of abstraction of the audio signal. Therefore, some systems apply *early fusion* and combine intermediate feature representations from different layers of the network as *multiscale features* [24,54,79]. Ensemble learning is a common *late fusion* technique where the prediction results of multiple classifiers are combined [5,9,19,23,34]. The predicted class scores can be averaged [81] or used as features for an additional classifier [82].

## 7. Open Challenges

This section discusses several open challenges which arise from deploying ASC algorithms to real-world application scenarios.

### 7.1. Domain Adaption

The performance of sound event classification algorithms often suffer from *covariate shift*, i.e., a distribution mismatch between training and test datasets. When being deployed to real-world application scenarios, ASC systems often face novel acoustic conditions which are caused by different recording devices or environmental influences. *Domain adaption* methods aim to increase the robustness of classification algorithms in such scenarios by adapting them to data from a novel target domain [83]. Depending on whether labels exist for the target domain data, supervised and unsupervised methods are distinguished. Supervised domain adaptation usually involves fine-tuning a model on a new target domain data after it was pre-trained on the annotated source domain data.

One unsupervised domain adaptation strategy is to alter the target domain data such that its distribution becomes closer to that of the source domain data. As an example, Kosmider use “spectral correction” to compensate for different frequency responses of the recording equipment. He estimates a set of frequency-dependent magnitude coefficients from the source domain data and used them for spectrogram equalization of the target domain data [84]. Mun and Shon perform an independent domain adaptation of both the source and target domain to an additional domain using factorized hierarchical variational autoencoder [85].

As a second unsupervised strategy, Gharib et al. use an adversarial training approach such that the intermediate feature mappings of an ASC model follow a similar distribution for both the source and target domain data [83]. This approach is further improved using the Wasserstein generative adversarial networks (WGAN) formulation [86].

### 7.2. Ambiguous Allocation between Sound Events and Scenes

Acoustic scenes often comprise multiple *sound events*, which are not class-specific but instead appear in a similar way in various scene classes [2,18]. As an example, sound recordings which were recorded in different vehicle types such as car, tram, or train often exhibit *prominent speech* from human conversations or automatic voice announcements. At the same time, class-specific sound components

like engine noises, road surface sounds, or door opening and closing sounds appear at a lower level in the background. Wu and Lee use the gradient-weighted class activation mappings (GradCAM) to show that CNN-based ASC models in general have the capability to ignore high-energy sound events and focus on quieter background sounds instead [32].

### 7.3. Model Interpretability

Despite their superior performance, deep learning based ASC models are often considered as “black boxes” due to their high complexity and large number of parameters. One main challenge is to develop methods which allow to better interpret the model predictions and internal feature representations. As discussed in Section 6.3, attention mechanisms allow neural networks to focus on relevant subsets of the input data. Wang et al. investigate an attention-based ASC model and demonstrate that only fractions of long-term scene recordings are relevant for its classification [72]. Similarly, Ren et al. visualize internal attention matrices obtained for different acoustic scenes [52]. The results confirm that either stationary and short-term signal components are most relevant for particular acoustic scenes.

Another common strategy to investigate the class separability in intermediate feature representations are dimension reduction techniques such as t-SNE [24]. Techniques such as Layerwise Relevance Propagation (LRP) [87] allow to interpret neural networks by investigating the pixel-wise contributions of input features to classification decisions.

### 7.4. Real-World Deployment

Many challenges arise when ASC models are deployed in smart city [88,89] or industrial sound analysis [90] scenarios. The first challenge is the *model complexity*, which is limited if data privacy concerns require the classification to be performed directly on mobile sensor devices. *Real-time processing requirements* often demand for fast model prediction with low latency. In a related study, Sigtia et al. contrast the performance of different audio event detection methods with the respective computational costs [91]. In comparison to traditional methods such as Support Vector Machines (SVM) and Gaussian Mixture Models (GMM), fully-connected neural networks achieve the best performance while requiring the lowest number of operations.

This often requires a *model compression* step, where trained classification models are reduced in size and redundant components need to be identified. Several attempts were made to make networks more compact by decreasing the number of operations and increasing the memory-efficiency. For instance, the MobileNetV2 architecture is based on a novel layer module, which mainly uses convolutional operations to avoid large intermediate tensors [92]. A similar approach is followed by Drossos et al. for AED [93]. The authors replace the common CNN-based frontend by depthwise separable convolutions and the RNN-backend with dilated convolutions to reduce the number of model parameters and required training time. In the MorphNet approach proposed by Gordon et al., network layers are shrunk and expanded in an iterative procedure to optimize network architectures in order to match given resource constraints [94]. Finally, Tan and Lee show in the EfficientNet approach that a uniform scaling of the network dimensions depth, width, and resolution of a convolutional neural network leads to highly effective networks [95].

A second challenge arises from the *audio recording devices* in mobile sensor units. Due to space constraints, microelectro-mechanical systems (MEMS) microphones are often used. However, scientific datasets used for training ASC models are usually recorded with high quality electret microphones [96]. As discussed in Section 7.1, changed recording conditions affect the input data distribution. Achieving robust classification systems in such a scenario requires the application of domain adaptation strategies.

## 8. Conclusion

In the research field of acoustic scene classification a rapid increase of scientific publications has been observed in the last decade. This progress was mainly stimulated by recent advances in the



field of deep learning such as transfer learning, attention mechanisms, and multitask learning as well as the release of public datasets. The DCASE community plays a major role in this development by organising annual evaluation campaigns on various machine listening tasks.

State-of-the-art ASC algorithms have matured and can be applied in context-aware devices such as hearables and wearables. In such real-world application scenarios, novel challenges need to be faced such as microphone mismatch and domain adaptation, open set classification, as well as model complexity and real-time processing constraints. The general demand of deep learning based classification algorithms for larger training corpora can be faced with novel techniques from unsupervised and self-supervised learning as it was shown in natural language processing, speech processing, and image processing.

**Funding:** This work has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 786993 and was supported by the German Research Foundation (AB 675/2-1).

**Acknowledgments:** The author would like to thank Hanna Lukashevich, Stylianos Mimilakis, David S. Johnson, and Sascha Grollmisch for valuable discussions and proof-reading.

**Conflicts of Interest:** The author declares no conflict of interest.

## References

1. Virtanen, T.; Plumbley, M.D.; Ellis, D., Eds. *Computational Analysis of Sound Scenes and Events*; Springer International Publishing, 2018. doi:10.1007/978-3-319-63450-0.
2. Mesaros, A.; Heittola, T.; Virtanen, T. Assessment of Human and Machine Performance in Acoustic Scene Classification: DCASE 2016 Case Study. *Proceedings of the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. New Paltz, NY, USA, 2017, pp. 319–323. 15–18 October.
3. Mesaros, A.; Heittola, T.; Benetos, E.; Foster, P.; Lagrange, M.; Virtanen, T.; Plumbley, M.D. Detection and Classification of Acoustic Scenes and Events: Outcome of the DCASE 2016 Challenge. *IEEE/ACM Transactions on Audio Speech and Language Processing* **2018**, *26*, 379–393. doi:10.1109/TASLP.2017.2778423.
4. Han, Y.; Park, J.; Lee, K. Convolutional Neural Networks with Binaural Representations and Background Subtraction for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
5. Mars, R.; Pratik, P.; Nagisetty, S.; Lim, C. Acoustic Scene Classification from Binaural Signals using Convolutional Neural Networks. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019, pp. 149–153. 25–26 October, doi:10.33682/6c9z-gd15.
6. Green, M.C.; Murphy, D. Acoustic Scene Classification using Spatial Features. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
7. Zieliński, S.K.; Lee, H. Feature Extraction of Binaural Recordings for Acoustic Scene Classification. *Proceedings of the Federated Conference on Computer Science and Information Systems (FedCSIS)*. Poznań, Poland, 2018, pp. 585–588. 9–12 September, doi:10.15439/2018F182.
8. Qian, K.; Ren, Z.; Pandit, V.; Yang, Z.; Zhang, Z.; Schuller, B. Wavelets Revisited for the Classification of Acoustic Scenes. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
9. Ren, Z.; Pandit, V.; Qian, K.; Yang, Z.; Zhang, Z.; Schuller, B. Deep Sequential Image Features for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
10. Li, Z.; Hou, Y.; Xie, X.; Li, S.; Zhang, L.; Du, S.; Liu, W. Multi-Level Attention Model with Deep Scattering Spectrum for Acoustic Scene Classification. *Proceedings of the IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*. Shanghai, China, 2019, pp. 396–401. 8–12 July, doi:10.1109/ICMEW.2019.00074.
11. Chen, H.; Zhang, P.; Bai, H.; Yuan, Q.; Bao, X.; Yan, Y. Deep convolutional neural network with scalogram for audio scene modeling. *Proceedings of the Annual Conference of the International Speech*

- Communication Association (INTERSPEECH). Hyderabad, India, 2018, pp. 3304–3308. 2–6 September, doi:10.21437/Interspeech.2018-1524.
12. Chen, H.; Liu, Z.; Liu, Z.; Zhang, P.; Yan, Y. Integrating the Data Augmentation Scheme with Various Classifiers for Acoustic Scene Modeling. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019. 25–26 October.
  13. Ye, J.; Kobayashi, T.; Toyama, N.; Tsuda, H.; Murakawa, M. Acoustic scene classification using efficient summary statistics and multiple spectro-temporal descriptor fusion. *Applied Sciences* **2018**, *8*, 1–12. doi:10.3390/app8081363.
  14. Li, Y.; Li, X.; Zhang, Y.; Wang, W.; Liu, M.; Feng, X. Acoustic Scene Classification Using Deep Audio Feature and BLSTM Network. *Proceedings of the 6th International Conference on Audio, Language and Image Processing (ICALIP)*. Shanghai, China, 2018, pp. 371–374. 16–17 July, doi:10.1109/ICALIP.2018.8455765.
  15. Bisot, V.; Essid, S.; Richard, G. HOG and Subband Power Distribution Image Features for Acoustic Scene Classification. *Proceedings of the 23rd European Signal Processing Conference (EUSIPCO)*. Nice, France, 2015, pp. 719–723. 31 August - 4 September, doi:10.1109/EUSIPCO.2015.7362477.
  16. Sharma, J.; Granmo, O.C.; Goodwin, M. Environment Sound Classification using Multiple Feature Channels and Deep Convolutional Neural Networks. *ArXiv pre-prints* **2019**, *14*, 1–11, [1908.11219].
  17. Moritz, N.; Schröder, J.; Goetze, S.; Anemüller, J.; Kollmeier, B. Acoustic Scene Classification using Time-Delay Neural Networks and Amplitude Modulation Filter Bank Features. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Budapest, Hungary, 2016. 3 September.
  18. Park, S.; Mun, S.; Lee, Y.; Ko, H. Acoustic Scene Classification Based on Convolutional Neural Network using Double Image Features. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
  19. Fonseca, E.; Gong, R.; Bogdanov, D.; Slizovskaia, O.; Gomez, E.; Serra, X. Acoustic Scene Classification by Ensembling Gradient Boosting Machine and Convolutional Neural Networks. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
  20. Maka, T. Audio Feature Space Analysis for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Surrey, UK, 2018. 19–20 November.
  21. Abidin, S.; Togneri, R.; Sohel, F. Enhanced LBP Texture Features from Time Frequency Representations for Acoustic Scene Classification. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. New Orleans, LA, USA, 2017, pp. 626–630. 5–9 March.
  22. Jiménez, A.; Elizalde, B.; Raj, B. DCASE 2017 Task 1: Acoustic Scene Classification using Shift-Invariant Kernels and Random Features. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
  23. Huang, J.; Lu, H.; Lopez-Meyer, P.; Maruri, H.A.C.; Ontiveros, J.A.d.H. Acoustic Scene Classification using Deep Learning-Based Ensemble Averaging. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019, pp. 94–98. 25–26 October.
  24. Singh, A.; Rajan, P.; Bhavsar, A. Deep Multi-View Features from Raw Audio for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019, pp. 229–233. 25–26 October.
  25. Chen, H.; Zhang, P.; Yan, Y. An Audio Scene Classification Framework with Embedded Filters and a DCT-Based Temporal Module. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Brighton, UK, 2019, pp. 835–839. 12–17 May.
  26. Amiriparian, S.; Freitag, M.; Cummins, N.; Gerczuk, M.; Pugachevskiy, S.; Schuller, B. A Fusion of Deep Convolutional Generative Adversarial Networks and Sequence to Sequence Autoencoders for Acoustic Scene Classification. *Proceedings of the 26th European Signal Processing Conference (EUSIPCO)*. Rome, Italy, 2018, pp. 977–981. 3–7 September, doi:10.23919/EUSIPCO.2018.8553225.
  27. Bisot, V.; Serizel, R.; Essid, S.; Richard, G. Feature Learning with Matrix Factorization Applied to Acoustic Scene Classification. *IEEE/ACM Transactions on Audio Speech and Language Processing* **2017**, *25*, 1216–1229. doi:10.1109/TASLP.2017.2690570.

28. Benetos, E.; Lagrange, M.; Dixon, S. Characterisation of Acoustic Scenes using a Temporally-Constrained Shift-Invariant Model. *Proceedings of the 15th International Conference on Digital Audio Effects (DAFx-12)*. York, UK, 2012, pp. 1–7. 17–21 September.
29. Seo, H.; Park, J.; Park, Y. Acoustic Scene Classification using Various Pre-Processed Features and Convolutional Neural Networks. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019, pp. 3–6. 25–26 October.
30. Wang, Y.; Getreuer, P.; Hughes, T.; Lyon, R.F.; Saurous, R.A. Trainable Frontend for Robust and Far-Field Keyword Spotting. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. New Orleans, LA, USA, 2017, pp. 5670–5674, [1607.05666]. 5–9 March, doi:10.1109/ICASSP.2017.7953242.
31. Lostanlen, V.; Salamon, J.; Cartwright, M.; McFee, B.; Farnsworth, A.; Kelling, S.; Bello, J.P. Per-channel energy normalization: Why and how. *IEEE Signal Processing Letters* **2019**, *26*, 39–43. doi:10.1109/LSP.2018.2878620.
32. Wu, Y.; Lee, T. Enhancing Sound Texture in CNN-based Acoustic Scene Classification. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Brighton, UK, 2019, pp. 815–819, [1901.01502]. 12–17 May, doi:10.1109/ICASSP.2019.8683490.
33. Rafii, Z.; Pardo, B. Music/Voice Separation using the Similarity Matrix. *Proceedings of the 13th International Society for Music Information Retrieval Conference (ISMIR)*. Porto, Portugal, 2012, pp. 583–588. 8–12 October.
34. Nguyen, T.; Pernkopf, F. Acoustic Scene Classification using a Convolutional Neural Network Ensemble and Nearest Neighbor Filters. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Surrey, UK, 2018. 19–20 November.
35. Mariotti, O.; Cord, M.; Schwander, O. Exploring Deep Vision Models for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Surrey, UK, 2018. 19–20 November.
36. Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A.C.; Fei-Fei, L. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* **2015**, *115*, 211–252, [1409.0575]. doi:10.1007/s11263-015-0816-y.
37. Gemmeke, J.F.; Ellis, D.P.W.; Freedman, D.; Jansen, A.; Lawrence, W.; Moore, R.C.; Plakal, M.; Ritter, M. Audio Set: An Ontology and Human-Labeled Dataset for Audio Events. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. New Orleans, LA, USA, 2017, pp. 776–780. 5–9 March.
38. Abeßer, J.; Mimilakis, S.I.; Gräfe, R.; Lukashevich, H. Acoustic Scene Classification By Combining Autoencoder-Based Dimensionality Reduction and Convolutional Neural Networks. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
39. Salamon, J.; Bello, J.P. Deep Convolutional Neural Networks and Data Augmentation for Environmental Sound Classification. *IEEE Signal Processing Letters* **2017**, *24*, 279–283, [1608.04363]. doi:10.1109/LSP.2017.2657381.
40. Xu, J.X.; Lin, T.C.; Yu, T.C.; Tai, T.C.; Chang, P.C. Acoustic Scene Classification Using Reduced MobileNet Architecture. *Proceedings of the IEEE International Symposium on Multimedia (ISM)*. Taichung, Taiwan, 2018, pp. 267–270. 10–12 December, doi:10.1109/ISM.2018.00038.
41. Koutini, K.; Eghbal-zadeh, H.; Widmer, G. Receptive-Field-Regularized CNN Variants for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019, pp. 124–128. 25–26 October.
42. Zhang, H.; Cisse, M.; Dauphin, Y.N.; Lopez-Paz, D. mixup: Beyond Empirical Risk Minimization. *Proceedings of the International Conference on Learning Representations (ICLR)*. Vancouver, Canada, 2018, [1710.09412]. 30 April - 3 May.
43. Park, D.S.; Chan, W.; Zhang, Y.; Chiu, C.C.; Zoph, B.; Cubuk, E.D.; Le, Q.V. SpecAugment: A simple data augmentation method for automatic speech recognition. *Proceedings of the Annual Conference of the International Speech Communication Association (INTERSPEECH)*. Graz, Austria, 2019, Vol. 2019-Sept, pp. 2613–2617, [1904.08779]. 2–15 November, doi:10.21437/Interspeech.2019-2680.

44. Zhong, Z.; Zheng, L.; Kang, G.; Li, S.; Yang, Y. Random Erasing Data Augmentation. *ArXiv pre-prints* **2017**, [1708.04896].
45. Lasseck, M. Acoustic bird detection with deep convolutional neural networks. Proceedings of the Detection and Classification of Acoustic Scenes and Events 2018 Workshop (DCASE). Surrey, UK, 2018, pp. 143–147. 19–20 November.
46. Goodfellow, I.J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Yoshua Bengio. Generative Adversarial Nets. *Advances in Neural Information Processing Systems (NIPS)*, 2014, pp. 2672–2680.
47. Mun, S.; Shon, S.; Kim, W.; Han, D.K.; Ko, H. Deep Neural Network Based Learning and Transferring Mid-Level Audio Features for Acoustic Scene Classification. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. New Orleans, LA, USA, 2017, pp. 796–800. 5–9 March, doi:10.1097/IOP.0000000000000348.
48. Mun, S.; Park, S.; Han, D.K.; Ko, H. Generative Adversarial Networks based Acoustic Scene Training Set Augmentation and Selection using SVM Hyperplane. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.
49. Kong, Q.; Xu, Y.; Iqbal, T.; Cao, Y.; Wang, W.; Plumbley, M.D. Acoustic Scene Generation with Conditional SampleRNN. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Brighton, UK, 2019, pp. 925–929. 12–17 May.
50. Xia, X.; Togneri, R.; Sohel, F.; Zhao, Y.; Huang, D. A Survey: Neural Network-Based Deep Learning for Acoustic Event Detection. *Circuits, Systems, and Signal Processing* **2019**, p. 3433?3453. doi:10.1007/s00034-019-01094-1.
51. Roletscheck, C.; Watzka, T.; Seiderer, A.; Schiller, D.; André, E. Using an Evolutionary Approach To Explore Convolutional Neural Networks for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019. 25–26 October.
52. Ren, Z.; Kong, Q.; Han, J.; Plumbley, M.D.; Schuller, B.W. Attention-based Atrous Convolutional Neural Networks: Visualisation and Understanding Perspectives of Acoustic Scenes. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Brighton, UK, 2019, pp. 56–60. 12–17 May, doi:10.1109/ICASSP.2019.8683434.
53. Koutini, K.; Eghbal-zadeh, H.; Widmer, G.; Kepler, J. CP-JKU Submissions to DCASE'19: Acoustic Scene Classification and Audio Tagging with REceptive-Field-Regularized CNNs. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019, pp. 1–5. 25–26 October.
54. Yang, L.; Chen, X.; Tao, L. Acoustic Scene Classification using Multi-Scale Features. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Surrey, UK, 2018. 19–20 November.
55. Cho, J.; Yun, S.; Park, H.; Eum, J.; Hwang, K. Acoustic Scene Classification Based on a Large-Margin Factorized CNN. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York, NY, USA, 2019, pp. 45–49, [1910.06784]. 25–26 October, doi:10.33682/8xh4-jm46.
56. Wang, C.Y.; Wang, J.C.; Wu, Y.C.; Chang, P.C. Asymmetric Kernel Convolution Neural Networks for Acoustic Scenes Classification. *Proceedings of the IEEE International Symposium on Consumer Electronics (ISCE)*. Kuala Lumpur, Malaysia, 2017, pp. 11–12. 14–15 November.
57. Basbug, A.M.; Sert, M. Acoustic Scene Classification Using Spatial Pyramid Pooling with Convolutional Neural Networks. *Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC)*. Newport, CA, USA, 30 January - 1 February, 2019, pp. 128–131. 18–20 November, doi:10.1109/ICOSC.2019.8665547.
58. Marchi, E.; Tonelli, D.; Xu, X.; Ringeval, F.; Deng, J.; Squartini, S.; Schuller, B. Pairwise Decomposition with Deep Neural Networks and Multiscale Kernel Subspace Learning for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Budapest, Hungary, 2016. 3 September.
59. Bisot, V.; Serizel, R.; Essid, S.; Richard, G. Nonnegative Feature Learning Methods for Acoustic Scene Classification. *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. Munich, Germany, 2017. 16–17 November.

- 544 60. Takahashi, G.; Yamada, T.; Ono, N.; Makino, S. Performance Evaluation of Acoustic Scene Classification  
545 using DNN-GMM and Frame-Concatenated Acoustic Features. Proceedings of the 9th Asia-Pacific Signal  
546 and Information Processing Association Annual Summit and Conference (APSIPA). Honolulu, Hawaii,  
547 USA, 2018, pp. 1739–1743. 2-15 November, doi:10.1109/APSIPA.2017.8282314.
- 548 61. Bae, S.H.; Choi, I.; Kim, N.S. Acoustic Scene Classification using Parallel Combination of LSTM and  
549 CNN. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE).  
550 Budapest, Hungary, 2016. 3 September.
- 551 62. Zöhrer, M.; Pernkopf, F. Gated Recurrent Networks Applied to Acoustic Scene Classification and Acoustic  
552 Event Detection. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop  
553 (DCASE). Budapest, Hungary, 2016. 3 September.
- 554 63. Jati, A.; Nadarajan, A.; Mundnich, K.; Narayanan, S. Characterizing dynamically varying acoustic scenes  
555 from egocentric audio recordings in workplace setting. submitted to IEEE International Conference on  
556 Acoustics, Speech and Signal Processing (ICASSP). Barcelona, Spain, 2020, [1911.03843].
- 557 64. Mesaros, A.; Heittola, T.; Virtanen, T. Acoustic Scene Classification in DCASE 2019 Challenge: Closed and  
558 Open Set Classification and Data Mismatch Setups. Proceedings of the Detection and Classification of  
559 Acoustic Scenes and Events Workshop (DCASE). New York, NY, USA, 2019, pp. 164–168. 25-26 October,  
560 doi:10.33682/m5kp-fa97.
- 561 65. Saki, F.; Guo, Y.; Hung, C.Y. Open-Set Evolving Acoustic Scene Classification System. Proceedings of the  
562 Detection and Classification of Acoustic Scenes and Events Workshop (DCASE). New York, NY, USA, 2019,  
563 pp. 219–223. 25-26 October.
- 564 66. Wilkinghoff, K.; Frank Kurth. Open-Set Acoustic Scene Classification with Deep Convolutional  
565 Autoencoders. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop  
566 (DCASE). New York, NY, USA, 2019, pp. 258–262. 25-26 October.
- 567 67. Lehner, B.; Koutini, K.; Schwarzmüller, C.; Gallien, T.; Widmer, G. Acoustic Scene Classification with  
568 Reject Option based on Resnets. Proceedings of the Detection and Classification of Acoustic Scenes and  
569 Events Workshop (DCASE). New York, NY, USA, 2019. 25-26 October.
- 570 68. McDonnell, M.D.; Gao, W. Acoustic Scene Classification Using Deep Residual Networks With Late Fusion  
571 of Separated High and Low Frequency Paths. Proceedings of the Detection and Classification of Acoustic  
572 Scenes and Events Workshop (DCASE). New York, NY, USA, 2019. 25-26 October.
- 573 69. Phayre, S.S.R.; Benetos, E.; Wang, Y. Subspectralnet - Using Sub-Spectrogram based Convolutional Neural  
574 Networks for Acoustic Scene Classification. Proceedings of the IEEE International Conference on Acoustics,  
575 Speech, and Signal Processing (ICASSP). Brighton, UK, 2019, pp. 825–829. 12-17 May.
- 576 70. Dang, A.; Vu, T.H.; Wang, J.C. Acoustic Scene Classification using Convolutional Neural Networks and  
577 Multi-Scale Multi-Feature Extraction. Proceedings of the IEEE International Conference on Consumer  
578 Electronics (ICCE). Hue City, Vietnam, 2018. 18-20 July, doi:10.1109/ICCE.2018.8326315.
- 579 71. Ren, Z.; Kong, Q.; Qian, K.; Plumbley, M.D.; Schuller, B.W. Attention-based Convolutional Neural  
580 Networks for Acoustic Scene Classification. Proceedings of the Detection and Classification of Acoustic  
581 Scenes and Events Workshop (DCASE). Surrey, UK, 2018. 19-20 November.
- 582 72. Wang, C.Y.; Santoso, A.; Wang, J.C. Acoustic Scene Classification using Self-Determination Convolutional  
583 Neural Network. Proceedings of the 9th Asia-Pacific Signal and Information Processing Association  
584 Annual Summit and Conference (APSIPA). Honolulu, Hawaii, USA, 2018, pp. 19–22. 2-15 November,  
585 doi:10.1109/APSIPA.2017.8281995.
- 586 73. Ruder, S. An Overview of Multi-Task Learning in Deep Neural Networks. *ArXiv pre-prints* **2017**,  
587 [1706.05098].
- 588 74. Bear, H.L.; Nolasco, I.; Benetos, E. Towards joint sound scene and polyphonic sound event  
589 recognition. Proceedings of the Annual Conference of the International Speech Communication Association  
590 (INTERSPEECH). Graz, Austria, 2019, Vol. 2019-Sept, pp. 4594–4598, [1904.10408]. 2-15 November,  
591 doi:10.21437/Interspeech.2019-2169.
- 592 75. Xu, Y.; Huang, Q.; Wang, W.; Plumbley, M.D. Hierarchical Learning for DNN-Based Acoustic Scene  
593 Classification. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop  
594 (DCASE). Budapest, Hungary, 2016. 3 September.
- 595 76. Nwe, T.L.; Dat, T.H.; Ma, B. Convolutional Neural Network with Multi-Task Learning Scheme for Acoustic  
596 Scene Classification. Proceedings of the 9th Asia-Pacific Signal and Information Processing Association



- Annual Summit and Conference (APSIPA). Honolulu, Hawaii, USA, 2018, pp. 1347–1350. 2–15 November, doi:10.1109/APSIPA.2017.8282241.
77. Boddapati, V.; Petef, A.; Rasmusson, J.; Lundberg, L. Classifying environmental sounds using image recognition networks. *Procedia Computer Science* **2017**, *112*, 2048–2056. doi:10.1016/j.procs.2017.08.250.
  78. Aytar, Y.; Vondrick, C.; Torralba, A. SoundNet: Learning Sound Representations from Unlabeled Video. *Advances in Neural Information Processing Systems (NIPS)* **2016**, pp. 892–900, [1610.09001].
  79. Singh, A.; Thakur, A.; Rajan, P.; Bhavsar, A. A Layer-Wise Score Level Ensemble Framework for Acoustic Scene Detection. Proceedings of the 26th European Signal Processing Conference (EUSIPCO). Rome, Italy, 2018, pp. 837–841. 3–7 September, doi:10.23919/EUSIPCO.2018.8553052.
  80. Kumar, A.; Khadkevich, M.; Fugen, C. Knowledge Transfer from Weakly Labeled Audio Using Convolutional Neural Network for Sound Events and Scenes. Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). Alberta, Canada, 2018, pp. 326–330. 15–20 April, doi:10.1109/ICASSP.2018.8462200.
  81. Zeinali, H.; Burget, L.; Cernocky, J. Convolutional Neural Networks and X-Vector Embeddings for DCASE2018 Acoustic Scene Classification Challenge. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE). Surrey, UK, 2018. 19–20 November.
  82. Weiping, Z.; Jiantao, Y.; Xiaotao, X.; Xiangtao, L.; Shaohu, P. Acoustic Scene Classification using Deep Convolutional Neural Networks and Multiple Spectrogram Fusions. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE). Munich, Germany, 2017. 16–17 November.
  83. Gharib, S.; Drossos, K.; Emre, C.; Serdyuk, D.; Virtanen, T. Unsupervised Adversarial Domain Adaptation for Acoustic Scene Classification. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE). Surrey, UK, 2018. 19–20 November.
  84. Kosmider, M. Calibrating Neural Networks for Secondary Recording Devices. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE). New York, NY, USA, 2019. 25–26 October.
  85. Mun, S.; Shon, S. Domain Mismatch Robust Acoustic Scene Classification Using Channel Information Conversion. Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). Brighton, UK, 2019, pp. 845–849. 12–17 May, doi:10.1109/ICASSP.2019.8683514.
  86. Drossos, K.; Magron, P.; Virtanen, T. Unsupervised Adversarial Domain Adaptation based on the Wasserstein Distance for Acoustic Scene Classification. Proceedings of the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA). New Paltz, NY, USA, IEEE, 2019, pp. 259–263. 20–23 October.
  87. Bach, S.; Binder, A.; Montavon, G.; Klauschen, F.; Müller, K.R.; Samek, W. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE* **2015**, *10*, 1–46. doi:10.1371/journal.pone.0130140.
  88. Bello, J.P.; Silva, C.; Nov, O.; DuBois, R.L.; Arora, A.; Salamon, J.; Mydlarz, C.; Doraiswamy, H. SONYC: A System for the Monitoring, Analysis and Mitigation of Urban Noise Pollution. *Communications of the ACM (CACM)* **2018**, *62*, [1805.00889].
  89. Abeßer, J.; Götze, M.; Clauß, T.; Zapf, D.; Kühn, C.; Lukashevich, H.; Kühnlenz, S.; Mimilakis, S. Urban Noise Monitoring in the Stadtlärm Project - A Field Report. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE). New York, NY, USA, 2019. 25–26 October.
  90. Grollmisch, S.; Abeßer, J.; Liebetrau, J.; Lukashevich, H. Sounding Industry: Challenges and Datasets for Industrial Sound Analysis (ISA). Proceedings of the 27th European Signal Processing Conference (EUSIPCO). A Corua, Spain, 2019, pp. 1–5. 2–6 September.
  91. Sigtia, S.; Stark, A.M.; Krstulović, S.; Plumbley, M.D. Automatic Environmental Sound Recognition: Performance Versus Computational Cost. *IEEE/ACM Transactions on Audio Speech and Language Processing* **2016**, *24*, 2096–2107. doi:10.1109/TASLP.2016.2592698.
  92. Sandler, M.; Howard, A.; Zhu, M.; Zhmoginov, A.; Chen, L.C. MobileNetV2: Inverted Residuals and Linear Bottlenecks. Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR). Salt Lake City, UT, USA, 2018, pp. 4510–4520, [1801.04381]. 18–23 June, doi:10.1109/CVPR.2018.00474.

93. Drossos, K.; Mimilakis, S.I.; Gharib, S.; Li, Y.; Virtanen, T. Sound Event Detection with Depthwise Separable and Dilated Convolutions. *ArXiv pre-prints* **2020**, [2002.00476].
94. Gordon, A.; Eban, E.; Nachum, O.; Chen, B.; Wu, H.; Yang, T.J.; Choi, E. MorphNet: Fast & Simple Resource-Constrained Structure Learning of Deep Networks. Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR). Salt Lake City, UT, USA, 2018, number 1, pp. 1586–1595, [1711.06798]. 18-23 June, doi:10.1109/CVPR.2018.00171.
95. Tan, M.; Le, Q.V. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. Proceedings of the 36th International Conference on Machine Learning (ICML). Long Beach, CA, USA, 2019, [1905.11946]. 9-15 June.
96. Mesaros, A.; Heittola, T.; Tuomas Virtanen. A Multi-Device Dataset for Urban Acoustic Scene Classification. Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE). Surrey, UK, 2018. 19-20 November.

© 2020 by the author. Submitted to *Appl. Sci.* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).