

# Acoustic Event Recognition Using Cochleagram Image and Convolutional Neural Networks

Roneel V Sharan<sup>1</sup> and Tom J Moir<sup>2</sup>

## Abstract

Convolutional neural networks (CNN) have produced encouraging results in image classification tasks and have been increasingly adopted in audio classification applications. However, in using CNN for acoustic event recognition, the first hurdle is finding the best image representation of an audio signal. In this work, we evaluate the performance of four time-frequency representations for use with CNN. Firstly, we consider the conventional spectrogram image. Secondly, we apply moving average to the spectrogram along the frequency domain to obtain what we refer as the smoothed spectrogram. Thirdly, we use the mel-spectrogram which utilizes the mel-filter, as in mel-frequency cepstral coefficients. Finally, we propose the use of a cochleagram image the frequency components of which are based on the frequency selectivity property of the human cochlea. We test the proposed techniques on an acoustic event database containing 50 sound classes. The results show that the proposed cochleagram time-frequency image representation gives the best classification performance when used with CNN.

## Keywords

Acoustic event recognition, cochleagram, convolutional neural network, mel-spectrogram, spectrogram

## 1. Introduction

Acoustic event recognition has many applications such as urban sound classification [1], fall detection [2], etc. As with other pattern recognition problems, a key challenge in acoustic event recognition is finding techniques for robust classification.

Convolutional neural networks (CNN) have produced encouraging results in image classification tasks such as on the ImageNet dataset [3], which has various image categories, and handwritten digit recognition [4]. This formed the inspiration to its extension to speech recognition tasks where it has shown to perform better than deep neural networks (DNN) [5]. Acoustic event recognition is a relatively new area of research and a number of techniques are inspired from other pattern recognition tasks, speech recognition, in particular. As such, some recent works have adopted CNN for sound

classification applications [6, 7].

CNN is an image classification technique and one of the major challenges in speech and acoustic event recognition has been how to best represent the audio signal using an image for this purpose. Two common approaches have been seen in addressing this problem. Firstly, the audio signal is converted to spectrogram images [6]. Secondly, a mel-filter, as used in computing mel-frequency cepstral coefficients (MFCC), is used to form an image-like representation [5]. We refer this as the mel-spectrogram. To ensure that all images are of an equal size, the audio signal is divided into a fixed number of frames [5] or image scaling options [7] are explored. In addition, techniques such as moving average have been applied to spectrograms [8].

The spectrogram offers equally spaced frequency components with equal bandwidth. This is not ideal for modeling the frequency characteristics of acoustic event recognition tasks where, depending on the application, most spectral energy lies in the lower frequency range. While the use of a mel-filter helps to some extent, in this work we propose the use of a cochleagram time-frequency representation for this purpose. The gammatone filter utilized in forming the cochleagram representation is modeled on the frequency selectivity property of the human cochlea. It offers narrow frequency components in the lower frequency range and wide frequency components in the upper frequency range. The finer frequency resolution in the lower frequency range helps reveal more spectral information without significantly losing spectral information in the upper frequency range [9]. The use of cochleagram has shown to be effective in audio classification and separation tasks [10, 11].

The rest of the paper is organized as follows. Section 2 provides an overview of the sound signal image representation techniques and the CNN architecture. Baseline methods for this work are presented in Section 3. In Section 4, we present the experimental setup and results followed by conclusion in Section 5.

## 2. CNN Method

An overview of the proposed input layer and the CNN architecture employed in this work is given in Fig. 1. The generic architecture of CNN has been discussed in detail in a number of literature, such as [5], therefore, we only

<sup>1</sup>School of Information Technology and Electrical Engineering, The University of Queensland, Brisbane, QLD 4072, Australia

<sup>2</sup>School of Engineering, Auckland University of Technology, Private Bag 92006, Auckland 1142, New Zealand

Email: r.sharan@uq.edu.au, tom.moir@aut.ac.nz

Corresponding Author: Roneel V Sharan (Email: r.sharan@uq.edu.au)

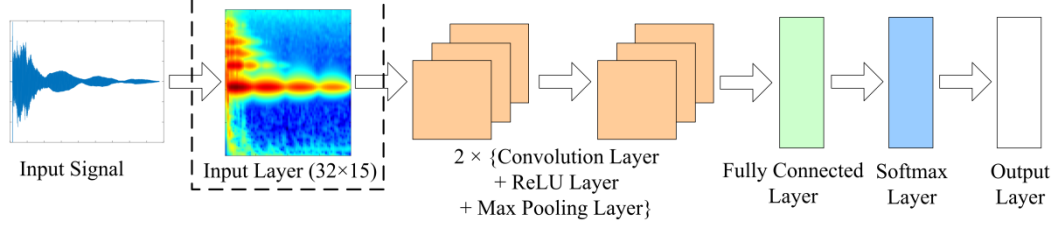


Figure 1: An overview of the proposed input layer and CNN architecture.

provide an overview of the CNN architecture and settings that we used. In this work, our primary focus is the formation of the input layer (image) of CNN for audio signal classification.

In this work, we consider spectrogram, smoothed spectrogram, mel-spectrogram, and cochleagram time-frequency images for the input layer. The dimension of the input layer image is chosen as  $32 \times 15$ , consistent with other similar work such as [5]. In addition, our sound database includes impulsive sounds and having more number of frames would make the frame size small, making it difficult to capture distinguishing frequency characteristics.

## 2.1. Spectrogram

In forming the spectrogram image, discrete Fourier transform (DFT) is applied to the windowed signal as

$$X(k, r) = \sum_{n=0}^{N-1} x(n)w(n)e^{-\frac{2\pi i k n}{N}}, \quad k=0, \dots, N-1 \quad (1)$$

where  $N$  is the length of the window,  $x(n)$  is the time-domain signal,  $X(k, r)$  is the  $k^{th}$  harmonic corresponding to the frequency  $f(k) = kF_s/N$  for the  $r^{th}$  frame,  $F_s$  is the sampling frequency, and  $w(n)$  is the window function.

The spectrogram values are obtained from log of the magnitude of the DFT values as

$$S(k, r) = \log |X(k, r)|. \quad (2)$$

To get the same time-frequency image resolution for all signals, each sound event signal is divided into 15 frames with 50% overlap between frames. A Hamming window is then applied to the frames and Fourier transform performed using 64 points so that the final image is  $32 \times 15$  dimensional.

## 2.2. Smoothed Spectrogram

Non-overlapping moving average computation is performed on the frequency components of the spectrogram in each frame to obtain a smoothed spectrogram. That is,

$$X_S(a, r) = \sum_{v=(a-1)W+1}^{aW} |X(k_v, r)|, \quad a=1, \dots, N/2W \quad (3)$$

where  $W$  is the length of the moving average window.

For the smoothed spectrogram image, each sound event signal is divided into 15 frames with 50% overlap between frames. A Hamming window is then applied and Fourier transform performed using 1024 points. For the smoothed spectrogram, the moving window length,  $W$ , is set to 16 to get 32 non-overlapping moving windows and a final image size of  $32 \times 15$ . The log values are then calculated as in Eq. (2).

## 2.3. Mel-Spectrogram

The mel-spectrogram image intensity values are computed similar to MFCC [12] but without applying the discrete cosine transform (DCT), that is, using the filter bank energies. The mel-filter bank output of the  $m^{th}$  filter can be determined as

$$E(m, r) = \sum_{k=0}^{\frac{N}{2}-1} V(m, k) |X(k, r)|, \quad m=1, 2, \dots, M \quad (4)$$

where  $E(m, r)$  is the filter bank energy of the  $m^{th}$  filter in the  $r^{th}$  frame,  $V(m, k)$  is the normalized filter response of the triangular filter banks which are equally spaced on the mel-scale [13], and  $M$  is the total number of mel-filters.

For the mel-spectrogram, the spectrogram image used in the smoothed spectrogram is computed and the number of mel-filters,  $M$ , is set to 32 to obtain a  $32 \times 15$  image. The log values are then calculated as in Eq. (2).

## 2.4. Cochleagram

In the cochleagram representation, the frequency components in the time-frequency image are based on the frequency selectivity property of the human cochlea and modeled by a gammatone filter as [14]

$$h(t) = At^{j-1}e^{-2\pi Bt} \cos(2\pi f_c t + \phi) \quad (5)$$

where  $A$  is the amplitude,  $j$  is the order of the filter,  $B$  is the bandwidth of the filter,  $f_c$  is the center frequency of the filter,  $\phi$  is the phase, and  $t$  is the time.

The equivalent rectangular bandwidth (ERB), a psychoacoustic measure of the auditory filter width at each point along the cochlea, is used to describe the bandwidth of each cochlea filter in [14]. In this work, we use the ERB filter model as described in [15] which was shown to produce the best results in [16].

After filtering the signal using the gammatone filter, the implementation of which can be found in [16, 17], a representation similar to the spectrogram is obtained by adding the energy in the windowed signal for each frequency channel as

$$C(g, r) = \sum_{n=0}^{N-1} |\hat{x}(g, n)| w(n), \quad g = 1, 2, \dots, G \quad (6)$$

where  $\hat{x}(g, n)$  is the gammatone filtered signal,  $C(g, r)$  is the  $g^{th}$  harmonic corresponding to the center frequency  $f_{cg}$  for the  $r^{th}$  frame, and  $G$  is the number of gammatone filters.

With the cochleagram representation, we set the number of gammatone filters to 32. The filtered signal is divided into 15 frames with 50% overlap between frames. The energy in each frame is added to obtain a cochleagram image of size  $32 \times 15$ . The log values are then calculated as in Eq. (2).

Illustration of spectrogram, mel-spectrogram, and cochleagram images of a sample sound signal, mapped to the jet colorspace for better visualization, are given in Fig. 2(a), (b), and (c), respectively. All three representations have the same frequency range from 0Hz to the Nyquist frequency of 22,050Hz. However, in the cochleagram representation, the finer resolution of the frequency components in the lower frequency range helps reveal more spectral information in the lower frequency range for the sound signals considered in this work without losing the spectral information in the upper frequency range. Refer to [9] for more information on cochleagram image formation and a detailed comparison between spectrogram and cochleagram representations.

## 2.5. CNN

The layout of the CNN was determined after a number of considerations and experimentations. The model is trained using stochastic gradient descent with momentum [18]. The network includes two convolution layers, each of which includes a rectified linear unit (ReLU) [19] and followed by a max pooling layer [20]. The filter size for both convolution layers is  $3 \times 3$ , stride  $1 \times 1$ , and padding  $1 \times 1$ . The number of filters in each layer is set to 16 after

experimenting with a number of filters. Similarly, the max pooling layer size is  $2 \times 2$ , stride  $1 \times 1$ , and padding  $1 \times 1$ . This is followed by a fully connected layer and a softmax layer [18] of size 50 and an output layer of size 50.

The settings for other parameters are as follows: image normalization = zero-centering, momentum = 0.4, initial learn rate = 0.01, learn rate schedule = piecewise, learn rate drop factor = 0.6, learn rate drop period = 6, L2 regularization = 0.05, mini batch size = 50, data shuffle = once, and max epochs = 20. The parameters were optimized based on the training/validation performance. The training stops after the maximum number of epochs is reached.

## 3. Baseline Method

We use three baseline feature sets (MFCC, time-frequency image feature extraction, and raw time-frequency image) and two baseline classifiers (K-nearest neighbor (KNN) and support vector machines (SVM)) to compare the performance of the proposed techniques.

### 3.1. MFCC

Firstly, we use MFCC as features which are computed as the DCT of the log compressed filter bank energies given as

$$c(l, r) = \sqrt{\frac{2}{M}} \sum_{m=1}^M \log(E(m, r)) \cos\left(\frac{\pi l}{M}(m - 0.5)\right) \quad (7)$$

which is evaluated from  $l = 1, 2, \dots, L$ , where  $L$  is the order of the cepstrum.

The number of mel-filters,  $M$ , was set to 32 to be consistent with the approach taken in forming the mel-spectrogram image for CNN classification. Two approaches were then taken for feature representation. In the first method, each signal is divided into frames of 1024 points with 50% overlap between frames. A Hamming window is then applied, Fourier transform performed using 1024 points, and MFCC values computed in each frame. The final feature vector is represented using the mean and standard deviation across each of the 32 dimensions resulting in a 64 dimensional feature vector. In the second method, the difference is that each signal is divided into 15 frames with 50% overlap between frames and the final feature vector represented using raw values of dimension 480 ( $32 \times 15$ ).

### 3.2. Time-Frequency Image Derived Features

Secondly, we use time-frequency image feature extraction. We use central moments to capture the spectral distribution in the spectrogram, a technique which

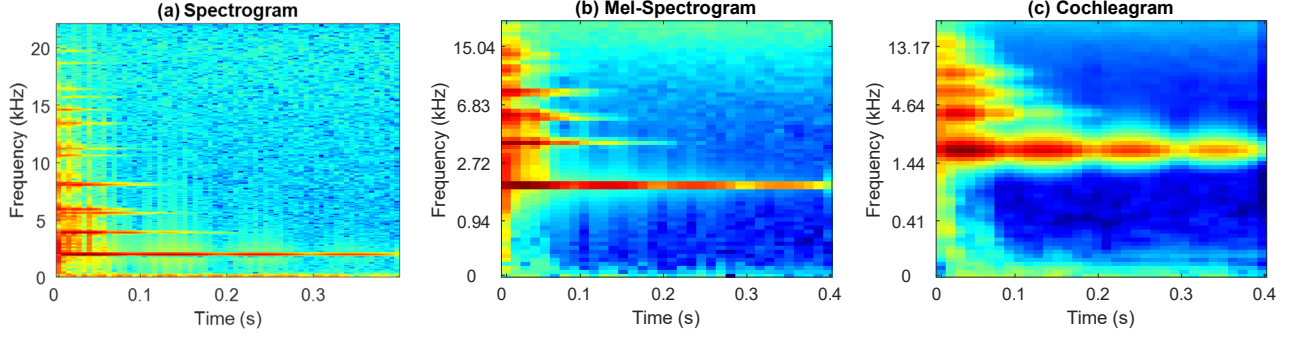


Figure 2: (a) Spectrogram, (b) mel-spectrogram, and (c) cochleagram image of a sample sound signal. The frequency range in each case is 0 – 22,050Hz.

produced encouraging results in [21]. The central moments are computed for the spectrogram and cochleagram representations and the corresponding features referred as spectrogram image feature (SIF) [21] and cochleagram image feature (CIF) [16], respectively.

In computing the SIF, we divide each signal into 1024 points with 50% overlap. Fourier transform is performed using 1024 points. The spectrogram image is divided into subbands and second and third central moments are extracted as features in each subband. Various number of subbands were experimented with but best results were achieved when using 64 subbands which results in a 128 dimensional feature vector.

For the cochleagram representation, each sound signal is filtered using 512 gammatone filters to have same number of frequency components as the spectrogram image. The filtered signal is divided into frames of size 1024 with 50% overlap and the energy in each frame computed as per Eq. 6. The same feature extraction procedure is then followed as the spectrogram image.

### 3.3. Raw Time-Frequency Image Feature

Finally, the raw time-frequency image ( $32 \times 15$ ) values were concatenated into a 480 dimensional feature vector for classification. The spectrogram, smoothed spectrogram, mel-spectrogram, and cochleagram representations, as described in section 2.1 – 2.4, were used for this purpose.

### 3.4. Baseline Classifiers

The Euclidean distance measure was used for the KNN classifier and the value of K was set to 1 after experimenting with a number of values. For the SVM classifier, nonlinear SVM with a Gaussian RBF kernel in a one-against-all strategy is used as it was found to give the best results. SVM parameters, the penalty parameter and the width of the Gaussian function [22], were tuned using Bayesian optimization [23] using 5-fold cross-validation

on the training data.

## 4. Experimental Evaluation

### 4.1. Dataset

The sound event files used in this work are taken from the Real World Computing Partnership (RWCP) Sound Scene Database (SSD) in Real Acoustical Environments [24]. The sound database has a total of 4000 manually segmented sound event files belonging to 50 classes, 80 files per class. All signals in the database have 16-bit resolution and a sampling frequency of 44,100 Hz.

Each file has one sound event the duration of which depends on the class of the sound. The sound event duration varies from a minimum of 12.6ms to a maximum of 3.84s.

### 4.2. Experimental Setup

For all experimentations, the classifier is trained and validated with 50 samples per class with the remaining 30 samples used for testing. As such, a total of 2500 samples are used for training and validating the classifier and the remaining 1500 samples used for testing the trained model. The classification accuracy is reported in percentage as *number of correctly classified test samples* divided by the *total number of test samples*. The classification accuracy in the case of CNN is averaged over 10 runs.

### 4.3. Baseline Results

The classification accuracy value using the baseline methods is given in Table I. The best baseline classification accuracy of 95.07% is achieved using the cochleagram image feature vector and SVM classifier. The cochleagram image feature vector or cochleagram image derived feature vector were seen to be more robust than features using other time-frequency representations.

Table 1 Results using baseline methods.

Method		Feature Dimension	Classification Accuracy
MFCC	MFCC-SVM (Method 1)	64	94.27
	MFCC-SVM (Method 2)	480	93.67
Time-frequency image features	SIF-KNN	128	70.27
	SIF-SVM	128	84.07
	CIF-KNN	128	82.40
	CIF-SVM	128	88.00
Raw time-frequency image feature vector (32 × 15)	Spectrogram-KNN	480	63.73
	Spectrogram-SVM		86.87
	Smoothed Spectrogram-KNN		85.00
	Smoothed Spectrogram-SVM		93.47
	Mel-Spectrogram-KNN		89.87
	Mel-Spectrogram-SVM		93.33
	Cochleagram-KNN		95.00
	Cochleagram-SVM		<b>95.07</b>

Table 2 Results using different time-frequency representations and CNN.

Method	Classification Accuracy
Spectrogram-CNN	93.46
Smoothed Spectrogram-CNN	96.34
Mel-Spectrogram-CNN	95.35
Cochleagram-CNN	<b>98.03</b>

#### 4.4. Results using CNN

In this subsection, spectrogram, smoothed spectrogram, mel-spectrogram, and cochleagram time-frequency representations are considered for classification using CNN. Results using the four time-frequency representations and CNN for classification are given in Table II. For spectrogram, smoothed spectrogram, mel-spectrogram, and cochleagram, the best classification accuracy of 93.46%, 96.34%, 95.35, and 98.03%, respectively, is achieved.

The classification accuracy using the smoothed spectrogram, mel-spectrogram, and cochleagram exceed the best baseline accuracy of 95.07%. The classification accuracy value using the mel-spectrogram is better than the spectrogram. Similarly, the classification accuracy value using smoothed spectrogram is observed to be better than the mel-spectrogram. The best classification accuracy of 98.03% is achieved using the cochleagram image. This shows the suitability of the proposed cochleagram image in acoustic event recognition.

#### 5. Conclusions

Four different time-frequency representations, spectrogram, smoothed spectrogram, mel-spectrogram,

and cochleagram, are considered for classification using CNN in an acoustic event recognition task. The performance with cochleagram image was determined to be better than spectrogram, smoothed spectrogram, and mel-spectrogram. A classification accuracy of 98.03% could be achieved which is an improvement over the best baseline classification accuracy of 95.07% using cochleagram-SVM. In summary, cochleagram time frequency representation is determined to be more suited for classification using CNN then the other time-frequency image representations for the acoustic event recognition task considered in this work.

#### References

- [1] J. Ye, T. Kobayashi, M. Murakawa, Urban sound event classification based on local and global features aggregation, *Applied Acoustics*, 117 (2017) 246-256.
- [2] S.M. Adnan, A. Irtaza, S. Aziz, M.O. Ullah, A. Javed, M.T. Mahmood, Fall detection through acoustic Local Ternary Patterns, *Applied Acoustics*, 140 (2018) 296-300.
- [3] A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, 2012, pp. 1097-1105.
- [4] D. Ciregan, U. Meier, J. Schmidhuber, Multi-column deep neural networks for image classification, in: *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 3642-3649.
- [5] O. Abdel-Hamid, A.R. Mohamed, H. Jiang, L. Deng, G. Penn, D. Yu, Convolutional neural networks for speech recognition, *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 22 (10) (2014) 1533-1545.
- [6] H. Zhang, I. McLoughlin, Y. Song, Robust sound event recognition using convolutional neural networks, in: *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 559-563.
- [7] I. Ozer, Z. Ozer, O. Findik, Noise robust sound event classification with convolutional neural network, *Neurocomputing*, 272 (2018) 505-512.
- [8] G. Kovács, L. Tóth, D. Van Compernelle, S. Ganapathy, Increasing the robustness of CNN acoustic models using autoregressive moving average spectrogram features and channel dropout, *Pattern Recognition Letters*, 100 (2017) 44-50.
- [9] R.V. Sharan, T.J. Moir, Subband time-frequency image texture features for robust audio surveillance, *IEEE Transactions on Information Forensics and Security*, 10 (12) (2015) 2605-2615.
- [10] B. Gao, W.L. Woo, L.C. Khor, Cochleagram-based audio pattern separation using two-dimensional non-negative matrix factorization with automatic sparsity adaptation, *The Journal of the Acoustical Society of America*, 135 (3) (2014) 1171-1185.
- [11] R.V. Sharan, T.J. Moir, Pseudo-color cochleagram image feature and sequential feature selection for robust acoustic event recognition, *Applied Acoustics*, 140 (2018) 198-204.
- [12] S. Davis, P. Mermelstein, Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences, *IEEE Transactions on*

- Acoustics, Speech and Signal Processing, 28 (4) (1980) 357-366.
- [13] D. O'Shaughnessy, Speech communication: human and machine, Addison-Wesley Pub. Co., 1987.
  - [14] R.D. Patterson, K. Robinson, J. Holdsworth, D. McKeown, C. Zhang, M. Allerhand, Complex sounds and auditory images, in: Y. Cazals, L. Demany, K. Horner (Eds.) Auditory physiology and perception, Pergamon, Oxford, 1992, pp. 429-446.
  - [15] D.D. Greenwood, A cochlear frequency-position function for several species - 29 years later, Journal of the Acoustical Society of America 87 (6) (1990) 2592-2605.
  - [16] R.V. Sharan, T.J. Moir, Cochleagram image feature for improved robustness in sound recognition, in: Proceedings of the IEEE International Conference on Digital Signal Processing (DSP), Singapore, 2015, pp. 441-444.
  - [17] M. Slaney, Auditory Toolbox for Matlab, Interval Research Corporation, Technical Report 1998-010, 1998.
  - [18] C.M. Bishop, Pattern recognition and machine learning, Springer, New York, 2006.
  - [19] V. Nair, G.E. Hinton, Rectified linear units improve restricted boltzmann machines, in: Proceedings of the 27th International Conference on International Conference on Machine Learning, Haifa, Israel, 2010, pp. 807-814.
  - [20] K. Jarrett, K. Kavukcuoglu, M. Ranzato, Y. LeCun, What is the best multi-stage architecture for object recognition?, in: 2009 IEEE 12th International Conference on Computer Vision, 2009, pp. 2146-2153.
  - [21] J. Dennis, H.D. Tran, H. Li, Spectrogram image feature for sound event classification in mismatched conditions, IEEE Signal Processing Letters, 18 (2) (2011) 130-133.
  - [22] C. Cortes, V. Vapnik, Support-vector networks, Machine Learning, 20 (3) (1995) 273-297.
  - [23] J. Snoek, H. Larochelle, R.P. Adams, Practical bayesian optimization of machine learning algorithms, in: Advances in Neural Information Processing Systems, 2012, pp. 2951-2959.
  - [24] S. Nakamura, K. Hiyane, F. Asano, T. Nishiura, T. Yamada, Acoustical sound database in real environments for sound scene understanding and hands-free speech recognition, in: Proceedings of the 2nd International Conference on Language Resources and Evaluation (LREC 2000), Athens, Greece, 2000, pp. 965-968.