# Blind Noise Reduction for Speech Enhancement by Simulated Auditory Nerve Representations

Anton Yakovenko[1(✉)], Aleksandr Antropov[1], and Galina Malykhina[1,2]

[1] Peter the Great St. Petersburg Polytechnic University, St. Petersburg, Russia
yakovenko_aa@spbstu.ru, malykhina@ftk.spbstu.ru
[2] Russian State Scientific Center for Robotics and Technical Cybernetics, St. Petersburg, Russia

**Abstract.** Background and environmental noises negatively affect the quality of verbal communication between humans as well as in human-computer interaction. However, this problem is efficiently solved by a healthy auditory system. Hence, the knowledge about the physiology of auditory perception can be used along with noise reduction algorithms to enhance speech intelligibility. The paper suggests an approach to noise reduction at the level of the auditory periphery. The approach involves an adaptive neural network algorithm of independent component analysis for blind source separation using simulated auditory nerve firing probability patterns. The approach has been applied to several categories of colored noise models and real-world acoustic scenes. The suggested technique has significantly increased the signal-to-noise ratio for the auditory nerve representations of complex sounds due to the variability in spatial positioning of sound sources and a flexible number of sensors.

**Keywords:** Speech enhancement · Noise reduction ·
Blind source separation · Independent component analysis ·
Machine hearing · Auditory periphery model ·
Auditory nerve responses · FastICA

## 1 Introduction

Background noises given by single or multiple sound sources are always present in the environment. In engineering practice, they have a significant impact on acoustic signal processing. However, human or mammalian auditory systems are less responsive to noise than computational systems, as they process sensory signals of the auditory periphery using high-level neuronal structures that form biological neural networks.

Humans can concentrate attention on a certain acoustic source, e.g. the speaker's voice, despite the variability of the sound environment. This allows for verbal communication in noisy environments, including conditions of multi-talker babble noise. This phenomenon of auditory perception is widely known as the "cocktail party effect", accentuated by E.C. Cherry back in 1953. It represents a unique hearing ability that enables extracting the necessary acoustic

signal source in the presence of varied background noises. In psychoacoustics, this feature is associated with auditory scene analysis (ASA) [1]. ASA is related to the problem of acoustic scene, event, or source recognition through the perceptual mechanisms of the auditory system. The principles of ASA underlie various biologically relevant computational studies, which examined the systems of practical acoustic signal processing and used one or two microphone recordings in the experimental setup. These studies are known as computational auditory scene analysis (CASA) [2].

However, there are certain conditions where technical systems of automatic speech signal processing and acoustic scene analysis may have an advantage over biological auditory systems. This advantage is due to the fact that microphone sensors used in the setup can be arbitrarily allocated and optimally positioned in space. Besides, unlike monaural or binaural hearing, a technical system can have multiple microphone sensors and channels, which allows improving the quality of results through information redundancy and appropriate signal processing [3–5]. Thus, if there are no restrictions on the number and relative position of microphones, the limitations of CASA can be circumvented.

To create machine hearing and audition systems, it is advisable to combine the advantages of auditory signal processing with technical capabilities. Auditory peripheral coding of an input acoustic signal in the form of neural responses provides a robust representation against background noises [6] due to neural phase-locking [7]. Furthermore, the representation and parametrization of speech signals based on the responses of auditory nerve (AN) fibers provides noise-robust features for automatic speech recognition, outperforming common mel-frequency cepstral coefficients under certain noise conditions [8–11].

Our study aims to develop a signal processing algorithm for noisy vowel phoneme representations in the form of simulated AN responses with the purpose of noise reduction. The approach described in the present paper imitates some features of biological neural processing. It employs a computational model of the auditory periphery and an artificial neural network for blind separation of AN responses. Three different spontaneous rates for signal and noise mixture were considered in the study.

## 2    Simulation of the AN Responses

A physiologically-motivated computer model of the auditory periphery by R. Meddis [9] was used to obtain neural responses of auditory nerve fibers. This model simulates the temporal fine structure of AN firing for three types of fibers corresponding to the input speech signal: low spontaneous rate (LSR)—less than 0.5 spikes/s, middle spontaneous rate (MSR)—0.5–18 spikes/s, and high spontaneous rate (HSR)—18–250 spikes/s [12]. In the present study, the model was set to generate a probabilistic firing rate pattern.

The model requires a digitalized speech signal in the WAV format, sampled at 44.1 kHz. The sound pressure level of the input signal was adjusted to 60 dB, as it must correspond to the preferred listening level for conversational

speech. Further, the signal passes through a processing cascade that simulates the functions of the outer, middle, and inner ear. The nonlinear mechanical behavior of the basilar membrane is modeled by a dual-resonance nonlinear filterbank (DRNL) [13]. Each segment of the basilar membrane provides the most pronounced response for a specific frequency of the acoustic stimulus, which is defined as best frequency (BF) for sounds near threshold. Thus, DRNL decomposes the signal into 41 frequency bands, logarithmically spaced from 250 to 8,000 Hz, corresponding to BFs in the most significant range for speech. The subsequent processing stages simulate stereocilia movement, inner hair cells transduction, synaptic exocytosis, and AN firing.

At the output, the auditory periphery model generates a signal encoded by the average firing rate of the auditory nerve fibers. The present study compares the results for three types of AN fibers as mentioned above. Figure 1 demonstrates a sequence of five English long vowels – clean (first column) and with additive white Gaussian noise at 0 dB SNR (second column). The duration of each vowel sound is 300 ms. The figure illustrates AN responses for LSR, MSR, and HSR fibers correspondingly. Every BF channel of the obtained AN firing probability pattern provides responses, which are then smoothed using a 20 ms Hann window and a 10 ms frame shift to extract feature data. Thus, each input signal is represented by its own multivariate data matrix consisting of 41 spectral features and an equal number of samples.
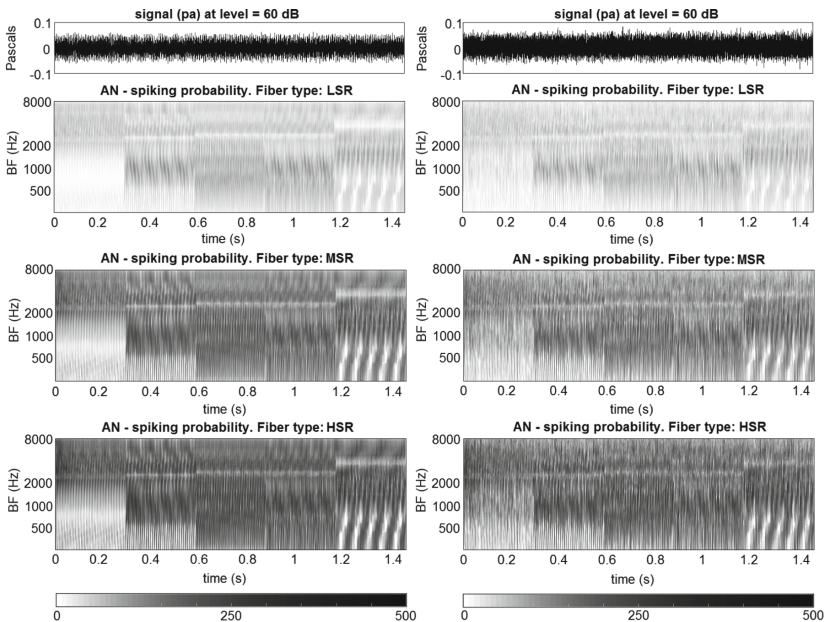


**Fig. 1.** AN firing probability patterns for three nerve fiber types for the vowel sequence: first column – clear speech signal, second column – signal corrupted by AWGN with 0 dB SNR.

## 3   Blind Signal Processing for AN Responses

Let us suppose that the sources of background noise are localized in environment about the target signal source. In that case, all sound signals are received by all sensors, but with different intensity, thus forming a linear mixture. Using the information about signal intensity difference on various sensors, we can solve the noise reduction problem for the target signal as a problem of blind source separation (BSS) [14]. Let us assume two mutually independent sound sources. The first source corresponds to the speech signal represented by a sequence of vowel phonemes. The second source is localized background environmental noise. In this case, the blind noise reduction problem [15] comes down to the task of independent component analysis (ICA).

The paper suggests using blind signal processing to solve the problem of noise reduction for a speech signal at the level of the auditory periphery. However, a technical system allows for arbitrary placement of multiple sensors in space, as distinct from the mammalian auditory system. Therefore, the intensity of signals may vary significantly, depending on the positions of sources in relation to sensors. Every sensor is represented by an auditory periphery model that encodes information in the form of stationary AN firing probability patterns. In this case, the mixing model of the speech signal and the background noise remains indeterminate, and source separation is based only on the AN responses on different sensors.

Let us consider a case where the two aforementioned sound sources are separated with the use of two biologically relevant sensors. The mixing model is a transformation of two AN output signals by a non-singular mixing matrix $\mathbf{H}$, the dimension of which depends on the number of mixed sound sources. If the sources of mixing are significantly different in amplitude or if the location of the sensors is chosen poorly, the mixing matrix is ill-conditioned. For a stable operation of the separation algorithm, it is advisable to perform a decorrelated transformation of the signal mixture in advance. The use of a decorrelation matrix makes it possible to present mixed signals in such a way that their correlation matrix is identity: $\mathbf{R}_{x_1 x_1} = E\left\{\mathbf{x}_1 \mathbf{x}_1^T\right\} = \mathbf{I}$. The mixing matrix will take the form $\mathbf{A} = \mathbf{Q}\mathbf{H}$, where $\mathbf{H}$ is the original unknown mixing matrix.

At the output of the auditory periphery model built-in each sensor, a mixture of sources is formed. Some of the signals represent the neural responses to the target signal, and others represent the responses to the noise caused by the sound environment: $\mathbf{s}(t) = [s_1(t), s_2(t)]^T$. The speech signal and the noise are mixed, and the additional mixture can affect the formed mixture to represent the intrinsic noise $\mathbf{n}(t) = [n_1(t), n_2(t)]^T$ of the system elements. The result of the conversion is the observed and measured signal $\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t) + \mathbf{v}(t)$, where $\mathbf{v}(t) = [v_1(t), v_2(t)]^T$. The task is reduced to the search for the separation matrix $\mathbf{W}$ of the observed signal vector $\mathbf{x}(t)$ by means of an artificial neural network. The matrix $\mathbf{W}$ should be such that the estimate $\mathbf{y}(t)$ of the unknown signal vector $\mathbf{s}(t)$ would be the result of applying the separation matrix to the measured signal: $\mathbf{y}(t) = \mathbf{W}\mathbf{x}(t)$. In other words, the BSS task for the AN firing rate pattern is reduced to estimating the original signal by searching for the inverse mixing operator.

We separate the stationary AN response patterns into components attributable to signal or noise, assuming the independence of neural responses to these two factors. The independence condition is determined by the minimum of information that the neural responses to signal and noise have in common. The transformation of two-dimensional signals of the AN output $\mathbf{X}_1\left(t\right), \mathbf{X}_2\left(t\right)$ into vectors is performed according to the equation: $x_{n(i-1)\times j} = x_{i,j}$. Therefore, the goal of source separation is to minimize the Kullback-Leibler divergence between the two distributions – the probability density function (PDF) $f\left(\mathbf{y}, \mathbf{W}\right)$, which depends on the coefficients of the matrix $\mathbf{W}$ and the factorial distribution:

$$f_1\left(\mathbf{y}, \mathbf{W}\right) = \prod_{i=1}^{m} f_{1,i}\left(y_i, \mathbf{W}\right) . \tag{1}$$

$$D_{f||f_1}\left(\mathbf{W}\right) = -h\left(\mathbf{y}\right) + \sum_{i=1}^{m} h_1\left(y_i\right) . \tag{2}$$

where $h\left(\mathbf{y}\right)$ is the entropy at the output of the separator, $h_1\left(y_i\right)$ is the entropy of the $i$-th element of the vector. For BSS, we used an approximation of probability density $f_{1,i}\left(y_i\right)$ by truncating the Gram-Charlier decomposition:

$$f_{1,i}\left(y_i\right) \approx N\left(y_i\right)\left[1 + \frac{\kappa_{i,3}}{3!}H_3\left(y_i\right) + \frac{\kappa_{i,2}}{4!}H_4\left(y_i\right) + \frac{\kappa_{i,6} + 10\kappa_{i,3}^2}{6!}H_6\left(y_i\right)\right] . \tag{3}$$

where $\kappa_{i,k}$ is the cumulant $k$-order of the variable $y_i$; $H_3\left(y_i\right) = y_i^3 - 3y_i, H_4\left(y_i\right) = y_i^4 - 6y_i^2 + 3, H_6\left(y_i\right) = y_i^6 - 15y_i^4 + 45y_i^2 - 15$ are Hermite polynomials; $N\left(y_i\right) = \frac{1}{\sqrt{2\pi}}exp\left(\frac{-y_i^2}{2}\right)$ is a PDF of a random quantity. The rule of weights correction when adapting a shared matrix is:

$$\mathbf{W}\left(n+1\right) = \mathbf{W}\left(n\right) + \mu\left(n\right)\left[\mathbf{I} - \varphi\left(\mathbf{y}\left(n\right)\right)\mathbf{y}^T\left(n\right)\right]\mathbf{W}^{-T}\left(n\right) . \tag{4}$$

where $\mu\left(n\right)$ is the convergence rate parameter, $\varphi\left(\mathbf{y}\left(n\right)\right) = \left[\varphi\left(y_1\left(n\right)\right), \varphi\left(y_2\left(n\right)\right)\right]^T$ is a vector consisting of activation functions, the form of which changes in the course of adaptation. The activation functions change in the learning process, since their magnitude depends on the observed values $y_i\left(n\right)$.

## 4 Results and Discussion

### 4.1 Experimental Setup

This study addresses the problem of blind noise reduction. A series of computational experiments was conducted in which noise with different spectral power distributions was removed from the signal. The study aimed to investigate the impact of noise on the stationary AN firing probability pattern distortion and

included different kinds of colored noises on the first stage: white, pink, red, blue, and violet. White Gaussian noise is a widespread noise model in robustness studies. In application areas, it is also important to remove pink (or flicker) noise, whose power spectral density is inversely proportional to frequency. The spectral density of red noise decreases in proportion to the square of the frequency. The spectral density of blue noise is specular with respect to pink noise, i.e. it increases with increasing frequency. Blue noise was synthesized using spectrum inversion. The spectral density of violet noise is inverted with respect to the red noise frequency spectrum.

On the second stage of blind noise reduction computational experiments, mixtures with real-world environmental noise were considered, including eight categories of urban acoustic scenes: airport, travelling by a bus, travelling by an underground metro, travelling by a tram, street with medium level of traffic, public square, metro station and indoor shopping mall. To obtain such categories of the environmental noises a TUT Urban Acoustic Scenes 2018 dataset [16] from DCASE Challenge was used.

Here is a summary of the experimental setup of our study. In accordance with the problem statement, we used two sensors and two sound sources. The first source was a clean speech signal. The second source was interference represented by one of the aforementioned noise types. On the first stage of computational experiments with colored noise interferences, the speech signal was a sequence of English long vowels represented by a multi-frequency complex tone that was synthesized as a sum of the first five formant frequencies – a speech-related model sound. Sound mixture had a duration of 1.5 s. On the second stage, the vowel sequence was pronounced several times by a male speaker. Sound mixture with real-world noise interferences had a duration of 10 s.

Each sensor received a sound mixture of two sources, with different mixing parameters specified in the mixing matrix. In this way, a certain spatial location of each sound source was simulated. Then, auditory peripheral representation was modelled for the sound mixtures in the form of AN average firing rate probability pattern. The output data matrices served as inputs for the FastICA algorithm of blind source separation [17]. The resulting data matrices describe the unmixed patterns of the corresponding sound sources. Finally, the impact on the blind noise reduction quality by the increase in the number of sensors from 2 to 8 has been considered.

## 4.2   Blind Noise Reduction Evaluation

The noise reduction performance for simulated AN fibers response patterns was evaluated through the signal-to-noise ratio (SNR) and noise intensity measurements. SNR allows estimating the ratio of target signal power to the power of background noise. For denoised response pattern $Y$ of AN fibers, SNR is defined as follows:

$$SNR_y = 10log_{10}\left(\frac{\|X\|^2}{\|Y - X\|^2}\right) .$$

(5)

where $X$ is the response pattern for the clean vowel sequence. Also, SNR was calculated for the response patterns for initial signal mixtures $S_1$ and $S_2$. The tables below demonstrate SNR estimation results corresponding to different spontaneous rate types of AN fibers and colored noise interferences. Table 1 presents the initial SNR for sound mixtures on two sensors provided by the mixing matrix, averaged for the vowel sequence. Table 3 presents the SNR values for the unmixed AN response pattern for the vowel sequence – a result of blind noise reduction.

**Table 1.** Initial SNR/dB for a mixture on two sensors by spontaneous rate

| Noise | White | | Pink | | Red | | Blue | | Violet | |
|---|---|---|---|---|---|---|---|---|---|---|
| Sensor | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd |
| LSR | 10.8 | 6.9 | 9.5 | 2.6 | 11.5 | 6.3 | 11.3 | 5.4 | 11.5 | 8.2 |
| MSR | 9.6 | 5.3 | 8.2 | 2.2 | 9.7 | 5.8 | 9.3 | 5.4 | 9.5 | 8.5 |
| HSR | 9.1 | 4.4 | 7.6 | 1.9 | 8.5 | 5.2 | 8.2 | 5.1 | 8.4 | 8.1 |

**Table 2.** Initial noise intensity for a mixture on two sensors by spontaneous rate

| Noise | White | | Pink | | Red | | Blue | | Violet | |
|---|---|---|---|---|---|---|---|---|---|---|
| Sensor | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd |
| LSR | 26.7 | 20.5 | 27.5 | 24.1 | 25.7 | 15.6 | 25.6 | 15.7 | 25.4 | 13.9 |
| MSR | 106.4 | 85.6 | 109.8 | 100.7 | 65.9 | 43.3 | 101.2 | 62.6 | 100.1 | 54.8 |
| HSR | 154.1 | 128.6 | 159.2 | 150.8 | 93.3 | 63.4 | 145.4 | 91.3 | 143.7 | 79.2 |

The noise intensity for each of 41 BF bands of the AN firing probability pattern can be approximated by the standard deviation. For the resultant response pattern $Y$, it can be defined as follows:

$$\sigma_y = \sqrt{\frac{1}{T}\sum_{t=1}^{T}\left(y(t) - \left[\frac{1}{T}\sum_{t=1}^{T}y(t)\right]\right)^2}. \tag{6}$$

where $T$ represents the total number of samples. Table 2 shows initial noise intensity values for the AN response pattern of the two-sensor sound mixture. Table 4 lists the resulting noise intensity values for the output AN response pattern obtained through blind noise reduction algorithm.

As can be seen, the results of noise reduction are most representative for HSR AN fibers. As for MSR and LSR fibers, the performance of blind noise reduction was poorer for colored noises. While the AN response pattern turned out to be less sensitive to red noise, the most distortion was made by violet noise. Let us consider the second stage of computational experiments. Table 5 summarizes the

**Table 3.** Resultant SNR/dB for mixture with colored noise

| Noise | White | Pink | Red | Blue | Violet |
|---|---|---|---|---|---|
| LSR | 10.9 | 10.1 | 15.8 | 12.7 | 12.6 |
| MSR | 11.4 | 11.4 | 16.2 | 11.9 | 11.7 |
| HSR | 11.5 | 12.1 | 16.2 | 11.4 | 11.3 |

**Table 4.** Resultant noise intensity for mixture with colored noise

| Noise | White | Pink | Red | Blue | Violet |
|---|---|---|---|---|---|
| LSR | 44.3 | 43.8 | 45.8 | 45.1 | 45.1 |
| MSR | 45.3 | 45.3 | 35.6 | 45.6 | 45.4 |
| HSR | 45.9 | 46.1 | 36.1 | 45.8 | 45.6 |

**Table 5.** Results of two-sensor blind noise reduction (SNR/dB) for HSR AN fibers: mixture with real-world environmental noise

| Noise | Airport | | Bus | | Metro | | Tram | | Traffic | | Square | | Station | | Mall | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Sensor | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd |
| Input | 6.3 | −4.4 | 6.2 | −4.4 | 5.5 | −6.5 | 6.1 | −5.1 | 4.6 | −7.4 | 4.7 | −7.3 | 7.4 | −4.6 | 4.5 | −7.6 |
| Average | 0.9 | | 0.9 | | −0.5 | | 0.5 | | −1.4 | | −1.3 | | 1.4 | | −1.5 | |
| Output | 7.6 | | 7.5 | | 6.6 | | 6.9 | | 7.2 | | 7.1 | | 6.2 | | 7.2 | |

blind noise reduction results in terms of SNR for a vowel sequence mixed with real-world environmental noises represented by eight categories of urban acoustic scenes. These noises largely overlap with the speech range, so their removal is the challenging task. The location of sound sources with respect to sensors was set by the mixing matrix so that the average SNR value for the sound mixtures was approximately 0 dB. As can be seen from the obtained results, the approach allowed us to improve the average value of SNR by 7 dB. This is a good result for an initial study.
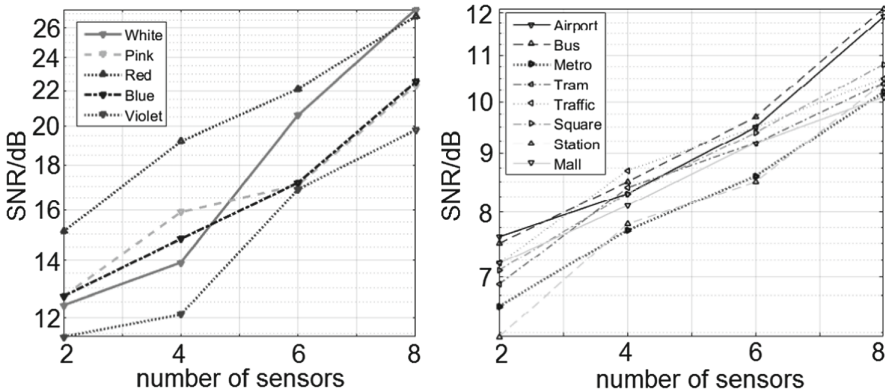


**Fig. 2.** Blind noise reduction performance for vowel sequence represented by HSR AN fibers, depending on the number of sensors: left panel – mixture with colored noises, right panel – mixture with real-world environmental noises.

As mentioned in the introduction, technical systems of speech signal processing allow the use of multiple microphone sensors. Therefore, we also evaluated blind noise reduction performance with increasing number of sensors for an HSR AN fiber response pattern. As seen from Fig. 2, performance was improved for the considered types of noise interference with increasing number of sensors, both for colored noise models and for real-world environmental noises.

## 5   Conclusions

The paper has suggested an approach to enhancement of noisy speech intelligibility by means of processing the signals of the auditory periphery. We have considered the task of designing a blind noise reduction system, which uses the information about the sound sources that is received by biologically relevant sensors distributed in space. The sensors simulates the processes of encoding information at the AN level of the auditory periphery. The speech signal, represented by a sequence of English long vowels, was separated from noise by means of independent component analysis of stationary AN firing probability patterns.

Two stages of computational studies were carried out – the first stage involved colored noise models, and the second dealt with background noises of real-world acoustic scenes. The quality of noise reduction largely depends on the mutual position of sound sources and sensors. In our case, arbitrary positions were chosen, modelled by a well-conditioned mixing matrix. The suggested approach has improved the SNR of the stationary AN firing activity pattern for colored and real-world noises. Besides, an increased number of sensors has demonstrated an improved quality of blind noise reduction.

An increase in SNR values can also be achieved through the optimization of quantity and relative placement of sensors in a given acoustic environment. Further elaboration of the approach will involve methods of blind signal extraction and real-time processing of dynamic AN firing activity patterns. The developed methodology can be used at the stage of pre-processing in machine hearing and biologically-inspired speech signal classification systems, such as [6,11,18]. Ultimately, it can become part of the new generation of neurocomputer interfaces and find use in cochlear implants [19].

## References

1. Bergman, A.S.: Auditory Scene Analysis: The Perceptual Organization of Sound. MIT Press, Cambridge (1994)
2. Wang, D.L., Brown, G.J.: Computational Auditory Scene Analysis: Principles, Algorithms, and Applications. Wiley-IEEE Press, Hoboken (2006)
3. Nugraha, A.A., Liutkus, A., Vincent, E.: Deep neural network based multichannel audio source separation. In: Makino, S. (ed.) Audio Source Separation. SCT, pp. 157–185. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-73031-8_7

4. Schwartz, O., David, A., Shahen-Tov, O., Gannot, S.: Multi-microphone voice activity and single-talk detectors based on steered-response power output entropy. In: 2018 IEEE International Conference on the Science of Electrical Engineering in Israel (ICSEE), pp. 1–4 (2018)

5. Bu, S., Zhao, Y., Hwang, M.Y., Sun, S.: A robust nonlinear microphone array postfilter for noise reduction. In: 2018 16th International Workshop on Acoustic Signal Enhancement (IWAENC), pp. 206–210 (2018)

6. Alam, M.S., Jassim, W.A., Zilany, M.S.A.: Neural response based phoneme classification under noisy condition. In: Proceedings of International Symposium on Intelligent Signal Processing and Communication Systems, pp. 175–179 (2014)

7. Miller, R.L., Schilling, J.R., Franck, K.R., Young, E.D.: Effects of acoustic trauma on the representation of the vowel "eh" in cat auditory nerve fibers. J. Acoust. Soc. Am. **101**(6), 3602–3616 (1997)

8. Kim, D.-S., Lee, S.-Y., Kil, R.M.: Auditory processing of speech signals for robust speech recognition in real-world noisy environments. IEEE Trans. Speech Audio Process. **7**(1), 55–69 (1999)

9. Brown, G.J., Ferry, R.T., Meddis, R.: A computer model of auditory efferent suppression: implications for the recognition of speech in noise. J. Acoust. Soc. Am. **127**(2), 943–954 (2010)

10. Jurgens, T., Brand, T., Clark, N.R., Meddis, R., Brown, G.J.: The robustness of speech representations obtained from simulated auditory nerve fibers under different noise conditions. J. Acoust. Soc. Am. **134**(3), 282–288 (2013)

11. Yakovenko, A., Sidorenko, E., Malykhina, G.: Semi-supervised classifying of modelled auditory nerve patterns for vowel stimuli with additive noise. In: Kryzhanovsky, B., Dunin-Barkowski, W., Redko, V., Tiumentsev, Y. (eds.) NEUROINFORMATICS 2018. SCI, vol. 799, pp. 234–240. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-01328-8_28

12. Liberman, M.C.: Auditory nerve response from cats raised in a low noise chamber. J. Acoust. Soc. Am. **63**(2), 442–455 (1978)

13. Lopez-Poveda, E., Meddis, R.: A human nonlinear cochlear filterbank. J. Acoust. Soc. Am. **110**, 3107–3118 (2001)

14. Houda, A., Otman, C.: Blind audio source separation: state-of-art. Int. J. Comput. Appl. **130**(4), 1–6 (2015)

15. Vorobyov, S., Cichocki, A.: Blind noise reduction for multisensory signals using ICA and subspace filtering, with application to EEG analysis. Biol. Cybern. **86**(4), 293–303 (2002)

16. Heittola, T., Mesaros, A., Virtanen, T.: TUT Urban Acoustic Scenes 2018, Development dataset [Data set]. Zenodo. https://doi.org/10.5281/zenodo.1228142

17. Miettinen, J., Nordhausen, K., Taskinen, S.: fICA: FastICA algorithms and their improved variants. R J. **10**(2), 148–158 (2018)

18. Yakovenko, A.A., Malykhina, G.F.: Bio-inspired approach for automatic speaker clustering using auditory modeling and self-organizing maps. Procedia Comput. Sci. **123**, 547–552 (2018)

19. Kokkinakis, K., Azimi, B., Hu, Y., Friedland, D.R.: Single and multiple microphone noise reduction strategies in cochlear implants. Trends Amplif. **16**(2), 102–116 (2012)