

Learning self-organized topology-preserving complex speech features at primary auditory cortex

Taesu Kim and Soo-Young Lee

*Brain Science Research Center and Departments of BioSystems & EECS
Korea Advanced Institute of Science and Technology, Daejeon, Republic of Korea*

Abstract

By applying independent component analysis (ICA) algorithm to auditory signals a computational model was developed for the speech feature extraction at the primary auditory cortex. Unlike the other ICA-based features with simple frequency selectivity at the basilar membrane and inner-hair-cells the learnt features represent complex signal characteristics at the primary auditory cortex such as onset/offset and frequency modulation in time. Also, the topology is preserved with the help of neighborhood coupling during the self-organization. The extracted complex features demonstrated good performance for the robust discrimination of speech phonemes.

Keywords: Independent component analysis, Auditory cortex, Complex speech features, Topology-preserving self-organization, Neural coding

1. Introduction

The computational models of human auditory nerve systems had attracted many attentions from both neuroscience and speech recognition communities. Although simple frequency-selective features are extracted at the basilar membrane and inner-hair-cells (IHCs) in the cochlea, physiological experiments prove that more complex features are extracted along the auditory pathway. Especially, some neurons at the primary auditory cortex respond to specific time-frequency signals such as onset/offset and frequency modulation in time. It is believed that this signal coding (or feature extraction) mechanism follows the information theory.

Recently, based on an information-theoretic theory of self-organization, i.e., independent component analysis (ICA), several researchers had come up with computational models of feature extraction from natural and speech sounds [1-4]. These self-organized speech features resemble the frequency-selectivity of the IHC, but do not show the complex time-dependent features at the primary auditory cortex.

In this paper, we report an ICA-based computational model for the complex speech features at the auditory cortex. The extracted features from human speeches imitate higher level auditory pathway than previous works [1-4]. Thus, in contrast to single frequency-selectivity in previous works, they are frequency transition components, on/off set components, and so on. Moreover, they provide the topology-preserving mapping of speech signals.

2. Computational Model of Auditory Nerve Systems

The feature extraction at cochlea is relatively well-understood [5][6], and we are interested in the mapping from the cochlear features to those at the auditory cortex only. Therefore, we use the existing model of cochlea to generate speech features at the inner-hair-cells, which is the input stage of the ICA-based mapping network.

We first modeled the outer and middle ear resonances as a simple linear high-pass filter as

$$s_{IE}(t) = s_{OE}(t) - 0.97 s_{OE}(t-1), \quad (1)$$

where $s_{OE}(t)$ is the speech signal at time step t and $s_{IE}(t)$ is the filtered output signal, which goes to inner ear and vibrates the tympanic membrane at the cochlea. Secondly, peripheral auditory frequency selectivity is modeled by a bank of bandpass filters with overlapping passbands [6]. Especially, we use a “gammatone” filterbanks [7] which have an impulse response of the following form :

$$g_i(t) = t^{n-1} \exp(-2\pi b_i t) \cos(2\pi f_i t + \varphi_i) H(t) \quad (1 \leq i \leq N) \quad (2)$$

where N is the number of filterbanks, n is the filter order and H is the unit step function that make the filter causal. For the i th filter bank, f_i is the center frequency (in Hz) of the filter, φ_i is the phase (in radians) and b_i determines the rate of decay of the impulse response, which is related to bandwidth. Here, we set $\varphi_i = 0$, because we set the phase doesn't make any differences.

Physiological studies of auditory nerve tuning curves [8] and psychophysical studies of critical bandwidth [9] indicate that auditory filters are distributed in frequency according to their bandwidths, which increase quasi-logarithmically with increasing center frequency. Here, we set the bandwidth of each filter according to its equivalent rectangular bandwidth (ERB), a psychophysical measurement of critical bandwidth in human subjects [9]

$$ERB(f) = 24.7(4.37 f / 1000 + 1) \quad (3)$$

Moreover, we define

$$b_i = 1.019 ERB(f_i) \quad (4)$$

and the center frequencies are equally distributed on the ERB scale between 80 Hz and 5kHz. Therefore, the outputs of each gammatone filterbanks are following:

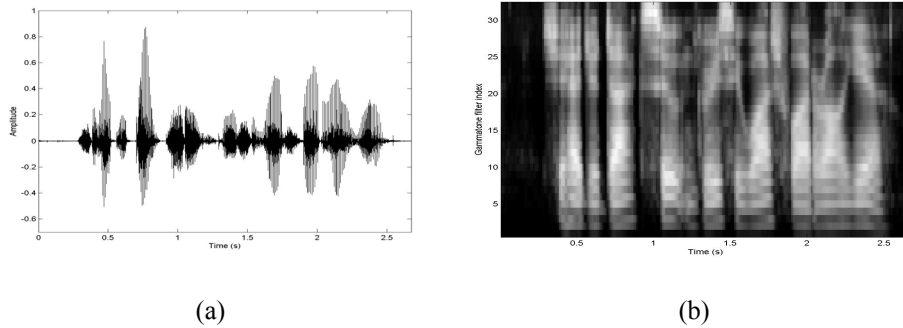


Figure 1: (a) Time-domain acoustic signal. (b) Average firing rate in the auditory nerve (IHC)

$$s_C^i(t) = \sum_{\tau} s_{IE}(\tau) g_i(t - \tau) \quad (5)$$

Finally, we calculate the power of each band in every time frame. To model the logarithmic sensitivity of the loudness, then, it goes through a logarithmic nonlinearity [5][6]

$$s_{AFR}(i, k) = \log \left(1 + \alpha \cdot \sum_{\text{all } t \text{ in the frame } k} \{s_C^i(t)\}^2 \right) \quad (6)$$

where α is scaling a factor and 1 is added to make the minimum value zero.

The results may be considered as average firing rates in the IHCs, which stimulate the higher auditory pathway. Figure 1 shows an example of the speech signal and corresponding average firing rates in time and frequency domain.

In both vision and auditory nerve systems the feature coding may be understood in the framework of the sparse coding or the ICA for super-Gaussian sources [10]. Also, an experiment on the dynamic functional connectivity in an auditory cortex shows the sparse activity [11]. Therefore, to get complex features at auditory cortex, we had applied an ICA algorithm for the IHC signals in Figure 1(b). It is worthy noting that the other ICA-based speech features were extracted from the time-domain speech signals in Figure 1(a), and resulted in frequency selectivity only. [1-4]

The basic idea of ICA is to find the hidden random variables \mathbf{v} , when the observed data (random variables) \mathbf{s} is given as a linear superposition of \mathbf{v} as

$$\mathbf{s} = \mathbf{A} \cdot \mathbf{v} + \mathbf{n}, \quad (7)$$

where \mathbf{A} is an $M \times N$ mixing matrix. [10][12] By maximizing log-likelihood of \mathbf{s} with a given \mathbf{A} , the learning rule can be derived as following:

$$\Delta \mathbf{A} \propto \mathbf{A} \mathbf{A}^T \frac{\partial}{\partial \mathbf{A}} \log P(\mathbf{s} | \mathbf{A}) = -\mathbf{A}(\mathbf{I} - \varphi(\hat{\mathbf{v}}) \cdot \hat{\mathbf{v}}^T), \quad (8)$$

where $\varphi(\hat{v}_i) = -\partial \log P_{v_i}(\hat{v}_i) / \partial \hat{v}_i$ is the score function, \mathbf{I} is the identity matrix, and $\mathbf{A} \mathbf{A}^T$ is used for the faster convergence [13]. $\hat{\mathbf{v}} = [\hat{v}_1 \ \hat{v}_2 \ \dots \ \hat{v}_N]^T$ is inferred hidden variable, which can be obtained by finding the maximum *a posteriori* value of \mathbf{v} :

$$\hat{\mathbf{v}} = \arg \max_{\mathbf{v}} P_{\mathbf{v}}(\mathbf{v} | \mathbf{s}, \mathbf{A}) = \arg \max_{\mathbf{v}} P_{\mathbf{s}}(\mathbf{s} | \mathbf{A}, \mathbf{v}) P_{\mathbf{v}}(\mathbf{v}). \quad (9)$$

When the mixing matrix is rectangular and the additive noise does not exist, the solution for \mathbf{s} can be found as $\hat{\mathbf{v}} = \mathbf{A}^{-1} \cdot \mathbf{s}$.

To implement topologically-preserving mapping from the IHC signals to those of auditory cortex, we also introduce a neighborhood function as [14]

$$h_{ij} = \exp \left[-\frac{(i-j)^2}{2\sigma^2} \right]. \quad (10)$$

where σ is the width of the neighborhood. Then complex cell output, $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_i]^T$, is defined as

$$y_i = \sum_j h_{ij} \cdot \hat{v}_j^2. \quad (11)$$

Then, the modified ICA learning rule becomes

$$\Delta \mathbf{A} \propto \mathbf{A} \mathbf{A}^T \frac{\partial}{\partial \mathbf{A}} \log P_{\mathbf{y}}(\mathbf{y} | \mathbf{A}) = -\mathbf{A}(\mathbf{I} - \psi(\hat{\mathbf{v}}) \cdot \hat{\mathbf{v}}^T), \quad (12)$$

where

$$\psi(\hat{s}_j) = \hat{s}_j \cdot \sum_i h_{ij} \cdot \phi(y_i), \quad (13)$$

$$\phi(y_i) = -\partial \log P_{y_i}(y_i) / \partial y_i. \quad (14)$$

For speech signal the probability density function of y_i can be assumed as supper-Gaussian. Thus, we used the following probability density function

$$p(y_i) = \alpha \cdot \exp(-\beta y_i^{1/2}). \quad (15)$$

Figure 2 shows the overall procedure of the method.

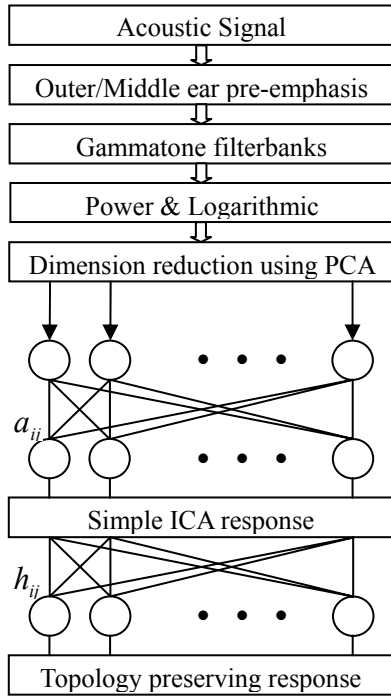


Figure 2: Overall procedure

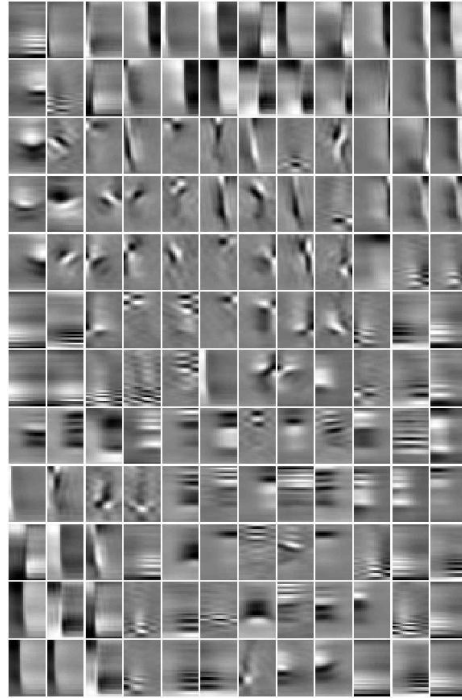


Figure 3: Learnt speech features

3. Speech Features at Auditory Cortex

In order to adapt the ICA network for self-organization, we used the TIMIT continuous speech corpus as the training data. It consists of 630 speakers, 438 males and 192 females. Ten sentences are pronounced by each speaker, and sampling frequency is 16 kHz. To calculate average firing rates in the inner-hair-cells using the gammatone filterbanks, we select 30ms frame length and 10ms shift length. For computational simplicity, 32 gammatone filters are selected. The time duration for basic speech signal is set to 20 frames, i.e., 200ms. Thus, the dimension of the data is $32 \times 20 = 640$. To remove the effect of the average component of the data, the local mean was subtracted. To remove noise and reduce computational complexity, the data dimension is reduced to 144 by principle component analysis (PCA). Therefore, the input dimension of ICA network is 144. For self-organizing topology-preserving mapping, we determine the neighborhood in a 12×12 rectangle. Thus,

Eq. (10) was modified to a 2-dimensional neighborhood function. The width, σ , of the neighborhood was set to 2, and then decreased for the fine tuning.

Figure 3 shows the learnt speech features based on self-organizing scheme. They are located topographically in the 2-dimensional map. The horizontal lines at the lower right corner represent time-independent frequency-maintaining components. Usually, in these cases, more than two frequency components are appeared concurrently. That can be interpreted as harmonics components. Figure 4 shows an example of the harmonics component. As shown in figure 4(b), 4 circled points indicate from 1st to 4th peak frequencies, which are approximately 350 Hz, 700 Hz, 1100 Hz, and 1500 Hz, respectively. 2nd to 4th ones, thus, are almost 2, 3, and 4 times of fundamental frequency 350 Hz.

In contrast to horizontal lines, vertical lines are appeared at the lower left and upper right corners of the map, and represent noise bursts or signal on- and off-set components. Around the middle and upper left corner of the map, lots of frequency-transitions are found. They are frequency-rising and frequency-falling components, which can be understood as frequency modulation (FM) components. For example, two features at 5th row and 2nd column, and 4th row and 3rd column in figure 3 are quite good examples of frequency-rising and frequency-falling component, respectively. In fact, human speech and animal sounds share the same three basic elements: steady-state harmonically related frequencies, FMs, and onset/offset bursts. Also, many auditory cortical areas are tonotopically. Therefore, our result fits those biological evidences well.

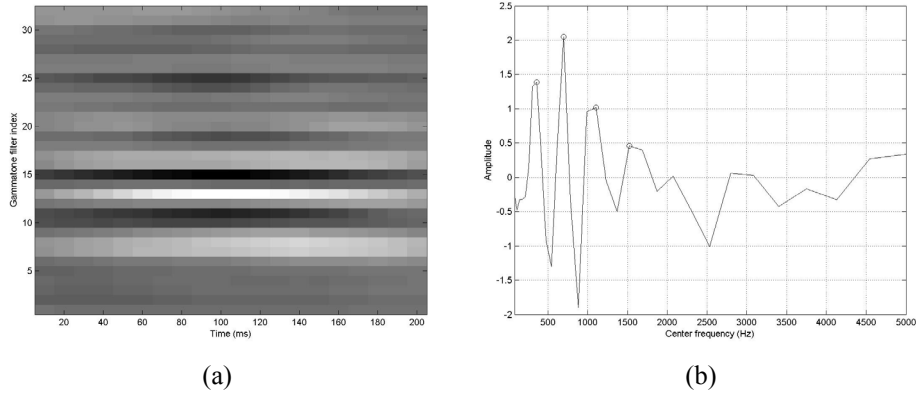


Figure 4: (a) A learnt speech feature at 8th row 11th column in figure 3. (b) Time averaged frequency of (a)

4. Conclusion

An information-theoretic model based on independent component analysis results in complex speech features such as onset/offset and frequency modulation in time, which are extracted at the primary auditory cortex. Also, similar to the auditory cortex, their topology is preserved. It demonstrates that the signal processing mechanism in human auditory nerve systems follows the information theory for the best efficiency in signal coding and feature extraction. These features are speaker independent, and may be applicable to the automatic speech recognition.

Acknowledgment: This work was supported by the Chung Moon Soul BioInformation and BioElectronics Center and also by the Korean Ministry of Science and Technology as a Brain Neuroinformatics Research Program.

References

- [1] A. J. Bell and T. J. Sejnowski, "Learning the higher-order structure of a natural sound," *Network : Computation in Neural Systems*, Vol. 6, pp. 261–266, 1996.
- [2] S. A. Abdallah and M.D. Plumbley, "If the independent components of natural images are edges, what are the independent components of natural sounds?," in *Proc. ICABSS*, 2001.
- [3] M. S. Lewicki, "Efficient coding of natural sounds," *Nature Neuroscience*, 2002.
- [4] J. H. Lee, T. W. Lee, H. Y. Jung, and S. Y. Lee, "On the efficient speech feature extraction based on independent component analysis," *Neural Processing Letters*, Vol. 15, pp.235–245, 2002.
- [5] W. A. Yost, "Fundamentals of hearing – An introduction," Academic Press, 2000.
- [6] B. C. J. Moore, "An introduction to the psychology of hearing," 4th edition, San Diego, CA: Academic, 1997.
- [7] R. D. Patterson, I. Nimmo-Smith, J. Holdsworth, and P. Rice, "APU Report 2341: An efficient auditory Filterbank based on the gammatone function," Cambridge: Applied psychology unit, 1988.
- [8] A. R. Palmer, "Physiology of the cochlear nerve and cochlear nucleus," in *Hearing*, M. P. Haggard and E. F. Evans, Eds., Edinburgh: Churchill Livingstone, 1987.
- [9] B. R. Glasberg and B. C. J. Moore, "Derivation of auditory filter shapes from notched-noise data," *Hearing research*, Vol. 47, pp. 103-138, 1990.
- [10] T. W. Lee, "Independent component analysis," Kluwer Academic Publishers, 1998.
- [11] R. S. Clement, W. Patrick, J. Rousche, and D. R. Kipke, "Functional connectivity in auditory cortex using chronic," *Neurocomputing* Vol. 26–27, pp. 347–354, 1999.
- [12] A. Hyvärinen, J. Karhunen, and E. Oja, "Independent component analysis," John Wiley & Sons, 2001.
- [13] S. Amari, A. Cichocki, and H.H. Yang, "A new learning algorithm for blind source separation," in *Proc. NIPS*, 1996.
- [14] A. Hyvarinen, P. O. Hoyer, and M. Inki, "Topographic independent component analysis," *Neural computation*, Vol 13, pp.1527–1558, 2001.



Taesu Kim was born in Korea, in 1978. He received the B.S. degree from the Hanyang University in 2001 and the M.S. degree from the Korea Advanced Institute of Science and Technology (KAIST) in 2003, both in electrical engineering. He is currently a Ph. D. candidate at the Department of BioSystems, KAIST. His research interests include biologically-motivated information processing and probabilistic method for unsupervised or reinforcement learning.



Soo-Young Lee received a Ph.D. degree from the Polytechnic Institute of New York in electrophysics in 1984. From 1982 to 1985 he had worked for General Physics Corporation at Columbia, MD, USA. In early 1986 he joined the Korea Advanced Institute of Science and Technology, and now is a Full Professor for the Department of BioSystems. His research interests have resided in artificial neural networks and auditory systems based on brain information processing mechanism. Dr. Lee had served as the President of Asia-Pacific Neural Network Assembly in 2001, and received Leadership Award and Achievement Award from International Neural Network Society in 1994 and 2001, respectively. Dr. Lee is the Editor-in-Chief of a new journal, *Neural Information processing – Letters and Reviews*.