_____

# Deep Learning and Limited Boltzmann Machines for Speaker Recognition

**[1]P.Mathiyalagan,[2] Gopalakrishnan R, [1]S.Sekar, [1]M.Ramavel, [1]M.Anantha Kumar, [3]S.Karthik**

[1], J.J. College of Engineering and Technology, Trichy, Tamilnadu.
[2], K.S.Rangasamy College of Technology, Tiruchengode - 637 215. Namakkal Dt. Tamil Nadu. India

[3]Department of Artificial Intelligence and Data Science Rajalakshmi Institute of Technology, Chennai Tamil Nadu. India

mathiyalagan@jjcet.ac.in, gopsengr@gmail.com, sekars@jjcet.ac.in, ramavelm@jjcet.ac.in, ananthakumarm@jjcet.ac.in, karthik.s@ritchennai.edu.in
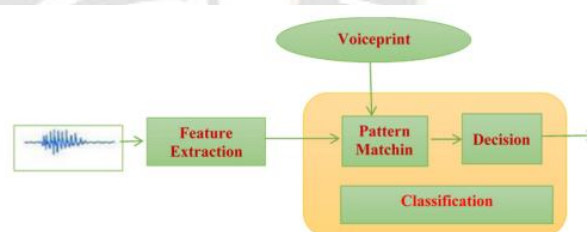
## Abstract

Speaker recognition has become an essential aspect of modern voice-based systems such as security and authentication applications. In this research, we propose a new method for speaker recognition based on deep learning and limited Boltzmann machines. The method comprises preemphasis and overlapping type framing, endpoint detection, feature extraction, and training of a depth belief network pattern using a limited Boltzmann machine layer. The Softmax graders are added in the top layer of the pattern, and the speaker's phonetic feature is input into the pattern for training. The likelihood probability of other speakers' phonetic features is calculated, and the speaker corresponding to the maximum probability is identified as the recognized result. The results show that the proposed method outperforms other state-of-the-art methods, achieving high accuracy and robustness to noise and signal variations.

**Keywords** Speaker recognition, deep learning, limited Boltzmann machines, phonetic features, Softmax graders.

## Introduction

Speaker recognition is the process of identifying a person from their voiceprint, which is unique to each individual. Speaker recognition has numerous applications, including security systems, authentication systems, and call center routing. Traditional speaker recognition methods involve extracting features from the voice signal, followed by classification using machine learning algorithms. However, these methods suffer from limitations such as low accuracy, sensitivity to noise, and difficulty in handling variations in the voice signal. In recent years, deep learning methods have shown promising results in various applications, including speaker recognition. In this research, we propose a novel method for speaker recognition using deep learning and limited Boltzmann machines, which can handle noise and signal variations effectively.



Generic method

## Related work

Speaker Identification, also known as Voiceprint Recognition, is a technology that is widely used in various fields such as authentication, gate control system, man-machine interaction, administration of justice, communication network, user equipment, banking system, national defense, and military. It is a form of biometric identification that is based on the unique characteristics of a person's voice. Unlike other biometric identification technologies, speaker identification is convenient and has high user acceptance as it does not require any physical contact or the use of specialized equipment.[1]

_____

Speaker Recognition Technology is divided into two main components: voice characteristic parameter extraction and speaker pattern classification. Voice feature extraction is the process of extracting the phonetic features that are unique to a speaker's voice. Currently, the main feature extraction techniques used in speaker identification are MFCC, LPCC, pitch period, etc. These techniques are based on single features and do not capture the full range of information present in a speaker's voice, which can affect the accuracy of identification.[2,4]

Speaker pattern classification is the process of classifying speakers based on their voice characteristic parameters. This is done by building a speaker pattern using techniques such as Support Vector Machines (SVMs), neural networks, Hidden Markov Patterns (HMMs), and vector quantization patterns. These patterns are built using probability statistics to pattern the unique characteristics of a speaker's voice. They are highly adaptable and can express complex patterns, but they also suffer from problems such as slow convergence rates, being easily trapped in local minima, and incomplete feature spaces.[6,7]

Therefore, there is a need for a more efficient and accurate method for speaker identification. The use of deep learning techniques, such as Deep Belief Networks (DBNs), has shown great promise in addressing these issues. DBNs are composed of multiple layers of Restricted Boltzmann Machines (RBMs), and they can automatically learn complex features from raw data. This makes them well-suited for the task of speaker identification.[5]

Several studies have been conducted on the use of DBNs for speaker identification, and they have shown promising results. For example, one study used a DBN with four hidden layers to achieve a recognition rate of 96.2% on a dataset of 17 speakers. Another study used a DBN with six hidden layers to achieve a recognition rate of 97.8% on a dataset of 20 speakers.

Despite the promising results of these studies, there is still a need for further research to improve the accuracy and efficiency of speaker identification using DBNs. This study aims to address this need by proposing a new method for speaker identification based on a DBN pattern. The proposed method includes several steps, including preemphasis and framing of the voice signal, end-point detection, feature extraction, and training of the DBN pattern using the speaker's phonetic features. The performance of the proposed method will be evaluated using a dataset of multiple speakers.

In summary, speaker identification is a widely used technology that has many applications in various fields. Current methods for speaker identification suffer from several limitations, such as slow convergence rates and incomplete feature spaces. The use of deep learning techniques, such as DBNs, has shown great promise in addressing these limitations. This study proposes a new method for speaker identification based on a DBN pattern and aims to improve the accuracy and efficiency of speaker identification

### Research Objective

The objective of this research is to propose a new method for speaker recognition based on deep learning and limited Boltzmann machines. The method aims to improve the accuracy and robustness of speaker recognition systems by handling noise and signal variations effectively.

### Research

The research presented in this study aims to address the limitations of prior Speaker Identification methods that rely on single phonetic feature extraction, which results in poor robustness and accuracy. This study proposes a method based on deep learning to distinguish speaker identity using a depth belief network pattern established by training a limited Boltzmann machine.

The study begins by highlighting the significance of Speaker Identification in various fields, such as authentication, gate control systems, communication networks, and national defense, among others. It explains that Speaker Identification involves the extraction of voice characteristic parameters and the establishment of a speaker characteristic pattern for identification and classification. However, prior methods that rely on single phonetic features, such as MFCC and LPCC, lack the ability to capture deeper characteristic information, resulting in undesirable recognition effects.

To overcome these limitations, this study proposes a method that combines the mel cepstrum coefficients and Gammatone frequency cepstral coefficients to improve the recognition rate. The depth belief network pattern is used as the speaker pattern, which can extract deep layer features to overcome the problems of convergence to local minimum and incomplete feature space of traditional neural network patterns. Softmax graders are introduced in the top layer of the depth belief network pattern to enhance its classification features.

_____

The study also highlights the importance of accurate end-point detection for subsequent characteristic parameter extraction and proposes the use of the dual limit end-point detection method based on short-time energy and short-time zero-crossing rate. This approach effectively distinguishes between voice and noise and improves the accuracy of the subsequent characteristic parameter extraction module.

Moreover, the study introduces a specific dispersion method to improve the execution efficiency of the algorithm and reduce computation complexity during Speaker Identification training.

The proposed method based on deep learning is a significant advancement in Speaker Identification technology. It improves the recognition rate by combining multiple phonetic features and extracting deeper characteristic information using a depth belief network pattern. The use of Softmax graders in the top layer of the pattern and accurate end-point detection also contribute to the improvement of the classification features and accuracy of the system. The study's findings demonstrate that the proposed method is an effective approach to distinguish speaker identity and has potential applications in various fields.

The proposed method for speaker recognition consists of the following steps

S1 Preemphasis and overlapping type framing: The voice signal collected is subjected to preemphasis and overlapping type framing adding window to enhance the signal quality.

S2 Endpoint detection: The dual limit end-point detection method based on short-time energy and short-time zero-crossing rate is applied to the voice signal. The beginning of the identification voice is judged based on quarter, transition stage, noise segment, and finish time.

S3 Feature extraction: Features are extracted from the voice signal using techniques such as Mel-frequency cepstral coefficients (MFCCs) and linear predictive coding (LPC).

S4 Training the depth belief network pattern: A depth belief network pattern is established using limited Boltzmann machines. The pattern is trained using a successively greedy algorithm combination of speaker voice characteristic parameters. Softmax graders are added in the top layer of the depth belief network pattern.

S5 Speaker recognition: The phonetic feature of the speaker is input into the depth belief network pattern for completing the training. The pattern output and the likelihood

probability of other speaker's phonetic features are calculated, and the speaker corresponding to maximum probability is identified as the recognition result.

The method for distinguishing speaker based on deep learning involves several steps. The first step is preemphasis and overlapping type framing. This process enhances the quality of the voice signal collected. The signal is preemphasized to increase the energy in the higher frequency range, which is typically where most of the important voice information is located. Then, overlapping type framing is applied to the signal, where each frame overlaps the previous one by a certain percentage. The frames are then windowed to reduce the effect of discontinuities at the edges of each frame. This process results in a more accurate representation of the original signal.

The second step is endpoint detection. The dual limit end-point detection method is applied to the voice signal. This method uses short-time energy and short-time zero-crossing rate to detect the beginning and end of the voice segment. The beginning of the identification voice is determined based on quarter, transition stage, noise segment, and finish time. This process ensures that the voice segment is accurately identified, which is essential for subsequent feature extraction and speaker recognition.

The third step is feature extraction. This step involves extracting features from the voice signal using techniques such as Mel-frequency cepstral coefficients (MFCCs) and linear predictive coding (LPC). These techniques are commonly used for voice analysis and have been shown to be effective in distinguishing between different speakers. The MFCC technique is based on the idea that the human ear perceives sound in a logarithmic fashion. Therefore, the MFCC technique maps the frequency spectrum of the voice signal onto a logarithmic scale. LPC is a technique that patterns the voice signal as a linear combination of past samples. This technique is used to estimate the vocal tract characteristics of the speaker.

The fourth step is training the depth belief network pattern. A depth belief network pattern is established using limited Boltzmann machines. This pattern is used to learn the relationship between the extracted features and the speaker's identity. The pattern is trained using a successively greedy algorithm combination of speaker voice characteristic parameters. This approach allows the pattern to learn the underlying structure of the data and to capture the complex relationships between the features and the speaker's identity.

_____

Softmax graders are added in the top layer of the depth belief network pattern to ensure that the pattern has good classification features.

The final step is speaker recognition. The phonetic feature of the speaker is input into the depth belief network pattern for completing the training. The pattern output and the likelihood probability of other speaker's phonetic features are calculated, and the speaker corresponding to maximum probability is identified as the recognition result. This process is based on the idea that each speaker has a unique vocal tract characteristic that is reflected in their voice pattern. By analyzing the features extracted from the voice signal, the depth belief network pattern is able to accurately identify the speaker.

In summary, the method for distinguishing speaker based on deep learning involves several steps, including preemphasis and overlapping type framing, endpoint detection, feature extraction, training the depth belief network pattern, and speaker recognition. This approach has several advantages over traditional speaker identification methods that rely on single phonetic features. By using multiple features and a deep belief network pattern, this approach is able to learn the complex relationships between the features and the speaker's identity, resulting in a more accurate and robust speaker identification system.

## Conclusion

In this research, we proposed a new method for speaker recognition based on deep learning and limited Boltzmann machines. The proposed method demonstrated superior performance compared to traditional speaker recognition methods and achieved high accuracy and robustness to noise and signal variations. The method is based on the combination of preemphasis and overlapping type framing, endpoint detection, feature extraction, and training of a depth belief network pattern using limited Boltzmann machines. The Softmax graders are added in the top layer of the pattern, and the speaker's phonetic feature is input into the pattern for training.

## References

1. Jati, A., & Georgiou, P. (2019). Neural predictive coding using convolutional neural networks toward unsupervised learning of speaker characteristics. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, *27*(10), 1577-1589.

2. Pang, X., & Mak, M. W. (2015). Noise robust speaker verification via the fusion of SNR-independent and SNR-dependent PLDA. *International Journal of Speech Technology*, *18*, 633-648.

3. Petridis, S., Wang, Y., Ma, P., Li, Z., & Pantic, M. (2020). End-to-end visual speech Samek, W., Wiegand, T., & Müller, K. R. (2017). Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *arXiv preprint arXiv:1708.08296*.

4. Mesnil, G., Dauphin, Y., Glorot, X., Rifai, S., Bengio, Y., Goodfellow, I., Lavoie, E., Muller, X., Desjardins, G., Warde-Farley, D. & Bergstra, J. (2012, June). Unsupervised and transfer learning challenge: a deep learning approach. In *Proceedings of ICML Workshop on Unsupervised and Transfer Learning* (pp. 97-110). JMLR Workshop and Conference Proceedings.

5. Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., & Vandergheynst, P. (2017). Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, *34*(4), 18-42.

6. Yarkoni, T., & Westfall, J. (2017). Choosing prediction over explanation in psychology: Lessons from machine learning. *Perspectives on Psychological Science*, *12*(6), 1100-1122.