



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 13 **Issue:** V **Month of publication:** May 2025

DOI: <https://doi.org/10.22214/ijraset.2025.71160>

www.ijraset.com

Call: ☎ 08813907089

E-mail ID: ijraset@gmail.com

Advancements in Speaker-Independent Speech Separation Using Deep Attractor Networks

Mrs. Asharani Chadchankar¹, Priya Shetty², Vaibhavi Dongare³, Samruddhi Bhandare⁴, Aliza Sayyad⁵

¹Assistant Professor, Information Technology, Marathwada Mitramandal's College of engineering, Sppu, Pune, India

^{2, 3, 4, 5}Final year student, Information Technology, Marathwada MitraMandal's College of engineering, Sppu, Pune, India

Abstract: *Speaker-independent speech separation, the task of isolating individual voices from a mixture without prior knowledge of the speakers, has gained significant attention due to its importance in various applications.*

However, challenges such as the arbitrary order of speakers and the unknown number of speakers in a mixture remain significant hurdles. This research paper analyzes Deep Attractor Networks (DANet), a novel deep learning framework designed to address these issues. DANet projects mixed speech signals into a high-dimensional embedding space where reference points, known as attractors, represent individual speakers. By encouraging time-frequency embeddings to cluster around their corresponding attractors, the network facilitates effective speech separation. This paper provides a comprehensive analysis of the DANet architecture, the methodologies for attractor formation, system analysis, potential enhancements, evaluation on standard datasets, and diverse applications, highlighting its potential in advancing the field of speech separation.

I. INTRODUCTION

The ability to segregate individual speech signals from a complex auditory scene, often referred to as speech separation, is a fundamental prerequisite for numerous speech-related technologies. Applications spanning automatic speech recognition (ASR), speaker recognition, and speech coding heavily rely on the accurate isolation of target voices from interfering sounds.¹ The "cocktail party problem," which describes the human capacity to focus on a single speaker amidst a multitude of voices, serves as a primary inspiration and benchmark for research in this domain. The increasing prevalence of voice-controlled devices and the growing need for effective communication in noisy environments further underscore the significance of robust speech separation techniques.

Speaker-independent speech separation presents a unique set of challenges compared to scenarios where prior information about the speakers is available. Two prominent obstacles are the permutation problem, arising from the arbitrary order in which speakers appear in a mixture, and the output dimension problem, caused by the unknown number of speakers present. These complexities necessitate the development of sophisticated models capable of generalizing across diverse speakers and adapting to varying numbers of sources.

Deep Attractor Networks (DANet) have emerged as a promising deep learning framework specifically engineered to confront the intricacies of speaker-independent speech separation. At its core, DANet employs a neural network to map mixed speech signals into a high-dimensional embedding space, where each speaker is represented by a reference point known as an Attractor. This approach encourages the embeddings corresponding to each speaker to form clusters around their respective attractors, thereby enabling the separation of individual voices. This paper aims to provide a comprehensive exploration of DANet, encompassing its architectural nuances, the various methodologies for determining Attractor points, a thorough system analysis, potential avenues for enhancement, a review of its evaluation on standard benchmarks, and a discussion of its wide-ranging applications. By delving into these aspects, this paper seeks to contribute to a deeper understanding and further advancement of speaker-independent speech separation technology.

II. BACKGROUND AND LITERATURE REVIEW

Historically, speech separation has been tackled using a variety of signal processing techniques. Methods such as spectral subtraction, which aims to reduce noise by estimating the noise spectrum and subtracting it from the mixed signal, and Wiener filtering, which uses statistical properties of the signal and noise to design optimal filters, have been employed. Computational Auditory Scene Analysis (CASA) represents another class of traditional approaches that attempts to mimic the human auditory system by using perceptual cues to segregate sound sources.

However, these traditional methods often struggle in complex real-world scenarios involving multiple unknown speakers and significant overlap, as they often rely on simplifying assumptions about the statistical properties of speech and noise. The field of speech processing has witnessed a paradigm shift with the advent of deep learning. Deep neural networks (DNNs) have demonstrated an exceptional ability to learn intricate, non-linear relationships directly from vast amounts of data. This capability has led to remarkable advancements across various speech processing tasks, including speech recognition, enhancement, and, notably, separation. The layered structure of deep networks enables the automatic extraction of relevant features from raw audio signals, surpassing the limitations of hand-engineered features used in traditional methods.

In the realm of deep learning-based speech separation, early efforts focused on utilizing DNNs to predict time-frequency masks, such as the ideal binary mask (IBM) or the ideal ratio mask (IRM), which could then be applied to the mixed signal to isolate individual speakers. Regression-based approaches also emerged, where DNNs are trained to directly estimate the clean speech features from the mixed input. Notably, Zhang and Wang proposed a deep ensemble learning method using multi-context networks for monaural speech separation, while Du et al. explored a regression approach employing high-resolution DNNs. These early applications of DNNs marked a significant step forward in data-driven speech separation.

Temporal Convolutional Networks (TCNs) have also proven highly effective in sequence modeling, including audio processing. Lea et al. introduced TCNs for action segmentation, highlighting their capacity to capture long-range temporal dependencies, a crucial aspect for understanding the evolution of speech over time. Bai et al. conducted an empirical evaluation demonstrating the superior performance of TCNs over recurrent neural networks (RNNs) in various sequence modeling tasks. The ability of TCNs to process sequential data in parallel and model long contexts efficiently makes them a valuable architecture for speech separation. Deep clustering represents another prominent deep learning technique for speaker-independent speech separation. This approach involves training a network to generate embeddings for each time-frequency bin of the mixed signal, such that embeddings belonging to the same speaker are located closer to each other in the embedding space. Deep clustering has been shown to effectively address both the permutation and output dimension problems inherent in speaker-independent separation. Wang et al. extended this concept to multi-channel audio by incorporating spatial information into the deep clustering framework. By learning discriminative embeddings, deep clustering offers a robust way to segment a mixture into its constituent speaker components. Deep Attractor Networks (DANet) build upon the success of embedding-based approaches and offer a distinct framework for speaker-independent speech separation. While sharing the concept of projecting the mixture into an embedding space with deep clustering, DANet introduces the idea of Attractor points that directly represent each speaker. This allows for an end-to-end trainable system that, in some variations, bypasses the explicit clustering step required at test time in deep clustering. By learning these Attractor points, DANet aims to provide a more direct and potentially more stable method for separating individual speakers from a mixed audio signal.

III. DEEP ATTRACTOR NETWORKS: ARCHITECTURE AND PRINCIPLES

A typical Deep Attractor Network (DANet) architecture often comprises several layers of Bi-directional Long Short-Term Memory (BLSTM) units, which serve as the primary feature extraction mechanism. These BLSTM layers are adept at capturing the temporal dynamics of speech by processing the input sequence in both forward and backward directions, allowing the network to consider both past and future context. Following the BLSTM layers, a fully connected feed-forward layer is commonly employed to project the extracted features into a high-dimensional embedding space. The input to this network is typically the time-frequency representation of the mixed audio signal, often in the form of a log magnitude spectrogram. To facilitate processing and training, the input audio is often segmented into smaller, non-overlapping chunks of frames.

The core of DANet lies in the concept of a high-dimensional embedding space. This latent space serves as a transformed representation of the complex audio mixture, where each time-frequency bin of the input is mapped to a corresponding point or vector. The primary objective of the neural network is to learn an embedding function that arranges these time-frequency bins in such a way that those originating from the same speaker are located closer to each other, while bins from different speakers are further apart. This creates a structured representation where the characteristics of individual speakers become more distinguishable. Within this embedding space, Attractor points play a crucial role. These reference points represent individual speakers present in the mixture, acting as anchors or prototypes for each voice. Typically, an Attractor point for a specific speaker is defined as the centroid or the weighted average of the embeddings of all the time-frequency bins that primarily belong to that speaker. During the training phase, these Attractor points are formed based on either the true speaker assignments (if available) or estimated assignments for each time-frequency bin. The network learns to position these attractors in the embedding space to effectively represent the unique characteristics of each speaker.

The speech separation process in DANet relies on the principle of encouraging the embeddings of each speaker to cluster around their corresponding Attractor points. For each time-frequency bin in the mixture, the network calculates its similarity or distance to each of the Attractor points. This measure of proximity is then used to generate soft separation masks, where bins closer to a particular speaker's Attractor are assigned a higher probability of belonging to that speaker. By applying these masks to the original mixture spectrogram, the separated speech signals for each speaker can be reconstructed.

DANet inherently addresses the challenges of the permutation and output dimension problems in speaker-independent speech separation. The permutation problem is tackled by establishing a direct link between the order in which the masks are generated and the order of the Attractor points, which is consistent with the target speaker order provided during training. As long as the target masks and the speaker assignment function used to form the attractors maintain the same order, the network learns a consistent mapping between attractors and speakers. The output dimension problem, arising from the unknown number of speakers, is handled by dynamically determining the number of attractors based on the speaker assignment function during the training phase. Since the number of attractors is a function of the speaker assignment, the network can adapt to mixtures with a variable number of speakers without requiring modifications to its architecture.

IV. METHODS FOR FINDING ATTRACTORS IN DANET

Several methods have been proposed for determining the Attractor points within the Deep Attractor Network framework. One approach involves using clustering algorithms, such as K-means, directly on the embeddings generated by the network. In this method, often referred to as DANet-Kmeans, the embeddings are grouped into clusters, and the centroids of these clusters are then treated as the Attractor points for generating the separation masks. This approach offers flexibility in handling mixtures with an unknown number of speakers, as the number of clusters can be set according to the estimated number of sources. Another strategy involves utilizing fixed Attractor points.¹ This method is based on the empirical observation that the locations of Attractor points in the embedding space tend to remain relatively stable across different speech mixtures. In this approach, Attractor points are first estimated from the training data across various mixtures.

Subsequently, the mean of these estimated attractors is used as a fixed set of attractors during the testing phase. A significant advantage of this method is that it eliminates the need for a clustering step during the testing phase, allowing for potentially faster and real-time implementations.

The Anchored Deep Attractor Network (ADANet) represents a third method for finding attractors. ADANet introduces a set of trainable reference points in the embedding space, referred to as "anchors," to estimate the speaker assignments in both the training and testing phases. For a mixture containing a certain number of speakers, ADANet considers all possible combinations of a predefined number of anchor points. The distance between the embeddings and the anchor points in each subset is calculated and used to estimate the speaker assignment using a SoftMax function. Attractors are then calculated for each anchor subset based on the estimated speaker assignments. Finally, the subset of attractors that exhibits the minimum in-set similarity (i.e., the largest distance between the attractors within the subset) is selected for mask estimation. This approach aims to address the "center mismatch problem" by not relying on true speaker assignments during training and enabling direct mask generation in both training and testing phases.

Each of these methods for finding attractors presents its own set of advantages and limitations.

Clustering-based methods offer adaptability to varying numbers of speakers but can be computationally intensive during testing. Fixed Attractor points provide computational efficiency but might lack the flexibility to adapt to significant variations in speaker characteristics. Anchored DANet offers a more robust approach by learning attractors in a less supervised manner, potentially improving generalization but introducing more trainable parameters. The choice of method often depends on the specific requirements of the application, balancing factors such as performance, computational resources, and the need for adaptability.

V. SYSTEM ANALYSIS AND POTENTIAL ENHANCEMENTS

The system proposed for speaker-independent speech separation using Deep Attractor Networks typically comprises two main modules: a training module and a testing module. In the training module, the system processes a dataset of voice recordings, extracts relevant features, and trains a DNN model for speech separation. This model learns to map mixed speech signals to an embedding space and to identify Attractor points representing individual speakers. The performance of the trained model is then evaluated using a separate validation dataset. In the testing module, users can input new, unseen mixed speech data, and the system performs feature extraction on this input.

Subsequently, the trained DNN model and one of the Attractor finding methods are used to perform target separation, and the separated speech signals are reconstructed from the estimated time-frequency masks. The final output consists of the isolated voices of the different speakers present in the original mixture.

Feature extraction plays a crucial role in the performance of the DANet model. A common technique employed is the Short-Time Fourier Transform (STFT), which converts the time-domain audio signal into a time-frequency representation known as a spectrogram. The log magnitude spectrogram is often used as the input features to the network. The choice of parameters for the STFT, such as the window length and hop size, can significantly impact the temporal and frequency resolution of the resulting spectrogram, thereby affecting the information available to the DANet model. Selecting appropriate parameters is crucial for capturing the essential characteristics of speech while minimizing noise and irrelevant information.

The field of deep learning is continuously evolving, and several potential enhancements could be explored to further improve the performance and efficiency of DANet. One promising direction involves investigating the use of different neural network architectures for generating the embeddings. For instance, Convolutional Neural Networks (CNNs) have shown success in capturing local spectral features, while Transformer networks, with their attention mechanisms, excel at modeling long-range dependencies. Exploring the integration of these architectures within the DANet framework could lead to improved embedding representations. Another potential enhancement involves using Bidirectional Gated Recurrent Units (BGRUs) instead of BLSTMs, which might offer comparable performance with reduced model complexity. Furthermore, alternative clustering algorithms, such as Gaussian Mixture Models (GMMs), could be investigated as a substitute for K-means in the clustering-based Attractor formation method. The development of time-domain DANets, which directly process the raw audio waveform, represents another interesting avenue for research. Finally, for scenarios involving multiple microphones, exploring the integration of multi-channel information with the DANet framework could leverage spatial cues to further enhance speech separation performance.

VI. EVALUATION AND RESULTS

The effectiveness of speaker-independent speech separation models, including DANet, is typically evaluated using standardized datasets. One of the most widely used datasets in this domain is the Wall Street Journal (WSJ0) corpus and its derived mixtures, such as WSJ0-2mix (two-speaker mixtures) and WSJ0-3mix (three-speaker mixtures).¹ These datasets provide a controlled environment for assessing the performance of different models under specific conditions, such as the number of speakers and the signal-to-noise ratio (SNR) of the mixture.

The performance of speech separation models is commonly quantified using objective evaluation metrics. Two widely adopted metrics are the Signal-to-Distortion Ratio (SDR) and the Scale-Invariant SDR (SI-SDR). SDR measures the overall quality of the separated speech by comparing it to the original clean speech, taking into account both interference and artifacts. SI-SDR is a more recently introduced metric that addresses some of the limitations of SDR, particularly its sensitivity to scaling differences between the estimated and reference signals. Higher SDR and SI-SDR values indicate better separation performance, signifying that the model has effectively isolated the target speech while minimizing unwanted noise and distortions.

Research on Deep Attractor Networks has demonstrated promising results on benchmark datasets. The initial work on DANet reported comparable or even better performance on the WSJ0 dataset compared to other state-of-the-art deep learning methods for speech separation.¹ Subsequent studies have further validated the effectiveness of DANet and its variants, reporting significant SI-SDR improvements on WSJ0-2mix and WSJ0-3mix.¹⁴ For instance, some reported SI-SDR_i values for DANet on WSJ0-2mix are in the range of 20-22 dB, indicating a substantial improvement over the mixed input.

Model	Dataset	Metric	Value (dB)	Reference Snippet(s)
Deep Attractor Network (DANet)	WSJ 0-2 mix	SI-S DRi	(Find value s)	1
Deep Attractor Network (DANet)	WSJ 0-3 mix	SI-S DRi	(Find value s)	1

Conv -Tas Net	WSJ 0-2 mix	SI-S DRi	(Find value s)	14
SepFormer	WSJ 0-2 mix	SI-S DRi	22.4	36
Gated Dual Path RNN	WSJ 0-2 mix	SI-S DRi	20.12	36
Gated Dual Path RNN	WSJ 0-3 mix	SI-S DRi	16.85	36
SepTDA (L=12)	WSJ 0-2 mix	SI-S DRi	24.0	30
SepTDA (L=12)	WSJ 0-3 mix	SI-S DRi	23.7	30

Note: The specific SI-SDRi values for the original DANet model on WSJ0-2mix and WSJ0-3mix need to be extracted from the research snippets for a complete table. The table includes some state-of-the-art models for comparison.

VII. APPLICATIONS OF SPEAKER-INDEPENDENT SPEECH SEPARATION

Speaker-independent speech separation holds significant potential for enhancing the performance of Automatic Speech Recognition (ASR) systems, particularly in challenging acoustic environments characterized by noise and the presence of multiple speakers. By effectively isolating individual speakers from a mixed audio signal, speech separation can serve as a crucial preprocessing step, providing cleaner audio input to the ASR system and thereby improving its robustness and accuracy. Studies have reported notable reductions in word error rate (WER) in ASR systems when employing speech separation techniques as a front-end. In the domain of speaker diarization, which involves identifying "who spoke when" in a multi-speaker recording, speaker-independent speech separation plays a vital role. The ability to separate the audio streams of individual speakers is a fundamental prerequisite for accurately determining speaker identities and their corresponding speaking turns. Continuous speech separation techniques are often employed for speaker segmentation, providing the isolated audio segments needed for subsequent speaker identification. Beyond ASR and speaker diarization, speaker-independent speech separation has numerous applications in broader audio processing and telecommunications. In telecommunications, it can significantly improve the clarity of audio communication in scenarios involving multiple participants, such as phone calls and video conferences. It also finds utility in multimedia analysis, where audio signals from various sources are intentionally mixed. Furthermore, this technology holds promise for enhancing hearing aids and assistive listening devices, enabling users to focus on a specific speaker in noisy environments. While the provided abstract broadly mentions applications in medical signal processing, further research could explore specific use cases, such as analyzing patient-doctor conversations or monitoring multiple speakers in healthcare settings.

VIII. CONCLUSION AND FUTURE DIRECTIONS

Deep Attractor Networks have demonstrated significant potential in addressing the challenging problem of speaker-independent speech separation. By projecting mixed speech into a high-dimensional embedding space and utilizing Attractor points to represent individual speakers, DANet offers an effective framework for isolating voices without prior knowledge of the speakers. The network's ability to handle the permutation and output dimension problems inherent in this task highlights its robustness and adaptability. Despite the advancements achieved by DANet, several challenges remain. The performance of current models can still be limited in highly noisy or reverberant acoustic environments, and separating mixtures with a large number of overlapping speakers continues to be a significant hurdle. Further research is needed to enhance the robustness and generalization capabilities of DANet across a wider range of acoustic conditions and speaker counts. Future research directions could explore the integration of novel network architectures, training strategies, and loss functions within the DANet framework. Investigating the use of attention mechanisms or Transformer networks for embedding generation could potentially improve the model's ability to capture complex speech dynamics. Continued efforts are also needed to enhance the handling of unknown numbers of speakers and to improve performance in extreme acoustic conditions. Finally, developing more efficient and lightweight DANet models is crucial for enabling real-time applications on resource-constrained devices. By pursuing these research avenues, the field of speaker-independent speech separation using Deep Attractor Networks can continue to advance, leading to more robust and versatile speech processing technologies.

REFERENCES

- [1] P. Rasane, H. Bhujbal, O. Dhore, M. Jagdale, and S. Sonkamble, "Speaker-Independent Speech Separation with Deep Attractor Network," *J. Emerg. Technol. Innov. Res. (JETIR)*, vol. 11, no. 4, pp. 1–7, 2024.
- [2] X.-L. Zhang and D. Wang, "A deep ensemble learning method for monaural speech separation," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 24, no. 5, pp. 967–977, 2016.
- [3] J. Du, Y. Tu, L.-R. Dai, and C.-H. Lee, "A regression approach to single-channel speech separation via high-resolution deep neural networks," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 24, no. 8, pp. 1424–1437, 2016.
- [4] C. Lea, M. D. Flynn, R. Vidal, A. Reiter, and G. D. Hager, "Temporal convolutional networks for action segmentation and detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 1003–1012.
- [5] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *arXiv preprint arXiv:1803.01271*, 2018.
- [6] Z. Wang, J. Le Roux, and J. R. Hershey, "Multi-channel deep clustering: Discriminative spectral and spatial embeddings for speaker-independent speech separation," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, 2018, pp. 1–5.
- [7] M. Maciejewski, G. Wichern, E. McQuinn, and J. L. Roux, "WHAMR!: Noisy and reverberant single-channel speech separation," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, 2020, pp. 696–700.
- [8] D. Michelsanti and Z.-H. Tan, "An overview of deep-learning-based audio-visual speech enhancement and separation," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 29, pp. 1–18, 2021.
- [9] S. Ansari, K. A. Alnajjaj, T. Khater, S. Mahmoud, and A. Hussain, "A robust hybrid neural network architecture for blind source separation of speech signals exploiting deep learning," *IEEE Access*, doi: 10.1109/ACCESS.2023.3313972.
- [10] D. Wang and J. Chen, "Supervised speech separation based on deep learning: An overview," *IEEE/ACM Trans. Audio Speech Lang. Process.*, doi: 10.1109/TASLP.2018.2842159.
- [11] Z. Wang, J. Le Roux, and J. R. Hershey, "Deep clustering with convolutional neural networks for large-scale speech separation," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, 2017, pp. 4846–4850.
- [12] Y. Luo, Z. Chen, and N. Mesgarani, "Speaker-independent speech separation with deep attractor network," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 26, no. 4, pp. 786–795, 2018.
- [13] H. Chen and P. Zhang, "Exploring the time-domain deep attractor network with two-stream architectures in a reverberant environment," *arXiv preprint arXiv:2007.00272*, 2020.
- [14] F. Jiang and Z. Duan, "Speaker attractor network: Generalizing speech separation to unseen numbers of sources," *IEEE Signal Process. Lett.*, vol. 27, pp. 1859–1863, 2020.
- [15] D. Michelsanti and Z.-H. Tan, "Online deep attractor network for real-time speaker-independent speech separation," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, 2020, pp. 666–670.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)