

¹ C Anna Palagan² N.Muthuvairavan Pillai³ K.Balamurugan⁴ T.A.Sundaravadivel

IMAR in Frequency Domain for Separating Speech Signals by STFT



Abstract: - Blind Source Separation is a statistical method for isolating signals coming from unidentified sources that have been picked up by a variety of sensors. In the fields of biological signal analysis and speech analysis technology, it is an active researcher. The quantity of noise that is present in the signal throughout the speech signal separation process is a disadvantage. This noise affects the signal that has been separated, and it is often loud melodic noise. For the purpose of separating speech signals from BSS, the suggested method makes use of a network that is configured with the parameters of Instantaneous Mixing in addition to an Auto Regressive model and a maximum-likelihood function. In order to separate the individual voice signal from the input mixed signals, the Back Propagation Network is put into operation. The Short Time Fourier Transform is then utilized for the purpose of dividing the spoken stream into extremely precise time periods. Next, Maximum-likelihood method is implemented to determine the optimum values of the IMAR model's W and G parameters. This step is done after the ideal values have been estimated. The suggested model was tested using a combination of voice signals and microphone signals, and the results demonstrated that it provided superior execution in comparison to other algorithms that are currently in use.

Keywords: Blind Source Separation, Maximum-likelihood function, Back Propagation Neural Network (BPN), Separation matrices, Instantaneous Mixing plus Auto Regressive.

I. INTRODUCTION

BSS's primary goal was to separate the framework's input signals from their mixed saw on the sensors. A consistent understanding of the predominant noise sources with respect to their contributions to the general clamor levels provides vital data for commotion control applications. For example, sonar and radar signal preparation, remote correspondence, geophysical examinations, biomedical flag management, discourse and image management and machine obligation assurance have all received extensive consideration for the BSS enhancement.

Single signals (sources) must be distinguished from several speakers' concurrent accounts (collectively known as mixing). In this signal partition issue, the sources and blending framework are referred to as the individual signals and commitments in the future blends. Source detachment is no stranger to higher-order metrics. These approaches are often linked to automated signals, which themselves belong to a different quantifiable class than spoken words. To resolve a mixed-drink party issue, for example, BSS can be used in a resonant manner. Currently, signal mixing and signal partitioning are the two approaches utilized to evaluate modular parameters with respect to their signal capabilities. There is a possibility of using BSS processes, which can divide signals, to examine the separated signal.

Both simultaneous and one-by-one extractions of the source signals are possible utilizing the IMAR model, which separates numerous signals into their individual components. In the literature, there are several systems for blind source separation that are mostly categorized based on classification techniques like BPN. Using STFT, the signal is segmented into smaller time intervals.

II. LITERATURE REVIEW

This system, which Iván Gómez has investigated, combines BSS with transmissibility-based approaches. Modular data from the recovered source signals was evaluated using BSS systems in conjunction with transmissibility-based approaches. Another technique was devised to describe a transmissibility work in order to achieve this combination. Recommendations on how to improve transmissibility were predicated on how well

¹Professor, Department of Electronics and Communication Engineering, Saveetha Engineering College, Chennai (T.N), India

²Associate Professor, Department of Computer Science and Business Systems, R.M.D. Engineering College, Chennai (T.N), India

³Professor, Department of Electronics and Communication Engineering, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Saveetha University, Chennai (T.N), India

⁴Assistant Professor, Department of Mechanical Engineering, R.M.D. Engineering College, Chennai (T.N), India
annapalagan@saveetha.ac.in¹, muthu.csbs@rmd.ac.in², bala237115@gmail.com³, tas.snh@rmd.ac.in⁴

*Correspondence Author: annapalagan@saveetha.ac.in

Copyright © JES 2024 on-line : journal.esrgroups.org

mixed-signal PSDs could be connected to single-source PSDs. Under these circumstances, it was necessary to prepare the numerical reactions of a truss structure to see how well it could discern modular parameters even with high levels of additional noise and clearly split modes.

In this situation, Bin Dong et al. 2016 examined two novel findings. It was initially established that "virtual" sources produced by PCA coincide with real ones; it demands that the sources of intrigue must be fragmented and spatially orthogonal. Non-covering bolster sets were encountered in a specific case of this scenario. With this information in mind, an orthogonal statistical and spatial rule was dismantled in order to systematically distinguish garbled sound sources transmitted from disjointed locations. Acoustic imaging techniques such as pillar framing or acoustical holography can easily incorporate this parameter into their algorithms for distinguishing sound sources of various origins.

Xingjie Wu proposed the use of MRI scans of the human cerebrum in light of acknowledged relationship examination. Group-level correlation inquiry was common in the human mind fMRI investigation. According to Canonical Correlation Analysis, all significant genuine signals are auto-connected, while white noise is not. White noise should be ignored in most cases. BSS-CCA was more casual than Independent Component Analysis because it only required that the second-order measurement be zero. A group BSS-CCA differentiating proof was found to be superior than a group ICA distinguishing proof for "sources" that were mostly covered by space.

It has recently been proposed for Blind Source Separation that Wei Zhao's group of criteria, collectively referred to as "referenced-based," consist of the cross-insights or cross-cumulants between evaluated outputs and reference signals. The associated quadratic optimization techniques for these contrast capacities have an appealing factor in the same way: the searched parameters are quadratic. With this reference-based strategy, a comparable contrast task is created by familiarizing the reference signals with negentropy, which then informs a new fast fixed-point method.

In the end, Abdelmalek Kouadri had proposed using Blind Source Separation Filters to identify and separate defects in a three-tank pilot plant. Blind recognizable evidence, rather than a statistical model, is used in this technique, which makes it extremely useful. The Independent Component Analysis (ICA) was used by each BSSF to extract signals from the procedure under consideration, based on the belief that the eliminated sources were statistically independent. Skewness is also important to keep in mind while diagnosing low-amplitude errors, as it has a high degree of sensitivity.

III. PROPOSED MODELE

The suggested method uses the IMAR model to split the signal and STFT to segment the frequency. Experimentation setups for DIR, SIR, and accuracy analysis are also provided. Separation of individual speech signals from a mixture of input signals is also accomplished through the use of BPN.

Fourier Transform in the Shortest Time (STFT) Non-stationary signals are well-suited to STFT's signal dispensation method. STFT measures the Fourier change of each part of the signal after dividing it into rigid time intervals. Scheduling systems for dealing with discourse begins with thinking about how to display speech in the STFT domain. Gaussian arbitrary factors are represented as zero-mean factually free variables in the Gaussian model of discourse motion development. It encourages the STFT space to develop a scientifically tractable structure for useful discourse upgrade computations. There are places where the Gaussian estimate can be extremely inaccurate, though. It includes the time frame index ($t = 0, 1, 2, \text{etc.}$).

$$X_{ik} = \sum_{m=0}^{N-1} x(n + tM)h(n)e^{-j(2\pi/k)nk} \dots\dots\dots (1)$$

and the frequency-bin index ($k = 0, 1, 2, \text{etc., } K-1$), as well as an analysis window ($h(n)$) of size ($K-1$) and a framing step (M) of size ($K-1$).

3.1. Back Propagation Neural Network (BPN)

For data and framework handling, speech signal segmentation and the recovery of distinct signals from a mixed signal is a testing ground. They've made it such that even if the mixing matrix is improperly shaped, they can nevertheless separate signals that are incredibly incapable or substantially scaled. The steepest plunge approach is used to refresh loads in the BPN, a multilayer controlled neural system. Figure 1 shows its design, which includes an information layer with nodes beginning with 'I,' an output layer with nodes beginning with 'K,' and a concealed layer with nodes beginning with 'J.' Hidden layer neurons each have their own information loads and the output of a neuron in one layer is sent to all neurons in the next layer.

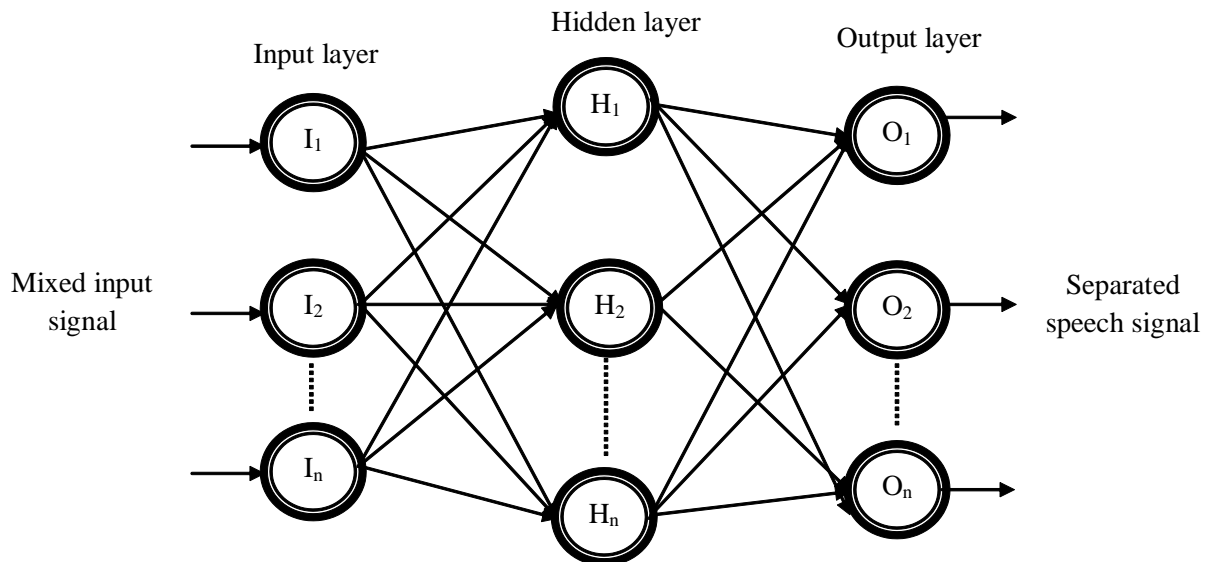


Fig: 1 Architecture of BPN

The mixed input is ignored by the input layer. It just sends the information tests to the buried layer's neurons. The sigmoid capacity of each hidden layer neuron is applied to its net input values in order to obtain the neuron's output. Each hidden neuron output feeds the yield layer's whole population of neurons. When neurons in the yield layer first calculate the net input value, they then apply nonlinear capacity to the net contribution and generate an output value It is a signal that separates speech from music.

3.2. Separation of speech and music signals

Based on past data on supply signal qualities, BSS techniques seek to achieve this goal. Convolutively mixing sounds from N sources with $m = 1...N$, according to sound intermixture physics. The recorded mixed signals at M sensors are designated by

$$b_i^t = \sum_{m=1}^N \sum_{e=1}^L F_{is}(e) S_m(t-e) \dots\dots\dots (2)$$

$10^3 - 10^4$ taps (each regulator's last second where ever Fs is that the sampling frequency) in a passing ordinary space is that the separate Green's perform of the world, jointly known as the Room Impulse Response (RIR). Given the following equation, Discrete Fourier Transform (DFT).

$$C_i(g, \delta) \approx \sum_{m=1}^N S_{is}(g) S_m(g, \delta) \dots\dots\dots (3)$$

Third, where is the unit DFT and is it used to choose frames? Due to the lack of regularity in d of DFT, the convolution does not exactly recapitulate the natural product.

The speech detachment of an inner area downside, the instant mix model, may even be proper because of the engendering delays, which are immaterial. The mixing of delayed and convolved sources is likely to occur in real-world situations because of considerable time-deferrals, which can occur in a variety of ways. Isolating speakers in space and implementing the instructional principle are the focus of this chapter. A feed forward system supported by the recurrence space design and pure polynomial channel arithmetic was used to implement partner casual technique. This pre-processing step was sufficient for signal separation under the most consolidated conditions. An audio recording was efficiently separated from the foundation value 'm' using these techniques.

3.3. IMAR model to generate microphone signals

BSS is improved by LTI filter and permutation in the IMAR model for generating microphone signals: The time-domain BSS progress is used in this instance. By using a combination of Associate in Nursing IM input signals and an IS output separation filter price on the newly discovered signal vector O, it is possible to estimate a blind source signal within the range of source signal vectors baccalaureate (n) (n). The order of the separation filter in the time domain BSS progress for extrication sound mixtures is based on a price that exceeds the reverberation time. For lengthy

reverberation times, the separation filter's order grows horrendously huge. Because the rate of convergence is low, the cost of calculation is quite high. A separation matrix is applied to the discovered spectral part vector to complete the estimation of supply spectral part vector using the frequency domain BSS technique.

When an input signal is applied to the observed signal vector, the output separation filter value is used to evaluate BSS in the structure of the source signal vector. Sound mixes can be extricated in the time domain BSS technique using values that extend beyond the room reverberation duration. The longer the reverberation duration, the greater the order of the separation filter. BSS frequency domain technique is depicted in Fig. 2.

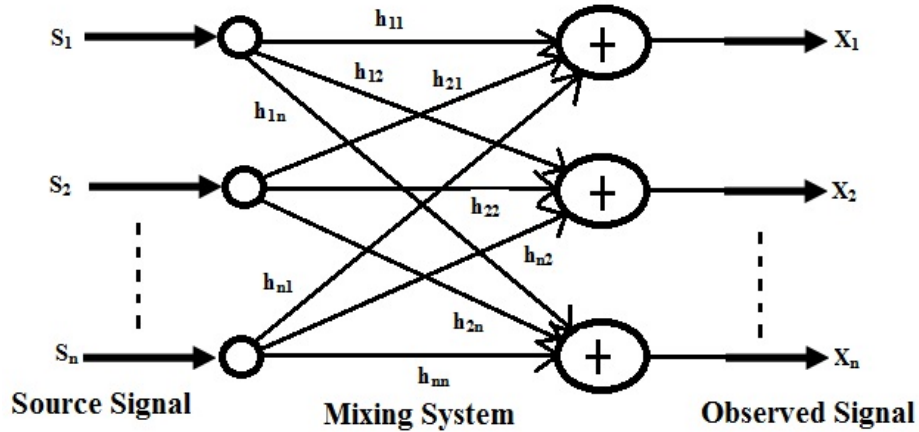


Fig: 2 Structure for frequency domain BSS approach

For any microphone, the prediction filter for the spectral component of the microphone as well as the prediction error that goes along with it are denoted as

$$Q_{r,s,t} = O_{r,s,t} - \sum_{s=K_i}^{L_t+M_i-1} W_{r,s,t} O_{r,s,t}, v \leq n \leq I_M \dots\dots\dots (4)$$

which we are familiar with, assumes that the bin indices 1 are identical for all frequencies in the set of existing values and so equalizes the output spectral component vector

$$Q_{r,s} = \sum_{s=K_i}^{L_t+M_i-1} H^G_{r,s} O_{u-r,s} + P_{r,s} \dots\dots\dots (5)$$

Therefore, any experimental details are needed in order to verify the Equation (5) assumption. As a result, the IMAR model achieves a high level of speech signal separation while still demonstrating that the potential assumption is at least partially correct.

IMAR model creation for the determined spectral element vector is represented by Equations (4) and (5). Equation (4) specified the spectral element of the sound supply as the sum of the spectral components of the various sound signals. Element vector is generated by mixing the remaining elements existing in element and the multichannel AR system with regression or prediction matrices as shown in Equation 4.

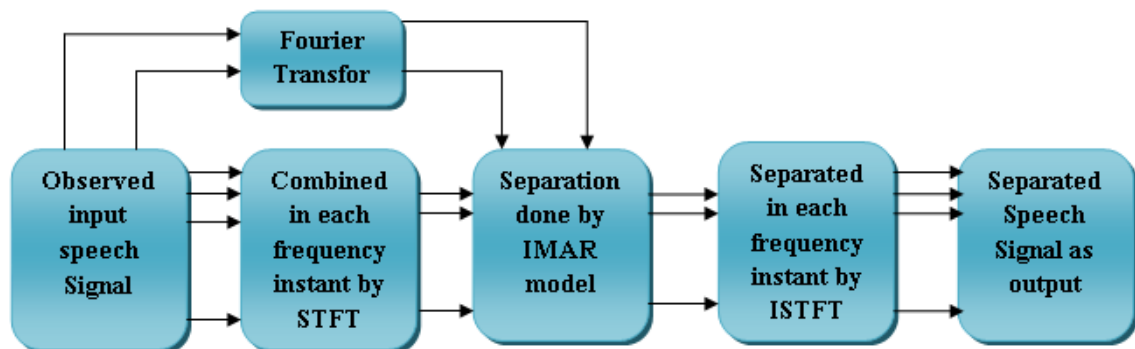


Fig: 3 Representation of IMAR model

The IMAR model is depicted in Fig. 3.

The parameters of the separation and prediction matrices utilized by the IMAR model across the whole frequency range are given by. Separation and prediction matrices are denoted by the following notation 6 and 7.

$$\Phi_w = \{W_v\} \quad 0 \leq k \leq k-1 \quad \dots\dots\dots (6)$$

$$\Phi_G = \{D_{s,v}\}_{L_1 \leq s \leq L_v + S_v - 1; \quad 0 \leq k \leq k-1 \quad \dots\dots\dots (7)$$

The optimum values of IMAR model parameters are estimated via maximum likelihood estimation. Equations (3) and estimate the parameter values of the source spectral component vectors (4).

$$X_{s,v} = \begin{cases} X_v & ; \text{if } s = 0 \\ 0 & ; \text{if } 1 \leq s \leq L_v \\ D_{s,v} X_v & ; \text{else} \end{cases} \quad \dots\dots\dots (8)$$

IMAR model parameters W and G are the best values for IMAR model parameters based on the most accurate probability estimation. Equations (4) and are used to calculate the parameter values for the supply spectral component vectors (5). Using the MIMO filter, the IMAR model is linked to the frequency domain by describing the matrix type for the IMAR model. The matrix type Xs,v expression given by equation (8)

IV. RESULT AND DISCUSSION

Different mixed speech signals were captured in the BSS. The BSS is used to mix the signals, and the IMAR model is used to separate the mixed signals. SIR and DRR charts illustrate how to compare various speech signals to the SIR.

4.1. Experimental setup

The voice signals will be tested using two sources and two microphones in this study. The 145 utterances in the TIMIT corpus were used in this study, which included 38 male speakers and 8 female speakers. The acoustic signals of these utterances are tested at a sampling frequency of 14 KHz, with a bandwidth restricted to 70 Hz to 4 KHz. Using data from 26 male speakers and 5 female speakers, utterances from male and female speakers can be compared to see if there are any similarities. 12 male speakers and 3 female speakers contributed the remaining information. In total, there were 55 male-male, 20 female-female, and 70 male-female utterances pairs. Based on past data on supply signal qualities, BSS techniques seek to achieve this goal. Convolutively mixing sounds from N sources with $m = 1 \dots N$, according to sound intermixture physics. The recorded mixed signals at M sensors are designated by

Table 4.1. Parameters using for Experimental verification

Parameters	Value
Number of sources	03
Number of Microphones	03
Number of Male Speakers	38
Number of Female Speakers	08
Number of Utterances	145
Sampled Frequency	1.4 KHz
Bandwidth	70 Hz to 4 KHz

These signals can then be combined with the room impulse responses obtained by microphone in the varechoic chamber to generate signals that can be picked up by speakers. The experimental setup should perform SIR and DRR analyses after each trial. Consider the following equation to explain the source of the microphones (6).

$$O_{IS,KM}(n) = \sum_{s \geq 0} b_{IS,KM}(s) S_{IS(m-s)} \quad \dots\dots\dots (9)$$

From the source speech signal to the microphone, an impulse response to the room is provided by. The most prominent part of the source voice signal is used to calculate the microphone's index value

$$mic(I_s) = \arg \max \{SIR_{IS,KM}\} \dots\dots\dots (10)$$

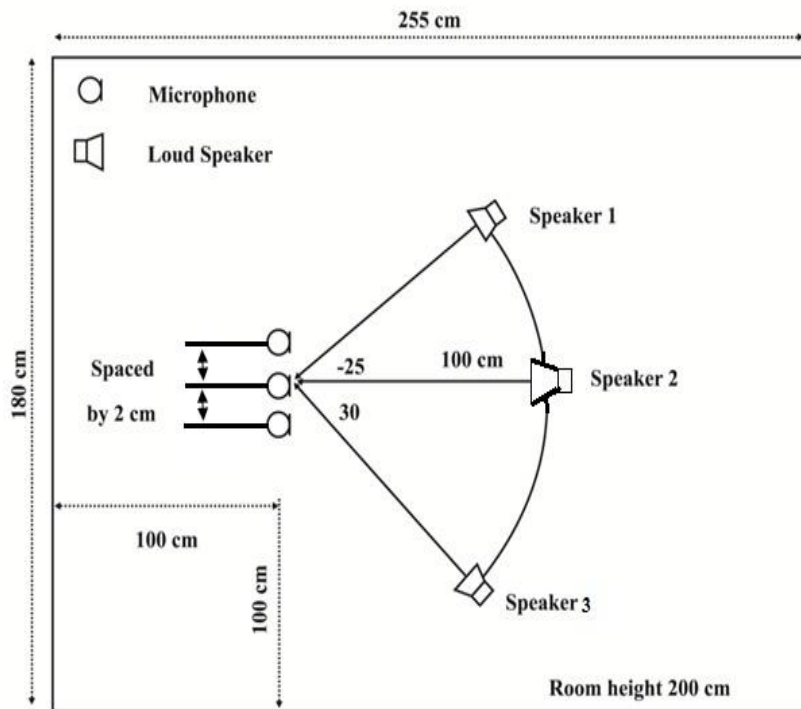


Fig: 4 Experimental Setup

The input SIR and DRR values are computed using the source signal as

$$SIR_{IS} = SIR_{IS,mic(IS)} \dots\dots\dots (11)$$

$$DRR_{IS} = 10 \log_{10} \frac{\sum_{r=0}^{N_j-1} O^{I_s,D}_{mic(IS)}(r)^2}{\sum_{r=0}^{N_j-1} O^{IS,R}_{mic(IS)}(r)^2} \dots\dots\dots(12)$$

Reverberation components (12) are derived from the experimental and estimated reverberation components accordingly. DRR can be characterized using the following experimental data in this way.

$$O_{KM}^{IS,D}(r) = \sum_{s=\Delta+1}^{\Delta} b_{IS,Km}(s) S_{IS}(m-s) \dots\dots\dots (13)$$

Output from microphones is depicted in Figure 4.

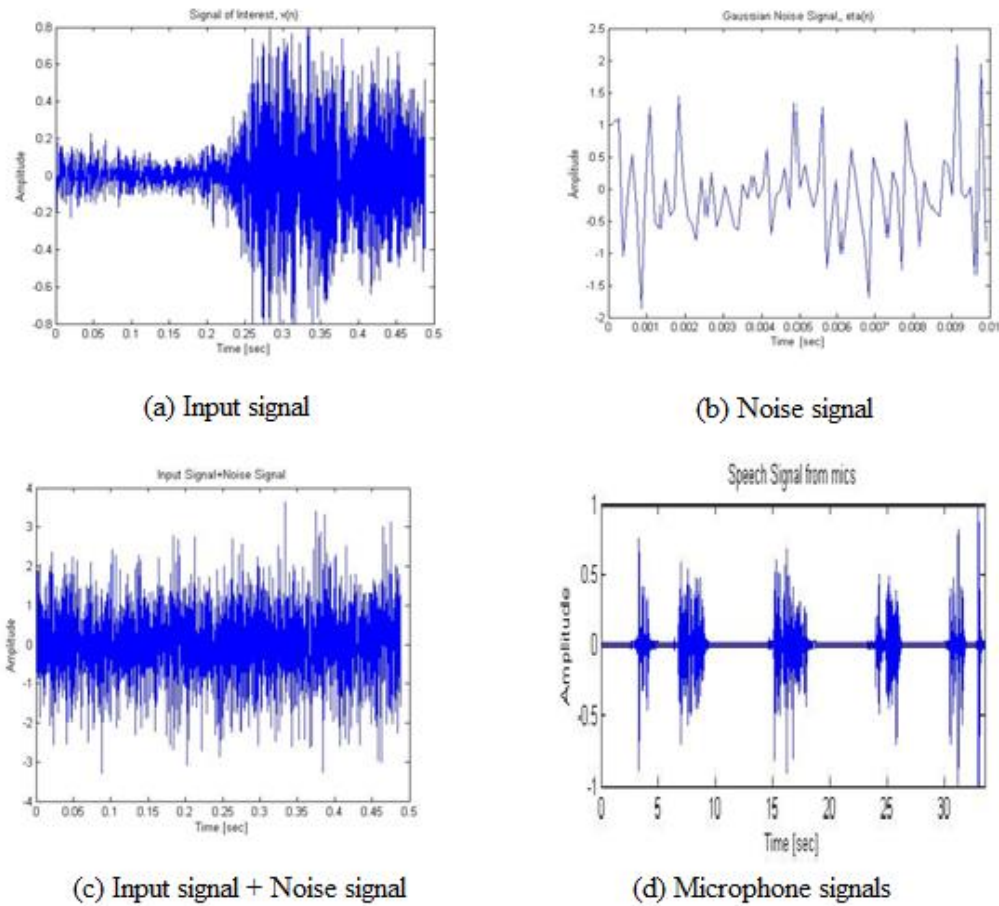


Fig: 4 Output of microphone signals

Figure 4 illustrates the output analysis of microphone voice sounds. For example, the input signals are depicted in figure (a) and the noise signals are depicted in figure (b), which we've divided into two independent figures. An explanation of the input and noise signals, as well as the associated STFT, may be seen in Fig. (c). Here, we've analyzed the results of two speeches. Figure (d) illustrates how the voice signal microphone signals are separated.

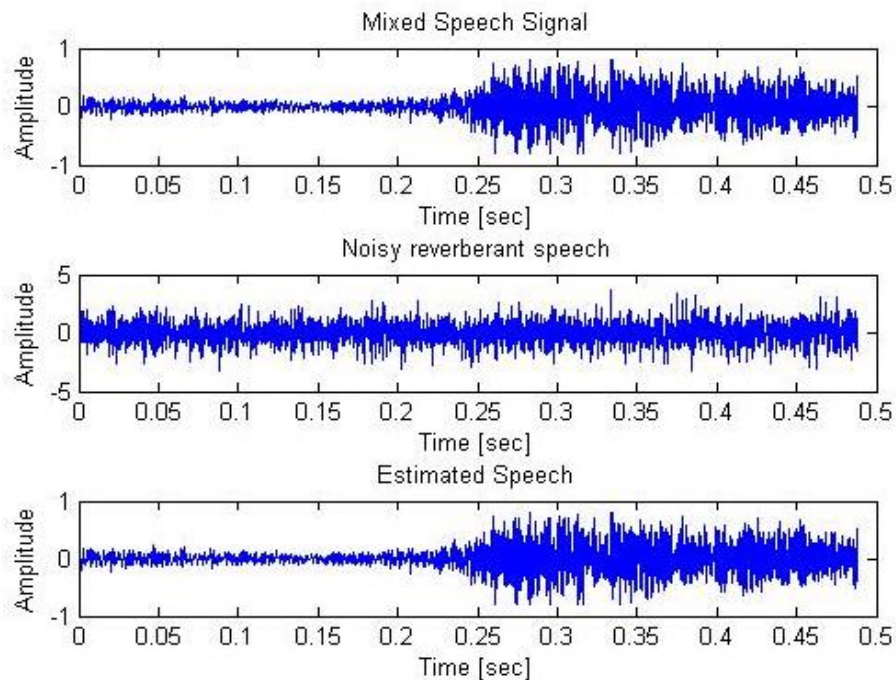


Fig:5 Predicted speech signal

For variously separable speech signals, Figure 5 illustrates the rapid blending of the samples using the IMAR model as a reference point and the separated speech signals of the blind sources. We have separated the noisy reverberant signal from the mixed speech signal using the IMAR model and then used it to separate the estimated speech signal from the mixed signal.

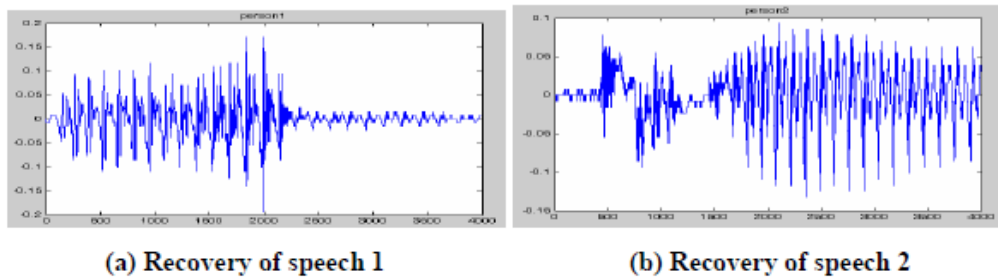


Fig: 6 Recovery of speech signal

Figure 6 depicts the results of applying the IMAR model to recover speech and speech 2. This model effectively separates the voice signal from noise. Because they keep iterating, BPN's separation is of poor quality. In terms of average SIR, the new IMAR model outperforms the old one in the frequency domain when it comes to blind source separation under both reverberation situations.

4.2. Comparison analysis for SIR

An example of SIR analysis using multiple techniques may be seen in Figure 7. We looked at the reverberation time at 0.3 and 0.5 seconds in this case. ICA, HMMs, RNN, and BPN are the algorithms being compared. At a reverberation time of 0.3 seconds, the ICA and HMMs gain 2.5 and 3 dB, respectively. In comparison to existing methods, the suggested IMAR model provides the highest SIR at 0.5 sec reverberation.

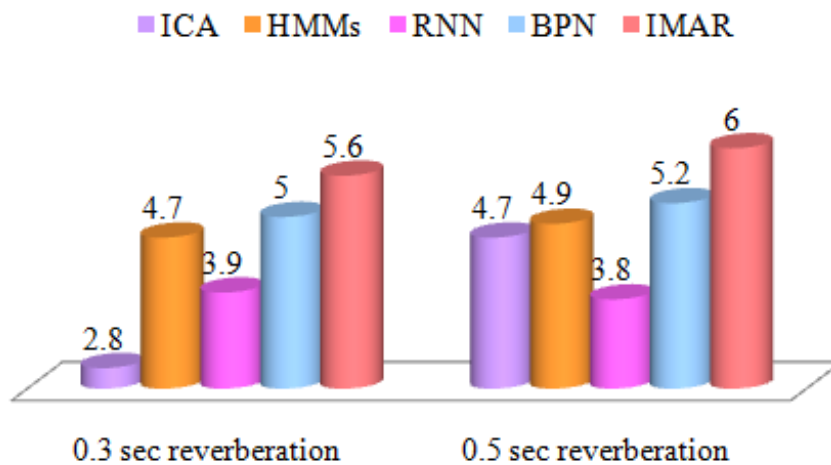


Fig: 7 SIR analysis

4.3. Comparison analysis for DRR

The DRR analysis is shown in Figure 8.

Various reverberation circumstances are depicted in Figure 8 to demonstrate the DRR analysis. x-axis denotes several algorithms, and y-axis represents the DRR in db, as shown in the figure. The proposed IMAR model reduces the DRR. DRR is increased by four times with the ICA method, three times with the HMM method, three times with the RNN method, and two times with the suggested technique at a reverb time of 0.3 seconds. For DRR, the IMAR model offers the lowest value at 0.5 sec reverberation Based on past data on supply signal qualities, BSS techniques seek to achieve this goal. Convolutively mixing sounds from N sources with $m = 1 \dots N$, according to sound intermixture physics. The recorded mixed signals at M sensors are designated by

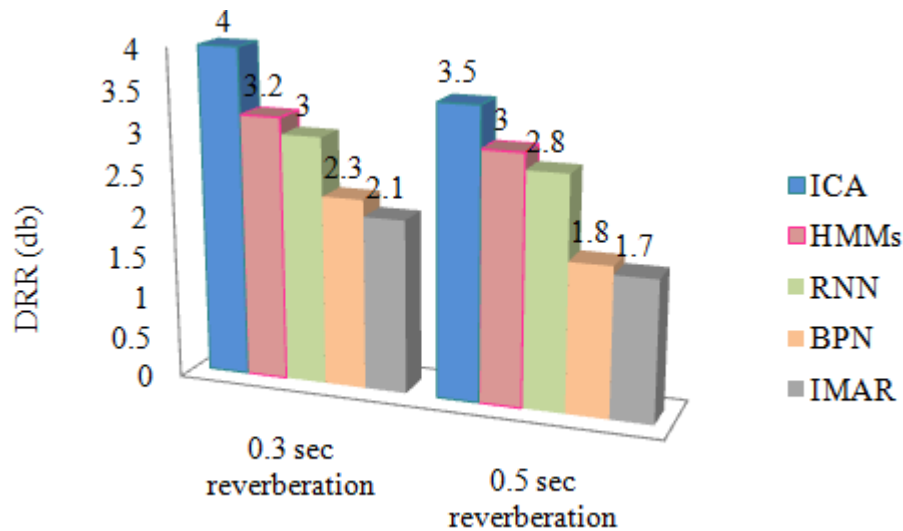


Fig: 8 DRR analysis

4.4. Comparison analysis for recognition accuracy

Figure 9 illustrates the accuracy of various methods. Here, a variety of algorithms are used to evaluate the accuracy of the results. In the 0.3 reverberation range, ICA has 93.2 accuracy, HMMs has 94.2, RNN has 91.23, and IMAR has the best accuracy. The IMAR model provides the most accurate results at 0.5 reverberation.

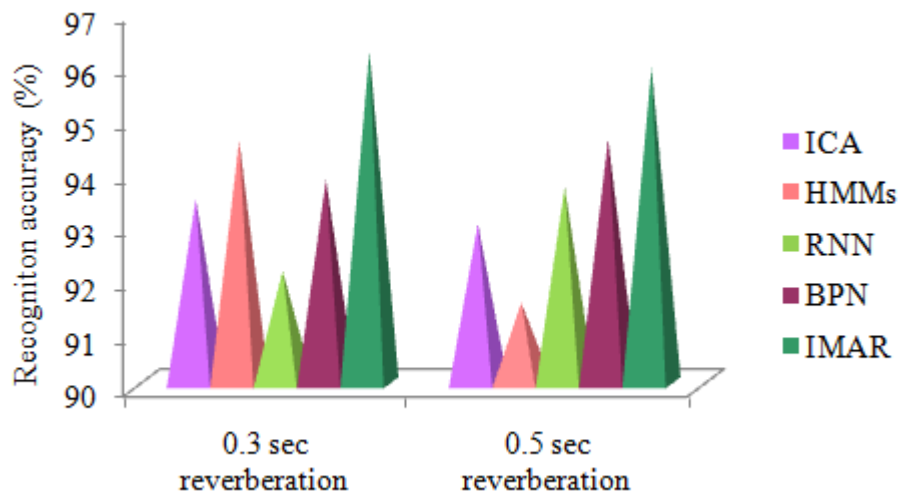


Fig: 9 Accuracy analysis

V. CONCLUSION

Optimized rapid IMAR model and also the most probable performance were used in this work to efficiently separate speech signals from the blind source. The quick IMAR methodology's major possibilities are the optimized speech signal separation and the resulting ability to perform a thought-based blind source separation approach. The SIR rate improves by nearly half a dozen decibels with the current methods. Once the reverberation time was zero, optimizing the IMAR approach yielded a suitable SIR and DIR ratio. Then, the IMAR technique is complete in providing a comprehensive instrument for the processing of mike array signals during a room impulse response.

REFERENCES

- [1] He, Pengju, Tingting She, Wenhui Li, and Weibiao Yuan. "Single channel blind source separation on the instantaneous mixed signal of multiple dynamic sources" *Mechanical Systems and Signal Processing* 113 (2018): 22-35.
- [2] Parimala Gandhi, A., and S. Vijayan. "Upgrading Sparse NMF algorithm for Blind Source Separation through Adaptive Parameterized Hybrid Kernel based approach" *Measurement* (2018).
- [3] Li, Jiong, Hang Zhang, and Pengfei Wang. "Blind separation of temporally correlated noncircular sources using complex matrix joint diagonalization" *Pattern Recognition* 87 (2019): 285-295.

- [4] Nasser Mourad. "Robust smoothing of one-dimensional data with missing and/or outlier values" *IET Signal Processing* (2021), 10.1049/sil2.12033.
- [5] Yang, Xingkai & Zhou, Peng & Zuo, Ming & Tian, Zhigang. "Normalization of gearbox vibration signal for tooth crack diagnosis under variable speed conditions " *Quality and Reliability Engineering International*. 38. 10.1002/qre.3029.
- [6] Sun, Linhui, Keli Xie, Ting Gu, Jia Chen, and Zhen Yang. "Joint dictionary learning using a new optimization method for Single-channel blind source separation" *Speech Communication* 106 (2019): 85-94.
- [7] Malathi P., Suresh G.R., Moorthi M., Shanker N.R. "Speech Enhancement via Smart Larynx of Variable Frequency for Laryngectomee Patient for Tamil Language Syllables Using RADWT Algorithm" *Circuits, Systems, and Signal Processing*, 2019, 38(9), pp. 4202–4228.
- [8] Nuño Ayón, José de Jesús, Julián Sotelo Castañón, and Carlos Alberto López de Alba. "Extracting Low-Frequency Spatio-Temporal Patterns in Ambient Power System Data Using Blind Source Separation" *Electric Power Components and Systems* 46, no. 2 (2018): 230-241.
- [9] Ryuichi Ashino a , Takeshi Mandai b , Akira Morimoto c & Fumio Sasaki," *Journal of Blind source separation of spatiotemporal mixed signals using time frequency analysis*", *Applicable Analysis: An International Journal*, pp.1-34, 2014
- [10] Rajmohan V., Shankar N., Suresh Kumar K. "Analysis of a High Competent Feedforward FFT Architecture", *Advances in Intelligent Systems and Computing*, 2021, 1187, DOI: https://doi.org/10.1007/978-981-15-6014-9_41.
- [11] Araújo, Iván Gómez, Jesús Antonio García Sánchez, and Palle Andersen. "Modal parameter identification based on combining transmissibility functions and blind source separation techniques" *Mechanical Systems and Signal Processing* 105 (2018): 276-293.
- [12] Dong, Bin, Jérôme Antoni, Antonio Pereira, and Walter Kellermann. "Blind separation of incoherent and spatially disjoint sound sources" *Journal of Sound and Vibration* 383 (2016): 414-445.
- [13] Wu, Xingjie, Ling-Li Zeng, Hui Shen, Ming Li, Yun-an Hu, and Dewen Hu. "Blind source separation of functional MRI scans of the human brain based on canonical correlation analysis" *Neurocomputing* 269 (2017): 220-225.
- [14] Deng, C, Wei, Y, Shen, Y. Zhao, W and Li, H. "RSNT-CFASTICA for complex-valued noncircular signals in wireless sensor networks" *KSII Transactions on Internet and Information Systemsthis link is disabled*, 2018, 12(10), pp. 4814–4834.
- [15] Fezai, Radhia & Abodayeh, Kamal & Mansouri, Majdi & Kouadri, Abdelmalek & HARKAT, Mohamed-Faouzi & Nounou, Hazem & Nounou, Mohamed & Messaoud, Hassani. (2020). *Reliable Fault Detection and Diagnosis of Large-Scale Nonlinear Uncertain Systems Using Interval Reduced Kernel PLS*. *IEEE Access*. 10.1109/ACCESS.2020.2989917.