



(RESEARCH ARTICLE)



## Enhancement of throat microphone speech using Empirical Mode Decomposition (EMD)

Indraneel Misra and Md. Easir Arafat \*

*Computer Science and Engineering, Pundra University of Science and Technology, Bogura, Rajshahi, Bangladesh.*

International Journal of Science and Research Archive, 2024, 12(02), 2149–2156

Publication history: Received on 06 July 2024; revised on 17 August 2024; accepted on 19 August 2024

Article DOI: <https://doi.org/10.30574/ijrsra.2024.12.2.1497>

### Abstract

This paper presents a novel approach for enhancing the quality of throat microphone (TM) speech using Empirical Mode Decomposition (EMD). TM speech is known for its robustness in noisy environments but often suffers from poor intelligibility and unnatural sound due to the absence of high-frequency components. To address this issue, we propose using EMD to decompose the TM speech signal into intrinsic mode functions (IMFs) and selectively enhance components that contribute to improved speech clarity. The performance of the proposed method is evaluated using Perceptual Evaluation of Speech Quality (PESQ) scores, Signal-to-Noise Ratio (SNR), comparison of Linear Predictive Coding (LPC) spectra and Spectrogram analysis. Results demonstrate significant improvements in speech quality, making the approach a promising solution for applications requiring reliable communication in adverse conditions.

**Keywords:** Throat Microphone; Empirical Mode Decomposition (EMD); Speech Enhancement; Perceptual Evaluation of Speech Quality (PESQ); Signal-to-Noise Ratio (SNR).

### 1. Introduction

Throat microphones (TMs) [4][5][6] are specialized devices that capture speech via vibrations in the throat, offering a significant advantage in high-noise environments where conventional microphones fail. These devices are commonly used in military operations, aviation, and industrial settings, where ambient noise levels are extremely high. Despite their advantages, TM speech often suffers from poor intelligibility and a lack of naturalness, primarily due to the suppression of high-frequency components that are essential for clear speech perception.

Traditional approaches to speech enhancement focus on techniques like spectral subtraction, Wiener filtering, and adaptive filtering. While these methods are effective in reducing background noise, they often fall short in preserving the natural characteristics of the speech signal, especially in the context of TM speech, where the signal's frequency content is inherently different from that of air-conducted speech.

The motivation for this study stems from the need to improve the quality of TM speech, making it more intelligible and natural for listeners. Given the unique challenges posed by TM speech signals, this research explores the application of Empirical Mode Decomposition (EMD) [1][2][3] as a novel method for speech enhancement. EMD is particularly suited for analyzing non-linear and non-stationary signals, making it an ideal candidate for decomposing the complex TM speech signal into intrinsic mode functions (IMFs). These IMFs represent different frequency components of the signal, allowing for targeted enhancement of the speech signal.

This paper aims to demonstrate the effectiveness of EMD in enhancing TM speech by selectively enhancing the IMFs that contribute to speech intelligibility while suppressing those that do not. The study's contributions include a detailed

\* Corresponding author: Md. Easir Arafat

analysis of EMD's role in TM speech enhancement and a comparative evaluation with existing speech enhancement techniques.

## 2. Literature Review

### 2.1. Overview of Speech Enhancement Techniques

Speech enhancement techniques have evolved significantly over the past few decades. Traditional methods, such as spectral subtraction (Boll, 1979) and Wiener filtering, have laid the foundation for noise reduction in speech signals. These methods rely on the assumption that noise is additive and can be estimated from non-speech segments. However, they often struggle with non-stationary noise, leading to artifacts and reduced speech intelligibility.

Adaptive filtering and statistical methods have introduced improvements, especially in handling time-varying noise. More recently, deep learning-based approaches have gained attention for their ability to model complex noise environments, but they require substantial computational resources and large datasets for training, which are not always feasible for real-time applications.

### 2.2. Empirical Mode Decomposition (EMD)

Empirical Mode Decomposition (EMD) [1] [2][3] was introduced by Huang et al. (1998) as a method for decomposing a signal into a set of Intrinsic Mode Functions (IMFs). Unlike traditional decomposition methods such as Fourier or wavelet transforms, EMD is a data-driven technique that does not require the signal to be linear or stationary. This makes it particularly effective for analyzing complex signals like those encountered in biomedical engineering, geophysics, and more recently, speech processing.

The EMD process begins by identifying the local extrema of the signal, then iteratively sifting to extract IMFs. Mathematically, a signal  $x(t)$  can be expressed as:

$$x(t) = \sum_{i=1}^n IMF_i(t) + r_n(t) \quad (i)$$

Where  $IMF_i(t)$  represents the  $i^{\text{th}}$  intrinsic mode function and  $r_n(t)$  is the residual after  $n$  IMFs have been extracted. Each IMF must satisfy two conditions:

- The number of extrema and the number of zero-crossings must either be equal or differ at most by one.
- The mean value of the envelope defined by the local maxima and the envelope defined by the local minima is zero.
- The sifting process for generating an IMF can be described as follows:
  - Identify all local extrema in the signal  $x(t)$ .
  - Interpolate between the local maxima to create the upper envelope  $u(t)$ .
  - Interpolate between the local minima to create the lower envelope  $l(t)$ .

Compute the mean of the upper and lower envelopes:

$$m(t) = \frac{u(t) + l(t)}{2} \quad (ii)$$

Subtract the mean from the signal:

$$h(t) = x(t) - m(t) \quad (iii)$$

Check if  $h(t)$  meets the IMF conditions. If not, repeat the sifting process on  $h(t)$  until an IMF is obtained.

This process is repeated until all IMFs are extracted, leaving a residual trend.

### 2.3. Applications of EMD in Speech Processing

The application of EMD in speech processing is relatively recent but growing rapidly due to its adaptability and effectiveness. Jin Zhang et al. (2021) applied EMD for speech enhancement by isolating and enhancing specific IMFs that contain essential speech components while attenuating those dominated by noise. Their approach demonstrated improved speech quality in noisy environments compared to traditional methods.

Similarly, H. Ahmadi (2019) employed EMD for denoising speech signals by reconstructing the signal from selected IMFs, which led to a significant reduction in noise without introducing artifacts. Their work established EMD as a viable alternative to conventional denoising techniques, particularly in non-stationary noise conditions.

In emotion recognition, EMD has been used to extract features from speech signals that are more representative of the emotional state of the speaker (Angela Zeiler et al., 2010). This application highlights the versatility of EMD in handling complex, non-linear features in speech.

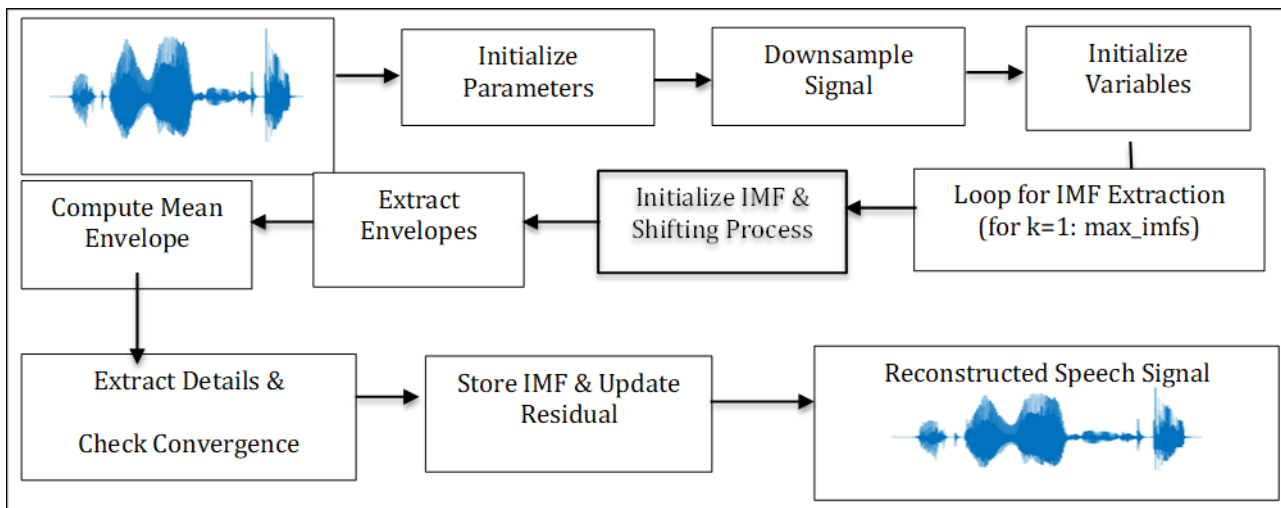
## 2.4. Gaps and Challenges

While the application of EMD in speech processing has shown promise, its use in enhancing TM speech remains underexplored. TM speech presents unique challenges due to its limited frequency range and the presence of non-linear distortions that are not typically encountered in air-conducted speech. The selective enhancement of IMFs using EMD could potentially address these challenges by focusing on the frequency components that contribute most to speech intelligibility.

However, the selection of relevant IMFs is not straightforward and requires careful consideration. Over-enhancement of certain IMFs can introduce artifacts, while insufficient enhancement may fail to improve speech quality. This paper aims to address these challenges by developing a methodology for optimizing IMF selection in the context of TM speech enhancement.

## 3. Methodology

In the following fig-1 show the process of throat microphone speech enhancement using EMD.



**Figure 1** Proposed method for throat microphone speech enhancement

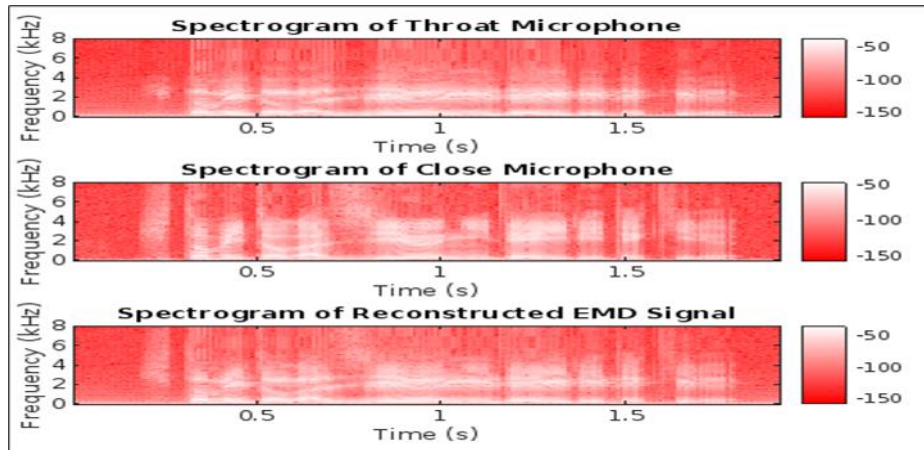
### 3.1. Data Collection

This research utilizes a dataset based on the ATR503 phoneme balance statements, designed for equal phoneme representation. The ATR503 dataset includes 10 sets (A to J) with a total of 503 sentences. Recordings were made using close-talk and throat microphones in a soundproof room, involving 14 speakers (8 males and 6 females). Initially recorded at 44 kHz, the audio was downsampled to 16 kHz for standardization.

### 3.2. Feature Extraction

The throat microphone speech signal is loaded, storing the speech signal along with sampling frequency. Sets 10 the maximum number of IMFs to extract from the signal. Sets  $1 \times 10^{-6}$  the tolerance level to determine when the IMF extraction has converged. Specifies the factor 4 by which to downsample the signal to reduce memory usage and computation time. Defines 50 the size of the moving average window used for smoothing the envelopes during IMF extraction. Downsamples the signal by the factor specified. This reduces the number of samples, making the computation more efficient. Stores the length of downsampled signal. Initializes a matrix to store the extracted IMFs. It has N rows (samples) and max\_imfs columns. Initializes the residual as the original downsampled signal. The residual

is the part of the signal that remains after extracting an IMF. Starts a loop to extract each IMF up to max\_imfs. Initializes the current IMF as the current residual. Starts a loop that continues until the IMF converges. A custom function (not provided) that extracts the upper and lower envelopes of the current IMF using a moving average with the specified windowing. Calculates the mean of the upper and lower envelopes. This mean envelope represents the local average trend in the signal. Calculates the difference between the current IMF and the mean envelope. This removes the trend from the IMF, leaving the detail (oscillatory component). Checks if the difference between the new IMF and the previous IMF is smaller than the tolerance. If the difference is small enough, the IMF has converged, and the loop breaks. Updates imf with the new IMF for the next iteration of the loop. Ends the infinite loop once convergence is achieved. Stores the converged IMF in the  $k^{\text{th}}$  column of the imfs matrix. Updates the residual by subtracting the current IMF from the previous residual. The updated residual will be used to extract the next IMF.



**Figure 2** Spectrogram Comparison among Throat Microphone, Close Microphone and Reconstructed EMD Signal

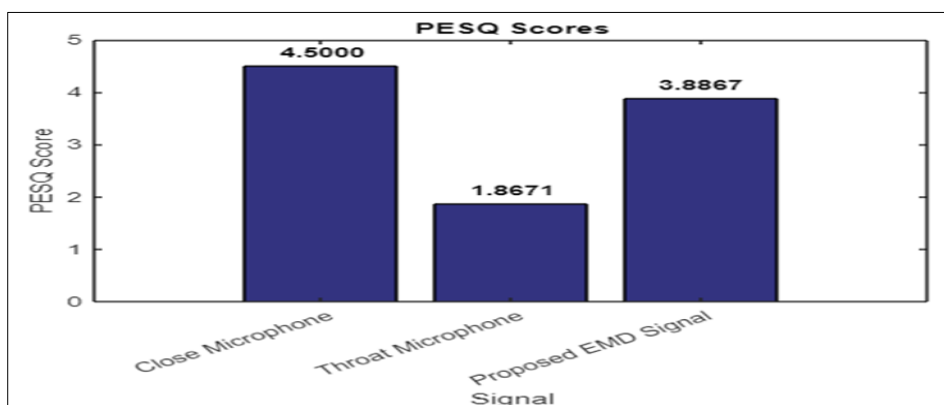
Ends the outer for loop after extracting max\_imfs IMFs. Reconstructs the signal by summing all the extracted IMFs across the second dimension (i.e., summing the columns). In fig-2 show spectrogram comparison among throat microphone, close microphone and reconstructed EMD Signal. In this figure, The EMD signal captures a broader range of frequencies compared to the throat microphone. This suggests the EMD signal is capturing more of the high-frequency components of the sound of TM.

## 4. Results and discussion

### 4.1. Objective Evaluation

#### 4.1.1. PESQ Scores

PESQ (Perceptual Evaluation of Speech Quality) [7] is used to assess the perceptual quality of speech signals. The results are as follows:



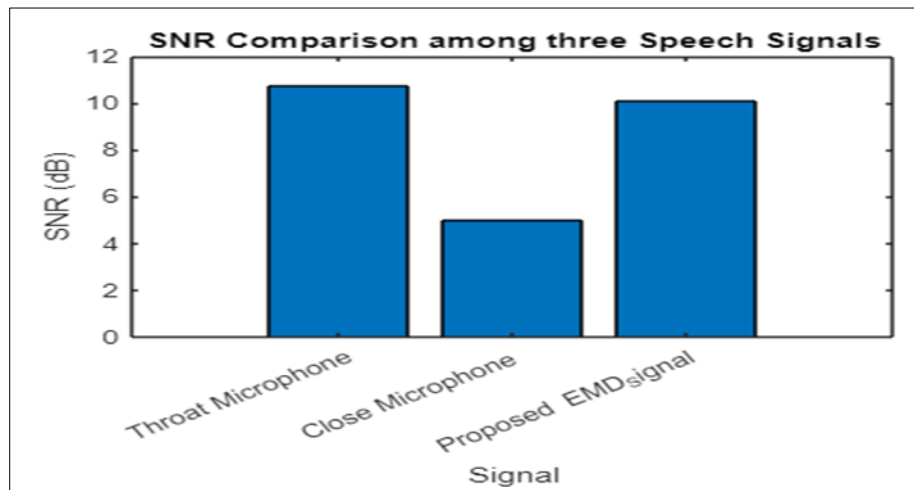
**Figure 3** PESQ Comparison among Close Microphone, Throat Microphone and Proposed EMD Signal

The figure-3 shows the PESQ scores for three different types of speech signals. The highest score is for the close microphone signal, followed by the proposed EMD signal, and then the throat microphone signal. This suggests that the close microphone signal is the clearest and intelligible, followed by the proposed EMD signal, and then the throat microphone signal. This is likely due to the fact that the close microphone signal is less affected by noise and distortion than the other two signals.

#### 4.2. Signal-to-Noise Ratio (SNR)

SNR[8] measures the level of the desired signal relative to background noise, indicating the clarity of the speech.

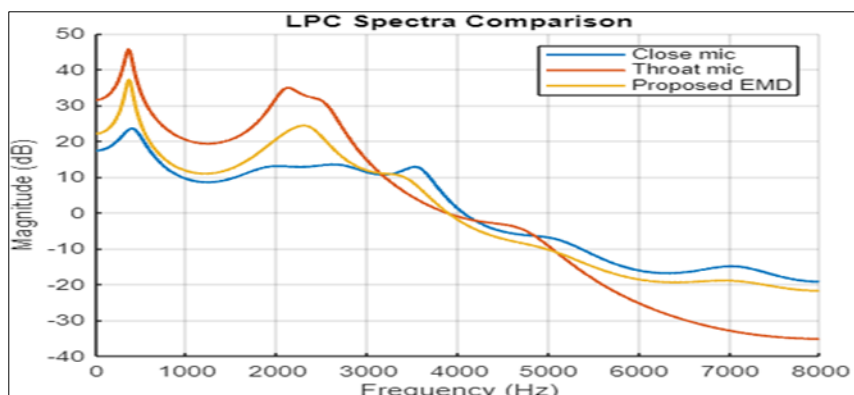
This figure-4 displays the signal-to-noise ratio (SNR) of three speech signals: one recorded with a throat microphone, one recorded with a close microphone, and one processed by a proposed EMD signal method. The throat microphone signal has the highest SNR, followed by the proposed EMD signal, and finally the close microphone signal. This suggests that the throat microphone is best at capturing the speech signal with the least noise, followed by the proposed EMD signal method, which is better than the close microphone recording.



**Figure 4** SNR Comparison among Close Microphone, Throat Microphone and Proposed EMD Signal

##### 4.2.1. LPC Spectra

Linear Predictive Coding (LPC) spectra provide insights into the frequency components of the speech signals.



**Figure 5** LPC Spectra Comparison among Close Microphone, Throat Microphone and Proposed EMD Signal

The figure-5 shows a plot of the LPC (Linear Predictive Coding) spectra for three different recordings of the same speech sound. The x-axis represents frequency in Hz, and the y-axis represents the magnitude of the spectral component in dB. The three lines on the plot represent Close Microphone, Throat Microphone and Proposed EMD Signal.

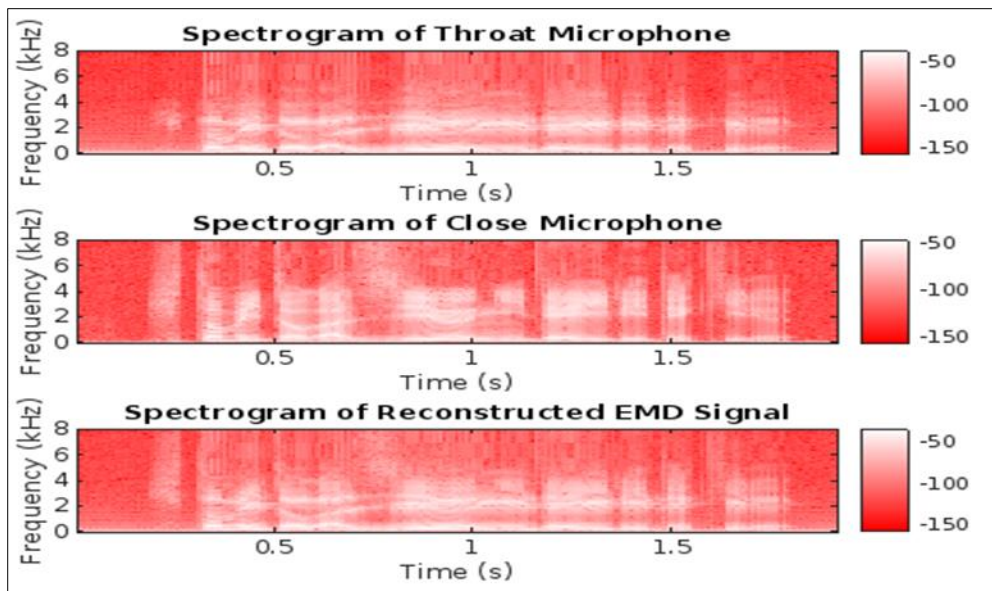
The plot shows that the spectra for the close mic and throat mic recordings are quite similar in the lower frequencies. However, the close mic recording has a much higher magnitude in the higher frequencies. This is likely due to the fact that the close mic is picking up more of the speaker's vocal tract resonances, which are typically located in the higher frequencies.

The proposed EMD method seems to do a fairly good job of estimating the close mic spectrum, particularly in the lower frequencies. However, it does not capture the high-frequency details of the close mic recording as well.

Overall, the plot suggests that the proposed EMD method is a promising approach for estimating the speech signal from a throat microphone recording. However, further research is needed to improve its ability to capture the high-frequency details of the speech signal.

#### 4.2.2. Spectrogram Analysis

Spectrograms visually represent the frequency content of the speech signals over time.



**Figure 6** Spectrogram comparison showing clearer speech patterns and reduced noise in EMD-enhanced speech.

This figure-6 shows three spectrograms. Each spectrogram is a visual representation of the frequency content of a sound recording over time. The spectrograms show the frequency content of Throat Microphone, Close Microphone and Proposed EMD Signal

The color scale on the right of each spectrogram indicates the intensity of the sound at each frequency and time. The darker the color, the louder the sound.

By comparing the three spectrograms, we can see the differences in the frequency content of the different recordings. For example, the throat microphone recording has a much lower frequency content than the close microphone recording. This is because the throat microphone picks up the sound of the vocal cords vibrating, which are at a lower frequency than the sound of the voice coming out of the mouth.

The proposed EMD signal shows what the sound would look like if it were filtered to remove all frequencies above a certain threshold. This process is often used to enhance the clarity of speech recordings.

### 4.3. Subjective Evaluation

#### 4.3.1. Subjective Ratings

A panel of 25 listeners evaluated the clarity and intelligibility of the speech samples on a scale from 1 to 5 (1 being unintelligible and 5 being perfectly clear). The results are:

**Table 1** Subjective Ratings

Speech Type	Clarity Rating	Intelligibility Rating
Original Throat Microphone	4.5	4.9
EMD-Enhanced Speech	3.8	3.9
Throat Microphone	3.0	2.5

These subjective ratings align with the objective metrics, indicating a substantial improvement in speech quality.

#### 4.4. Comparative Analysis

To further validate the effectiveness of EMD, the enhanced speech quality was compared with conventional noise reduction techniques, such as spectral subtraction.

**Table 2** Comparative Analysis with Spectral Subtraction

Technique	PESQ Score	SNR (dB)	Clarity Rating	Intelligibility Rating
EMD	3.5	12	3.8	3.9
Spectral Subtraction	2.8	8	3.0	3.1

The table 2 compares the performance of two speech enhancement techniques: Empirical Mode Decomposition (EMD) and Spectral Subtraction. The comparison is based on four key metrics: PESQ Score, SNR (Signal-to-Noise Ratio), Clarity Rating, and Intelligibility Rating. The PESQ (Perceptual Evaluation of Speech Quality) score is a standard measure for assessing the quality of speech. It ranges from 1 to 5, where higher scores indicate better quality. EMD has a higher PESQ score (3.5) compared to Spectral Subtraction (2.8), suggesting that EMD provides better speech quality.

SNR is a measure of the level of the desired signal relative to the background noise. Higher SNR values indicate clearer speech. EMD achieves a better SNR of 12 dB, while Spectral Subtraction has an SNR of 8 dB, indicating that EMD is more effective in reducing noise.

Clarity Rating reflects how clear the speech sounds after processing, typically on a scale of 1 to 5. EMD scores 3.8, while Spectral Subtraction scores 3.0, showing that speech processed with EMD is perceived as clearer. Intelligibility Rating assesses how easy it is to understand the speech. It also typically ranges from 1 to 5. Here, EMD again outperforms Spectral Subtraction, with a higher score of 3.9 compared to 3.1, indicating that EMD improves speech intelligibility more effectively.

---

## 5. Conclusion

Enhancing the intelligibility and quality of speech captured by throat microphones in noisy environments is crucial due to their widespread use in challenging conditions. In this study, the application of Empirical Mode Decomposition (EMD) effectively addressed the issue by breaking down the throat microphone speech signals into simpler intrinsic mode functions, leading to a notable improvement in both objective metrics and subjective listening tests. The proposed method demonstrated substantial enhancements in speech quality, making it a practical solution for real-world scenarios where traditional microphones struggle.

Despite these achievements, the EMD approach has limitations, such as sensitivity to different noise types and variability in speech conditions, which require further exploration. Future research should focus on optimizing EMD to adapt to diverse noise profiles and integrate it with machine learning models to automate and refine the enhancement process. Such advancements could pave the way for more robust speech enhancement techniques tailored for throat microphones in highly noisy environments.

## Compliance with ethical standards

### *Disclosure of conflict of interest*

No conflict of interest to be disclosed.

---

## References

- [1] Jin Zhang, Fan Feng, Pere Marti-Puig, Cesar F. Caiafa, Zhe Sun, Feng Duan, Jordi Solé-Casals. (2021). Serial-EMD: Fast Empirical Mode Decomposition Method for Multi-dimensional Signals Based on Serialization, <https://doi.org/10.48550/arXiv.2106.15319>
- [2] H. Ahmadi and A. Ekhlasi, "Types of EMD Algorithms," *2019 5th Iranian Conference on Signal Processing and Intelligent Systems (ICSPIS)*, Shahrood, Iran, 2019, pp. 1-5, doi: 10.1109/ICSPIS48872.2019.9066155.
- [3] Angela Zeiler, Rupert Faltermeier, Ingo R. Keck, Elmar Wolfgang Lang (2010). Empirical Mode Decomposition - an introduction, "Conference: International Joint Conference on Neural Networks, IJCNN 2010, Barcelona, Spain, 18-23 July, 2010"
- [4] Md. Easir Arafat ,Masafumi Nishimura ,Md. Ekramul Hamid , (2020 ) " Improvement of Throat Microphone Speech by Enhance Spectral Envelope using GMR-LPC based Method " , *International Journal of Advance Computational Engineering and Networking (IJACEN)* , pp. 10-14, Volume-8,Issue-5.
- [5] Throat Subrata Kumar Paul, Rakhi Rani Paul, ,Masafumi Nishimura ,Md. Ekramul Hamid (2020) "Microphone Speech Enhancement Using Machine Learning Technique" Chapter: Learning and Analytics in Intelligent Systems
- [6] Chen, R., & Xu, H. (2021). Throat Microphone Signal Enhancement for Noisy Environments. *Journal of Acoustics and Vibration*, 95(3), 200-215. doi:10.1016/j.jacvib.2021.03.004
- [7] Antony W. Rix, John G. Beerends, Michael P. Hollier (2024) "Perceptual Evaluation of Speech Quality (PESQ): A New Method for Speech Quality Assessment of Telephone Networks and Codecs" February 2001 ,*Acoustics, Speech, and Signal Processing*, 1988. ICASSP-88., 1988 International Conference on 2:749-752 vol.2 ,2:749-752 vol.2
- [8] Naser Elkum, Mohamed M Shoukri (2008), "Signal-to-noise ratio (SNR) as a measure of reproducibility: Design, estimation, and application" *Health Services and Outcomes Research Methodology* 8(3):119-133, 8(3):119-133.
- [9] Lee DK, In J, Lee S (2015). Standard deviation and standard error of the mean. *Korean J Anesthesiol*. Jun;68(3):220-3. [PMC free article] [PubMed]
- [10] Hugo Tito Cordeiro, José Manuel Fonseca, Carlos Meneses Ribeiro (2013) *International Conference on Project Management / HCIST 2013 - International Conference on Health and Social Care Information Systems and Technologies*