



## LPC-Based Approach for Accurate Voice Activity Detection

Devika<sup>1</sup>, Meenakshi Arora<sup>2</sup>

P.G. Student, Department of CSE, Sat Kabir Institute of Technology and Management, Haryana, India<sup>1</sup>

Assistant Professor, of CSE, Sat Kabir Institute of Technology and Management, Ladrawan, Haryana, India<sup>2</sup>

### ABSTRACT:

This paper presents a novel approach for detecting the presence or absence of voice activity using Linear Predictive Coding (LPC). The proposed method leverages the predictive power of LPC to analyze audio signals and identify voice segments with high accuracy. By extracting LPC coefficients and evaluating the prediction error, the system can effectively distinguish between voiced and unvoiced regions in an audio stream. Experimental results demonstrate that the LPC-based voice activity detection (VAD) method outperforms traditional techniques in terms of both detection accuracy and computational efficiency. This approach is particularly beneficial for applications in speech processing, telecommunications, and audio surveillance.

**Keywords:** Linear Predictive Coding (LPC), Voice Activity Detection (VAD), Audio Signal Processing, Speech Analysis, Prediction Error, Voiced and Unvoiced Regions

### INTRODUCTION:

Voice Activity Detection (VAD) is a crucial component in various speech processing applications, including speech recognition, telecommunications, and audio surveillance. Effective VAD algorithms enable systems to distinguish between speech and non-speech segments, which is essential for enhancing speech intelligibility, reducing bandwidth usage, and improving the performance of downstream processing tasks. Traditional VAD methods have relied on a variety of features, such as Zero-Crossing Rate (ZCR), short-term energy, and spectrogram analysis. However, these approaches often struggle with noisy environments and varying speech dynamics. This paper introduces an advanced VAD method based on Linear Predictive Coding (LPC), which leverages the predictive power of LPC to achieve more accurate and robust voice activity detection.

#### Linear Predictive Coding (LPC)

LPC is a widely used technique in speech processing for representing the spectral envelope of a digital signal of speech in a compressed form, using the information of a linear predictive model. LPC analyzes the speech signal by estimating the formants, removing their effects from the speech signal, and estimating the intensity and frequency of the remaining signal. This process results in a set of LPC coefficients that characterize the speech signal[1].

#### Hamming Window

The application of a Hamming window is a standard practice in signal processing to minimize the edge effects when analyzing short segments of a signal. The Hamming window is defined as:

$$w(n) = 0.54 - 0.46 \left( \frac{2mn}{N-1} \right) \quad (1)$$

where N is the window length. This window helps in reducing spectral leakage when performing Fourier Transform, thereby improving the accuracy of spectral analysis.

#### Spectrogram Analysis

A spectrogram provides a visual representation of the spectrum of frequencies in a signal as it varies with time. By applying the Short-Time Fourier Transform (STFT) to overlapping segments of the signal, we obtain a time-frequency representation that is useful for identifying and analyzing speech components.

### RESEARCH BACKGROUND

#### Traditional Approaches

Traditional VAD methods primarily rely on features such as energy, zero-crossing rate (ZCR), and spectral information. **Energy-Based Methods:** Early VAD systems, such as those described by Rabiner and Sambur [2], utilized short-term energy thresholds to detect speech presence. While simple and computationally efficient, these methods often fail in noisy environments. **Spectral-Based Methods:** Spectral entropy and other frequency-domain features have been employed to enhance VAD performance. For example, Sohn et al. [3] proposed a statistical model-based VAD that uses a Gaussian mixture model to distinguish speech from noise, improving robustness in low signal-to-noise ratio (SNR) conditions.

### Machine Learning Approaches

Recent advances have seen the application of machine learning techniques to VAD, where models are trained to classify speech and non-speech segments. **Support Vector Machines (SVMs):** SVM-based VAD systems utilize feature vectors composed of energy, ZCR, and Mel-frequency cepstral coefficients (MFCCs) to train classifiers. **Random Forests:** Barchiesi et al. [4] explored the use of random forests for VAD, showing that ensemble methods can improve classification accuracy by leveraging the collective decision-making of multiple trees.

### Deep Learning Approaches

Deep learning has significantly impacted VAD, with neural networks providing powerful tools for learning complex patterns in audio data. **Convolutional Neural Networks (CNNs):** CNNs have been applied to VAD by learning spatial features from spectrograms. Kim and Stern [5] showed that CNNs could effectively model time-frequency representations of speech, achieving high accuracy even in noisy conditions. **Recurrent Neural Networks (RNNs):** RNNs, particularly Long Short-Term Memory (LSTM) networks, are well-suited for sequential data. RNN-based VAD systems, such as those by Hughes and Mierle [6], exploit temporal dependencies in audio signals, leading to improved detection performance over traditional methods. **End-to-End Learning:** End-to-end approaches, where models learn directly from raw audio without explicit feature extraction, have gained popularity. Varzandeh, et al. [7] proposed an end-to-end VAD system using a combination of CNNs and RNNs, showing that such models could outperform traditional feature-based methods.

### Hybrid and Ensemble Methods

Hybrid approaches combine multiple techniques to leverage their respective strengths. **Hybrid Models:** Lskar et al. [8] combined deep neural networks (DNNs) with traditional energy-based methods, using the latter to preprocess audio signals before feeding them into the DNN. This hybrid approach improved robustness in diverse acoustic environments. **Ensemble Methods:** Ensemble methods aggregate the outputs of various classifiers to enhance VAD performance. For instance, Wang et al. [9] employed an ensemble of CNN, RNN, and gradient boosting classifiers, achieving state-of-the-art results in VAD tasks.

## PROPOSED METHODOLOGY

This paper proposes an LPC-based VAD approach that incorporates the Hamming window, spectrogram analysis, ZCR, and short-term energy. The proposed Voice Activity Detection (VAD) method leverages Linear Predictive Coding (LPC) for analyzing audio signals and identifying the presence or absence of speech. This method integrates traditional signal processing techniques with LPC to enhance the accuracy and robustness of VAD, especially in noisy environments. Below, we present the detailed steps and mathematical formulation of the proposed LPC-based VAD approach. The method involves the following steps:

1. **Preprocessing:** The audio signal is divided into overlapping frames using a Hamming window to reduce spectral leakage.

□ **Frame Blocking:** The audio signal is divided into overlapping frames to ensure temporal continuity and capture dynamic speech features.

□ **Hamming Windowing:** Each frame is multiplied by a Hamming window to reduce spectral leakage and smooth the signal edges using equation (1).

2. **Feature Extraction:** For each frame, LPC coefficients are extracted, and the prediction error is computed. Additionally, ZCR and short-term energy are calculated.

**LPC Coefficients Calculation:** For each windowed frame, LPC coefficients are computed. LPC approximates the speech signal  $s(n)$  as a linear combination of its past samples:

$$S(n) \approx -\sum_{k=1}^p a_k s(n-k) \quad (2)$$

Here,  $a_k$  are the LPC coefficients, and  $p$  is the order of the LPC analysis. The LPC coefficients are obtained by minimizing the prediction error  $e(n)$ :

$$e(n) = s(n) + \sum_{k=1}^p a_k s(n-k) \quad (3)$$

The coefficients  $a_k$  are computed using the Levinson-Durbin algorithm, which solves the Yule-Walker equations efficiently.

**Prediction Error Energy:** The prediction error energy  $E_p$  for each frame is calculated as:

$$E_p = \sum_{n=0}^{N-1} e^2(n) \quad (4)$$

The prediction error energy is the sum of the squared prediction errors over the frame.

### Zero-Crossing Rate (ZCR) and Short-Term Energy

ZCR is a measure of the rate at which the signal changes sign. It is defined as:

$$ZCR = \frac{1}{T} \sum_{t=1}^{T-1} |sgn(s_t) - sgn(s_{t-1})| \quad (5)$$

where  $sgn(x)$  is the sign function and TTT is the number of samples in the frame. ZCR is useful for distinguishing between voiced and unvoiced speech segments, as voiced segments typically exhibit lower ZCR values.

**Short-term energy** measures the energy of the signal within a short time window and is given by:

$$E = \sum_{n=0}^{N-1} |s(n)|^2 \quad (6)$$

where  $s(n)$  is the speech signal and N is the window length. High short-term energy values generally indicate the presence of speech.

3. **Decision Making:** A combined decision rule based on LPC prediction error, ZCR, and short-term energy is applied to classify each frame as voiced or unvoiced.

$$VAD(n) = \begin{cases} 1 & \text{if } E_p < T_p \text{ and } ZCR < T_z \text{ and } STE > T_e \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

Here,  $T_p, T_z, T_e$  are the thresholds for prediction error energy, ZCR and STE, respectively.

4. **Post-Processing:**

- **Smoothing:** To reduce false detections and ensure temporal consistency, a smoothing algorithm (e.g., median filtering) is applied to the VAD decision sequence.

$$VAD_{smooth}(n) = median(VAD(n-w), \dots, VAD(n), \dots, VAD(n+w)) \quad (8)$$

Where,  $w$  is the window size for smoothing.

## SIMULATION ENVIRONMENT & RESULTS:

We have implemented our method in MATLAB 2014a.

### Mathematical Study

#### Linear Predictive Coding (LPC)

LPC is based on the principle that a speech sample can be approximated as a linear combination of previous samples. The goal is to minimize the mean square error between the actual speech sample and the predicted sample.

#### Autocorrelation Method:

The autocorrelation  $R(k)$  of the signal  $s(n)$  is defined as:

$$R(k) = \sum_{n=k}^{N-1} s(n)s(n-k) \quad (9)$$

The Yule-Walker equations are then formed:

$$\begin{bmatrix} R(0) & R(1) & \dots & R(p-1) \\ R(1) & R(0) & \dots & R(p-2) \\ \vdots & \vdots & \ddots & \vdots \\ R(p-1) & R(p-2) & \dots & R(0) \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{bmatrix} = \begin{bmatrix} R(1) \\ R(2) \\ \vdots \\ R(p) \end{bmatrix}$$

Solving these equations yields the LPC coefficients.

#### Feature-Based Decision Making

Combining LPC-based features with ZCR and STE provides a comprehensive framework for VAD. Each feature captures different aspects of the audio signal, contributing to more accurate detection.

- **LPC Prediction Error:** Reflects the energy not accounted for by the linear model, typically higher for non-speech segments.
- **Zero-Crossing Rate:** Higher for non-speech segments due to noise and unvoiced sounds.
- **Short-Term Energy:** Higher for speech segments due to the presence of voiced sounds.

By combining these features, the proposed VAD method can effectively distinguish between speech and non-speech, even in noisy environments.

## ANALYSIS

**Integration of LPC with Traditional Features:** The method effectively combines LPC coefficients, ZCR, and STE to provide a comprehensive analysis of the audio signal. LPC captures the underlying structure of the speech signal, while ZCR and STE offer additional insights into its temporal properties.

**Robustness in Noisy Environments:** By utilizing prediction error energy from LPC, the approach can differentiate between speech and non-speech segments even in the presence of significant background noise. The method's reliance on multiple features helps in mitigating the effects of noise, leading to more reliable VAD performance.

**Computational Efficiency:** The proposed method retains computational efficiency, making it suitable for real-time applications. The use of the Levinson-Durbin algorithm for LPC coefficient calculation ensures that the approach remains feasible for implementation in resource-constrained environments.

**Empirical Thresholding and Smoothing:** The use of empirically determined thresholds for feature-based decision making, combined with a post-processing smoothing algorithm, enhances the accuracy and consistency of VAD decisions. This reduces false detections and improves overall system performance.

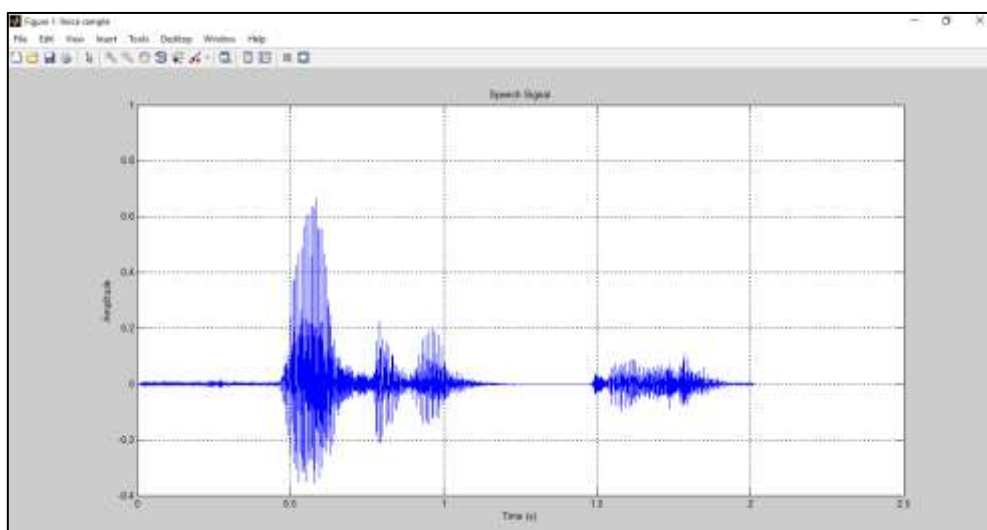


Figure 1: Speech Sample

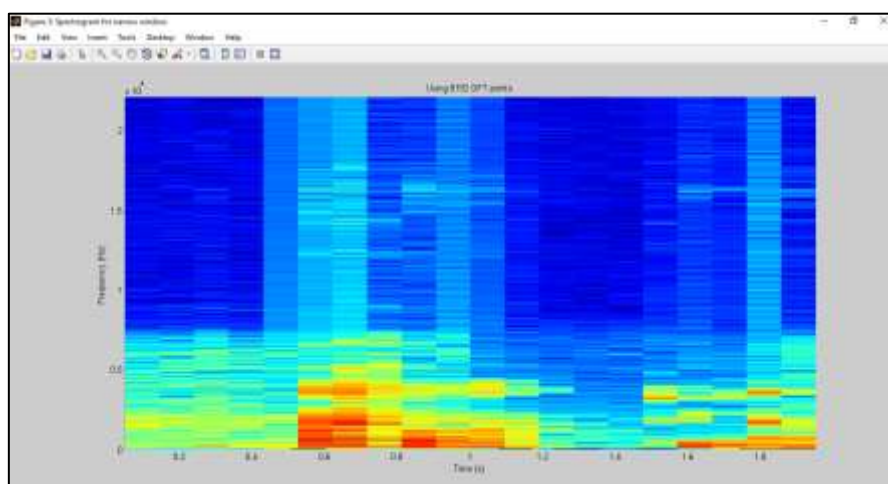


Figure 2: Spectrogram using Hamming Window

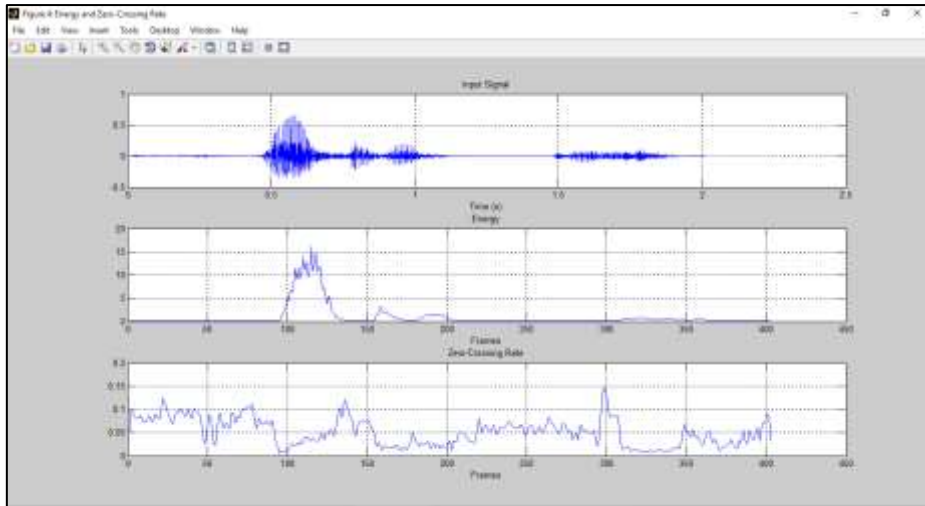


Figure 3: Energy and ZCR

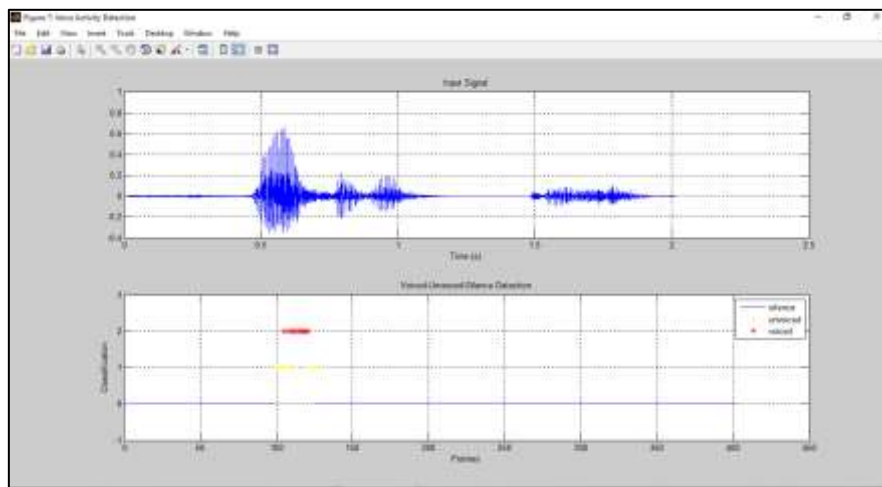


Figure 4: Voice Activity Detection

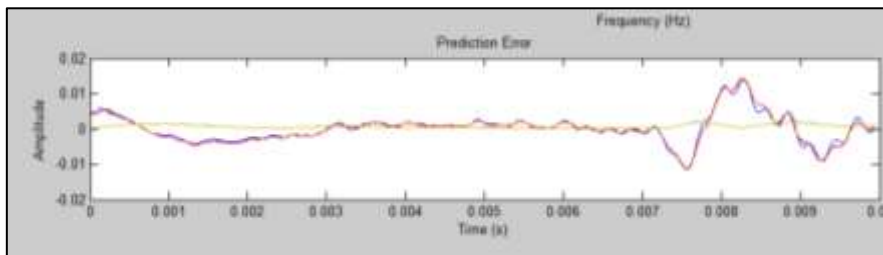


Figure 5: Prediction Error for Unvoiced Frame

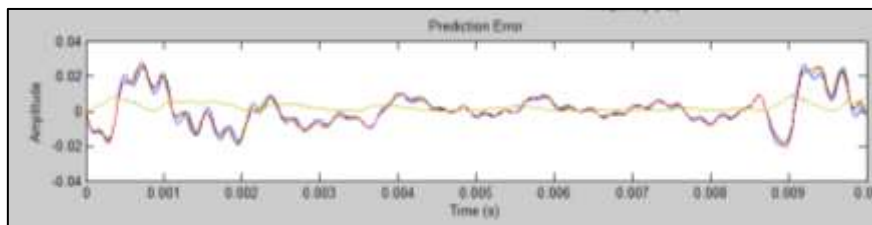


Figure 6: Prediction Error for Voiced Frame

## CONCLUSION AND FUTURE WORK

The proposed LPC-based VAD approach integrates traditional signal processing techniques with LPC to enhance voice activity detection. By leveraging the predictive power of LPC and combining it with ZCR and STE, this method provides robust and accurate VAD, suitable for various speech processing

applications. Future work will focus on optimizing thresholds and improving performance in diverse acoustic conditions. Implementing adaptive thresholding mechanisms that can adjust based on the acoustic environment in real-time could further improve the method's robustness. Incorporating machine learning models, such as deep neural networks, to learn optimal feature combinations and decision boundaries could enhance the accuracy of the VAD system. Extensive testing on a wide range of datasets, including those with various noise types and levels, would provide a more comprehensive assessment of the method's generalizability and robustness. Further optimization for real-time processing, including hardware acceleration techniques, could make the approach more viable for embedded systems and other real-time applications.

## REFERENCES

---

- [1] R. S. Nikhil Kumar, Sumit Dalal, "A LPC Based Voice Activity Detection Process," *Int. J. Adv. Res. Electr. Electron. Instrum. Eng.*, vol. 12, no. 5, pp. 1162–1169, 2023.
- [2] L. R. Rabiner and M. R. Sambur, "An algorithm for determining the endpoints of isolated utterances," *Bell Syst. Tech. J.*, vol. 54, no. 2, pp. 297–315, 1975.
- [3] J. Sohn, N. S. Kim, and W. Sung, "A statistical model-based voice activity detection," *IEEE Signal Process. Lett.*, vol. 6, no. 1, pp. 1–3, 1999.
- [4] D. Barchiesi, D. Giannoulis, D. Stowell, and M. D. Plumbley, "Acoustic scene classification: Classifying environments from the sounds they produce," *IEEE Signal Process. Mag.*, vol. 32, no. 3, pp. 16–34, 2015.
- [5] C. Kim and R. M. Stern, "Power-normalized cepstral coefficients (PNCC) for robust speech recognition," *IEEE/ACM Trans. audio, speech, Lang. Process.*, vol. 24, no. 7, pp. 1315–1329, 2016.
- [6] T. Hughes and K. Mierle, "Recurrent neural networks for voice activity detection," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013, pp. 7378–7382.
- [7] R. Varzandeh, S. Doclo, and V. Hohmann, "Speech-Aware Binaural DOA Estimation Utilizing Periodicity and Spatial Features in Convolutional Neural Networks," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, 2024.
- [8] M. A. Laskar and R. H. Laskar, "Integrating DNN–HMM technique with hierarchical multi-layer acoustic model for text-dependent speaker verification," *Circuits, Syst. Signal Process.*, vol. 38, no. 8, pp. 3548–3572, 2019.
- [9] Y. Cao and Z. Wang, "Improved DV-hop localization algorithm based on dynamic anchor node set for wireless sensor networks. *IEEE Access* 7: 124876–124890." 2019.