

VOP Detection for Read and Conversation Speech using CWT Coefficients and Phone Boundaries

Kumud Tripathi · K. Sreenivasa Rao

Received: date / Accepted: date

Abstract In this paper, we propose a novel approach for accurate detection of the vowel onset points (VOPs). VOP is the instant at which the vowel begins in the speech signal. Precise identification of VOPs is important for various speech applications such as speech segmentation and speech rate modification. The existing methods detect the majority of VOPs within 40 ms deviation, and it may not be appropriate for the above speech applications. To address this issue, we proposed a two-stage approach for accurate detection of VOPs. At the first stage, VOPs are detected using continuous wavelet transform coefficients, and the position of the detected VOPs are corrected using the phone boundaries in the second stage. The phone boundaries are detected by the spectral transition measure method. Experiments are done using TIMIT and Bengali speech corpora. Performance of the proposed approach is compared with two standard signal processing based methods. The evaluation results show that the proposed method performs better than the existing methods.

Keywords Vowel onset point (VOP) · continuous wavelet transform (CWT) · spectral transition measure (STM) · phone boundary · read and conversation modes.

1 Introduction

Broadly, speech can be divided into two modes, namely, read and conversation modes [1–3]. In read mode, an individual utters in restricted conditions, for

Kumud Tripathi, K. Sreenivasa Rao
Department of Computer Science and Engineering,
Indian Institute of Technology Kharagpur,
Kharagpur-721302, West Bengal, India
E-mail: kumudtripathi.cs@gmail.com

K. Sreenivasa Rao
E-mail: ksrao@iitkgp.ac.in

instance, reading of books or TV news. However, in conversation mode, two or more individuals are communicating in an unrestricted condition. In general, conversation speech is spontaneous, informal, unstructured, and unorganized. On the other hand, read speech is relatively less expressive and more organized. Due to high dissimilarity, it's important to study the impact of both the modes on the vowel onset point (VOP) detection.

Vowel onset point represents the start of a vowel in a speech signal. VOP is utilized for various speech applications such as (i) spotting consonant-vowel (CV) units in a speech signal [4,5], (ii) identifying other speech events such as formant transition, burst, aspiration, which help significantly in speech recognition [6] (iii) dividing speech into vowel and non-vowel like regions [7], etc. In general, the speech signals are processed at sub-word level such as monophone and triphone for speech recognition. Gangashetty et al. [8], have shown that syllables are the relevant sub-word units for speech recognition in Indian languages. Syllable represents a group of consonants (C) and vowel (V) in the form of C^aVC^x , where a and x indicate the count of consonants before and after the vowel, respectively. Among all the C^aVC^x combinations, the CV units are the most common (about 90%) existing syllables in Indian languages [8]. The CV units can be identified by accurate detection of VOPs in a continuous speech [9]. Therefore, the performance of the speech recognition system will be affected by the accuracy of the VOP detection method [4]. In literature [6,8,10–14], traditional vowel onset point detection methods are developed for read mode of speech. However, in a realistic scenario, conversation speech is more frequently observed than the read speech [15,16]. In terms of acoustic and linguistic characteristics, conversation mode has significant variations than the read mode [15,16]. Thus, the existing VOP detection methods may lead to spurious detection as well as missing VOPs for conversation speech. Hence, the performance of the traditional speech recognition system will be drastically reduced for conversation speech. Therefore, it is required to accurately detect the VOPs in conversation and read modes of speech for achieving the better recognition accuracy. So, the current study is motivated by the recognition of speech for Indian languages in read and conversation modes.

In previous studies, various methods are explored for VOP detection based on statistical modeling and signal processing approaches. Different statistical modeling methods are explained in [8,10,17] which utilizes multilayer feed-forward neural network, hierarchical neural network, and auto-associative neural network for VOP detection. These networks are developed using the speech features extracted from both sides of the VOPs. The predicted frame types are used for detecting VOPs in a speech signal. On the other hand, the detection of VOPs using signal processing methods is implemented by deriving various speech features. In [6], VOPs are detected by finding the locations of rapid growth in the vowel intensity. The change in the energy of each peak and valley of a speech signal is representing the vowel intensity. Prasanna et al., [13] proposed a VOP detection method based on the fusion of evidence from the excitation source, spectral peak energy, and modulation spectrum. The performance of the combined approach is better than the individual methods

for VOP detection. Vuppala et al. [14] utilized spectral energy present in glottal closure regions of speech signal for VOP detection. The computed spectral energy is robust and high around the glottal closure instants (GCIs).

Mostly, the statistical modelling methods and signal processing techniques may falsely detect the VOPs in the presence of diphthongs and semivowel-vowel transitions [18–20]. This is due to the similar acoustic characteristics of the semivowels and vowels. Hence, in the recent works, various statistical modelling [21, 22], and signal processing [7, 18] methods have been evolved for detecting the vowel-like region onset points (VLROP) instead of vowel onset point. The VLROPs represents the start of the vowel, semivowel, and diphthong speech regions [18]. But these methods may not be suitable for speech applications where only vowel regions are needed to be identified, such as consonant-vowel recognition, speech-rate modification, speaker recognition, and so on.

The existing methods based on statistical modelling techniques depend on a huge amount of training data. In these methods, at the first step, a classifier is trained for detecting the vowel regions, and then the VOPs are detected by locating the instant at which detected vowels are started. The accuracy of the detected VOPs will depend on the performance of the vowel detection algorithm. However, the signal processing methods can be directly applied to speech signals for identifying VOPs as compared to statistical modelling methods. The signal processing methods follow simple and less number of steps, then the statistical modelling methods which follow the complex process. In terms of accuracy, both methods are providing almost similar results. Hence, in this work, we have proposed a signal processing based method for accurate detection of VOPs. As the proposed method is signal oriented, so; the state of the art signal processing methods [13, 14] are included for performance comparison. The existing methods [13, 14, 18–20] detect the majority of the VOPs within 40 ms deviation. Therefore, attaining a better accuracy for VOP detection at lower deviation is the primary goal of the proposed approach.

In this work, we have proposed a novel method to accurately detect the VOPs in a speech signal. The proposed method is performed at two stages for robust detection of VOPs. At the first stage, continuous wavelet transform (CWT) is explored for determining the VOP evidence. Continuous wavelet transform [23] is capable of detecting the instants of sharp transitions, including steady regions in a speech signal. This is the motivation behind choosing CWT for VOP detection. At the second stage, a new approach is explored based on phone boundary information for correcting the positions of detected VOPs. Spectral transition measure (STM) method [24–26] is applied for detecting the phone boundaries. Dusan et al. [25] have shown that the STM is accurately detecting 90% of phone boundaries under 20 ms deviation. In this work, it is analyzed that the majority of VOPs detected using CWT coefficients are within 40 ms deviation. Therefore, to improve the accuracy of the proposed method at low deviation, the location of the detected VOPs are corrected with the help of detected phone boundaries. The proposed method is significant for segregating the vowel onset points from the remaining speech

regions. To validate this fact, the proposed method is compared with signal processing techniques reported in [13, 14] using TIMIT corpus. In addition to that, the importance of the proposed method is shown by detecting VOPs for read and conversation modes of Bengali speech.

The organization of the paper is as follows. Section 2 describes the baseline VOP detection methods. The description of the proposed method for accurate VOP detection is presented in Section 3. The performance and significance evaluations of the proposed method using TIMIT and Bengali (read and conversation modes) speech corpora are presented in Section 4. Section 5, includes the conclusions of the current study and works that need to be explored in the future.

2 Baseline VOP Detection Methods

Performance of the proposed approach is compared with two standard signal processing based methods. The first method combines the evidence from the excitation source, spectral peaks energy, and modulation spectrum [13], and the second method is based on spectral energy around glottal closure regions [14]. The detailed description of these methods are presented below.

2.1 Combined evidences from excitation source, spectral peaks, and modulation spectrum for VOP Detection

In this method, evidence from excitation source, spectral energy, and modulation spectrum are combined at frame level for detecting VOPs. The Hilbert envelope of LP residual contains the information about excitation source. The sum of 10 major peaks of the DFT computed for each frame, represents the energy of spectral peaks. The modulation spectrum corresponds to the gradually changing temporal envelope of speech. These methods contain different information for VOP detection and thus can be combined. The combined method leads to better performance than the excitation source, spectral energy, and modulation spectrum methods, respectively. This method is titled as COMB-ESM for the rest of the paper.

2.2 Spectral energy around glottal closure regions for VOP Detection

This method detects VOPs in a glottal closure regions of the speech signal using evidence from the spectral energy. The spectral energies are more prominent at GCIs. Therefore, the spectral energy is computed for the frames present in the 30% of the glottal cycle around the GCIs. The zero frequency filter is applied for detecting the glottal closure instants in a given speech sequence. The spectral energies in the range of 500-2500 Hz are considered for VOP detection. The spectral energy signal was smoothed over the window of 50 ms to reduce the inconstancies. Further, the smoothed spectral signal is enhanced by

computing the slope using first-order difference. In the enhanced signal, prominent variations (peaks) are extracted by convolving with first-order Gaussian difference operator of size 100 ms. The peaks in the convolved signal were representing the vowel onset points. This method is named as SE-GCI for the rest of the paper.

3 Proposed VOP Detection Method

In the proposed approach, continuous wavelet transform is explored along with spectral transition measure to enhance the accuracy of the VOPs. Continuous wavelet transform can predict smooth signal features as well as abrupt transitions [27]. For some phonemes such as /ax/, /axr/ and /ux/, CWT may fail to predict VOPs under 40 ms because of the very short and devoiced vowel. However, STM will accurately provide 97% of phone boundaries within 40 ms deviation. For that reason, we have incorporated the information carried by STM along with CWT for further improving the performance of VOP detection. The details about VOP detection using CWT is included in Section 3.1. The detailed description of phone boundary detection using STM is provided in Section 3.2. The combined model for improving the performance of detected VOPs is described in Section 3.3.

3.1 VOP Detection using CWT

CWT gives a complete representation of a signal by varying the scale value of the wavelets repeatedly. Mathematically, the CWT of a speech signal $x(t)$ can be represented as:

$$C_x(p, q) = \frac{1}{\sqrt{q}} \int_{-\infty}^{\infty} x(t) \phi^* \left(\frac{t-p}{q} \right) dt \quad (1)$$

where $\phi(t)$ is the mother wavelet and $\phi^*(t)$ is the complex conjugate of $\phi(t)$. $C_x(p, q)$ is representing the wavelet coefficient for scale parameter q ($q > 0$) and translation parameter p . The CWT coefficients computed in Eq. (1) can be observed as the product of signal $x(t)$ and wavelet (shifted and scaled) $\phi(t) : \phi_{p,q}(t) = (1/\sqrt{q})\phi((t-p)/q)$. In this work, VOPs are detected from the mean signal derived using CWT coefficients. The mean signal can be computed as follows:

$$A_c(p) = \frac{1}{N} \sum_{q \in q_s} C_x(p, q) \quad (2)$$

where N is the number of scales and q_s is the set of chosen scale. In rest of the paper, the mean signal derived from CWT coefficients is named as ‘‘mean-signal’’.

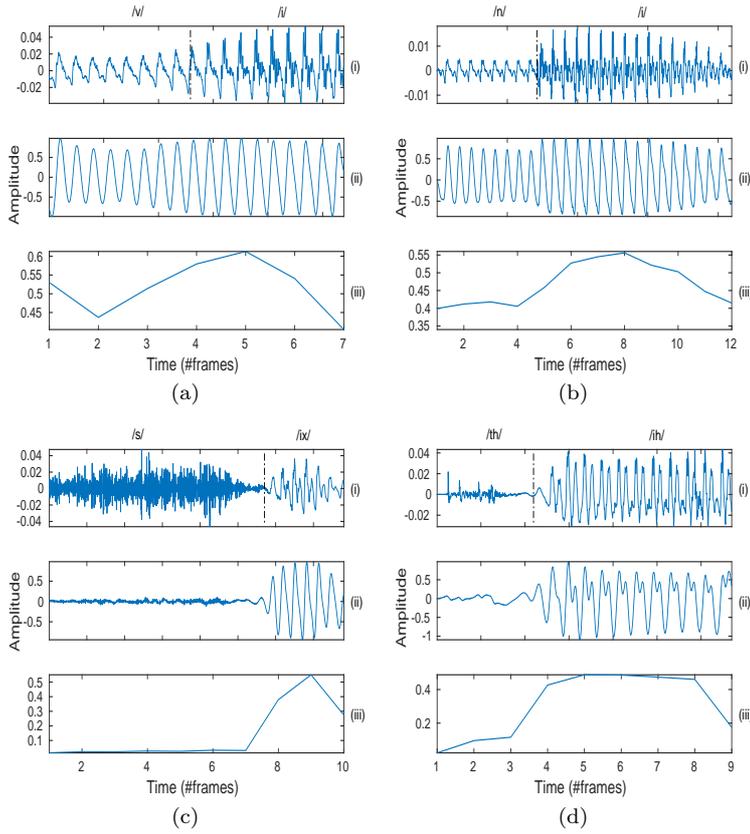


Fig. 1: Illustration of mean-signal and AAM of mean-signal. (a) Semi-vowel to vowel transition, (b) nasal to vowel transition, (c) fricative to vowel transition and (d) unvoiced to vowel transition.

For detecting the VOPs, mean-signal is segmented into frames of 20 ms with 10 ms overlap. For every segment, we have calculated the average absolute magnitude (AAM) of the mean-signal. Inconsistencies in the AAM of the mean-signal are flattened by applying mean-smoothing for 40 ms segment. For determining VOPs, an optimal threshold (th_1) is fixed at 15% of the maximum of smoothed AAM. This threshold is decided by estimating the errors for VOP detection. The VOP detection errors are representing the percentage of missed and spurious VOPs. In this work, 5 distinct threshold values are examined between 11% to 20% at the interval of 2% as displayed in Table 1. These experiments are performed on a subset of TIMIT speech corpus [28]. The first column of the table contains the different threshold values, and the second column represents the percentage of missed VOPs. The third column indicates the percentage of spurious VOPs. The results shown in columns 2 and 3 are computed within 40 ms deviation. It can be observed from Table 1

that the threshold value of 15% of the maximum of smoothed AAM is providing an appropriate decision for VOPs. In a similar manner, experiments are performed with Bengali read and conversation speech signals and analyzed that $th_1=15\%$ is giving the minimum error. As per the experimental results, threshold value of 15% is providing the global minimum error within given data. Therefore, in the proposed method, 15% of the maximum of smoothed AAM is considered for determining the VOPs in TIMIT and Bengali speech signals.

Table 1: VOP detection errors for various threshold values.

Threshold Values (%)	Miss Rate	Spurious Rate
11	15	28
13	15	25
15	15	20
17	17	21
19	20	23

It is noticed from the literature survey that there is no work related to CWT for VOP detection. This motivated us to explore vowel discriminative characteristics of CWT for predicting VOPs in a speech signal. The well-known difficulties in VOP detection are finding false VOPs in case of semi-vowels, nasals, and fricatives as they are periodic in nature [18–20]. The AAM of mean-signal for semi-vowel, nasal, fricative and unvoiced speech regions are illustrated in Figure 1. The speech waveform and mean-signal of given speech regions are displayed in Figures 1(a,b,c,d)(i and ii), respectively. However, for given speech regions, the AAM of mean-signals are displayed in Figures 1(a,b,c,d)(iii). It can be observed from the figure that the mean-signals are generally periodic in vowels, semi-vowels, and nasals, and almost zero in fricatives and unvoiced speech segments. This represents that the speech signal and mother wavelet have the least correlation at all the scales in the fricative and unvoiced regions. However, the CWT coefficients are higher and, shows a periodic shape at all scales for speech regions such as vowels, semi-vowels, and nasals. As we can see in Figure 1(a,b)(iii) that the peak amplitudes of vowel regions are significantly higher than other speech regions. Thus, the unwanted peaks can be easily removed using an optimal threshold parameter as well as using mean-smoothing. Hence, it is clear that the CWT can significantly identify the VOPs in vowel regions by suppressing the remaining speech regions.

The steps for the VOP detection using CWT based method are summarized as follows:

1. Compute the CWT coefficients of the speech signal.
2. Derive the mean-signal.
3. Determine AAM for each frame of the mean-signal where the length of a frame is 20 ms and frame shift is 10 ms.

4. Inconsistencies in the AAM of the mean-signal are flattened by applying mean-smoothing of frame size 40 ms.
5. Detect local peaks of the smoothed AAM of the mean-signal.
6. For removing the undesirable peaks, an optimal threshold is fixed at 15% of the maximum of smoothed AAM. Further, undesired peaks are removed if two consecutive peaks are present within 50 ms window. On this basis, the smaller amplitude peak is removed.
7. After removing the undesired peaks, frames whose AAM is larger than or equal to the threshold value are chosen as VOP frames.

Figure 2 demonstrates the evidence of VOP detected using CWT for an utterance /“*She had your dark suit in greasy wash*”/ from the TIMIT corpus sampled at 16 kHz. Figure 2(a) displays the waveform of the speech signal. Mean-signal derived from CWT is represented in Figure 2(b). The average absolute magnitude (AAM) of the mean-signal is shown in Figure 2(c). The local peaks are represented by the circle (o) symbol. Figure 2(d), shows the smoothed AAM where mean-smoothing is applied for removing the fluctuations present in the AAM of mean-signal. The undesirable peaks in Figure 2(d) are omitted by applying the threshold value, which is 15% of the maximum of smoothed AAM. The threshold value is empirically chosen based on several experiments on a large volume of data. In addition to that if two consecutive peaks are reported within 50 ms; the peak with smaller magnitude will be omitted. This relies on the hypothesis that there will be only one VOP within the window of 50 ms [14]. The detected peaks in Figure 2(e) after removing the undesired peaks are representing the desired locations of VOPs.

3.2 Phone Boundary Detection using STM

The spectral transition measure provides an unsupervised way of detecting phone boundaries in a speech signal. This is the key motivation for exploring STM in this study. The 13-dimensional Mel Frequency Cepstral Coefficients (MFCCs) along with Δ and $\Delta\Delta$ coefficients are considered for deriving spectral details of the speech signal. The Δ and $\Delta\Delta$ coefficients corresponds to the first and second order derivative of MFCCs, respectively. The spectral details are extracted by considering a frame size of 25 ms and a frame shift of 10 ms using the Hamming window.

The implementation of STM in this work is the same as that described in [25]. STM can be understood as the degree of variation in the spectral value of a speech signal. The spectral variation at the phone transition is maximum compared to steady speech regions. Such higher spectral variations represent peaks and these peaks are considered as detected phone boundaries in a speech sequence. It can be observed from the STM contour (see Figure 3(c)) that the phone boundaries have more spectral deviation. That is responsible for producing high Mean Square Error (MSE) in linear regression [24]. The STM,

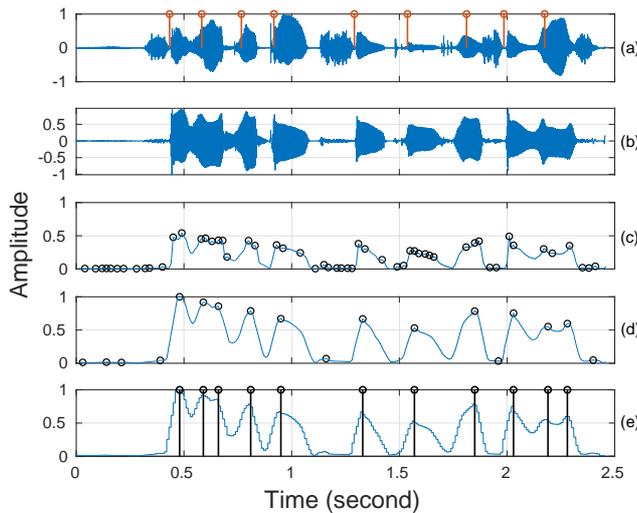


Fig. 2: VOP detection using CWT for an utterance /“*She had your dark suit in greasy wash*”/. (a) Speech waveform with actual VOPs, (b) Mean-signal, (c) AAM of mean-signal with peak locations, (d) Smoothed AAM with peak locations, (e) CWT detected VOPs at 15% of the optimal threshold of the maximum of smoothed AAM. Actual VOPs are marked with the red line and detected VOPs are marked with the black line.

at frame g , can be computed as a mean-squared value [26], i.e.,

$$S_g = \frac{1}{D} \sum_{i=1}^D r_i^2(g) \quad (3)$$

where S_g represents the STM at frame g , D is the dimension of the spectral feature vector (39 in this case) and $r_i(g)$ shows the rate of variation in spectral details $MFCC_i$ and defined as [25],

$$r_i(g) = \frac{\sum_{n=-I}^I MFCC_i(n+g) * n}{\sum_{n=-I}^I n^2} \quad (4)$$

where n shows the frame index, i is the coefficient index, and I displays the number of frames (on each side of the current frame) utilized for computing the regression coefficients. The considered value of I is 2 for calculating STM [26]. The value of I greater than 2, result in missing desired phone locations, whereas I smaller than 2, generate many unwanted phone locations.

In this work, phone boundaries are extracted from the STM contour of the speech signal. Figure 3 illustrates the phone boundary detection for an utterance /“*She had your dark suit in greasy wash*”/. Figure 3(a) represents the speech signal with actual phone boundaries. The STM contour of the speech

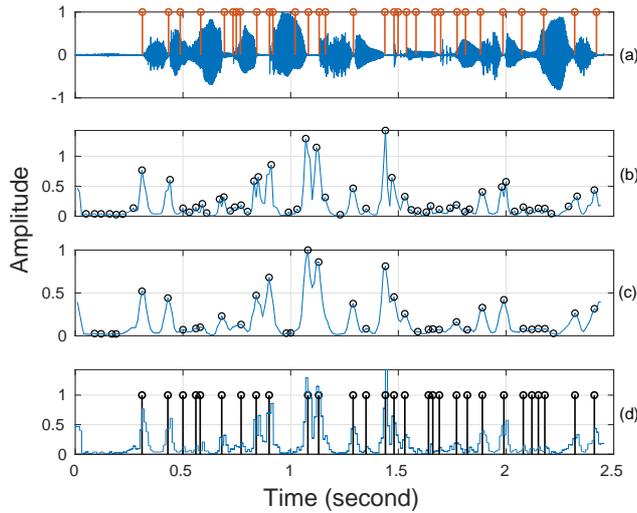


Fig. 3: Phone boundary detection using STM for an utterance /“*She had your dark suit in greasy wash*”/. (a) Speech waveform with actual boundaries, (b) STM contour with peak locations, (c) Smoothed STM contour with peak locations, (d) Detected phone boundaries at 12% of the optimal threshold of the maximum of STM contour amplitude. Actual phone boundaries are marked with the red line and detected phone boundaries are marked with the black line.

signal is shown in Figure 3(b). The local peaks are marked as a circle (o). The inconsistencies in the STM contour are flattened by applying mean-smoothing for 20 ms frame represented in Figure 3(c). The noisy peaks present in Figure 3(c) are removed by implementing a threshold (th_2) of 12% of the maximum of STM contour amplitude. The threshold value is decided after analyzing the STM contour on a subset of TIMIT corpus. The performance of the detection task is dependant on the selection of the optimal threshold. It is noted that the threshold value smaller than th_2 generate more spurious boundaries, whereas the threshold value higher than th_2 have missed phone boundaries. In Figure 3(b), it can be seen that the spectral value is varying significantly at phone transition. This resulted in peaks which are considered as phone boundaries in this study. The detected phone boundaries after eliminating the unwanted peaks are detected at a fixed threshold ($th_2 = 12\%$) is shown in Figure 3(d). The process of phone boundary detection can be summarized as follows:

1. Extract 39-dimensional MFCC, Δ , and $\Delta\Delta$ feature, with a frame size of 25 ms and frame shift of 10 ms,
2. Compute STM for each frame of a given signal,
3. Remove inconsistencies in the STM contour by applying mean-smoothing of frame size 20 ms,

4. Detect local peaks of the smoothed STM contour,
5. For eliminating the undesirable peaks, an optimal threshold is set at 12% of the maximum of smoothed STM contour amplitude.
6. After removing the false peaks, frames whose amplitude is larger than or equal to the threshold value are chosen as frames for phone boundary.

3.3 Two-Stage Method for VOP Detection

The proposed method for VOP detection is based on the evidence of two different methods discussed in Section 3.1 and 3.2. In the first method, VOPs are hypothesized from AAM of CWT derived mean-signal. In the second method, phone boundaries are detected from STM contour for correcting the position of CWT detected VOPs. The block diagram of the proposed VOP detection method is shown in Figure 4. Figure 5 demonstrates the VOP detection for an utterance /*“She had your dark suit in greasy wash”*/. The speech waveform with actual VOPs is shown in Figure 5(a). The detected VOPs using CWT as well as actual VOPs are depicted in Figure 5(b). It can be analyzed that the CWT detected VOPs are deviated from the actual VOPs. Due to this at low deviation, some actual VOPs will be accounted as missed, and some of the detected VOPs will be accounted as spurious. This will become the reason for generating missed and spurious VOPs using CWT. In addition to that, CWT will detect false VOPs in case of high energy voiced consonants; for example, in Figure 5(b) the detected VOPs such as 3rd and 11th are noisy.

Table 2: Performance of VOP detection using CWT, and STM on TIMIT corpus.

VOP Detection Method	VOPs Detected within ms (%)				Spurious VOPs (%)
	10	20	30	40	
CWT	52	65	78	91	20
STM	90	91	96	98	70

Therefore, STM detected phone boundaries in Figure 5(c) are utilized for eliminating noisy VOPs as well as for reducing the deviation between actual and predicted VOPs. Intuitively, it is reported for a subset of TIMIT speech corpus that detected VOPs are mostly occurring right side of the actual VOPs. It can also be visualized in Figure 5(b) that each detected VOPs have unevenly deviated towards the right side of actual VOPs. Hence, to correct the location, the detected VOPs need to be relocated to its left side. One way to achieve this is by shifting each VOP with the fixed length. However, this is not a feasible solution because the VOPs have unevenly deviated. Hence, a method is required to automatically adjust the detected VOPs to their accurate locations. In this work, we have explored a new approach based on phone boundary details for

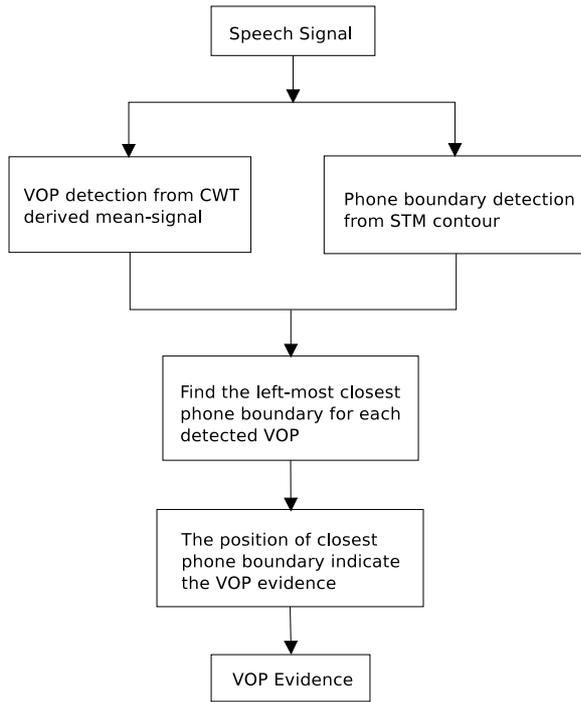


Fig. 4: Block diagram of proposed VOP detection method.

correcting the detected VOPs. Here, the STM method is explored for phone boundary detection. Table 2 represents the accuracy of vowel onset point detection using CWT and STM based methods. For performing this experiment, 50 utterances are randomly selected from TIMIT speech corpus. It is observed from Table 2 that the percentage of detected VOPs using STM is better than the CWT. However, the percentage of spurious VOPs is huge using STM than the CWT. In this work, STM is evolved for phone boundary detection. Thus, at the time of VOP detection using STM, vowel boundaries are considered as VOPs, and remaining phone boundaries are considered as spurious VOPs. For that reason, the percentage of spurious VOPs are overestimated in STM based method. From experiments, it is noted that the performance of the proposed two-stage method does not affected by the STM detected spurious phone boundaries. Further, it can be seen from Table 2 that the percentage of STM detected VOPs within 10 ms and 20 ms is much higher than the CWT detected VOPs. This explains that STM can detect vowel boundaries (VOPs) better than CWT. Therefore, the STM based vowel boundary details are incorporated with CWT based VOP detection, for improving the performance of the proposed method at a smaller deviation. The detected VOPs in Figure 5(d) after removing the spurious VOPs are depicting the desired location of VOPs. These VOPs are within 10 ms deviation from the actual VOPs. Hence, the

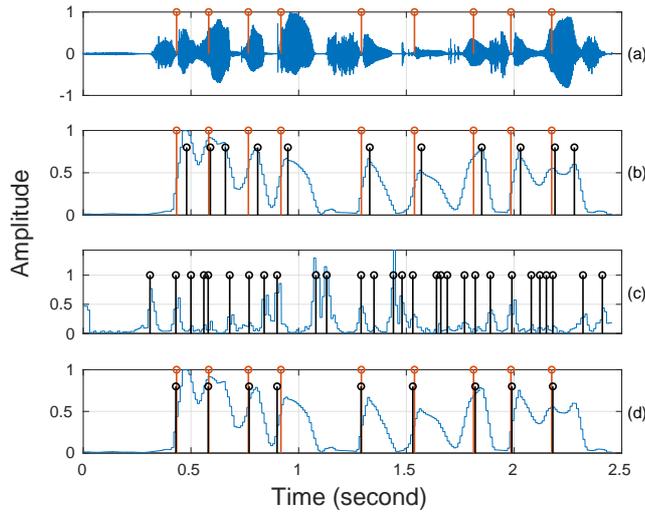


Fig. 5: VOP detection for an utterance /“She had your dark suit in greasy wash”/. (a) Speech waveform with actual VOPs, (b) Actual VOPs and detected VOPs using CWT coefficients, (c) Located phone boundaries using STM contour, (d) Actual VOPs and detected VOPs using the proposed method. Actual phone boundaries are marked with the red line and detected phone boundaries are marked with the black line.

proposed method is suitable for speech applications where accurate VOPs are required. The steps involved in correcting the detected VOPs using proposed two-stage method are:

1. Detect the VOPs using the CWT based method.
2. Detect the phone boundaries using the STM based method.
3. For each detected VOP, find the left-most closest phone boundary (see in Figure 5(c)).
4. The position of the detected phone boundary is marked as the location of modified VOP, as shown in Figure 5(d).

4 Performance Evaluation

In this work, the performance of the proposed method is compared with two existing methods based on COMB-ESM and SE-GCI. Here, TIMIT corpus is considered for evaluating the performance of VOP detection methods. However, Indian speech corpora in read and conversation modes are considered for depicting the significance of proposed VOP detection method.

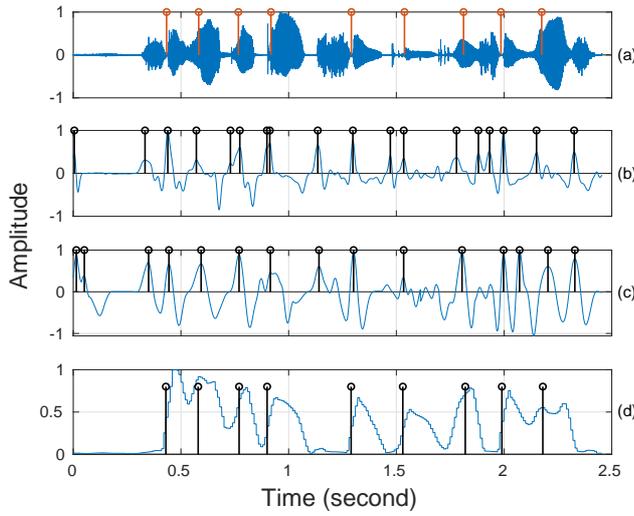


Fig. 6: VOP detection using proposed and existing methods for an utterance /“*She had your dark suit in greasy wash*”/. (a) Speech waveform with actual VOPs, (b) Detected VOPs using COMB-ESM method, (c) Detected VOPs using SE-GCI method, (d) Detected VOPs using proposed method. Actual phone boundaries are marked with red line and detected phone boundaries are marked with black line.

4.1 VOP Detection on TIMIT Corpus

Experiments are conducted on TIMIT speech corpus for performance analysis of proposed and two existing VOP detection methods. About 120 randomly selected utterances sampled at 16 kHz are used for analyzing the performance of the explored VOP detection methods. The metrics considered for measuring the performance of the various methods are identification rate (IR), average deviation (AD), missing rate (MR), and spurious rate (SR). The percentage of actual VOPs that correspond to the detected VOPs within the considered (10-40 ms) time-resolutions is known as identification rate. The average deviation (in ms) is demonstrating the average time difference between the actual and predicted VOPs. The percentage of actual VOPs that are undetected within the considered deviation is termed as missed rate. The percentage of detected VOPs other than the actual VOPs is termed as spurious rate.

VOP detection accuracy of proposed, COMB-ESM and SE-GCI methods in terms of IR, AD, MR, and SR is demonstrated in Table 3. The first column contains the list of methods involved in analyzing the VOP detection performance. Columns second to fifth represent the IR (%) within the mentioned deviations. The sixth column specifies the AD (in ms) with respect to the actual VOPs. Seventh and eighth columns, represent the missed and spurious

Table 3: Performance of VOP detection using proposed method, COMB-ESM and SE-GCI on TIMIT corpus.

VOP Detection Method	VOPs Detected within ms (%)				Average Deviation (\approx ms)	Missed VOPs (%)	Spurious VOPs (%)
	10	20	30	40			
COMB-ESM	51	59	74	90	18	10	6
SE-GCI	62	79	86	91	13	9	5
Proposed	82	88	91	92	7	8	3

rates, respectively. It is observed from the table that the overall performance of the proposed method is better than the existing (COMB-ESM, and SE-GCI) methods for VOP detection. The average deviation in the proposed method (7 ms) is significantly smaller than the COMB-ESM (18 ms) and SE-GCI (13 ms) methods. The rate of missed and spurious VOPs is relatively higher in COMB-ESM and SE-GCI methods (see Figure 6). In both COMB-ESM and SE-GCI, the detection of VOPs relies upon the spectral energy and its enhancement. In these methods, the spectral energy of a speech signal is enhanced by computing its slope value using first-order derivative (FOD). Further, the enhanced features are convolved with first order Gaussian difference (FOGD) operator for locating the VOP evidences. Here, enhanced feature help in improving the identification rate but at the same time it highlighted the peaks for periodic non-vowel regions such as semi-vowels, and nasals, which leads to spurious detection of VOPs. The proposed method outperformed the existing methods in case of identification and spurious rates. It can be seen in Figure 6 that the existing spurious VOPs in COMB-ESM and SE-GCI methods are removed in the proposed method. Additionally, it is noticed that the IR for the proposed method is almost 30% higher within 10 ms as compared to other existing methods. However, the proposed method is shown the significantly better performance within 20 ms than the 10 ms deviation. This is due to the cases where high energy voiced consonants are preceded by the vowels, which resulted in the deviation of detected VOPs with respect to the genuine VOPs.

4.2 VOP Detection in Read and Conversation Modes of Bengali Speech Corpora

In this work, the significance of the proposed method is demonstrated by detecting the VOPs in read and conversation modes of Bengali speech. The Bengali speech dataset is collected as part of consortium project titled *Prosodically guided phonetic engine for searching speech databases in Indian languages* supported by DIT, Govt. of India [29]. In this study, read speech is collected from news reading, and the conversational speech is collected from casual talks. The speech signals are sampled at a rate of 16 kHz with the precision of 16 bits per sample. Altogether, 20 utterances are collected from 5 distinct speakers, where 3 male and 2 female speakers are considered for each mode. About 100

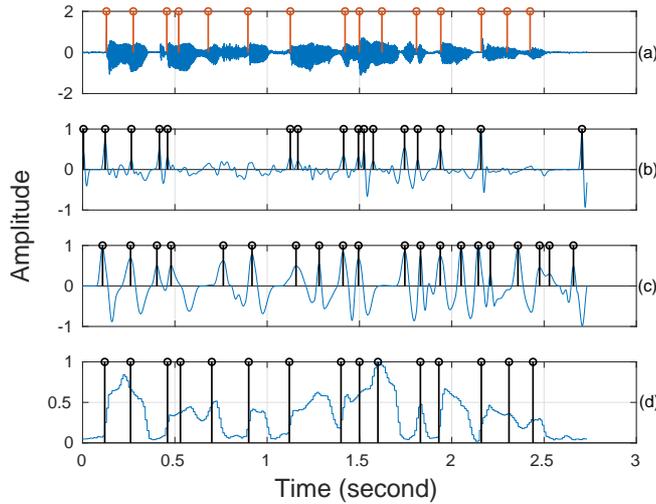


Fig. 7: VOP detection for Bengali sentence /“Tomake abosi fael gulo muche felte hobe”/ uttered in read mode. (a) Speech waveform with actual VOPs, (b) Located VOPs using COMB-ESM method, (c) Located VOPs using SE-GCI method, (d) Detected VOPs using the proposed method. Actual phone boundaries are marked with the red line and detected phone boundaries are marked with the black line.

utterances from each mode are selected for evaluating the performance of the proposed and existing methods.

Mostly, VOP detection methods are studied for read mode of speech [6, 8, 10, 12–14]. However, speech can be broadly divided into two modes, such as read and conversation. The acoustic and linguistic characteristics of these modes are very different. As conversation speech is a type of spontaneous and unconstrained communication between two or more than two people. However, the read speech includes planning before reading in constrained conditions such as news reading. The conversation mode includes higher variations in activity of vocal folds than the read mode. Due to the aforementioned variations, read, and conversation modes are examined in this work for evaluating the accuracy of VOP detection methods.

Table 4 demonstrates the performance of VOPs detected in read and conversation modes of Bengali speech using the existing and proposed methods. The first column shows the modes of speech considered for detecting VOPs. The second column signifies various methods used in the analysis of VOP detection. Third, to sixth columns represent the IR (%) within the given deviations (10 to 40 ms). Seventh and eighth columns are specifying the average deviation and spurious rates, respectively.

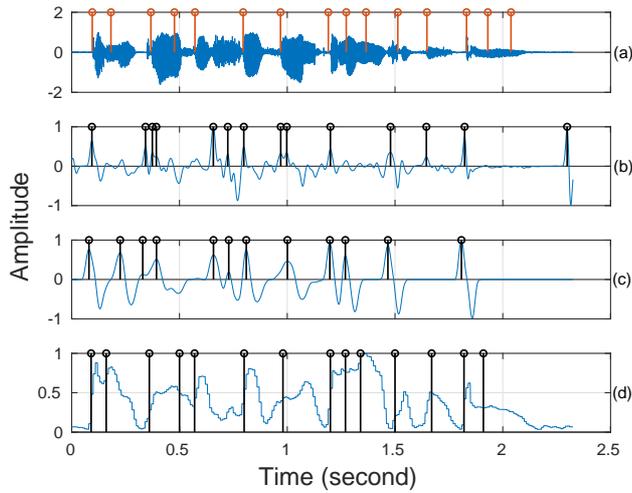


Fig. 8: VOP detection for Bengali sentence /“Tomake abosi fael gulo mucche felle hobe”/ uttered in conversation mode. (a) Speech waveform with actual VOPs, (b) Located VOPs using COMB-ESM method, (c) Located VOPs using SE-GCI method, (d) Detected VOPs using the proposed method. Actual phone boundaries are marked with the red line and detected phone boundaries are marked with the black line.

Table 4: Performance of VOP detection using COMB-ESM, SE-GCI, and proposed methods on read and conversation (Conv) modes of Bengali speech corpora.

Speech Mode	VOP Detection Method	VOPs Detected within ms (%)				Average Deviation (\approx ms)	Spurious VOPs (%)
		10	20	30	40		
Read	COMB-ESM	44	63	78	87	20	8
	SE-GCI	59	70	82	89	14	5
	Proposed	75	81	88	89	13	4
Conv	COMB-ESM	40	46	61	73	21	15
	SE-GCI	45	58	72	78	19	13
	Proposed	72	79	83	85	15	7

It is noticed from Table 4 that the performance of detected VOPs is higher under read mode compared to conversation mode. In conversation mode, percentage of detected VOPs is reduced by 14%, 11% and 4% (within 40 ms deviation) as compared to read mode using COMB-ESM, SE-GCI and proposed method, respectively. It is important to note that the performance reduction in conversation mode is minimum for the proposed method than the existing methods. Hence, this result states that the proposed method is least influenced by the speech mode discriminative characteristics. This analysis can be visu-

alized in Figures 7 and 8 where, a Bengali sentence /“*Tomake abosi fael gulo mache felte hobe*”/ is spoken in read and conversation modes, respectively. For both the modes, the same female speaker has uttered the given sentence. The speech recorded in conversation mode includes fear expression. However, uttered speech in read mode contains a neutral expression. Here, Figures 7(a) and 8(a) show the speech waveform with actual VOPs in read and conversation modes, respectively. The detected VOPs using COMB-ESM method in read and conversation modes are respectively depicted in Figures 7(b) and 8(b). Figures 7(c) and 8(c) show the detected VOPs using SE-GCI method in read and conversation modes, respectively. Figures 7(d) and 8(d) represent the identified VOPs using the proposed method in read and conversation modes, respectively. It is observed from the Figures 7 and 8, the energy variation is highly dynamic in conversation mode than the read mode. This observation is justified intuitively by analyzing the Bengali sentences in given modes. Further, it is claimed that the spectral energy is significantly varying in conversation mode than the read mode due to the presence of emotions. The involvement of emotions leads to a small percentage of clean speech in conversation mode than the read mode of speech. Therefore, the percentage of detected VOPs are improved in read mode than the conversation mode. In addition to that, while expressing emotions, some voiced consonants are also got emphasized in conversation mode, which results in spurious VOPs. It is observed that the duration of vowels is smaller in conversation mode as compared to read mode. This is the reason for missing VOPs in conversation mode. The paralinguistic aspects of the speech, such as gasp, sigh, and mhm are more often present in conversation speech than the read speech. This resulted in spurious detection of VOPs in conversation mode. For all these reasons, the overall performance of VOP detection methods is reduced in conversation mode as compared to read mode.

It can be noted from Table 4 that proposed method is performing better than the existing methods for spotting VOPs in case of read and conversation modes. In read mode, the proposed method is detecting about 17% more VOPs under 10-20 ms deviation as compared to COMB-ESM and SE-GCI. Similarly, in conversation mode, the proposed method is extracting almost 35% more VOPs within 10-20 ms deviation than the COMB-ESM and SE-GCI. The complete result represents that the performance of the proposed method is relatively similar for read and conversation speech. The instants of VOPs represent sharp energy transitions in both read and conversation speech. The ability of CWT to confine these sharp energy transitions, help in accurately detecting VOPs even in the presence of conversation speech. Hence, the proposed method is providing better identification rate for both the conversation and read modes of speech. Further, it can be noticed that the proposed method has a significant reduction in the average deviation and spurious rates. Percentage of spurious VOPs is reduced around 4% and 8% in the proposed method as compared to COMB-ESM for read and conversation modes, respectively. Similarly, the percentage of spurious VOPs using the proposed method as compared to SE-GCI is reduced around 1% and 6%

in read and conversation modes, respectively. This is because, at first step in the proposed method, AAM of mean-signal obtained from CWT coefficients is enhancing the voiced region and suppressing the unvoiced region of speech. At second step, the detected VOPs are corrected using phone boundaries derived from STM contour. The use of STM in proposed method helps in removing the spurious VOPs as well as in reducing the deviation between actual and predicted VOPs. The overall result indicates that combining CWT and STM methods into a single framework can accurately detect the VOPs present in a speech utterance spoken in any mode.

5 Conclusion

In this work, a novel method is proposed for accurate VOP detection. The proposed method consists of two-stages. At the first stage, VOPs are detected by using the AAM of mean-signal derived from continuous wavelet transform (CWT) coefficients. At the second stage, the evidence of identified VOPs is corrected with the presence of the nearest phone boundary detected using spectral transition measure (STM) method. In the proposed method, CWT and STM are utilized to obtain the sharp energy transitions around the VOPs. VOP detection experiments are carried out with TIMIT corpus (read speech) and Bengali corpus (read and conversation speech). Performance of the proposed approach is compared with two standard methods: COMB-ESM and SE-GCI. The proposed method was demonstrated to be significantly better in eliminating spurious VOPs and for accurately detecting the VOPs within 10 ms deviation as compared to COMB-ESM and SE-GCI methods. The efficiency of the proposed method is demonstrated by detecting VOPs in two acoustically and linguistically different speech modes such as read and conversation modes. The results achieved for read and conversation modes signify that the performance of the proposed method is insignificantly affected by the acoustic variation among the modes and achieved almost similar performance for both the modes. As the proposed approach demonstrates accurate computation of VOP locations in a speech signal, this can be utilized for consonant-vowel recognition, speech rate modification, voiced-unvoiced classification, and so on. Further, the robustness of the proposed method can be explored for noisy speech corpora. In this work, we have explored two broad modes of speech for examining the significance of the proposed method; one can explore other emotional modes of speech like anger, happy, sad, etc. The proposed method is explored for VOP detection and in future, it may be examined for detecting the vowel end points (VEPs) in read and conversation modes.

References

1. A. Batliner, R. Kompe, A. Kießling, E. Nöth, H. Niemann, Can you tell apart spontaneous and read speech if you just look at prosody?, in: *Speech Recognition and Coding*, Springer, 1995, pp. 321–324.

2. E. Blaauw, Phonetic characteristics of spontaneous and read-aloud speech, in: *Phonetics and Phonology of Speaking Styles*, 1991.
3. V. Dellwo, A. Leemann, M.-J. Kolly, The recognition of read and spontaneous speech in local vernacular: The case of Zurich German, *Journal of Phonetics* 48 (2015) 13–28.
4. S. M. Prasanna, S. V. Gangashetty, B. Yegnanarayana, Significance of vowel onset point for speech analysis, in: *Proceedings of International Conference on signal processing and communications*, Citeseer, 2001, pp. 81–88.
5. B. D. Sarma, S. M. Prasanna, P. Sarmah, Consonant-vowel unit recognition using dominant aperiodic and transition region detection, *Speech Communication* 92 (2017) 77–89.
6. D. J. Hermes, Vowel-onset detection, *The Journal of the Acoustical Society of America* 87 (2) (1990) 866–873.
7. S. Deb, S. Dandapat, Emotion classification using segmentation of vowel-like and non-vowel-like regions, *IEEE Transactions on Affective Computing*.
8. S. V. Gangashetty, C. C. Sekhar, B. Yegnanarayana, Detection of vowel onset points in continuous speech using autoassociative neural network models, in: *Proceedings of INTERSPEECH*, 2004, pp. 1081–1084.
9. S. V. Gangashetty, C. C. Sekhar, B. Yegnanarayana, Spotting multilingual consonant-vowel units of speech using neural network models, in: *International Conference on Nonlinear Analyses and Algorithms for Speech Processing*, Springer, 2005, pp. 303–317.
10. S. V. Gangashetty, C. C. Sekhar, B. Yegnanarayana, Extraction of fixed dimension patterns from varying duration segments of consonant-vowel utterances, in: *Proceedings of International Conference on Intelligent Sensing and Information Processing*, IEEE, 2004, pp. 159–164.
11. S. M. Prasanna, B. Yegnanarayana, Detection of vowel onset point events using excitation information, in: *Proceedings of INTERSPEECH*, 2005, pp. 1133–1136.
12. K. S. Rao, B. Yegnanarayana, Duration modification using glottal closure instants and vowel onset points, *Speech communication* 51 (12) (2009) 1263–1269.
13. S. M. Prasanna, B. S. Reddy, P. Krishnamoorthy, Vowel onset point detection using source, spectral peaks, and modulation spectrum energies, *IEEE Transactions on audio, speech, and language processing* 17 (4) (2009) 556–565.
14. A. K. Vuppala, J. Yadav, S. Chakrabarti, K. S. Rao, Vowel onset point detection for low bit rate coded speech, *IEEE Transactions on Audio, Speech, and Language Processing* 20 (6) (2012) 1894–1903.
15. S. Furui, Recent advances in spontaneous speech recognition and understanding, in: *Proceedings of ISCA & IEEE workshop on spontaneous speech processing and recognition*, 2003.
16. M. Nakamura, K. Iwano, S. Furui, Differences between acoustic characteristics of spontaneous and read speech and their effects on speech recognition performance, *Computer Speech & Language* 22 (2) (2008) 171–184.
17. J.-F. Wang, C.-H. Wu, S.-H. Chang, J.-Y. L. Lee, A hierarchical neural network model based on a c/v segmentation algorithm for isolated Mandarin speech recognition, *IEEE Transactions on Signal Processing* 39 (9) (1991) 2141–2146.
18. S. M. Prasanna, G. Pradhan, Significance of vowel-like regions for speaker verification under degraded conditions, *IEEE transactions on audio, speech, and language processing* 19 (8) (2011) 2552–2565.
19. A. Kumar, S. Shahnawazuddin, G. Pradhan, Exploring different acoustic modeling techniques for the detection of vowels in speech signal, in: *Proceedings of Twenty Second National Conference on Communication (NCC)*, IEEE, 2016, pp. 1–5.
20. A. Kumar, S. Shahnawazuddin, G. Pradhan, Improvements in the detection of vowel onset and offset points in a speech sequence, *Circuits, systems, and signal processing* 36 (6) (2017) 2315–2340.
21. B. D. Sarma, S. M. Prasanna, Analysis of spurious vowel-like regions (vLRS) detected by excitation source information, in: *Proceedings of Annual IEEE India Conference (INDICON)*, IEEE, 2013, pp. 1–5.
22. B. K. Khonglah, B. D. Sarma, S. Prasanna, Exploration of deep belief networks for vowel-like regions detection, in: *Proceedings of Annual IEEE India Conference (INDICON)*, IEEE, 2014, pp. 1–5.
23. M. Stephane, A wavelet tour of signal processing, *The Sparse Way*.

-
24. M. Madhavi, H. Patil, B. B. Vachhani, Spectral transition measure for detection of obstruents, in: Proceedings of 23rd European Signal Processing Conference (EUSIPCO), IEEE, 2015, pp. 330–334.
 25. S. Dusan, L. Rabiner, On the relation between maximum spectral transition positions and phone boundaries, in: Proceedings of INTERSPEECH, 2006, p. 1317–1320.
 26. S. Furui, On the role of spectral transition for speech perception, The Journal of the Acoustical Society of America 80 (4) (1986) 1016–1025.
 27. M. K. Reddy, K. S. Rao, Robust pitch extraction method for the hmm-based speech synthesis system, IEEE Signal Processing Letters 24 (8) (2017) 1133–1137.
 28. J. S. Garofalo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, D. S. Pallett, N. L. Dahlgren, The DARPA TIMIT acoustic-phonetic continuous speech corpus cdrom, Linguistic Data Consortium.
 29. S. S. Kumar, K. S. Rao, D. Pati, Phonetic and prosodically rich transcribed speech corpus in Indian languages: Bengali and Odia, in: Proceedings of International Conference on Oriental COCODSA held jointly with Conference on Asian Spoken Language Research and Evaluation (O-COCOSDA/CASLRE), Gurgaon, India, 2013, pp. 1–5.