

Automatic Speech Segmentation for Myanmar Language

msec. Each frame overlap by half. To reduce the edge effect of each frame segment windowing is done.

C. Speech Feature Extraction

After windowing, compute the short-time energy features and spectral centroid features of each frame of the speech signal. These features have been discussed in detail in Section 3. In this step, median filtering of these feature sequences also computed.

D. Speech Segment Detection

After computing speech feature sequences, a simple dynamic threshold-based algorithm is applied in order to detect the speech word segments.

- Compute the Mean or average values of smoothed feature sequences.
- Find the local maxima of histogram.
- If at least two maxima M_1 and M_2 have been found, then:

Threshold,

$$T = \frac{W * M_1 + M_2}{W + 1} \quad (6)$$

Otherwise,
Threshold,

$$T = \frac{Mean}{2} \quad (7)$$

Where W is a user-defined weight parameter [5], Large values of W obviously lead to threshold values closer to M_1 , Here, $W=10$.

The above process is applied for both feature sequences and finding two thresholds: T_1 based on the energy sequences and T_2 based on the spectral centroid sequences. After computing two thresholds, the speech word segments are formed by successive frames for which the respective feature values are larger than the computed thresholds (for both feature sequences).

E. Post Processing

As a post-processing step, the detected speech segments are lengthened by 5 short term window. Finally, successive segments are merged.

IV. EXPERIMENTAL STUDY

All the techniques and algorithms discussed in this paper have been implemented in Matlab. In this implement, various speech sentences in Myanmar Language have been recorded, analyzed and segmented by using time-domain and frequency-domain features with dynamic threshold technique. Fig.2 shows the filtered short-time energy and spectral centroid features of the Myanmar sentence. Fig.3 shows the segmented sub-word. The percentage of accuracy rate and failure rate of segmentation had been calculated using the following eq:

$$\text{segmentation accuracy} = \frac{\text{No.of correct words segmented by the system}}{\text{No.of words in the sentence}} * 100\% \quad (8)$$

In this experiment, example of 10 spoken sentences was allowed to segments by the system. The system output was word and sub-word. An example of Myanmar continuous sentences are given in Table 2 and the detailed segmentation results were shown in Table 3.

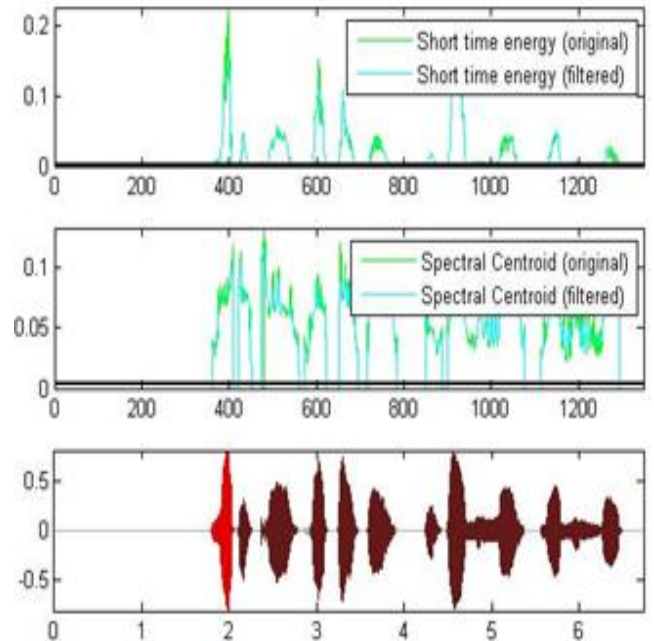


Fig.2. The first sub figure shows the sequence of the signal's energy. In the second sub figure the spectral centroid sequence is presented. In both cases, the respective thresholds are also shown. The third figure presents the whole audio signal. Red color represents the detected speech segments.

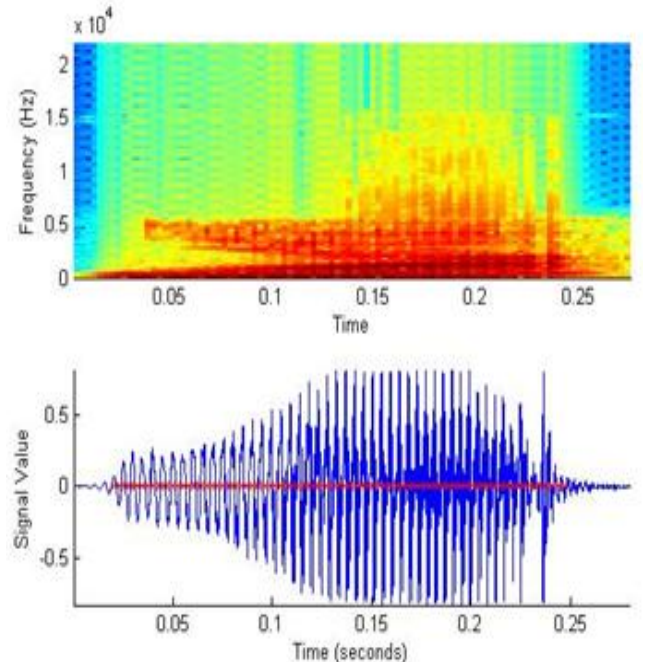


Fig.3. The first sub figure shows the spectrogram of the segment “□” and the second sub figure shows the segmented sub-word “□”.

TABLE II: Example of Myanmar Continuous Sentences

Sentence ID	Myanmar continuous sentence
S1	████████████████████ ██
S2	████████████████████ ████████████████████
S3	████████████████████ ███ █████
S4	██████████ ██████████ ████████████████████ ██████████
S5	████████████████████ ██
S6	██████ ██████████ ████████████████████
S7	████████████████████ ████████████████████
S8	████████████████████ ████████████████████ ██████████
S9	████████████████████ ██████████ ███ ████████████████████ ███
S10	████████████████████ ████████████████████

TABLE III: Segmentation Result

Sentence ID	No. of expected words in the sentence	No. of correct words segmented by the system	Accuracy rate
S1	13	11	84.6%
S2	19	17	89.4%
S3	14	12	85.7%
S4	25	22	88%
S5	14	13	92%
S6	18	16	88.9%
S7	17	16	94.1%
S8	24	20	83.3%
S9	22	20	90.9%
S10	18	15	83.3%
Total	184	162	88%

V. CONCLUSION

It has presented a simple speech features extraction approach for segmenting continuous speech into word/sub-words in a simple and efficient way. From the experiments, it was observed that some of the words were not segmented properly. This is due to different causes: (i) the utterance of words and sub-words differs depending on their position in the sentence, (ii) the pauses between the words or sub-words are not identical in all cases because of the

variability of the speech signal and (iii) the non-uniform articulation of speech. Also, the speech signal is very much sensitive to the speaker’s properties such as age, sex, and emotion. This reduces the memory requirement and computational time in any automatic speech recognition system.

VI. ACKNOWLEDGMENT

I am very grateful to Dr. K Zin Lin for fruitful discussion during the preparation of this paper and also specially thank to Rector, Professors and colleagues from Technology University (Yatanarpon Cyber City), Myanmar.

VII. REFERENCES

[1] C. T. Zhang and C. J. Kuo, “Hierarchical classification of audio data for archiving and retrieving”, In International Conference on Acoustics, Speech and Signal Processing, volume VI, pages 3001–3004. IEEE, 1999.
 [2] R. Niederjohn and J. Grotelueschen, “The enhancement of speech intelligibility in high noise level by high-pass filtering followed by rapid amplitude compression”, IEEE Transactions on Acoustics, Speech and Signal Processing, 24(4), pp277-282, 1976.
 [3] T. Giannakopoulos, “Study and application of acoustic information for the detection of harmful content and fusion with visual information” Ph.D. dissertation, Dept.of Informatics and Telecommunications, University of Athens, Greece, 2009.
 [4] C.W. James and T.W. John. “An algorithm for the machine calculation of complex Fourier series”, Mathematics of Computation: Journal Review, 19: 297–301, 1965.
 [5] T. Giannakopoulos, A. Pikrakis and S. Theodoridis “A Novel Efficient Approach for Audio Segmentation”, Proceedings of the 19th International Conference on Pattern Recognition (ICPR2008), December 8-11 2008, Tampa, Florida, USA.

Author’s Profile:

Ingyin Khaing received Bachelor of Computer Technology from University of Computer Studies, Yangon in Myanmar. She completed the Master Course from University of Computer Studies since 2008 and especially studied and finished the thesis by Digital Signal Processing. She is working now in University of Computer Studies, Maubin as at tutor under Hardware Department. Now, she is a Ph.D student in University of Technology (Yatanarpon Cyber City) near PyinOoLwin, Upper Myanmar. Her fields of interest are Digital Signal Processing, speech recognition.