

Statistical Tests for Voice Activity Detection

J.M. Górriz, C.G. Puntonet, J. Ramírez, and J.C. Segura

Facultad de Ciencias, Universidad de Granada
Fuentenueva s/n, 18071 Granada, Spain
gorriz@ugr.es

Abstract. A robust and effective voice activity detection (VAD) algorithm is proposed for improving speech recognition performance in noisy environments. The approach is based on filtering the input channel to avoid high energy noisy components and then the determination of the speech/non-speech bispectra by means of third order autocumulants. This algorithm differs from many others in the way the decision rule is formulated (detection tests) and the domain used in this approach. Clear improvements in speech/non-speech discrimination accuracy demonstrate the effectiveness of the proposed VAD. It is shown that application of statistical detection test leads to a better separation of the speech and noise distributions, thus allowing a more effective discrimination and a tradeoff between complexity and performance. The algorithm also incorporates a previous noise reduction block improving the accuracy in detecting speech and non-speech. The experimental analysis carried out on the AURORA databases and tasks provides an extensive performance evaluation together with an exhaustive comparison to the standard VADs such as ITU G.729, GSM AMR and ETSI AFE for distributed speech recognition (DSR), and other recently reported VADs.

1 Introduction

Speech/non-speech detection is an unsolved problem in speech processing and affects numerous applications including robust speech recognition [1], discontinuous transmission [2, 3], real-time speech transmission on the Internet or combined noise reduction and echo cancellation schemes in the context of telephony [4]. The speech/non-speech classification task is not as trivial as it appears, and most of the VAD algorithms fail when the level of background noise increases. During the last decade, numerous researchers have developed different strategies for detecting speech on a noisy signal [5, 6] and have evaluated the influence of the VAD effectiveness on the performance of speech processing systems. Most of them have focussed on the development of robust algorithms with special attention on the derivation and study of noise robust features and decision rules [7–9]. The different approaches include those based on energy thresholds [7], pitch detection [10], spectrum analysis [9], zero-crossing rate [3], periodicity measure [11], higher order statistics in the LPC residual domain [12] or combinations of different features [3, 2].

This paper explores a new alternative towards improving speech detection robustness in adverse environments and the performance of speech recognition systems. The proposed VAD proposes a noise reduction block that precedes the VAD, and uses Bispectra of third order cumulants to formulate a robust decision rule. The rest of the paper is organized as follows. Section II reviews the theoretical background on Bispectra analysis and shows the proposed signal model. Section III analyzes the motivations for the proposed algorithm by comparing the speech/non-speech distributions for our decision function based on bispectra and when noise reduction is optionally applied. Section IV describes the experimental framework considered for the evaluation of the proposed endpoint detection algorithm. Finally, section V summarizes the conclusions of this work.

2 Model Assumptions

Let $\{x(t)\}$ denote the discrete time measurements at the sensor. Consider the set of stochastic variables y_k , $k = 0, \pm 1 \dots \pm M$ obtained from the shift of the input signal $\{x(t)\}$:

$$y_k(t) = x(t + k) \quad (1)$$

where k is the differential delay (or advance) between the samples. This provides a new set of $2 \cdot M + 1$ vector variables $\mathbf{y}_j = \{y_j(t_1), \dots, y_j(t_N)\}$ by selecting $i = 1 \dots N$ samples of the input signal. It can be represented using the associated Toeplitz matrix:

$$T_{x(t_0)} = \begin{pmatrix} y_{-M}(t_0) & \dots & y_{-M}(t_N) \\ y_{-M+1}(t_0) & \dots & y_{-M+1}(t_N) \\ \dots & \dots & \dots \\ y_M(t_0) & \dots & y_M(t_N) \end{pmatrix} \quad (2)$$

Using this model the speech-non speech detection can be described by using two essential hypothesis(re-ordering indexes):

$$H_o = (\mathbf{y}_0 = n_0; \mathbf{y}_{\pm 1} = n_{\pm 1}; \dots; \mathbf{y}_{\pm M} = n_{\pm M}) \quad (3)$$

$$H_1 = (\mathbf{y}_0 = s_0 + n_0; \mathbf{y}_{\pm 1} = s_{\pm 1} + n_{\pm 1}; \dots; \mathbf{y}_{\pm M} = s_{\pm M} + n_{\pm M}) \quad (4)$$

where s_k 's/ n_k 's are the speech/non-speech (any kind of additive background noise i.e. gaussian) signals, related themselves with some differential delay (or advance, see section 2.1). All the process involved are assumed to be jointly stationary and zero-mean. Consider the third order cumulant function $C_{y_k y_l}$ defined as $C_{y_k y_l} \equiv E[\mathbf{y}_0 \mathbf{y}_k \mathbf{y}_l]$, and the two-dimensional discrete Fourier transform (DFT) of $C_{\mathbf{y}_k \mathbf{y}_l}$, the bispectrum function:

$$C_{\mathbf{y}_k \mathbf{y}_l}(\omega_1, \omega_2) = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} C_{\mathbf{y}_k \mathbf{y}_l} \cdot \exp(-j(\omega_1 k + \omega_2 l)) \quad (5)$$

Sampling the equation 5 and assuming a finite number of samples, the biespectrum estimate can be written as:

$$\hat{C}_{\mathbf{y}_k\mathbf{y}_l}(n, m) = \sum_{k=-M}^M \sum_{l=-M}^M C_{\mathbf{y}_k\mathbf{y}_l} \cdot w(k, l) \cdot \exp(-j(\omega_n k + \omega_m l)) \quad (6)$$

where $\omega_{n,m} = \frac{2\pi}{M}(n, m)$ with $n, m = -M, \dots, M$ are the sampling frequencies, $w(k, l)$ is the window function (to get smooth estimates [13]) and $C_{\mathbf{y}_k\mathbf{y}_l} = \sum_{i=0}^{N-1} y_0(t_i)y_k(t_i)y_l(t_i) = \mathbf{y}_0\mathbf{y}_k\mathbf{y}_l|_{t_0}$.

2.1 A Model for Speech / non Speech

The voice detection is achieved applying biespectrum function to the set of new variables detailed in the previous section. Then the essential difference between speech (s_k) and non-speech (n_k) (i.e. noise) will be modelled in terms of the value of the spectral frequency coefficients. We also assume that the noise sequences (n_k) are statistically independent of s_k with vanishing biespectra. Of course the third order cumulant sequences of all process satisfy the summability conditions retailed in [13].

The sequence of cumulants of the voiced speech is modelled as a sum of coherent sine waves:

$$C_{y_k y_l} = \sum_{n,m=1}^K a_{nm} \cos[kn\omega_0^1 + lm\omega_0^2] \quad (7)$$

where a_{nm} is amplitude, $K \times K$ is the number of sinusoids and ω is the fundamental frequency in each dimension. It follows from [12] that a_{nm} is related to the energy of the signal $\mathcal{E}_s = E\{s^2\}$. The VAD proposed in the later reference only works with the coefficients in the sequence of cumulants and is more restrictive in the model of voiced speech. Thus the Biespectrum function associated to this sequence is the DFT of equation 7 which consists in a set of Dirac's deltas in each excitation frequency $n\omega_0^1, m\omega_0^2$. Our algorithm will detect any high frequency peak on this domain matching with speech frames, that is under the above assumptions and hypotheses, it follows that on H_0 , $C_{y_k y_l}(\omega_1, \omega_2) \equiv C_{n_k n_l}(\omega_1, \omega_2) \simeq 0$ and on H_1 $C_{y_k y_l}(\omega_1, \omega_2) \equiv C_{s_k s_l}(\omega_1, \omega_2) \neq 0$. Since $s_k(t) = s(t+k)$ where $k = 0, \pm 1 \dots \pm M$, we get:

$$C_{s_k s_l}(\omega_1, \omega_2) = \mathcal{F}\{E[s(t+k)s(t+l)s(t)]\} \quad (8)$$

The estimation of the bispectrum is deep discussed in [14] and many others, where conditions for consistency are given. The estimate is said to be (asymptotically) consistent if the squared deviation goes to zero, as the number of samples tends to infinity.

2.2 Detection Tests for Voice Activity

The decision of our algorithm is based on statistical tests including the Generalized Likelihood ratio tests (GLRT) [15] and the Central χ^2 -distributed test statistic under H_0 [16]. We will call them GLRT and χ^2 tests. The tests are based on some asymptotic distributions and computer simulations in [17] show that the χ^2 tests require larger data sets to achieve a consistent theoretical asymptotic distribution.

GRLT: Consider the complete domain in bispectrum frequency for $0 \leq \omega_{n,m} \leq 2\pi$ and define P uniformly distributed points in this grid (m, n) , called coarse grid. Define the fine grid of L points as the L nearest frequency pairs to coarse grid points. We have that $2M + 1 = P \cdot L$. If we reorder the components of the set of L Bispectrum estimates $\hat{C}(n_l, m_l)$ where $l = 1, \dots, L$, on the fine grid around the bifrequency pair into a L vector β_{ml} where $m = 1, \dots, P$ indexes the coarse grid [15] and define P -vectors $\phi_i(\beta_{1i}, \dots, \beta_{Pi})$, $i = 1, \dots, L$; the generalized likelihood ratio test for the above discussed hypothesis testing problem:

$$H_0 : \mu = \mu_n \quad \text{against} \quad H_1 : \eta \equiv \mu^T \sigma^{-1} \mu > \mu_n^T \sigma_n^{-1} \mu_n \quad (9)$$

where $\mu = 1/L \sum_{i=1}^L \phi_i$ and $\sigma = 1/L \sum_{i=1}^L (\phi_i - \mu)(\phi_i - \mu)^T$ are the maximum likelihood gaussian estimates of vector $\mathcal{C} = (\mathcal{C}_{\mathbf{y}_k \mathbf{y}_l}(m_1, n_1) \dots \mathcal{C}_{\mathbf{y}_k \mathbf{y}_l}(m_P, n_P))$, leads to the activity voice detection if:

$$\eta > \eta_0 \quad (10)$$

where η_0 is a constant determined by a certain significance level, i.e. the probability of false alarm. Note that:

1. We have supposed independence between signal s_k and additive noise n_k ¹ thus:

$$\mu = \mu_n + \mu_s; \quad \sigma = \sigma_n + \sigma_s \quad (11)$$

2. The right hand side of H_1 hypothesis must be estimated in each frame (it's a-priori unknown). In our algorithm the approach is based on the information in the previous non-speech detected intervals.

The statistic considered here η is distributed as a central $F_{2P, 2(L-P)}$ under the null hypothesis. Therefore a Neyman-Pearson test can be designed for a significance level α .

χ^2 tests: In this section we consider the χ_{2L}^2 distributed test statistic[16]:

$$\eta = \sum_{m,n} 2M^{-1} |\Gamma_{\mathbf{y}_k \mathbf{y}_l}(m, n)|^2 \quad (12)$$

¹ Observe that now we do not assume that n_k $k = 0 \dots \pm M$ are gaussian

where $\Gamma_{\mathbf{y}_k\mathbf{y}_l}(m, n) = \frac{|\hat{\mathcal{C}}_{\mathbf{y}_k\mathbf{y}_l}(n, m)|}{[S_{\mathbf{y}_0}(m)S_{\mathbf{y}_k}(n)S_{\mathbf{y}_l}(m+n)]^{0.5}}$ which is asymptotically distributed as $\chi_{2L}^2(0)$ where L denotes the number of points in interior of the principal domain. The Neyman-Pearson test for a significant level (false-alarm probability) α turns out to be:

$$H_1 \quad \text{if} \quad \eta > \eta_\alpha \quad (13)$$

where η_α is determined from tables of the central χ^2 distribution. Note that the denominator of $\Gamma_{\mathbf{y}_k\mathbf{y}_l}(m, n)$ is unknown a priori so they must be estimated as the bispectrum function (that is calculate $\hat{\mathcal{C}}_{\mathbf{y}_k\mathbf{y}_l}(n, m)$). This requires a larger data set as we mentioned above in this section.

2.3 Noise reduction block

Almost any VAD can be improved just placing a noise reduction block in the data channel before it. The noise reduction block for high energy noisy peaks, consists of four stages (spectrum smoothing, noise estimation, Wiener filtering, and frequency domain filtering) and was first developed in [1].

2.4 Some remarks about the algorithm

In order to observe the potential of the proposed method we first propose an approximated decision based on an average of the components of the bispectrum (the absolute value of them). In this way we define η as:

$$\eta = \frac{1}{L \cdot N} \sum_{i=1}^L \sum_{j=1}^N |\hat{\mathcal{C}}(i, j)| = \frac{1}{L} \sum_{i=1}^L |\eta(i)| \quad (14)$$

where L, N defines the selected grid (high frequencies with noteworthy variability). We also include long term information (LTI) in the decision VAD in the on-line algorithm [18] which essentially improves the efficiency of the proposed method. According to [1], using a noise reduction block previous to endpoint detection together with a long-term measure of the noise parameters, reports important benefits for detecting speech in noise since misclassification errors are significantly reduced.

Fig. 1 shows the operation of the proposed VAD on an utterance of the Spanish SpeechDat-Car (SDC) database [19]. The phonetic transcription is: [“siete”, “θinko”, “dos”, “uno”, “otSo”, “seis”]. Fig 1(b) shows the value of η versus time. Observe how assuming η_0 the initial value of the magnitude η over the first frame (noise), we can achieve a good VAD decision. It is clearly shown how the detection tests yield improved speech/non-speech discrimination of fricative sounds by giving complementary information. The VAD performs an advanced detection of beginnings and delayed detection of word endings which, in part, makes a hang-over unnecessary. In Fig 2 we display the differences between noise and speech in general and in figure we settle these differences in the evaluation of η on speech and non-speech frames.

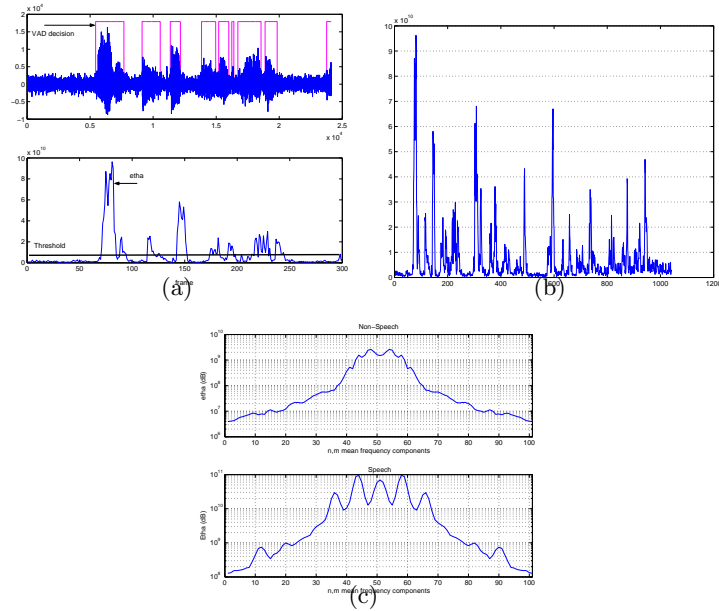


Fig. 1. Operation of the VAD on an utterance of Spanish SDC database. (a) Evaluation of η and VAD Decision. (b) Evaluation of the test hypothesis on an example utterance of the Spanish SpeechDat-Car (SDC) database [19]. (c) Speech/non-Speech η_i values for Speech-Non Speech Frames.

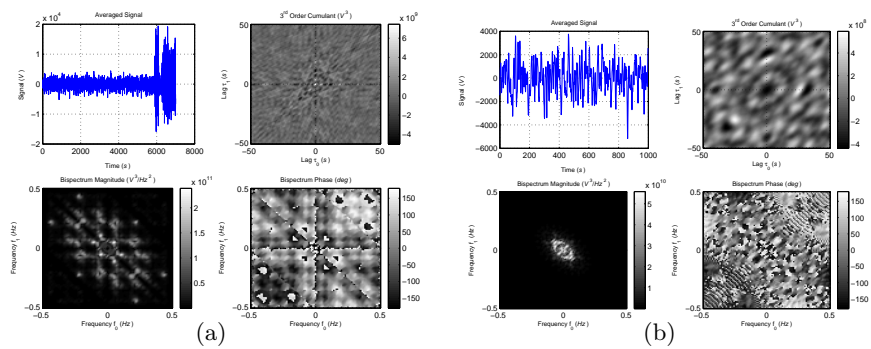


Fig. 2. Different Features (Biespectrum, magnitude and phase and 3th order cumulant over the set \mathbf{y}_k) allowing voice activity detection. (a) Features of Speech Signal. (b) Features of non Speech Signal.

3 Experimental framework

The ROC curves are frequently used to completely describe the VAD error rate. Only the AURORA subset of the original Spanish SpeechDat-Car (SDC) database [19] was used in this analysis for space reasons. The files are categorized into three noisy conditions: quiet, low noisy and highly noisy conditions, which represent different driving conditions with average SNR values between 25dB, and 5dB. The non-speech hit rate (HR0) and the false alarm rate (FAR0=100-HR1) were determined in each noise condition. These noisy signals represent the most probable application scenarios for telecommunication terminals (sub-urban train, babble, car, exhibition hall, restaurant, street, airport and train station). Fig. 3 shows the ROC curves of frequently referred algorithms [7–9, 5] for recordings from the distant microphone in quiet, low and high noisy conditions. The working points of the G.729, AMR and AFE VADs are also included. If we compare the two test discussed above we can conclude that GRLT prevails over χ^2 tests. Fig. 3 shows the ROC curves of the two proposed test varying the confidence level α (we actually vary the parameter η_α)

The results show improvements in detection accuracy over standard VADs and similarities over a representative set VAD algorithms [7–9, 5] in high noise scenario. The benefits are especially important over G.729 and over the Li’s algorithm. On average, it improves Marzinzik’s VAD that tracks the power spectral envelopes, and the Sohn’s VAD. These results clearly demonstrate that there is no optimal VAD for all the applications. Each VAD is developed and optimized for specific purposes. Hence, the evaluation has to be conducted according to the specific goal of the VAD. Frequently, VADs avoid losing speech periods leading to an extremely conservative behavior in detecting speech pauses (for instance, the AMR1 VAD). Thus, in order to correctly describe the VAD performance, both parameters have to be considered. On average the results are conclusive (see table 1).

Table 1. Average speech/non-speech hit rates for SNRs between 25dB and 5dB. Comparison of the proposed BSVAD to standard and recently reported VADs.

(%)	G.729	AMR1	AMR2	AFE (WF)	AFE (FD)
HR0	55.798	51.565	57.627	69.07	33.987
HR1	88.065	98.257	97.618	85.437	99.750
(%)	Woo	Li	Marzinzik	Sohn	χ^2 /GLRT
HR0	62.17	57.03	51.21	66.200	66.520/68.048
HR1	94.53	88.323	94.273	88.614	85.192/90.536

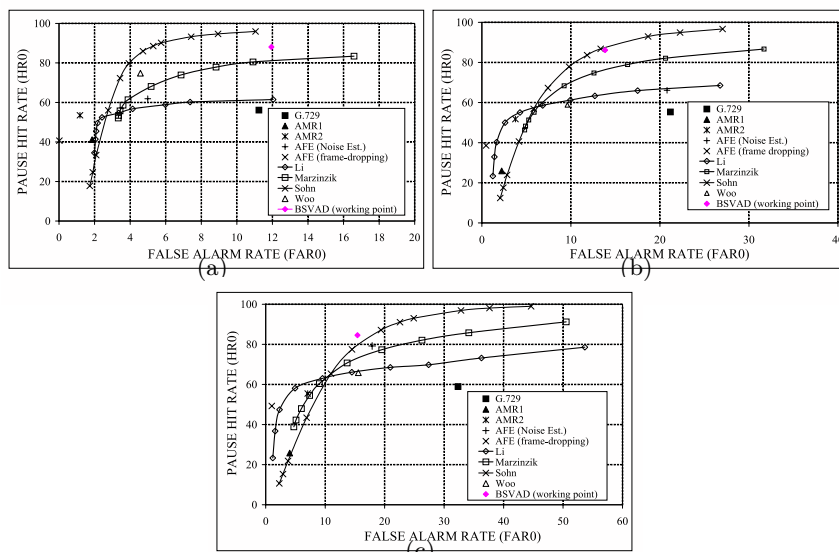


Fig. 3. ROC curves obtained for different subsets of the Spanish SDC database at different driving conditions: (a) Quiet (stopped car, motor running, 12 dB average SNR). (b) Low (town traffic, low speed, rough road, 9 dB average SNR). (c) High (high speed, good road, 5 dB average SNR).

4 Conclusion

This paper presented a new VAD for improving speech detection robustness in noisy environments. The approach is based on higher order spectra analysis employing noise reduction techniques and statistic tests for the formulation of the decision rule. The VAD performs an advanced detection using the estimated components of the Bispectrum function and robust statistical tests GLRT and χ^2 over the set of vector variables \mathbf{y}_k . As a result, it leads to clear improvements in speech/non-speech discrimination especially when the SNR drops. With this and other innovations, the proposed algorithm outperformed G.729, AMR and AFE standard VADs. We think that it also will improve the recognition rate when it was considered as part of a complete speech recognition system. The major benefit of the proposed algorithm is robustness and simplicity of the decision rule as well as the potential inclusion of the recently reported approaches for endpoint detection.

References

1. J. Ramírez, J. Segura, C. Benítez, A. delaTorre, and A. Rubio, “An effective sub-band osf-based vad with noise reduction for robust speech recognition,” *In press IEEE Transactions on Speech and Audio Processing*, vol. X, no. X, pp. X–X, 2004.
2. ETSI, “Voice activity detector (VAD) for Adaptive Multi-Rate (AMR) speech traffic channels,” *ETSI EN 301 708 Recommendation*, 1999.

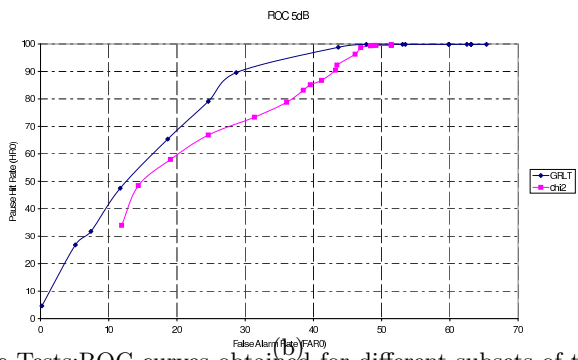
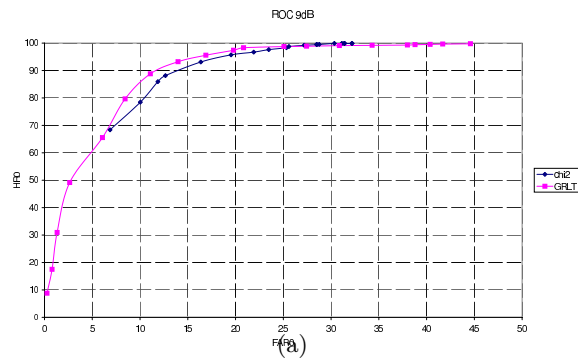


Fig. 4. Statistic Tests: ROC curves obtained for different subsets of the Spanish SDC database at different driving conditions: (a) Low (town traffic, low speed, rough road, 9 dB average SNR). (b) High (high speed, good road, 5 dB average SNR).

3. ITU, "A silence compression scheme for G.729 optimized for terminals conforming to recommendation V.70," *ITU-T Recommendation G.729-Annex B*, 1996.
4. S. Gustafsson, R. Martin, P. Jax, and P. Vary, "A psychoacoustic approach to combined acoustic echo cancellation and noise reduction," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 5, pp. 245–256, 2002.
5. J. Sohn, N. S. Kim, and W. Sung, "A statistical model-based voice activity detection," *IEEE Signal Processing Letters*, vol. 16, no. 1, pp. 1–3, 1999.
6. Y. D. Cho and A. Kondoz, "Analysis and improvement of a statistical model-based voice activity detector," *IEEE Signal Processing Letters*, vol. 8, no. 10, pp. 276–278, 2001.
7. K. Woo, T. Yang, K. Park, and C. Lee, "Robust voice activity detection algorithm for estimating noise spectrum," *Electronics Letters*, vol. 36, no. 2, pp. 180–181, 2000.
8. Q. Li, J. Zheng, A. Tsai, and Q. Zhou, "Robust endpoint detection and energy normalization for real-time speech and speaker recognition," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 3, pp. 146–157, 2002.
9. M. Marzinzik and B. Kollmeier, "Speech pause detection for noise spectrum estimation by tracking power envelope dynamics," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 6, pp. 341–351, 2002.
10. R. Chengalvarayan, "Robust energy normalization using speech/non-speech discriminator for German connected digit recognition," in *Proc. of EUROSPEECH 1999*, Budapest, Hungary, Sept. 1999, pp. 61–64.
11. R. Tucker, "Voice activity detection using a periodicity measure," *IEE Proceedings, Communications, Speech and Vision*, vol. 139, no. 4, pp. 377–380, 1992.
12. E. Nemer, R. Goubran, and S. Mahmoud, "Robust voice activity detection using higher-order statistics in the lpc residual domain," *IEEE Trans. Speech and Audio Processing*, vol. 9, no. 3, pp. 217–231, 2001.
13. C. Nikias and A. Petropulu, *Higher Order Spectra Analysis: a Nonlinear Signal Processing Framework*. Prentice Hall, 1993.
14. D. Brillinger and M. Rosenblatt, *Spectral Analysis of Time Series*. Wiley, 1975, ch. Asymptotic theory of estimates of kth order spectra.
15. T. S. Rao, "A test for linearity of stationary time series," *Journal of Time Series Analysis*, vol. 1, pp. 145–158, 1982.
16. J. Hinich, "Testing for gaussianity and linearity of a stationary time series," *Journal of Time Series Analysis*, vol. 3, pp. 169–176, 1982.
17. J. Tugnait, "Two channel tests for common non-gaussian signal detection," *IEE Proceedings-F*, vol. 140, pp. 343–349, 1993.
18. J. Ramírez, J. C. Segura, M. C. Benítez, A. de la Torre, and A. Rubio, "Efficient voice activity detection algorithms using long-term speech information," *Speech Communication*, vol. 42, no. 3-4, pp. 271–287, 2004.
19. A. Moreno, L. Borge, D. Christoph, R. Gael, C. Khalid, E. Stephan, and A. Jeffrey, "SpeechDat-Car: A Large Speech Database for Automotive Environments," in *Proceedings of the II LREC Conference*, 2000.