### 2. Manuscript (2nd revision)

Postprint of: Szwoch G., Czyżewski A., Kulesza M., A low complexity double-talk detector based on the signal envelope, Signal Processing, Vol. 88, Iss. 11 (2008), pp. 2856-2862, DOI: 10.1016/j.sigpro.2008.05.013

© 2008. This manuscript version is made available under the CC-BY-NC-ND 4.0 license

https://creativecommons.org/licenses/by-nc-nd/4.0/

A low complexity double-talk detector based on the signal envelope Grzegorz Szwoch<sup>1,2</sup> Andrzej Czyżewski<sup>1</sup> Maciej Kulesza<sup>1</sup> <sup>1</sup> Gdansk University of Technology, Multimedia Systems Department, Gdansk, Poland

<sup>2</sup> Corresponding author; e-mail: greg@sound.eti.pg.gda.pl, phone: +48 58 347 23 98 Gdansk University of Technology, Multimedia Systems Department, Narutowicza 11/12, 80-952 Gdansk, Poland

# ABSTRACT

A new algorithm for double-talk detection, intended for use in the acoustic echo canceller for voice communication applications, is proposed. The communication system developed by the authors required the use of a double-talk detection algorithm with low complexity and good accuracy. The authors propose an approach to double-talk detection based on the signal envelopes. For each of three signals: the far-end speech, the microphone signal and the echo estimate, an envelope is detected. Next, using these envelopes, a detection function is determined and compared to the threshold. Additionally, a dynamic threshold is introduced in order to improve the accuracy of the algorithm. The results of the simulations presented in the paper proved that the accuracy of doubletalk detection obtained using the proposed algorithm is higher than in the Geigel algorithm and comparable to the correlation-based methods, while the computational complexity of the proposed method remains at an acceptable level. The double-talk detection algorithm presented here may be used in voice communication systems having limited resources, allowing for accurate double-talk detection and, as a consequence, efficient acoustic echo cancellation.

## **KEYWORDS**:

Double Talk Detection; Acoustic Echo Cancellation; Speech Communication; Speech Signal Processing

10 pages + 1 page with references 5 figures 1 table No list of symbols necessary

#### INTRODUCTION

The popularity of Internet-based voice communication systems has increased significantly during the recent years. Since many users of VoIP systems prefer to use a loudspeaker instead of headphones during the communication, developers of VoIP applications have to solve the problem of an acoustic echo, occurring when a speech signal flows from the loudspeaker to the microphone and back to the sender [1]. This type of the echo is not easy to remove due to a constantly changing properties of an acoustic feedback path and the reverberations added to the signal. The task of acoustic echo cancellation is usually performed by means of an adaptive filter that simulates the acoustic feedback path and subtracts an echo estimate from the microphone signal [2,3]. The main problem of this approach is that the adaptation of filter coefficients must be stopped when the 'near speech' is present, otherwise the filter diverges and the signal quality is deteriorated significantly. Therefore, a double-talk detector (DTD) is needed in order to check whether a double-talk is present or not and to switch the filter adaptation off and on accordingly [2,4].

The authors of this paper were involved in the development of a VoIP system with the focus on speech quality improvement. Achieving this goal required using complex and computationally expensive speech coding algorithms [5,6]. In order to maintain transmission delays at a reasonable level, a low-complexity acoustic echo canceller (AEC) was needed, providing good accuracy and efficiency at the same time. The complexity and accuracy of the AEC depends greatly on the implemented DTD algorithm. Since none of the currently available DTDs met the requirements of the system, the authors decided to propose their own algorithm, details of which are presented in the The AEC based on the adaptive filter proved to be the most efficient method of acoustic echo cancellation in voice communication systems [2]. The diagram of the AEC used in the system developed by the authors is presented in Fig. 1. The far-end speech signal x(n) is distorted and attenuated by the acoustic feedback path having an unknown impulse response. Next, this signal is mixed with the near-end speaker signal and the background noise, forming the microphone signal d(n). The adaptive filter processes the far-end signal in order to produce an estimate y(n) of the

Fig. 1

acoustic echo. After subtracting the echo estimate from d(n), an error signal e(n) is created. The coefficients of the adaptive filter are modified using the values of e(n) in order to obtain a more accurate estimate of the acoustic feedback path [3]. A nonlinear processor (NLP) suppresses a residual echo. The task of the DTD is to stop the adaptation process when the near-end speech is present and restart it when there is no double-talk. A normalized least-mean square (NLS) algorithm was selected for the filter adaptation, because it provides a good trade-off between the complexity and the accuracy of the algorithm [3]. It was assumed that the AEC described above will fulfill the low complexity and good accuracy requirements, provided that the DTD will also fulfill the same requirements. Therefore, a proper choice of the DTD algorithm is crucial for achieving the desired goal.

#### DOUBLE-TALK DETECTOR

Two main approaches are typically used for the task of double-talk detection in AEC systems [7]. The first approach is based on a comparison of energy of the far-end speech and the microphone signal, while the second approach examines the correlation between the far-end speech and the microphone signal. The most popular representative of energy-based DTDs is the Geigel algorithm [8]. It is based on an observation that the energy of echo is typically much smaller than the energy of far-end speech (due to attenuation of the signal in the acoustic feedback path). Therefore, if the near-end speech is present, the energy of microphone signal increases. In the Geigel algorithm, a double-talk detection function is calculated as:

$$\xi(n) = \frac{|d(n)|}{\max\{|x(n-1)|, \dots, |x(n-L)|\}}$$
(1)

where x(n) is the far-end speech, d(n) is the microphone signal, L is a constant that determines a number of past samples of the far-end speech signal that are used for the double-talk detection. A final decision is made by comparing  $\xi(n)$  with a constant threshold value T. If  $\xi(n) > T$ , the double-talk is declared. The Geigel DTD is a low complexity algorithm, but its accuracy in case of the

acoustic echo and unstable acoustic feedback path, resulting in introduction of reverberations to the signal, was proved to be low [9]. The main problem of this approach is that it is not possible in the general case to set the constant threshold *T* that results in a low number of both false-positive and false-negative decisions. Moreover, the detection function  $\xi(n)$  changes rapidly, which makes the detection task even more problematic.

Some alternative approaches to the double-talk detection were proposed in the literature [10,11,12,13]. For example, the algorithm based on a cross-correlation utilizes the observation that the echo signal is highly correlated with the far-end speech while the correlation between the near-speech and the far-speech is low. Therefore, the double-talk detection is possible by means of calculation of a cross-correlation vector between the microphone signal and the far-end speech [12]. This approach fulfills one of the requirements (high accuracy), but it's computational cost is too high to implement this method in the system developed by the authors of this paper. Therefore, a novel algorithm for double-talk detection was proposed by the authors.

#### DTD ALGORITHM BASED ON SIGNAL ENVELOPES

The algorithm for the double-talk detection developed by the authors, shown in Fig. 2, is based on a comparison of the microphone signal energy with the energy of far-end speech, similarly to the Geigel algorithm. However, the Geigel algorithm is based on a comparison of the absolute sample values which results in a detection function that changes rapidly. More accurate results may be obtained by using the signal energy instead of the absolute sample values. The approach proposed by the authors of this paper is based on a calculation of the signal energy calculation. Various methods of the envelope detection may be used, one example is a low-pass filtering of the signal. In the proposed DTD algorithm, the envelope  $v_x(n)$  is calculated from the absolute values of x(n) samples, using the formula:

 $v_x(n) = \alpha v_x(n-1) + (1-\alpha) | x(n) | (2)$ 

Fig. 2

where  $\alpha$  is a 'forgetting factor' that defines how quickly the envelope detector reacts to rapid changes in the signal amplitude. In order to obtain an accurate estimation of the signal energy, the value of  $\alpha$  should be slightly less than one. In the experiments,  $\alpha = 0.99$  provided a sufficiently accurate energy estimation.

In order to form the decision function, two envelope detectors are needed:  $v_x(n)$  for the far-end speech signal x(n) and  $v_d(n)$  for the microphone signal d(n). The detection function is given by the formula:

$$\xi(n) = \frac{v_d(n)}{v_x(n) + \gamma}$$
(3)

The  $\gamma$  parameter is used in order to limit the values of detection function during parts of the signal containing only the noise, when values of both envelopes are low. The  $\gamma$  value should be small, in the experiments the highest detection accuracy was obtained for  $\gamma = 0.05$ .

It is assumed that the echo energy is low compared to the far-end speech energy while the nearspeech energy is significantly higher than the echo energy. Therefore, the value of the detection function will increase significantly when the double-talk occurs and it will decrease when the double-talk ends. This concept is similar to the one utilized in the Geigel algorithm. However, the detection function in the proposed method is much more smooth and slowly changing, which allows for easier detection of the signal parts that contain the double-talk. If the value of the detection function exceeds the selected threshold value T, the double-talk is detected.

The algorithm described above suffers from the same drawback as the Geigel algorithm: a choice of the threshold value allowing for an accurate double-talk detection for a wide range of signals having different energy is problematic. If the threshold is set too low, it will result in many false-positive decisions, if the threshold is set too high, it will cause many false-negative results. In order to avoid the dependency of the DTD accuracy on the threshold selection, the authors proposed a method of dynamic threshold setting, at the cost of introducing one additional envelope detector to the algorithm. The concept of the DTD with the dynamic threshold is based on an observation that in

the presence of double-talk, the energy of microphone signal is considerably larger than the echo energy. In order to make use of this observation, the estimated echo signal produced by the converged adaptive filter in the AEC may be used as an approximation of a real echo. The envelope of the echo estimate y(n) is compared to the envelope of the microphone signal d(n). If the envelope of d(n) is significantly higher than envelope of y(n), it indicates the presence of the double-talk. However, the DTD based directly on this idea proved to be inaccurate in the experiments. A modified approach that provided much more accurate double-talk detection was based on using the estimated echo signal envelope  $v_y(n)$  in the calculation of the dynamic detection threshold, according to the formula:

$$T(n) = \frac{v_y(n)}{v_y(n) + \gamma} + \beta$$
(4)

which is almost identical to the formula (3), except that the envelope of the echo estimate replaces the envelope of microphone signal. The parameter  $\beta$  is used in order to leave some margin for the detection error and it should be a small positive value.

The modification described above removes the need of setting the constant threshold value for the detection. The dynamic threshold adapts to the changes in echo signal envelope in real time. If the double-talk is present, the envelope of microphone signal is greater than the envelope of estimated echo. Therefore, the double-talk may be detected by comparing the current value of detection function to the current dynamic threshold value. The far-end speech envelope  $v_x(n)$  may be treated as a normalization term that improves the detection accuracy comparing to the original version of the algorithm described earlier.

It has to be noted that the accuracy of the presented DTD depends on the accuracy of echo estimation produced by the adaptive filter, therefore the calculation of dynamic threshold should not be performed until the filter convergence is finished. Therefore, the threshold is initially set to a constant value  $T_{init}$ . All three envelopes are calculated starting from the first sample. When the filter adaptation is finished, the dynamic threshold is calculated and used instead of  $T_{init}$ . Additionally, it is convenient to set lower and upper bounds on the dynamic threshold ( $T_{min}$  and  $T_{max}$ ) that help to avoid detection errors in parts of the signal containing only noise. To summarize the proposed algorithm, a double-talk is detected if the initial filter convergence period has passed and:

$$\xi(n) > T, \text{ where } T = \begin{cases} T(n), & T_{\min} < T(n) < T_{\max} \\ T_{\min}, & T(n) < T_{\min} \\ T_{\max}, & T(n) > T_{\max} \end{cases}$$
(5)

where T(n) is the dynamic threshold given by (4).

#### EVALUATION OF THE DOUBLE-TALK DETECTOR PERFORMANCE

The proposed DTD based on signal envelopes, described in the previous section, was implemented in the AEC based on the NLS algorithm, and its performance was evaluated. For the assessment of the DTD accuracy, a method proposed by Cho et al [9] was used. However, Cho et al tested various DTD algorithms using a set of studio recordings with artificially added reverberation, delay and noise. The authors of this paper used speech signals recorded in an actual Internet-based speech communication system, in order to measure the DTD performance in conditions similar to those existing in real communication systems. The delay between the echo and the far-end speech in the recordings was about 3.375 ms and the delay on the far-end of the system was about 273.25 ms, which resulted in an audible echo effect. For the double-talk detection, the near-end speech was artificially added to the recorded echo signal.

Fig. 3 shows an example of double-talk detection results obtained using the DTD presented in this paper. It can be seen that identification of the double-talk section in the signal was correct. There is a noticeable delay between the beginning of the actual double-talk section and the beginning of the detected double-talk. However, this delay is not long enough to cause a divergence of the adaptive filter. The dynamic threshold function, shown in Fig. 3, allows for an accurate detection of both the beginning and the end of the double-talk. The experiment was repeated for other test signals and similar results were obtained.

Fig. 3

A quantitative assessment of the envelope-based DTD was performed using the method proposed by Cho et al. [9], with some modifications. Two main parameters that describe the accuracy of DTD are a probability of miss  $P_m$  (a false negative decision) and a probability of false alarm  $P_f$  (a false positive decision). Of these two parameters, a low value of  $P_m$  is much more crucial, because a missed detection of the double-talk causes a divergence of the adaptive filter and a deterioration of the signal quality. In order to measure the DTD performance, averaged values of  $P_m$  and  $P_f$  were calculated for different ratios of the near-end speech energy to the far-end speech energy (near-tofar ratio, NFR) and for different test signals. Additionally, the same procedure was applied to the standard Geigel DTD and to the DTD based on the normalized cross-correlation (NCC) [9,12] in order to do a comparative analysis of all three DTD algorithms. For each DTD algorithm, 'optimal' values of the parameters were found during the experiments so that the best trade-off between  $P_m$ and  $P_f$  is obtained. For the envelope-based DTD, value  $\beta = 0.02$  was selected. The results of the assessment of all DTDs are plotted in Figures 4 and 5, showing probability of miss and probability of false alarm, respectively, as a function of NFR.

In order to obtain more detailed information on the performance of the envelope-based algorithm, three parameters were calculated for each tested DTD, as well as for the AEC without a DTD and for a signal not processed by the AEC. The test set contained four studio recordings of Polish sentences, two by a male speaker and another two by a female speaker, having length of 80 to 120 seconds. The averaged results of this experiment, shown in Table 1, allow for comparison of the proposed algorithm performance with other DTDs. A signal to noise ratio (SNR) was calculated as a ratio of the power of the near-end speech signal section to the power of the signal part containing only the echo. An echo return loss enhancement (ERLE) was calculated as a ratio of the power of the averaged signal (containing the echo) to the power of the signal processed by the AEC (after the echo cancellation). A PESQ (Perceptual Evaluation of Speech Quality) score is a measure of quality of the double-talk part of the signal, ranging from 4.5 (the highest possible quality) to 0 (the worst quality). Signal distortions caused by an incorrect double-talk detection decrease the PESQ

value. This measure was computed using the Opticom OPERA software, a clean near-end speech was used as a reference signal.

#### DISCUSSION

The accuracy of the DTD is measured in terms of the probability of miss and the probability of false detection. Both values should be as low as possible. The results of the simulations show that the proposed algorithm yields a significant improvement in the double-talk detection accuracy over the Geigel DTD. Introduction of the dynamic threshold that adapts to the signal removes the most important problem of the Geigel algorithm – the dependency of the DTD accuracy on the threshold setting. Using slowly changing envelopes instead of actual sample values in the calculation of detection functions simplifies the task of automatic selection of the dynamic threshold. Two parameters,  $\beta$  and  $\gamma$ , have to be set in the envelope-based DTD. They can be seen as tune-up settings that allow for decreasing  $P_f$  at the cost of some increase in  $P_m$ .

From Figs 4 and 5 it can be clearly seen that the accuracy of the envelope-based DTD for all NFR values is higher than in the Geigel algorithm and comparable to the NCC method. Although  $P_f$  increases with the NFR in the envelope DTD, it still remains at a sufficiently low level. Moreover, higher  $P_f$  values do not have as significant impact on DTD performance as higher  $P_m$  values. The values of the parameters shown in Table 1 indicate that the quality of the signal processed with the AEC and the envelope DTD is comparable to the quality obtained using the NCC method and higher than in both the Geigel algorithm and the unprocessed signal.

Regarding the computational complexity, the envelope-based DTD with the dynamic threshold requires a calculation of three envelopes, three absolute values, the detection function and the dynamic threshold, as well as a comparison of the detection function with the threshold. Each envelope detector requires two multiplications and two addition/subtraction operations. In total, the proposed DTD requires eight multiply/divide operations and nine add/subtract operations. Therefore, the computational complexity of the envelope-based DTD is higher than in the Geigel

DTD, but it is still acceptable. In contrast to that, the original NCC method requires the computation of two correlation vectors and an autocorrelation matrix that has to be inverted [9], therefore this method is impractical due to a high computational cost. A simplified method [11] was used in the experiments described in this paper, however, this 'fast' method still requires (5L + 4) multiplications, where *L* is the length of the analysis window (L = 1024 in the simulations). Therefore, the computational cost of the proposed algorithm is considerably lower than in the NCC method.

From Fig. 3 it can be seen that there is a delay between the end of the actual double-talk section and the end of the double-talk detected by the envelope-based DTD. This effect is caused by the envelope detector that is not able to react to rapid decreases in signal amplitude without a delay. However, this effect does not necessarily decrease the accuracy of echo cancellation, because in most AEC systems, a period of latency is introduced so that the filter adaptation is restarted only if the double-talk was not found during this period. In the proposed algorithm, there is no need to introduce this additional latency, because the delay caused by the long decay time of the envelope detector takes care of this problem. Therefore, the slow reaction time for a rapid decrease of the signal amplitude in the proposed DTD does not reduce the accuracy of double-talk detection in a significant way.

One more important thing to note is that the accuracy of double-talk detection in the proposed algorithm depends on the accuracy of echo estimation in the adaptive filter. If the echo estimate is inaccurate, the dynamic threshold value will be inappropriate, which may result in wrong DTD decision. In return, a wrong DTD decision may cause an unneeded switching of the adaptation filter state, which results in wrong echo estimation. Consequently, this feedback loop may cause the divergence of the filter and the deterioration of the signal quality. However, this effect was not observed during the experiments. The reaction time of the DTD on the double-talk occurrence is short enough and the risk of filter divergence due to a wrong DTD decision is low as long as the  $\beta$  parameter is set properly.

The results of the comparative analysis presented in Table 1 show that the accuracy of echo cancellation in the AEC utilizing the proposed DTD algorithm is comparable to the cross-correlation algorithm and significantly higher than in the Geigel algorithm. In terms of quality of the processed signal, the results of the PESQ test prove that the proposed algorithm provides an acceptable trade-off between the low complexity of the Geigel algorithm and the high accuracy of the cross-correlation method. Hence, the aim of the research was successfully achieved.

# CONCLUSIONS

The goal of the work was to develop a double-talk detector suitable for use in the AEC in speech communication systems, that provides both low computational complexity and good accuracy of the double-talk detection. The approach presented in this paper is based on the detection of signal envelopes of the far-end speech, microphone and echo estimate signals. It can be seen as an expansion of the original Geigel algorithm. However, using a signal envelope instead of absolute values of signal samples, together with the introduction of the dynamic threshold, resulted in an improved accuracy of double-talk detection.

The simulations of the acoustic echo cancellation, performed using the test signals recorded in the real speech communication system, proved that the proposed DTD algorithm provides good accuracy of double-talk detection, comparable to the correlation-based method and higher than in the Geigel algorithm. It is also possible to tune up the performance of the algorithm by changing the value of  $\beta$  parameter. The computational complexity of the proposed approach, although higher than in the Geigel DTD, is still low, allowing for implementation of this algorithm in the developed communication system. The complexity of the presented method is considerably lower than in the NCC algorithm, while the accuracy of both methods and the processed signal quality are similar. As a conclusion, it may be stated that the DTD algorithm presented here fulfills the requirements of the voice communication system that is being developed by the authors. Further experiments will focus

on the implementation and testing of the proposed DTD algorithm in the physical system and also on further modifications of the algorithm, improving its performance.

# ACKNOWLEDGEMENTS

Research funded by the Polish Ministry of Education and Science within the Grant No. 3 T11D 004

REFERENCES

International Telecommunication Union, Acoustic echo controllers, ITU-T Recommendation
G.167, 1993.

[2] S.M. Kuo, B.H. Lee, W. Tian, Real-Time Digital Signal Processing: Implementations and Applications, John Willey & Sons, 2006, Chapter 10: Adaptive Echo Cancellation.

[3] S. Haykin, Adaptive filter theory, 4th edition, Prentice Hall Inc, 2002.

[4] S.V. Vaseghi, Advanced Digital Signal Processing and Noise Reduction, John Willey & Sons,2000, Chapter 14: Echo cancellation.

[5] M. Kulesza, A. Czyżewski, Speech Codec Enhancements Utilizing Time Compression and Perceptual Coding, Proc. 112th Audio Eng. Conv, Vienna, 5-8 May 2007, preprint No. P-3.

[6] G. Szwoch, M. Kulesza, A. Czyżewski, Transient detection for speech coding applications, Int.Journ. of Computer Science and Network Security 6 (12) (2006) 320-325.

[7] T. Vu, H. Ding, M. Bouchard, A Survey of Double-talk Detection Schemes for Echo Cancellation Applications, Canadian Acoustics 32 (3) (2004) 144-145.

[8] A. Adrian, Voice over Internet Acoustic Echo Cancellation, DFS Deutsche Flugsicherung, 2004, Internet: http://home.arcor.de/andreadrian/echo\_cancel/.

[9] J. Cho, D. Morgan, An Objective Technique for Evaluating Doubletalk Detectors in Acoustic Echo Cancelers, IEEE Trans. Speech Audio Processing 7 (6) (1999) 718-724.

[10] T. Gaensler, M. Hansson et al., A double-talk detector based on coherence, IEEE Trans. on Communications, 44 (11) (1996) 1421-1427.

[11] T. Gaensler, J. Benesty, The fast normalized cross-correlation double-talk detector, Signal Processing, 86 (6) (2006) 1124-1139.

[12] J. Benesty, D.R. Morgan, J.H. Cho, A New Class of Doubletalk Detectors Based on Crosscorrelation, IEEE Trans. Speech Audio Processing 8 (3) (2000) 168-172.

[13] H. Ye, B. Vu, A new double-talk detection algorithm based on the orthogonality theorem, IEEETrans. on Communications 39 (11) (1991) 1542-1545.

# TABLE CAPTIONS

Table 1. Comparison of the averaged quality measures of the signal processed by the AEC based on the NLS adaptive filter, with various DTD algorithms

## FIGURE CAPTIONS

Fig. 1. Block diagram of the acoustic echo canceller with the double-talk detector (DTD)

Fig. 2. Block diagram of the acoustic echo canceller with the proposed double-talk detector based on signal envelopes

Fig. 3. The example of the double-talk detection obtained in the simulations using the proposed algorithm The dashed vertical lines mark the beginning and the end of the inserted double-talk signal. The light gray background marks the signal part detected as a double-talk by the proposed algorithm. The thin line is the detection function, the thick line is the dynamic threshold function

Fig. 4. The measurement results of the probability of miss  $(P_m)$  as a function of NFR, for the proposed (EnvDTD), Geigel and normalized cross-correlation (NCC) algorithms

Fig. 5. The measurement results of the probability of false alarm ( $P_f$ ) as a function of NFR, for the proposed (EnvDTD), Geigel and normalized cross-correlation (NCC) algorithms

Table 1. Comparison of the averaged quality measures of the signal processed by the AEC based on the NLS adaptive filter, with various DTD algorithms

| DTD algorithm                      | SNR [dB] | ERLE [dB] | PESQ score |
|------------------------------------|----------|-----------|------------|
| Proposed algorithm (envelope)      | 68.12    | 32.26     | 3.71       |
| Normalized cross-correlation (NCC) | 69.98    | 34.02     | 3.86       |
| Geigel                             | 44.59    | 8.38      | 3.63       |
| No double-talk detection           | 36.63    | 7.06      | 2.91       |
| No echo cancellation               | 36.36    | 0         | 3.29       |















