

## NOISE REDUCTION BY SPECTRAL SUBTRACTION

Authors: Felix ALBU, Neculai DUMITRIU, Livia Doina STANCIU  
Institution: University "POLITEHNICA" Bucharest  
Faculty of Electronics and Telecommunications  
Address: B-dul Iuliu Maniu nr 1-3, sector 6, Bucuresti  
Tel : 410 54 00

Spectral subtraction is performed over frames by obtaining the short-term magnitude spectrum of the noisy speech, subtracting an estimated noise magnitude spectrum from the noisy spectrum, and inverse transforming this spectral amplitude using the phase of the noisy speech. The techniques presented in this paper differ by the window used for frame analysis, the power exponent, the weighted subtraction coefficient. The results have been compared by informal listening and SNR computations.

### Introduction

In practice, recorded or transmitted speech contain an amount of noise. The noise passes a disproportionate amount of linguistic information. Because vowels pass larger amounts of energy, broadband noise degradation tends to mask consonant sections more than voiced, thus causing decreased intelligibility. This uncorrelated background noise has a spectral density depending on the recording conditions during recording. However, in most cases the additive noise is simulated as white. The background noise causes a degradation of speech which leads to total unintelligibility. Two approaches dominate in the literature: the first is based on the use of spectral subtraction, and the second on appropriate filters for noise reduction [2].

Spectral subtraction is a family of frequency-domain noise reduction techniques based on direct estimation of the short-term spectral magnitude. In this approach, speech is modeled as a random process to which uncorrelated random noise is added. It is assumed that the noise is short term stationary, with second-order statistics estimated using silent frames. The estimated noise power spectrum is subtracted from the transformed noisy input signal. A generalized estimator is given by:

$$\hat{S}_s(\omega, m) = \left[ |S_y(\omega, m)|^2 - k \cdot |\hat{S}_d(\omega, m)|^2 \right]^{1/2} \cdot e^{j\varphi_s(\omega, m)} \quad (1)$$

where  $|S_y(\omega, m)|$  and  $\varphi_s(\omega, m)$  are obtained from the short-term Discrete Time Fourier transform (stDTFT) of the noisy speech frame;  $|\hat{S}_d(\omega, m)|$  is the short-term magnitude spectrum of the noise that must be updated during the absence of the speech;  $\hat{S}_s(\omega, m)$  is

the estimated stDFT of the frame of speech from the spectral subtraction coefficient, the power exponent and  $k$  is the weighted subtraction coefficient

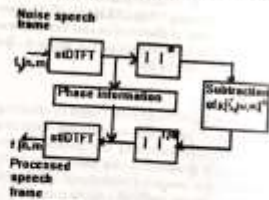


Fig. 1. Generalized spectral subtraction.

The resultant spectrum is converted to the speech signal by using the phase information of the original degraded speech.

The speech waveform was sampled at a rate of 8000 Hz. The 180 samples of the current frame were overlapped with the 76 trailing samples of the previous frame through trapezoidal windowing (Fig. 2).

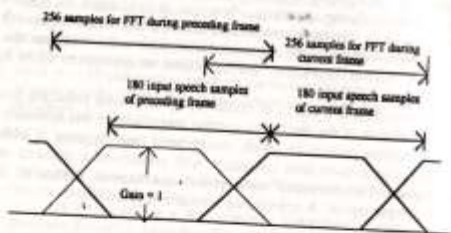


Fig. 2. Frame overlapping.

We chose a frame of 76 samples because the resulting 256 samples permit the use of a standard FFT for the time-to-frequency transformation. The speech samples are overlapped so that the sum of amplitude weights of the overlapped windows is unity. As indicated in Fig. 2, a portion of the speech samples is used twice in the time-to-frequency transform; the trailing end of the preceding FFT and the leading edge of the current FFT frame. Then, a two-way transform involving FFT and inverse FFT

undesirable audio effects are generated such as: clicking, popping, and distortion of the speech signal. We propose another window (composed of two overlapping sigmoidal windows) given by:

$$\begin{aligned}
 & 1 \leq n \leq L \\
 & L+1 \leq n \leq N-L \\
 & N-L+1 \leq n \leq N
 \end{aligned} \quad (2)$$

where  $L$  is the length of the window,  $p$  is the slope of the sigmoidal function, and  $L$  is the overlap between the two windows. A normalized magnitude in dB for the composed sigmoidal windows ( $p = 0.2$ ) is shown in Fig. 3b. For any window there are two desirable features:

- 1. narrow mainlobe width;
  - 2. high attenuation for the sidelobes;
- The window has a relatively narrow mainlobe width (Fig. 3b) which decreases as the slope  $p$  increases. The sidelobe attenuation is not dependent on the slope  $p$  and the overlap  $L$ , and the sidelobe attenuation is not dependent on the overlap  $L$ .

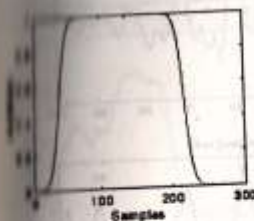


Fig. 3a. Composed sigmoidal window ( $N=256$ ,  $L=76$ ,  $p=0.2$ ).

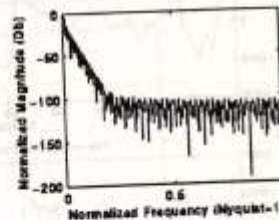


Fig. 3b. Normalized magnitude in dB for the composed sigmoidal window ( $N=256$ ,  $L=76$ ,  $p=0.2$ ).

From (1) it can be noticed that the estimated speech magnitude spectrum is not guaranteed to be positive. Different systems remedy this by performing half-wave rectification or full-wave rectification. We used half-wave rectification (i.e. negative magnitudes are set to zero). Forcing negative spectral magnitude values to zero, however,

introduce a "musical" tone artifact in the reconstructed speech. The major limitation of spectral subtraction techniques is the improvement at low SNR (Fig. 5) but introduced a "musical" reconstructed speech.

Both windows (trapezoidal window and composed sigmoidal window) show comparable performance in removing of undesirable audio effects in terms of SNR for presented windows shows comparable results with original for voiced frames in case of composed sigmoidal window with low SNR. Our results showed that if the weighted subtraction coefficient  $k$  and maximum spectral floor based on the estimated input SNR. The optimal must be between 1 and 3, with lower value for augmented segmental SNR, the performance of the system decreases.

The results showed that spectral subtraction might increase the processed signal for a large range of the parameter  $a$ . The processed signal equal to 1 or 0.5 sounded "less noisy" at relatively high SNR. The results show an improvement of SNR especially for unvoiced frames.

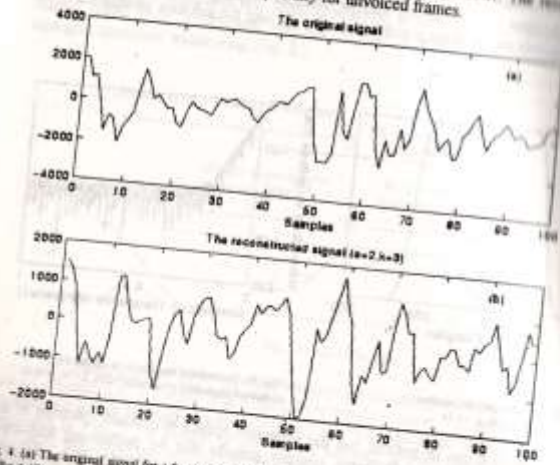


Fig. 4. (a) The original signal for a frame of the phoneme "A" (b) The reconstructed signal ( $a=2, k=3, p=1, 2, \text{global SNR} = 0 \text{ dB}$ )

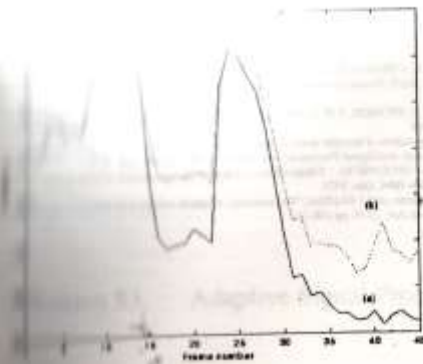


Fig. 5. The segmental SNR for the various frames of the processed Romanian word "SASE" with trapezoidal window (solid line) and composed sigmoidal window (dashed line) ( $k=1, p=1.0, \text{Global SNR} = 0 \text{ dB}$ )

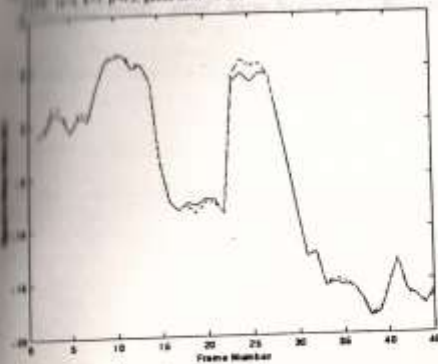


Fig. 6. The segmental SNR for the various frames of a processed Romanian word "SASE" with trapezoidal window (solid line) and composed sigmoidal window (dashed line) ( $k=1, p=1.0, \text{Global SNR} = 0 \text{ dB}$ )



## REFERENCES

- [ 1 ] G. S. KANG and L. J. FRANSEN, "Quality Improvement of LPC-Processed Noisy Speech by Spectral Subtraction", *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 37, pp. 1113-1120, June, 1989
- [ 2 ] J. R. DELLER, J. G. PROAKIS, J. H. L. HANSEN, "Discret-Time Processing of Speech Signals", Wiley, New York, 1993
- [ 3 ] S. F. BOLL, "Suppression of acoustic noise in speech using spectral subtraction", *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. ASSP-27, pp. 113-120, Apr. 1979
- [ 4 ] J. S. LIM and A. V. OPPENHEIM, "Enhancement and bandwidth compression of noisy speech", *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 67, pp. 1586-1604, Dec. 1979
- [ 5 ] M. Berouti, R. Schwartz, and J. Makhoul, "Enhancement of speech corrupted by acoustic noise", *ICASSP Record*, Apr. 1979, pp 208-211.