# Analysis of speech signal restored from the pitch harmonics

Snejana G. Pleshkova-Bekyarska[1] and Michail B. Momcshedjikov[2]

*Abstract*— **The most known speech coding method using decoding as synthese of the speech signal with the coded speech characteristics. One of these characteristics is the pith frequency of the voice sounds. It is very important to transmit the values of the pitch frequency and its harmonics with the minimum numbers of bits. The goal of this article is to make an analysis of the restored speech signal with the pitch harmonics and to make an estimation of the distortion in the restored speech signal.**

*Keywords*— **Speech analysis, speech coding, speech synthesis, speech coding distortions**

## I. INTRODUCTION

The pitch frequency determination can be made in the time or frequency domain [1]. Each of these methods give some precision for the current values of the pitch frequency. In this article the goal is to analyze the restoration of speech signal from the pitch harmonics and not the goal to study the pitch frequency determination methods.

## II. SPEECH REPRESENTATION WITH PITCH HARMONICS

It is chosen in this article to work with the pitch frequency determination method in frequency domain [2]:

$$E(\omega_0) = \frac{1}{2\pi} \int_{-\pi}^{+\pi} \left| S_\omega(\omega) - S_\omega(\overset{\wedge}{\omega}, \omega_0) \right|^2 d\omega, \qquad (1)$$

where:

$\omega_0$ is the pitch frequency;

$S_\omega(\omega)$ and $S_\omega(\omega,\omega_0)$ - the spectrum of the original and restored speech signals, respectively, with the use of window function $\omega$.

The determination of the pitch frequency $\omega_0$ is used in the speech coding methods with restoration of the speech signal from the harmonics of the pitch frequency [3]:

$$A_l(\omega) = \frac{\sum\limits_{n=a_l}^{b_l} S_\omega(n) . W_R(n)}{\sum\limits_{n=a_l}^{b_l} |W_R(n)|^2}, \qquad (2)$$

[1]Snejana G. Pleshkova-Bekyarska is with the Faculty of Communications and Communications Technologies, Technical University, Kliment Ohridski 8, 1000 Sofia, Bulgaria, E-mail: snegpl@vmei.acad.bg

[2]Michael B. Momschedjikov is with the Faculty of Communications and Communications Technologies, Technical University, Kliment Ohridski 8, 1000 Sofia, Bulgaria, E-mail: snegpl@vmei.acad.bg

where:

$\grave{A}_l(\omega)$ is $l^{th}$ harmonics of the pitch frequency $\omega_0$;

$S_\omega(n)$ – the spectrum of a frame of speech signal;

$W_R(n)$ – the window function for separating the $l^{th}$ harmonic of the pitch frequency from of the whole spectrum $S_\omega(n)$;

$n=a_l \div b_l$ – frequency interval around the $l^{th}$ harmonic of the pitch frequency, chosen very narrow to determine with a good precision the average amplitude of the $l^{th}$ harmonic of the pitch frequency.

Using the values $\grave{A}_l(\omega_0)$ of harmonics of pitch frequency from (2) in the synthesis of restored speech signals in the decoders must be considered with the characteristics of the speech signals and the chosen method for coding characteristics determination in the domain of a sequence of frames with $20 \div 30mS$ duration [1]. It is not real to decide that the pitch frequency is constant even if there is a case of a speech signal for a separation voiced letter, with the same values for the frequency and phase in the two adjacent frames.

## III. SPEECH FRAME TRANSITION CASES

It is necessary to consider in this analysis of the distortions the following cases of transition between two frames j-1 and j in accordance with values of logical signal v/uv for voiced (v/uv=1) and unvoiced (v/uv=0) speech signals:
the transition between two frames, belonging to unvoiced speech:

$$v/uv(j-1) = 0, \quad v/uv(j) = 0; \qquad (3)$$

- the transition between two frames, belonging to voiced speech:

$$v/uv(j-1) = 1, \quad v/uv(j) = 1; \qquad (4)$$

the transition between two frames, first frame belonging to unvoiced speech and the second – to voiced speech:

$$v/uv(j-1) = 0, \quad v/uv(j) = 1; \qquad (5)$$

or vice versa:

$$v/uv(j-1) = 1, \quad v/uv(j) = 0. \qquad (6)$$

For the case (4) the transition between two frames, belonging to voiced speech, can be define two conditions:
the speech signal for one voiced letters:

$$\omega_0(j-1) - \omega_0(j) \leq \Delta\omega_0; \qquad (7)$$

the speech signal for two difference voiced letters:

$$\omega_0(j-1) - \omega_0(j) \geq \Delta\omega_0, \qquad (8)$$

where:

$\Delta\omega_0$ is a predefined value of the difference between $\omega_0(j-1)$ and $\omega_0(j)$ in frame (j-1) and j, respectively.

## IV. RESTORED SPEECH DISTORTION ANALYSIS

In the condition (7) is satisfied, then the synthesis of the speech signal with pitch harmonics must be done with the concordance of the amplitudes $a_l(n)$ and phases $\theta_l(n)$ of the pitch harmonics in the two adjacent frames. It is used very often the methods of interpolation between amplitudes $M_l(j-1)$ and $M_l(j)$ and phases $\varphi_l(j-1)$ and $\varphi_l(j)$ for two adjacent frames:

$$a_l(n) = f_a(M_l(j-1), M_l(j)) \qquad (9)$$

$$\theta_l(n) = f_\theta(\varphi_l(j-1), \varphi_l(j)), \qquad (10)$$

where:

- $f_a$ and $f_\theta$ are respectively the functions chosen for interpolation of amplitudes and phases of each $l^{th}$ harmonic in two adjacent frames ;

- $M_l(j-1)$, $M_l(j)$, $\varphi_l(j-1)$ and $\varphi_l(j)$ – the amplitudes and phases of $l^{th}$ pitch harmonic in the transition between two adjacent frames (j-1) and j.

The amplitudes $M_l(j-1)$ and $M_l(j)$ and phases $\varphi_l(j-1)$ and $\varphi_l(j)$ can be determined from $\grave{A}_l(\omega)$ (2) – the complex value of $l^{th}$ pitch harmonic in frame (j-1) and j, respectively.

The most known speech coding methods, used the linear interpolation for the functions $f_a$ and $f_\theta$ in (9) and (10):

$$a_l(n) = M_l(j-1) + [M_l(j) - M_l(j-1)]\frac{n}{N} \qquad (11)$$

$$\begin{aligned} \theta_l(n) = \varphi_l(j-1) + [l\omega_0(j-1) + \Delta\omega_l(j)]n + \\ + [\omega_0(j) + \omega_0(j-1)]\frac{l.n^2}{N}, \end{aligned} \qquad (12)$$

where:

$\Delta\omega_l(j)$ is the variation of the frequency:

$$\Delta\omega_l(j) = \frac{1}{N}\left(\Delta\Phi_l(j) - 2\pi\left[\frac{\Delta\Phi_l(j) + \pi}{2\pi}\right]\right) \qquad (13)$$

$\Delta\Phi_l(j)$ – the variation of the phase angle:

$$\Delta\Phi_l(j) = \Phi_l(j) - \Phi_l(j-1) - [\omega_0(j-1) + \omega_0(j)]\frac{l.N}{2}. \quad (14)$$

The expressions (11), (12), (13) and (14) guarantee the continuos variations of the pitch harmonics in the order of two adjacent frame (j-1) and j.

For the case (7) the restored signal for the voiced speech $S_v(n)$ is:

$$S_v(n) = \sum_{l=1}^{L} a_l(n)\cos[\theta_l(n)], \qquad (15)$$

where:

L is the chosen maximum number of pitch harmonics for the synthesis of speech signal.

For the case (8) it is necessary to use the overlap-add method for concordance of the amplitudes and phases in two adjacent frames (j-1) and j. This mean that in the place of transition the restored speech signal is the sum of two signals $S_{v1}(n)$ and $S_{v2}(n)$, respectively for the pitch harmonics of $\omega_0(j-1)$ the voiced speech in the (j-1) frame and for the pitch harmonics of $\omega_0(j)$ for the voiced speech in the j frame:

$$S_v(n) = S_{v1}(n) + S_{v2}(n), \qquad (16)$$

where:

$$S_{v1}(n) = \sum_{l=1}^{L} \omega_s(n) M_l(j-1)\cos[\omega_0(j-1)nl + \Phi_l(j-1)] \qquad (17)$$

$$S_{v2}(n) = \sum_{l=1}^{L} \omega_s(n-N) M_l(j)\cos[\omega_0(j)(n-N) + \Phi_l(j)], \qquad (18)$$

where:

$\omega_s(\ldots)$ – window function for the synthesized speech.

It is seen from the expressions (15) and (16), that in the place of transition between two adjacent frames of speech signal, it is possible to have the distortions from the variations of the phase of synthesized speech for the case of the single variations of pitch frequency $\omega_0(j-1) \approx \omega_0(j)$ in two adjacent frames or from the variations of the amplitude and phase of the synthesized speech signal for the strong variations of the pitch $\omega_0(j-1) \neq \omega_0(j)$ in two adjacent frames. The quantity estimation of these of these distortions must be done from the expression (11) for the amplitude distortions (13) for the frequency distortions and (12,14) for the phase distortions.

## V. SIMULATION RESULTS

Some of the simulation results for different methods for speech signals synthesis in the place of transition between two adjacent frames are shown for visual impression from Fig.1 to Fig.6.

On the all figures are shown respectively the time and spectrum graphics in the following sequence: original and restored signal. The time interval in the time graphics is chosen such, that to display more then one frame for analysis or synthesis of speech signal. With an arrow mark it is shown the place where there, are the distortion in the restored speech signal and the respectively place in the spectrum of the restored speech signal, where there are the differences with the spectrum of the original speech signal.

## VI. CONCLUSION

The time and spectrum graphics in Fig.1, 2, 3, 4, 5 and 6 are different each to an other with the choice of speech synthesis method. The comparison of the graphics for these three methods shown, that the choice of method of speech synthesis is important to the value of amplitude distortion of the speech signals and their corresponding difference of spectrum in respect of original speech signals.
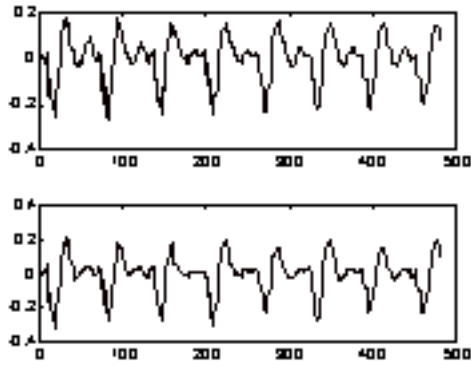
Fig. 1. The time graphics in the following sequence: original and restored signal of method 1.
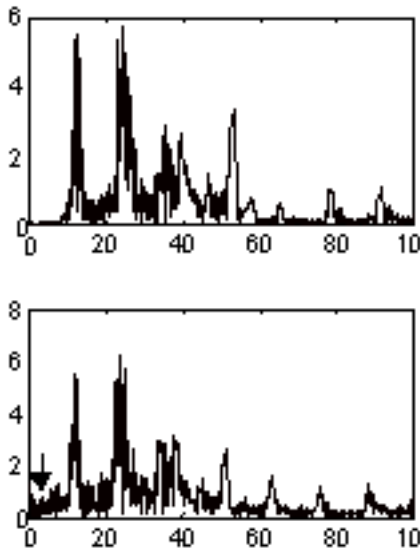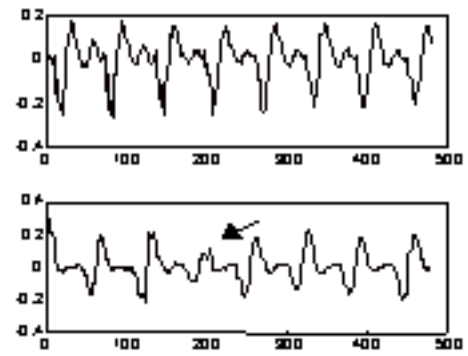


Fig. 3. The time graphics in the following sequence: original and restored signal of method 2.



Fig. 2. The spectrum graphics in the following sequence: original and restored signal of method 1.



Fig. 4. The spectrum graphics in the following sequence: original and restored signal of method 2.

REFERENCES

[1] Kondoz A.M, *Digital Speech*, John Wiley and Sons LTD, New Jork, 1994.
[2] Griffin D.W. and J.S.Lim, "Multi-Band Exitation Vocoder", *IEEE Transactions an ASSP*, 36(8) August 1988, pp.664-678.
[3] Spanias A.S, "Speech coding:A Tutorial Review", *Proceedings of the IEEE*, Volume 82, No10, October 1994..

Fig. 5. The time graphics in the following sequence: original and restored signal of method 3.
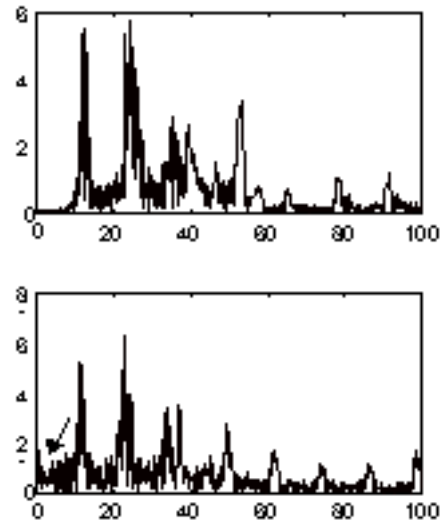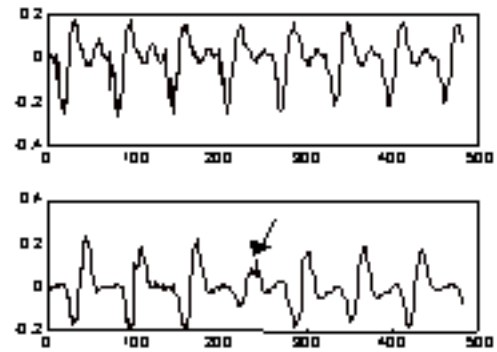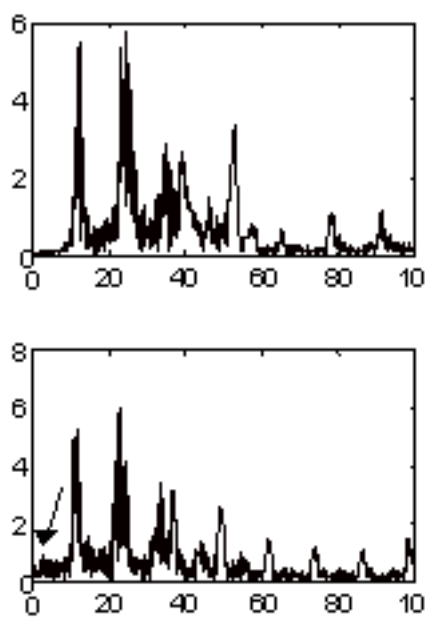
Fig. 6. The spectrum graphics in the following sequence: original and restored
signal of method 3.