

Estimation of the Fundamental Frequency of the Speech Signal Compressed by MP3 Algorithm

Zoran N. MILIVOJEVIĆ⁽¹⁾, Darko BRODIĆ⁽²⁾

⁽¹⁾ *Technical College*

Aleksandra Medvedeva 20, 18000 Nis, Serbia; e-mail: zoran.milivojevic@vtsnis.edu.rs

⁽²⁾ *Technical Faculty in Bor, University of Belgrade*

Vojske Jugoslavije 12, 19210 Bor, Serbia; e-mail: dbrodic@tf.bor.ac.rs

(received September 28, 2012; accepted May 17, 2013)

The paper analyzes the estimation of the fundamental frequency from the real speech signal which is obtained by recording the speaker in the real acoustic environment modeled by the MP3 method. The estimation was performed by the Picking-Peaks algorithm with implemented parametric cubic convolution (PCC) interpolation. The efficiency of PCC was tested for Catmull-Rom, Greville, and Greville two-parametric kernel. Depending on MSE, a window that gives optimal results was chosen.

Keywords: fundamental frequency, speech compression, speech processing, signal representation, MP3.

Notations

α_{opt} – optimal kernel parameter,
 α – kernel parameter,
 x – audio or speech signal,
 X – spectrum,
 w – window function,
 r – interpolation kernel,
 p – interpolation function,
 n – time lag, $0 \leq n \leq N-1$,
 N – window length,
 M – number of points between two samples in spectrum,
 L – kernel length, $4 \leq L \leq 8$,
 k – spectrum lag, $0 \leq k \leq N-1$,
 K – number of harmonics,
 f_s – sampling frequency,
 f_{max} – maximum of the interpolated frequency,
 f_e – estimated fundamental frequency,
 f – frequency,
 β_{opt} – optimal kernel parameter,
 β – kernel parameter,
 s – MP3 coded audio or speech signal.

13818-3, 1994). For instance, audio record rate in stereo technique at the sampling frequency $f_s = 44.1$ kHz is 10.584 MB/min. Transferring that number of bits is a very slow process even in very fast communication media. Hence, the development of compressing techniques is mandatory. A number of algorithms for audio signal compressing has been appeared. Most of them used MP3 algorithm with a compression degree of 1:12. Such compression ratio enables archiving a digitalized audio signal as well as transferring it by multimedia systems. Accordingly, MP3 became especially popular in internet applications (HACKER, 2000; MCCANDLESS, 1999). MP3 is a shortened name for coding algorithm derived from the standard MPEG-1, Layer III, developed by the German Technology Group. It was standardized by International Standards Organization (ISO) (ISO/IEC, 1992). MP3 does the compression tasks eliminating redundancy. It is similar to zip algorithm in accordance with a psycho-acoustic model that describes mechanisms of the human sound perception. Technically, MPEG-1 Layer III and MPEG-2 Layer III are declared as MP3 standard. MPEG-1 Layer III is used for 32 kHz, 44.1 kHz, and 48 kHz of sampling frequency, while MPEG-2 Layer III is used for 16 kHz, 22.05 kHz, and 24 kHz of sampling frequency. The standard broadening with a sign MPEG 2.5 is used for 8 kHz and 11 kHz (HACKER, 2000). MP3 compres-

1. Introduction

The rising trend of multimedia communications has imposed the need for archiving and transferring the audiovisual information. The amount of data which is archived or transferred, is very large (BRANDENBURG *et al.*, 1992; ISO/IEC, 1992; ISO/IEC

sion algorithm is based on the combination of several techniques the function of which is to maximize the relation between the perceived quality and the necessary file size. Spectrum of an audio signal is divided into 32 equally spaced frequency sub-bands. After that, a Modified Discrete Cosine Transformation (MDCT) is applied (BRITANAK, 2011). Precision of MDCT coefficient is reduced by the process of quantization. Further, the signal is processed according to the psychoacoustic model. This model emulates the human perception, i.e. the masking effects, which represent auditory and temporal masking (HACKER, 2000). After the signal processing according to the psychoacoustic model, Huffman's coding is performed. This coding additionally performs reduction of file size for 20%. In the name of copyright, the algorithms for inserting audio watermarks have been developed (YEO, KIM, 2003; WANG, HONG, 2006; DHAR, ECHIZEN, 2011). Latest advances in MP3 incorporate DFT-based MP3 multi-channel audio system (MOON, 2012). The parametric multi-channel audio coding concept enables the legacy system to reproduce stereo audio as well as the advanced system to reproduce multi-channel audio.

A number of the old music and speech records are digitalized and compressed by MP3 algorithm. However, there was a need for re-recording the significant historical and musical materials previously made by analog medium (magnetic tapes, vinyl records 78 rpm, LP, ...). The main deficiency of the analog sound recording is a high level of noise. In vinyl recordings, degradation effects come from imperfections and subsequent mechanical damage to the recording medium, and manifest themselves as clicks, sputtering, and noise from scratches. There is a need for this kind of processing as well as for restoration of the audio signals (AVILA, BISCAINHO, 2012).

In many multimedia applications, it is necessary to process audio records in order to improve the quality, intelligibility of speech, verification of the speaker, etc. A typical example is the quality improvement of the speech signal by reducing dissonant frequencies (JOEN *et al.*, 2003; KANG, 2004; KANG, KIM, 2006). Besides analyzing the trajectories of fundamental frequency, it is possible to classify the emotional state of a man (sadness, anger, joy, ...) (AYADI *et al.*, 2011), evaluate health status, and the conditions of hypoxia, which is manifested as a decrease in the concentration of oxygen in the blood (due to incidents during a flight, working in the mines, tunnels, etc.) (MILIVOJEVIC *et al.*, 2012). In processing of music and speech signal, it is necessary to determine the fundamental frequencies. Music signals are characterized by a fundamental frequency and the series of harmonic components that are integer multiples of fundamental frequency, i.e. partials.

In musical strings instruments, harmonic shifts occurrence leads to inharmonicity of an instrument. It

is defined through the inharmonicity coefficient. Determining the inharmonicity coefficient requires an accurate estimation of the fundamental frequency (BARBANCHO *et al.*, 2012). Digital processing of the music signal is possible for string instruments (e.g. guitar). Accordingly, it plays an estimate note, in which the string is played (E, H, G, D, A, E) as a fret. The more complex algorithms can detect the chords and form the score music (FRAGOULIS *et al.*, 2006).

The authors of this paper asked themselves: "What is the degradation of the fundamental frequency for MP3 encoding and decoding speech signals?"

In order to answer this question, the authors have conducted a number of experiments by applying the algorithms to estimate the fundamental frequency (F_0) in the frequency domain. After the calculation of DFT, the Picking-Picks is made. The highest peak represents the fundamental frequency. Particular attention is devoted to the application of parametric cubic convolution (PCC) algorithm in order to increase the precision of fundamental frequency estimation, when it is located between the spectral components on which DFT is calculated. The experiments were based on the time-domain (application of the window functions) and the frequency-domain processing, which implemented the cubic convolution kernels (Catmull-Rom, Greville, and Greville two-parametric kernel). Retrieval of the maximum position in the continuous convolution interpolation function is a mathematically complex and time-consuming process. Analytical expressions for the calculation of F_0 according to Keys kernel is proposed in (PANG *et al.*, 2000), while for the calculation of F_0 for Greville and Greville two-parametric (G2P) kernel is given in (MILIVOJEVIC, BRODIC, 2011).

In this paper, the authors present the results of the fundamental frequency assessment for:

- a) mathematically generated sine signal proposed in (PANG *et al.*, 2000) and
- b) real speech signals recorded in a real environment proposed in (MILIVOJEVIC, BRODIC, 2011).

The results will be analyzed by the mean square error (MSE) method. Finally, comparative analysis of the estimation accuracy F_0 by MP3 algorithm SYMPES (YARMAN *et al.*, 2006; MILIVOJEVIC, MIRKOVIC, 2009) and G.3.721 (MILIVOJEVIC, BRODIC, 2011) will be made.

This paper is organized as follows: Sec. 2 presents the previous works in the field. Section 3 describes the PCC algorithm. Subsection 3.1 defines the interpolation kernels. Subsection 3.2 presents the algorithm for determination of the optimal kernel parameters. Subsection 3.3 defines the test signals. Section 4 presents MSE results for the fundamental frequency estimation of the real speech signal modeled by the MP3 method. Section 5 shows the comparative analysis as well as the optimal kernel and window function selection. Section 6 gives the conclusion.

2. Previous works

The estimation of the fundamental frequency has received immense interest from different speech research areas, such as speech segregation, speech synthesis, speech coding, speech and speaker recognition, and speech articulation training for the deaf (GRIFFIN, LIM, 1988; ATAL, 1972; KAWAHARA *et al.*, 1999). A number of algorithms for determining the fundamental frequency has been developed. Theirs processing is performed in the time-domain (TD) and frequency-domain (FD) methods (KAWAHARA, 2002; SEKHAR, SREENIVAS, 2004; HUSSAIN, BOASHASH, 2002; KACHA, BENMAHAMMED, 2005; VEPREK, SCORDILIS, 2002; KLAPURI, 2003). In TD methods, one or more speech features (the fundamental harmonic, a quasi-periodic time structure, an alternation of high and low amplitudes, and points of discontinuities in the speech waveform) are identified first, and then the pitch markers or epochs are obtained in a pitch synchronous manner. In FD methods, a short-time frame or block of speech samples is transformed into spectral or frequency-domain in order to enhance the periodicity information contained in the speech. These methods determine an average pitch from several contiguous periods in the analysis frame. The performance of TD methods compared to FD methods depends more on the shape of the time waveform of speech (RESCH *et al.*, 2007). The autocorrelation function (ACF) (RABINER, 1977) and the average magnitude difference function (AMDF) (ROSS *et al.*, 1974) have been commonly employed for pitch estimation. In (KAWAHARA, 2002), an estimator named YIN has been proposed, where a series of modifications (a difference function formulation, normalization, and parabolic interpolation) has been introduced to decrease the error rates in pitch estimation from a clean speech (SHAHNAZ *et al.*, 2012).

The widespread method for determination of the fundamental frequency is based on Picking Peaks of the amplitude characteristic in the specific frequency range. This method is used for analyzing the signal values in the spectrum at frequencies on which the Discrete Fourier Transform (DFT) was calculated. Usually, the real value of the fundamental frequency is not there at the frequencies where DFT is calculated. In contrast, it lies between the two spectrum samples. That causes the frequency estimation error that lies in the interval $[-(f_s/(2N)) \text{ Hz}, (f_s/(2N)) \text{ Hz}]$, where f_s is the sampling frequency and N is the DFT window size. One way of reducing the error is determination of the interpolation function and estimation of the spectrum characteristics in the interval between two samples. This procedure gives the reconstruction of the spectrum on the base of DFT. The spectrum parameters are then determined by analytic procedures (differentiation, integration, extreme values, etc).

The calculation of the interpolation function by using PCC was represented in (KEYS, 1981; PARK,

SCHOWENGERDT, 1983). The special case of PCC interpolation applied in computer graphics has been called the Catmull-Rom interpolation (MEIJERING, UNSER, 2003). PANG *et al.* (2000) give detailed analysis of the fundamental frequency estimation and show the advantage of PCC interpolation. The application of PCC interpolation for determining the fundamental frequency in specific conditions is presented in (MILIVOJEVIC *et al.*, 2004). The efficiency of the algorithm for the evaluation of the fundamental frequency is determined by the simulation. As a quality measure of the algorithm, the mean square error (MSE) has been used. The best results were shown by the algorithm with the implemented Blackman window. The analysis of the algorithm efficiency where the signal-to-noise relation (SNR) is changeable according to the presence of the important harmonics in the fundamental function, is shown in (MIRKOVIC *et al.*, 2004). It confirmed the efficiency of the algorithm with the Blackman window. In (MIRKOVIC *et al.*, 2006), an analysis of PCC interpolation algorithm efficiency is made for the case where Greville two-parametric cubic convolution kernel (G2P) was implemented. The window was determined and the kernel parameters (α, β) were calculated where the minimum MSE was generated (in relation to Catmull-Rom kernel the error was smaller by 58.1%). The new method of speech signal modeling called “A Novel Systematic Procedure to Model Speech Signals via Predefined Envelope and Signature Sequences” (SYMPESES) is presented in the paper (YARMAN *et al.*, 2006). The results of the fundamental frequency estimation of the speech signal modeled by SYMPESSES method are shown in (MILIVOJEVIC, MIRKOVIC, 2009). Furthermore, the results of the fundamental frequency estimation for the speech signal coded by G.3.721 method are shown in (MILIVOJEVIC, BRODIC, 2011).

3. Proposed algorithms

Algorithm for the estimation of the fundamental frequency, based on the algorithm from (PANG *et al.*, 2000), is presented in Fig. 1.

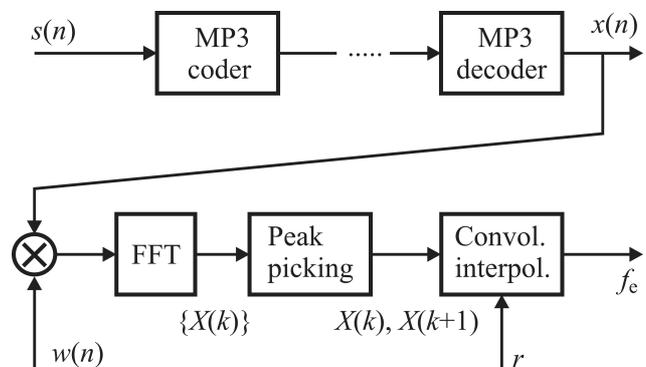


Fig. 1. Algorithm for the estimation of the fundamental frequency.

This algorithm is realized as follows:

Step 1: Audio or speech signal $s(n)$ is coded by MP3 coder.

Step 2: Coded signal is decoded by MP3 decoder and formed as signal $x(n)$.

Step 3: Window $w(n)$, length of which is N , is applied to decoded signal $x(n)$.

Step 4: Spectrum $X(k)$ is calculated by using DFT:

$$X(k) = \text{DFT}(x(n)). \quad (1)$$

The spectrum is calculated in discrete points $k = 0, \dots, N-1$, where N is the length of DFT. The real spectrum of signals $x(n)$ is continuous, whereas DFT defines the values of the spectrum at some discrete points.

Step 5: The maximum of the real spectrum that is between k -th and $(k+1)$ -th samples is determined by using the Picking-Peak algorithm. The values $X(k)$ and $X(k+1)$ are the highest in the specified domain.

Step 6: The maximum of the spectrum is calculated by PCC interpolation. The reconstructed function is:

$$X_r(f) = \sum_{i=k-L/2+1}^{k+L/2} p_i \cdot r(f-i), \quad k \leq f \leq k+1, \quad (2)$$

where $p_i = X_r(i)$, $r(f)$ is the kernel of interpolation, and L is the number of samples that participate in the interpolation.

Step 7: By differentiation $X_r(f)$ and zero adjustment, the position of the maximum is determined. It represents the estimated fundamental frequency f_e .

The quality of the algorithm for the fundamental frequency estimation can be also expressed by MSE:

$$\text{MSE} = \overline{(f - f_e)^2}, \quad (3)$$

where f is true fundamental frequency and f_e is estimated fundamental frequency.

3.1. Interpolation kernel

The definitions of the interpolation kernels, which are tested in this paper, are:

a) Keys interpolation kernel (KEYS, 1981; PARK, SCHOWENGERDT, 1983):

$$r(f) = \begin{cases} (\alpha+2)|f|^3 - (\alpha+3)|f|^2 + 1, & |f| \leq 1, \\ \alpha|f|^3 - 5\alpha|f|^2 + 8\alpha|f| - 4\alpha, & 1 < |f| \leq 2, \\ 0, & \text{otherwise.} \end{cases} \quad (4)$$

For $L = 4$, from Eq. (2), the position of maximum is determined:

$$f_{\max} = \begin{cases} k - \frac{c}{2b}, & a = 0, \\ k + \frac{-b - \sqrt{b^2 - ac}}{a}, & a \neq 0, \end{cases} \quad (5)$$

where

$$\begin{aligned} a &= 2(\alpha p_{k-1} + (\alpha + 2)p_k - (\alpha + 2)p_{k+1} - \alpha p_{k+2}), \\ b &= -2\alpha p_{k-1} - (\alpha + 3)p_k + (2\alpha + 3)p_{k+1} + \alpha p_{k+2}, \\ c &= -\alpha p_{k-1} - \alpha p_{k+1}, \end{aligned} \quad (6)$$

b) Greville interpolation kernel (MEIJERING, UNSER, 2003):

$$r(f) = \begin{cases} \left(\alpha + \frac{3}{2}\right)|f|^3 - \left(\alpha + \frac{5}{2}\right)|f|^2 + 1; & \text{if } 0 \leq |f| \leq 1, \\ \frac{1}{2}(\alpha - 1)|f|^3 - \left(3\alpha - \frac{5}{2}\right)|f|^2 + \left(\frac{11}{2}\alpha - 4\right)|f| - (3\alpha - 2); & \text{if } 1 \leq |f| \leq 2, \\ -\frac{1}{2}\alpha|f|^3 + 4\alpha|f|^2 - \frac{21}{2}\alpha|f| + 9\alpha; & \text{if } 2 \leq |f| \leq 3, \\ 0; & \text{if } 3 \leq |f|. \end{cases} \quad (7)$$

For $L = 6$, from Eqs. (2) and (7), the position of maximum is determined according to Eq. (5), where

$$\begin{aligned} a &= -\frac{3}{2}\alpha p_{k-2} + \frac{3}{2}(\alpha - 1)p_{k-1} + 3\left(\alpha + \frac{3}{2}\right)p_k - 3\left(\alpha + \frac{3}{2}\right)p_{k+1} - \frac{3}{2}(\alpha - 1)p_{k+2} + \frac{3}{2}\alpha p_{k+3}, \\ b &= -2\alpha p_{k-2} + (-3\alpha + 2)p_{k-1} - (2\alpha + 5)p_k + 4(\alpha + 1)p_{k+1} - p_{k+2} - \alpha p_{k+3}, \\ c &= -\frac{1}{2}\alpha p_{k-2} + \left(\alpha - \frac{1}{2}\right)p_{k-1} - \left(\alpha - \frac{1}{2}\right)p_{k+1} + \frac{1}{2}\alpha p_{k+2}, \end{aligned} \quad (8)$$

c) Greville two-parametric cubic convolution kernel (G2P) (MEIJERING, UNSER, 2003):

$$r(f) = \begin{cases} \left(\alpha - \frac{5}{2}\beta + \frac{3}{2} \right) \cdot |f|^3 - \left(\alpha - \frac{5}{2}\beta + \frac{5}{2} \right) \cdot |f|^2 + 1; & 0 \leq |f| \leq 1, \\ \frac{1}{2}(\alpha - \beta - 1) \cdot |f|^3 - \left(3\alpha - \frac{9}{2}\beta - \frac{5}{2} \right) \cdot |f|^2 + \left(\frac{11}{2}\alpha - 10\beta - 4 \right) \cdot |f| - (3\alpha - 6\beta - 2); & 1 \leq |f| \leq 2, \\ -\frac{1}{2}(\alpha - 3\beta) \cdot |f|^3 + \left(4\alpha - \frac{25}{2}\beta \right) \cdot |f|^2 - \left(\frac{21}{2}\alpha - 34\beta \right) \cdot |f| + (9\alpha - 30\beta); & 2 \leq |f| \leq 3, \\ -\frac{1}{2}\beta \cdot |f|^3 + \frac{11}{2}\beta \cdot |f|^2 - 20\beta \cdot |f| + 24\beta; & 4 \leq |f|. \end{cases} \quad (9)$$

For $L = 8$, from Eqs. (2) and (9), the position of maximum is determined according to Eq. (5), where

$$\begin{aligned} a &= -\frac{3}{2}\beta p_{k-3} - \frac{3}{2}(\alpha - 3\beta) p_{k-2} \\ &+ \frac{3}{2}(\alpha - \beta - 1) p_{k-1} + 3 \left(\alpha - \frac{5}{2}\beta + \frac{3}{2} \right) p_k \\ &- 3 \left(\alpha - \frac{5}{2}\beta + \frac{3}{2} \right) p_{k+1} - \frac{3}{2}(\alpha - \beta - 1) p_{k+2} \\ &+ -\frac{3}{2}(\alpha - 3\beta) p_{k+3} + \frac{3}{2}\beta p_{k+4}; \\ b &= -2\beta p_{k-3} - (2\alpha - 7\beta) p_{k-2} \\ &+ (-3\alpha + 6\beta + 2) p_{k-1} - \left(2\alpha - 5\beta + \frac{5}{2} \right) p_k \\ &+ (4\alpha - 10\beta + 1) p_{k+1} + (3\beta - 1) p_{k+2} \\ &+ (-\alpha + 2\beta) \alpha p_{k+3} - \beta p_{k+4}; \\ c &= -\frac{1}{2}\beta p_{k-3} + \left(-\frac{1}{2}\alpha + 2\beta \right) p_{k-2} \\ &+ \left(\alpha - \frac{5}{2}\beta - \frac{1}{2} \right) p_{k-1} - \left(\alpha + \frac{5}{2}\beta + \frac{1}{2} \right) p_{k+1} \\ &+ \left(\frac{1}{2}\alpha - 2\beta \right) p_{k+2} + \frac{1}{2}\beta p_{k-3}. \end{aligned} \quad (10)$$

In Eqs. (4)–(10), there are α and β parameters. The optimal values of these parameters will be determined by the minimum value of MSE, for Keys, Greville, and G2P kernel. For the first two of them

$$\alpha_{\text{opt}} = \arg \min_{\alpha} (\text{MSE}), \quad (11)$$

and for the G2P kernel,

$$(\alpha_{\text{opt}}, \beta_{\text{opt}}) = \arg \min_{\alpha, \beta} (\text{MSE}). \quad (12)$$

The detailed analysis in (PANG *et al.*, 2000; MILIVOJEVIC *et al.*, 2004; 2006; MIRKOVIC *et al.*, 2004; YARMAN *et al.*, 2006; MILIVOJEVIC, MIRKOVIC, 2009; MILIVOJEVIC, BRODIC, 2011) showed that the minimum value of MSE depends on the application of window by which signal processing $x(n)$ is carried out in time domain. MSE will be defined for: (a) Hamming, (b) Hanning, (c) Blackman, (d) Rectangular, (e) Kaiser, and (f) Triangular window.

3.2. Interpolation kernel parameters

The algorithm for determination of interpolation kernel parameters α and β is realized as follows:

Step 1: Signal $x(n)$, which was previously coded and decoded by MP3 algorithm, is modified by the window function $w(n)$, length of which is N .

Step 2: Spectrum $X(k)$ is determined by application of DFT.

Step 3: Reconstruction of the continual function that represents spectrum $X(f)$ is performed by application of PCC interpolation.

Step 4: MSE is calculated for various values of parameters α and β depending on the implemented window.

Step 5: α_{opt} and β_{opt} are determined for which the minimum value of MSE is obtained.

3.3. Test signals

PCC algorithm of the fundamental frequency estimation will be applied to:

- a) simulation sine test signal and
- b) real speech test signal.

Simulation sine signal for testing of PCC algorithm is defined in (Pang *et al.*, 2000):

$$s(t) = \sum_{i=1}^K \sum_{g=0}^M a_i \sin \left(2\pi i \left(f_o + g \frac{f_s}{NM} \right) t + \theta_i \right), \quad (13)$$

where f_o is fundamental frequency, θ_i and a_i are phase and amplitude of the i -th harmonic, respectively, K is the number of harmonics, and M is the number of points between the two samples in spectrum where PCC interpolation is being made. The real speech test signal is obtained by recording of a speaker in the real acoustic environment.

PCC algorithm will be applied to:

- uncoded real speech test signals and
- real speech test signals coded and decoded by MP3 algorithm.

The results will be summarized and comparative analysis will be established in accordance to MP3 algorithm applied to the sine test signal.

4. Experimental results and discussion

4.1. Testing parameters

In the simulation process, f_0 and θ_i are random variables with uniform distribution in the range [G2 (97.99 Hz), G5 (783.99 Hz)] and $[0, 2\pi]$ with sine and real speech test signals. Signal frequency of sampling is $f_s = 8$ kHz, and the length of window is $N = 256$, which assures the analysis of subsequences that last 32 ms. Furthermore, the results will relate to $f_0 = 125\text{--}140.625$ Hz (frequencies between the 8-th and 9-th DFT components). Number of frequencies in the specified range, for which the estimation is done, is $M = 100$. The sine test signal is with $K = 10$ harmonics. All further analyses will relate to: (a) Hamming, (b) Hanning, (c) Blackman, (d) Rectangular, (e) Kaiser, and (f) Triangular window.

4.2. Experimental results

4.2.1. Keys kernel

By applying the algorithm for determination of Keys interpolation kernel parameters, some diagrams $MSE(\alpha)$ are drawn (Fig. 2 and Fig. 3), the minimum value $MSE_{K_{\min}}$ is determined, and on the base of it, the optimum value of Keys kernel α_{opt} is determined for: (a) Hamming, (b) Hanning, (c) Blackman, (d) Kaiser, and (e) Triangular window functions. Values $MSE_{K_{\min}}$ and α_{opt} are presented in Table 1 (uncoded sine test signal $MSE_{K_{\min}}$, MP3 coded sine test

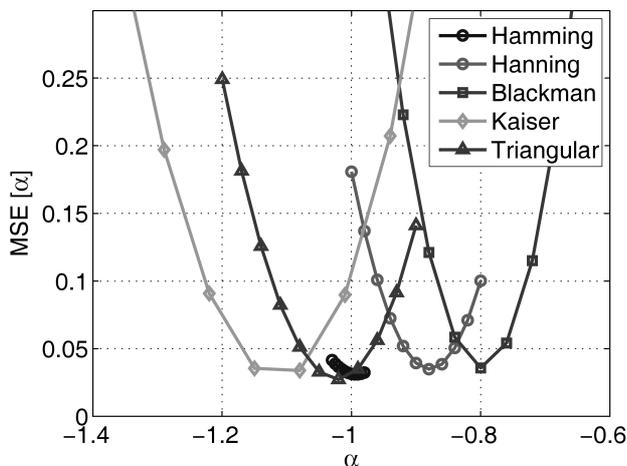


Fig. 2. $MSE(\alpha)$ for Keys kernel and uncompressed real speech test signal.

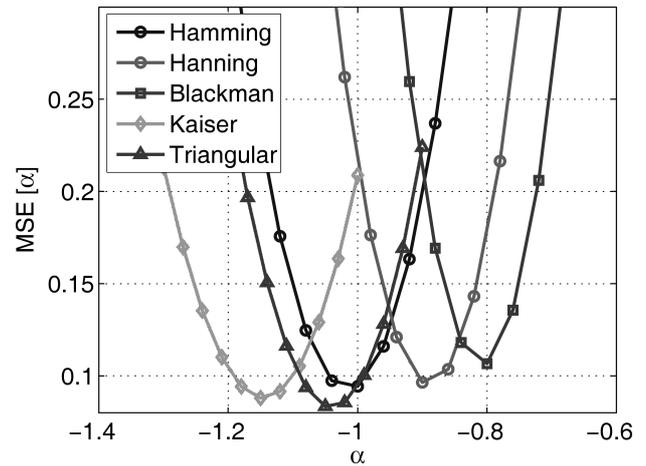


Fig. 3. $MSE(\alpha)$ for Keys kernel and MP3 compressed real speech test signal.

signal $MSE_{K_{\text{MP3min}}}$) and Table 2 (real speech test signal $MSE_{KSP_{\min}}$, MP3 coded real sine test signal $MSE_{KSP_{\text{MP3min}}}$).

Table 1. Minimum MSE and α_{opt} for sine test signal (Keys kernel).

	Uncoded signal		Signal coded by MP3 algorithm	
	α_{opt}	$MSE_{K_{\min}}$	α_{opt}	$MSE_{K_{\text{MP3min}}}$
Hamming	-1.005	0.023	-1.0100	0.0320
Hanning	-0.885	0.004	-0.8825	0.0031
Blackman	-1.801	0.001	-0.8024	0.0028
Rectangular	-2.61	0.515	-2.5500	0.4388
Kaiser	-1.125	0.02	-1.1250	0.0203
Triangular	-1.028	0.0028	-1.0280	0.0068

Table 2. Minimum MSE and α_{opt} for real speech test signal (Keys kernel).

	Uncoded signal		Signal coded by MP3 algorithm	
	α_{opt}	$MSE_{KSP_{\min}}$	α_{opt}	$MSE_{KSP_{\text{MP3min}}}$
Hamming	-0.995	0.0310	-1	0.0943
Hanning	-0.880	0.0349	-0.9000	0.0965
Blackman	-0.800	0.0358	-0.8000	0.1067
Rectangular	-2.400	0.4323	-2.3000	0.7011
Kaiser	-1.080	0.0339	-1.1500	0.0880
Triangular	-1.030	0.0277	-1.0500	0.0835

According to the results presented in Tables 1 and 2, it is obvious that:

- At sine test signal, the greatest precision of fundamental frequency estimation is when Blackman window ($MSE_{K_{\min}} = 0.001$) is applied. At MP3 coded sine test signal, the greatest precision of estimation is in Blackman ($MSE_{K_{\text{MP3min}}} = 0.0028$) window. When MP3 coding is applied, the pre-

cision of the fundamental frequency estimation is $MSE_{K_MP3min}/MSE_{KGmin} = 0.0028/0.001 = 2.8$ times lower.

- b) At real speech test signal, the greatest precision is in triangular window ($MSE_{KSPmin} = 0.0277$). At MP3 coded real speech signal, the greatest precision is in triangular window ($MSE_{KSP_MP3min} = 0.0835$). When coding is applied, the precision of the fundamental frequency estimation is $MSE_{KSP_MP3min}/MSE_{KSPmin} = 0.0835/0.0277 = 3.0144$ times lower.
- c) At coded real speech signal in relation to coded sine signal, the non-precision of the fundamental frequency estimation is $MSE_{KSP_MP3min}/MSE_{K_MP3min} = 0.0835/0.0028 = 29.821$ times higher.

4.2.2. Greville kernel

By applying the algorithm for determination of Greville interpolation kernel parameters, some diagrams $MSE(\alpha)$ are drawn (Fig. 4 and Fig. 5), minimum value MSE_{Gmin} is determined, and on the base

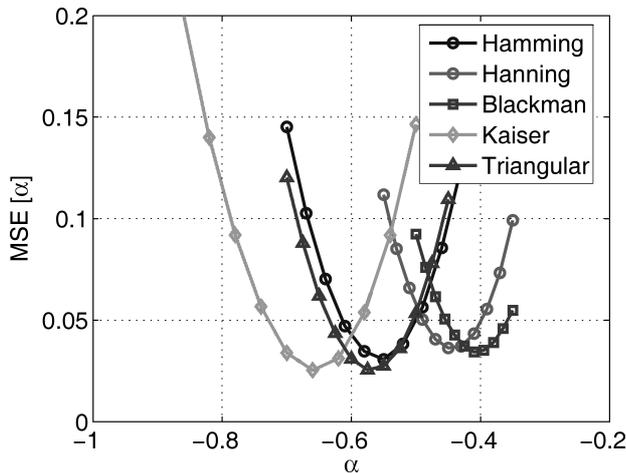


Fig. 4. $MSE(\alpha)$ for Greville kernel and uncompressed real speech test signal.

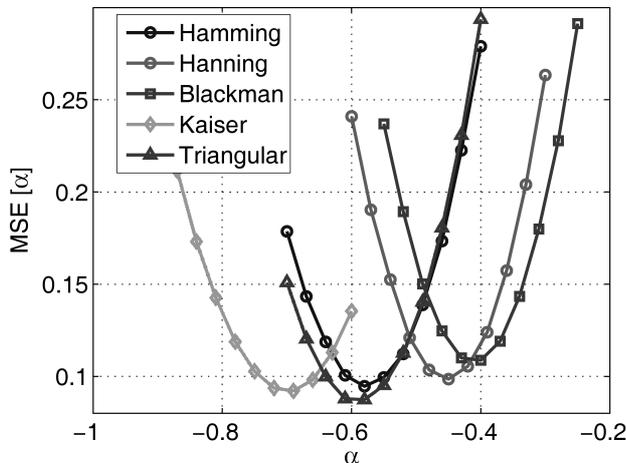


Fig. 5. $MSE(\alpha)$ for Greville kernel and MP3 compressed real speech test signal.

of it, the optimum value of Greville kernel parameters α_{opt} is determined for: (a) Hamming, (b) Hanning, (c) Blackman, (d) Kaiser, and (e) Triangular window. Values MSE_{min} and α_{opt} are presented in Table 3 (uncoded sine test signal MSE_{Gmin} , coded sine test signal MSE_{G_MP3min}) and Table 4 (real speech test signal MSE_{GSPmin} , coded real speech test signal MSE_{GSP_MP3min}).

Table 3. Minimum MSE and α_{opt} for sine test signal (Greville kernel).

	Uncoded signal		Signal coded by MP3 algorithm	
	α_{opt}	MSE_{Gmin}	α_{opt}	MSE_{G_MP3min}
Hamming	-0.57	0.0175	-0.5750	0.0272
Hanning	-0.449	0.0027	-0.4500	0.0032
Blackman	-0.415	0.0009	-0.4200	0.0037
Rectangular	-2.254	0.4054	-2.2000	0.3966
Kaiser	-0.6676	0.0124	-0.6600	0.0207
Triangular	-0.575	0.002	-0.5750	0.0064

Table 4. Minimum MSE and α_{opt} for real speech test signal (Greville kernel).

	Uncoded signal		Signal coded by MP3 algorithm	
	α_{opt}	MSE_{GSPmin}	α_{opt}	MSE_{GSP_MP3min}
Hamming	-0.560	0.0310	-0.5800	0.0947
Hanning	-0.450	0.0363	-0.4500	0.0986
Blackman	-0.410	0.0344	-0.4000	0.1088
Rectangular	-2.100	0.2016	-2.2000	0.3481
Kaiser	-0.660	0.0255	-0.6900	0.0922
Triangular	-0.575	0.0256	-0.5800	0.0874

According to the results presented in Tables 3 and 4, it is obvious that:

- a) At sine test signal, the greatest precision of fundamental frequency estimation is when Blackman ($MSE_{Gmin} = 0.0009$) window is applied. At MP3 coded sine test signal, the greatest precision of estimation is in Hanning window ($MSE_{G_MP3min} = 0.0032$). When coding is applied, the precision of the fundamental frequency estimation is $MSE_{G_MP3min}/MSE_{Gmin} = 0.0032/0.0009 = 3.55$ times lower.
- b) At real speech test signal, the greatest precision is in Kaiser window ($MSE_{GSPmin} = 0.0255$). At coded real speech signal, the greatest precision is in triangular window ($MSE_{GSP_MP3min} = 0.0874$). When coding is applied, the precision of the fundamental frequency estimation is $MSE_{GSP_MP3min}/MSE_{GSPmin} = 0.0874/0.0255 = 3.427$ times lower.
- c) At coded real speech signal in relation to coded sine signal, the non-precision of the fundamen-

tal frequency is $\text{MSE}_{\text{GSP_MP3min}}/\text{MSE}_{\text{G_MP3min}} = 0.0874/0.0032 = 27.31$ times higher.

4.2.3. G2P kernel

By applying the algorithm for determination of Greville two-parametric interpolation kernel parameters, some diagrams $\text{MSE}(\alpha)$ are drawn and minimum values $\text{MSE}_{\text{G2Pmin}}$ are determined for windows with the smallest MSE. Three-dimensional $\text{MSE}(\alpha, \beta)$ graphics are drawn for uncompressed real speech test signal (Fig. 6a), the shift of minimum MSE_{min} in (α, β) level (Fig. 6b) for Blackman window, and real speech test signal coded by MP3 (Fig. 7a), the shift of minimum MSE_{min} in (α, β) level (Fig. 7b) for Kaiser window. In Figs. 6b and 7b, positions of $\text{MSE}_{\text{min}} = \text{MSE}(\alpha_{\text{opt}}, \beta_{\text{opt}})$ minimum in (α, β) plane for Greville (point **A**) and G2P (point **B**) interpolation kernel are shown. Vector **AB** shows the position change of the minimum ($\text{MSE}(\alpha_{\text{opt}}, \beta_{\text{opt}})$). The determined parameters α_{opt} and β_{opt} are presented in Table 5.

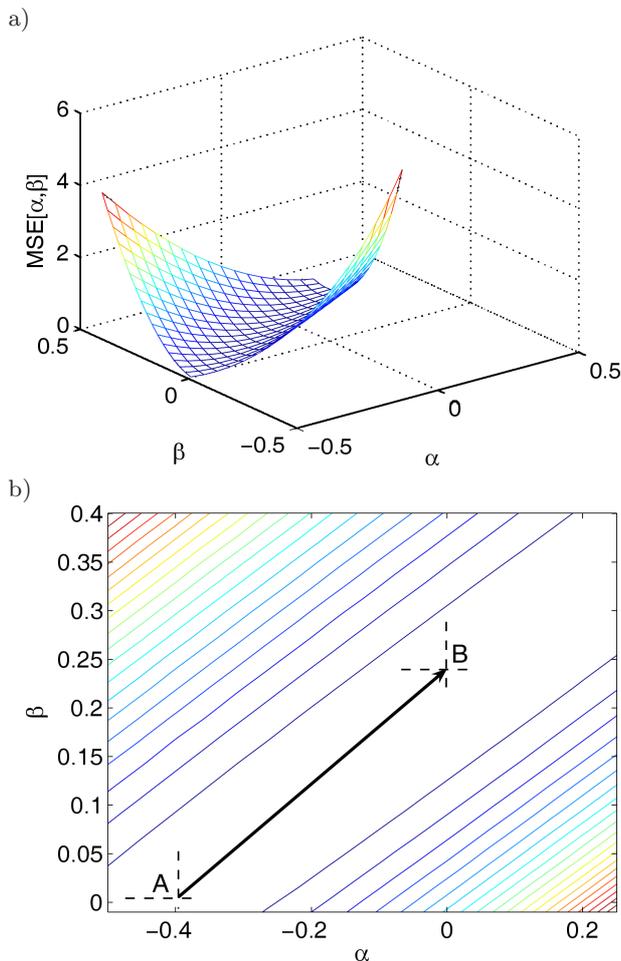


Fig. 6. Real speech test signal with the application of Blackman window without compression: a) $\text{MSE}(\alpha, \beta)$ for the application of G2P PCC interpolation; b) positions of $\min(\text{MSE}(\alpha_{\text{opt}}, \beta_{\text{opt}}))$ in plane (α, β) for Greville (point **A**) and G2P PCC (point **B**) interpolation.

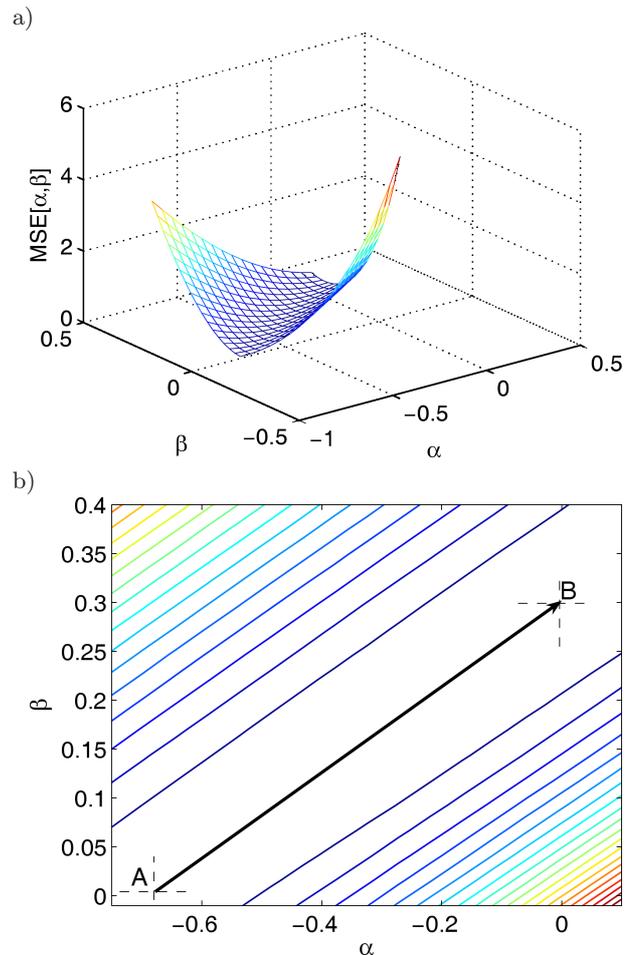


Fig. 7. Real speech test signal with the application of Kaiser window with MP3 compression: a) $\text{MSE}(\alpha, \beta)$ for the application of G2P PCC interpolation; b) positions of $\min(\text{MSE}(\alpha_{\text{opt}}, \beta_{\text{opt}}))$ in plane (α, β) for Greville (point **A**) and G2P PCC (point **B**) interpolation.

According to the results presented in Table 5, it is obvious that:

- At real speech test signal with G2P kernel comparing with Greville kernel, the precision of the fundamental frequency estimation ($\text{MSE}_{\text{GSP_MP3min}}/\text{MSE}_{\text{G2PSP_MP3min}}$) is: (a) 1.57 (Hamming), (b) 1.76 (Hanning), (c) 1.83 (Blackman), (d) 1.75 (Kaiser), and (e) 1.5 (Triangular) times higher.
- At real speech test signal, the greatest precision is in Blackman window ($\text{MSE}_{\text{GSPmin}} = 0.0009$). At MP3 coded real speech signal, the greatest precision is in Kaiser window ($\text{MSE}_{\text{G2PSP_MP2min}} = 0.0524$). When MP3 coding with G2P kernel is applied, the precision of the fundamental frequency estimation is $\text{MSE}_{\text{G2PSP_MP3min}}/\text{MSE}_{\text{GSPmin}} = 0.0524/0.0009 = 58.22$ times lower.
- At MP3 coded real speech with Greville kernel, the greatest precision is in Triangular window ($\text{MSE}_{\text{GSPmin}} = 0.0874$). At MP3 coded real speech

Table 5. Minimum MSE, α_{opt} , and β_{opt} (G2P kernel).

Sine test signal (uncoded)			
Window	α_{opt}	β_{opt}	MSE_{G2Pmin}
Hamming	-0.55	0.03	0.0046
Hanning	0.5	0.015	0.0018
Blackman	-0.42	0.002	0.000377
Kaiser	-0.681	0.001	0.0096
Triangular	0.6	-0.001	0.001
Sine test signal (coded by MP3 algorithm)			
Window	α_{opt}	β_{opt}	MSE_{G2P_MP3min}
Hamming	-0.5900	-0.0060	0.0270
Hanning	-0.4600	-0.0060	0.0029
Blackman	-0.4200	-0.0020	0.0022
Kaiser	-0.6600	0.0060	0.0178
Triangular	-0.5680	0.0030	0.0060
Real speech test signal (uncoded)			
Window	α_{opt}	β_{opt}	$MSE_{G2PSPmin}$
Hamming	0.1	0.2975	0.0072
Hanning	0.1531	0.2719	0.0025
Blackman	0.0625	0.2463	0.001
Kaiser	-0.1	0.2463	0.0075
Triangular	-0.3219	0.1181	0.0016
Real speech test signal (coded by MP3 algorithm)			
Window	α_{opt}	α_{opt}	MSE_{G2PSP_MP3min}
Hamming	0.0687	0.2788	0.0600
Hanning	0.0625	0.2375	0.0559
Blackman	-0.0625	0.1738	0.0592
Kaiser	-0.0062	0.2975	0.0524
Triangular	-0.3500	0.1194	0.0581

signal with G2P kernel, the greatest precision is in Kaiser window ($MSE_{G2PSP_MPmin} = 0.0524$). When MP3 coding with G2P kernel is applied, the precision of the fundamental frequency estimation is $MSE_{GSPmin} / MSE_{G2PSP_MP3min} = 0.0874 / 0.0524 = 1.66$ times higher.

5. Comparative analysis

The comparative analysis of the estimated fundamental frequency for the sine test signal and the real speech test signal, without and with MP3 compression, will be performed on the base of MSE minimum values. The minimum value of MSE is determined on the base of the diagram in the Figs. 2 and 3 (Keys), Figs. 4 and 5 (Greville), and Figs. 6 and 7 (G2P). It is presented in Table 1 ($MSE_{Kmin}, MSE_{K_MP3min}$), Table 2 ($MSE_{KSPmin}, MSE_{KSP_MP3min}$), Table 3 ($MSE_{Gmin}, MSE_{G_MP3min}$), Table 4 ($MSE_{GSPmin}, MSE_{GSP_MP3min}$), and Ta-

ble 5 ($MSE_{G2Pmin}, MSE_{G2P_MP3min}, MSE_{G2PSPmin}, MSE_{G2PSP_MP3min}$), respectively.

Comparing the values MSE_{min} from Tables 1–5, it can be concluded that:

- a) The optimum choice for sine test signal is Blackman window for all interpolation kernels. G2P interpolation kernel, which generates MSE by 60.05% less than Keys and 55.55% less than Greville kernel, showed the best results.
- b) The optimum choice for real speech test signal is G2P kernel with Blackman window, which generates MSE by 96.387% less than Keys kernel (Triangular window), and 96.07% less than Greville kernel (Kaiser window).
- c) The optimum choice for sine test signal coded by MP3 algorithm is G2P interpolation kernel with Blackman window, which generates MSE by 21.43% less than Keys (Blackman window) and 31.25% less than Greville kernel (Kaiser window), showed the best results.
- d) The optimum choice for real speech test signal coded by MP3 algorithm is G2P interpolation kernel with Kaiser window, which generates MSE by 37.27% less than Keys (Triangular window) and 43.16% less than Greville kernel (Kaiser window), showed the best results.
- e) Comparing MSE for G2P kernel for uncoded real speech test signal (Blackman window, $MSE_{G2PSPmin} = 0.0022$) and MP3 coded real speech test signal (Kaiser window, $MSE_{G2PSP_MP3min} = 0.0524$), relation $MSE_{G2PSP_MP3min} / MSE_{G2PSPmin} = 0.0524 / 0.0022 = 23.818$ has been obtained.

Comparison of the estimation of the fundamental frequency for the signal coded with SYMPES algorithm (MILIVOJEVIC, MIRKOVIC, 2009), G.723.1 algorithm (MILIVOJEVIC, BRODIC, 2011), and MP3 algorithm with real speech test signal given in this paper shows that the proposed algorithm has the least MSE values. Accordingly, the obtained results recommend the use of PCC algorithm with G2P kernel in preprocessing signals which are compressed by MP3 method. Hence, it is recommended for further processing by algorithms that require a precise determination of the fundamental frequency (automatic verification of a speaker, recognition of the speech, etc.).

6. Conclusions

This paper presents the comparative analysis of the fundamental frequency estimation for the real speech signal modeled by MP3 method. The estimation of the fundamental frequency has been made by the Picking-Peaks algorithm with implemented PCC interpolation. Experiments have been performed with Keys, Greville, and Greville two-parametric kernels. In order

to minimize MSE, different windows have been implemented. The detailed analysis has shown that the optimal choice is Greville two-parametric kernel and the Kaiser window implemented in PCC algorithm. The optimum choice for real speech test signal coded by MP3 algorithm is G2P interpolation kernel with Kaiser window, which generates MSE by 37.27% less than Keys (Triangular window) and 43.16% less than Greville kernel (Kaiser window). Comparing these results with the results of the estimation of the fundamental frequency in the real speech signal that is not modeled by MP3 method, the relation of minimum MSEs 23.818 has been obtained. Comparison between algorithms shows MSE for SYMPES ($MSE_{\min} = 3.174$) and G.723.1 algorithm ($MSE_{\min} = 0.2898$), and for the proposed MP3 algorithm ($MSE_{\min} = 0.0524$). These results prove the quality of the proposed solution. Hence, the obtained results recommend the use of PCC algorithm with G2P kernel in preprocessing of signals compressed by MP3 method for further processing by algorithms which require a precise determination of the fundamental frequency.

References

1. ATAL B. (1972), *Automatic speaker recognition based on pitch contours*, Journal of the Acoustical Society of America, **52**, 6, 1687–1697.
2. AVILA F., BISCAINHO L. (2012), *Bayesian Restoration of Audio Signals Degraded by Impulsive Noise Modeled as Individual Pulses*, IEEE Transactions On Audio, Speech, And Language Processing, **20**, 9, 2470–2481.
3. AYADI M., KAMEL M., KARRAY F. (2011), *Survey on speech emotion recognition: Features, classification schemes, and databases*, Pattern Recognition, **44**, 572–587.
4. BARBANCHO I., TARDON L., SAMMARTINO S., BARBANCHO A. (2012), *Inharmonicity-Based Method for the Automatic Generation of Guitar Tablature*, IEEE Transactions On Audio, Speech, And Language Processing, **20**, 6, 1857–1868.
5. BRANDENBURG K., STOLL G., DEHERY Y.F., JOHNSTON J.D., KERKHOF L.V., SCHROEDER E.F. (1992), *The ISO/MPEG Audio Codec: A Generic Standard for Coding of High Quality Digital Audio*, 92nd. AES-convention, preprint 3336, Vienna.
6. BRITANAK V. (2011), *A survey of efficient MDCT implementations in MP3 audio coding standard: Retrospective and state-of-the-art*, Signal Processing, **91**, 624–672.
7. DHAR P.K., ECHIZEN I. (2011), *Robust FFT Based Watermarking Scheme for Copyright Protection of Digital Audio Data*, Seventh International Conference on Intelligent Information Hiding and Multimedia Signal Processing (IIH-MSP), 181–184.
8. FRAGOULIS D., PAPAODYSSSEUS C., EXARHOS M., ROUSSOPOULOS G., PANAGOPOULOS T., KAMAROTOS D. (2006), *Automated classification of piano-guitar notes*, IEEE Transactions On Audio, Speech, And Language Processing, **14**, 3, 1040–1050.
9. GRIFFIN D., LIM J. (1988), *Multiband excitation vocoder*, IEEE Transactions On Audio, Speech, And Language Processing, **36**, 8, 1223–1235.
10. HACKER S. (2000), *MP3: The Definitive Guide*, O'Reilly & Associates, Sebastopol, CA 95472.
11. HUSSAIN Z.M., BOASHASH B. (2002), *Adaptive instantaneous frequency estimation of multicomponent signals using quadratic time-frequency distributions*, IEEE Transaction on Signal Processing. **50**, 8, 1866–1876.
12. ISO/IEC (1992), *Information Technology – Coding of Moving Pictures and Associated Audio for Digital Storage Media at up to about 1.5 Mbit/s, Part3: Audio*, ISO/IEC JTC1/SC29/WG11 MPEG, International Standard 11172-3 (MPEG-1).
13. ISO/IEC 13818-3 (1994), *Information Technology. Generic Coding of Moving Pictures and Associated Audio: Audio*. ISO/IEC JTC1/SC29/WG11 MPEG, International Standard 13818-3(MPEG-2), 1994.
14. JOEN B., KANG S., BAEK S.J., SUNG K.M. (2003), *Filtering of a Dissonant Frequency Based on Improved Fundamental Frequency Estimation for Speech Enhancement*, IEICE Trans. Fundamentals, E86-A, **8**, 2063–2064.
15. KACHA F., BENMAHAMMED G.K. (2005), *Time-frequency analysis and instantaneous frequency estimation using two-sided linear prediction*, IEEE Signal Processing, **85**, 491–503.
16. KANG S. (2004), *Dissonant frequency filtering technique for improving perceptual quality of noisy speech and husky voice*, IEEE Signal Processing, **84**, 431–433.
17. KANG S., KIM Y. (2006), *A Dissonant Frequency Filtering for Enhanced Clarity of Husky Voice Signals*, Lecture Notes Comp Science, **4188**, Berlin Springer, 517–522.
18. KAWAHARA H., KATSUSE I., CHEVEIGNE A. (1999), *Restructuring speech representations using a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based f0 extraction: Possible role of a repetitive structure in sounds*, Speech Communication, **27**, 3–4, 187–207.
19. KAWAHARA C.H. (2002), *YIN, a fundamental frequency estimator for speech and music*, Journal of the Acoustical Society of America, **111**, 4, 1917–1930.
20. KEYS R.G. (1981), *Cubic convolution interpolation for digital image processing*, IEEE Transaction on Acoustics, Speech & Signal Processing, **29**, 6, 1153–1160.
21. KLAPURI A. (2003), *Multiple fundamental frequency estimation based on harmonicity and spectral smoothness*, IEEE Transactions On Audio, Speech, And Language Processing, **11**, 6, 804–816.
22. MCCANDLESS M. (1999), *The MP3 revolution*, IEEE Intelligent Systems and their Applications, **14**, 3, 8–9.

23. MEIJERING E., UNSER M. (2003), *A Note on Cubic Convolution Interpolation*, IEEE Transaction on Image Processing, **12**, 4, 447–479.
24. MILIVOJEVIC Z., MIRKOVIC D. (2009), *Estimation of the fundamental frequency of the speech signal modeled by the SYMPES method*, International Journal of Electronics and Communications (AEU), **63**, 200–208.
25. MILIVOJEVIC Z., MIRKOVIC M., RAJKOVIC P. (2004), *Estimating of the fundamental frequency by the using of the parametric cubic convolution interpolation*, Proceedings of International Scientific Conference UNITECH '04, 138–141, Gabrovo, Bulgaria.
26. MILIVOJEVIC Z., BRODIC D. (2011), *Estimation of the Fundamental Frequency of the Speech Signal Compressed by G.723.1 Algorithm Applying PCC Interpolation*, Journal of Electrical Engineering, **62**, 4, 181–189.
27. MILIVOJEVIC Z., MIRKOVIC M., MILIVOJEVIC S. (2006), *An Estimate of Fundamental Frequency Using PCC Interpolation – Comparative Analysis*, Information Technology and Control, **35**, 2, 131–136.
28. MILIVOJEVIC Z., MILIVOJEVIC M., BRODIC D. (2012), *The Effects of the Acute Hypoxia to the Fundamental Frequency of the Speech Signal*, Advances in Electrical and Computer Engineering, **12**, 2, 57–60.
29. MIRKOVIC M., MILIVOJEVIC Z., RAJKOVIC P. (2004), *Performances of the system with the implemented PCC algorithm for the fundamental frequency estimation*, Proceedings of XII Telecommunications Forum TELFOR '04, Section 7, Signal processing, Beograd.
30. MOON H. (2012), *A Low-Complexity Design for an MP3 Multi-Channel Audio Decoding System*, IEEE Transactions On Audio, Speech, and Language Processing, **20**, 1, 314–321.
31. PANG H.S., BAEK S.J., SUNG K.M. (2000), *Improved Fundamental Frequency Estimation Using Parametric Cubic Convolution*, IEICE Trans. Fund., **E83-A**, 12, 2747–2750.
32. PARK K.S., SCHOWENGERDT R.A. (1983), *Image reconstruction by parametric cubic convolution*, Computing, Vision, Graphics & Image Processing, **23**, 258–272.
33. RABINER L. (1977), *On the use of autocorrelation analysis for pitch detection*, IEEE Transactions On Acoustic, Speech, Signal Processing, ASSP-25, 1, 24–33.
34. RESCH B., NILSSON M., EKMAN A., KLEIJN W. (2007), *Estimation of the instantaneous pitch of speech*, IEEE Transactions On Audio, Speech, And Language Processing, **15**, 3, 813–822.
35. ROSS M., SCHAFER H., COHEN A., FREUDBERG R., MANLEY H. (1974), *Average magnitude difference function pitch extractor*, IEEE Transactions On Acoustic, Speech, Signal Processing, ASSP-22, 5, 353–362.
36. SEKHAR S.C., SREENIVAS T.V. (2004), *Effect of interpolation on PWVD computation and instantaneous frequency estimation*, IEEE Signal Processing, **84**, 107–116.
37. SHAHNAZ C., ZHU W., AHMAD M. (2012), *Pitch Estimation Based on a Harmonic Sinusoidal Autocorrelation Model and a Time-Domain Matching Scheme*, IEEE Transactions On Audio, Speech, And Language Processing, **20**, 1, 310–323.
38. VEPREK P., SCORDILIS M. (2002), *Analysis, enhancement and evaluation of five pitch determination techniques*, Speech Communication, **37**, 3–4, 249–270.
39. WANG X., HONG H. (2006), *A Novel Synchronization Invariant Audio Watermarking Scheme Based on DWT and DCT*, IEEE Transactions on Signal Processing, **54**, 12, 4835–4840.
40. YARMAN B., GUZ U., GURKAN H. (2006), *On the comparative results of SYMPES: A new method of speech modeling*, International Journal of Electronics and Communications (AEU), **60**, 421–427.
41. YEO I., KIM H.J. (2003), *Modified patchwork algorithm: a novel audio watermarking scheme*, IEEE Transactions on Speech and Audio Processing, **11**, 4, 381–386.