

(12)

EUROPEAN PATENT APPLICATION

<div>(43) Date of publication:</div> <div>10.10.2001 Bulletin 2001/41</div>	<div>(51) Int Cl.7: G10L 11/04</div>
<div>(21) Application number: 00610034.1</div>	
<div>(22) Date of filing: 06.04.2000</div>	
<div>(84) Designated Contracting States:</div> <div>AT BE CH CY DE DK ES FI FR GB GR IE IT LI LU MC NL PT SE</div> <div>Designated Extension States:</div> <div>AL LT LV MK RO SI</div>	<div>(72) Inventors:</div> <ul style="list-style-type: none"> <li>Brandel, Cecilia 224 72 Lund (SE)</li> <li>Johannisson, Henrik 212 74 Malmö (SE)</li> </ul>
<div>(71) Applicant: Telefonaktiebolaget L M Ericsson (Publ)</div> <div>126 25 Stockholm (SE)</div>	<div>(74) Representative: Boesen, Johnny Peder et al</div> <div>Hofman-Bang Zacco A/S</div> <div>Hans Bekkevolds Allé 7</div> <div>2900 Hellerup (DK)</div>

(54)

Estimating the pitch of a speech signal using an intermediate binary signal

(57)

A method of estimating the pitch of a speech signal (2) comprises the steps of sampling the speech signal to obtain a series of samples, dividing the series of samples into segments, each segment having a fixed number of consecutive samples, calculating for each segment a conformity function, and detecting peaks in the conformity function. The method further comprises the steps of providing an intermediate signal derived from the speech signal, converting the intermediate signal to a binary signal, which is set to logical "1" where the intermediate signal exceeds a pre-selected threshold and to logical "0" where the intermediate signal does not exceed the pre-selected threshold, calculating the autocorrelation of the binary signal, and using the distance between peaks in the autocorrelation of the binary signal as an estimate of the pitch. The large amount of operations needed in prior art algorithms is thus avoided. A similar device is also provided.

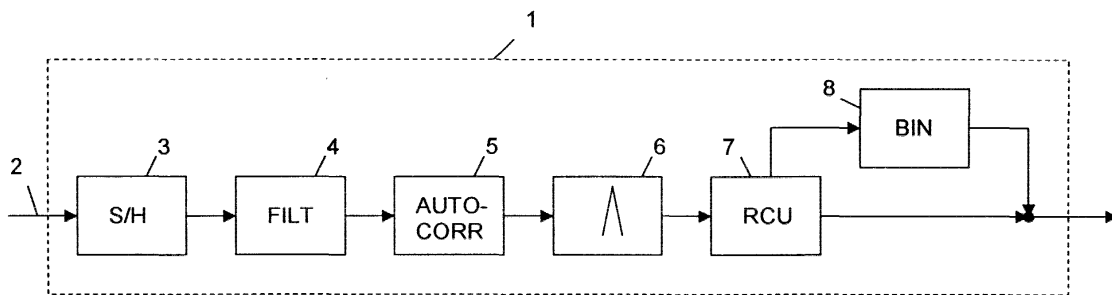


Fig. 1

## Description

**[0001]** The invention relates to a method of estimating the pitch of a speech signal, said method being of the type where the speech signal is divided into segments, a conformity function for the signal is calculated for each segment, and peaks in the conformity function are detected. The invention also relates to the use of the method in a mobile telephone. Further, the invention relates to a device adapted to estimate the pitch of a speech signal.

**[0002]** In many speech processing systems it is desirable to know the pitch period of the speech. As an example, several speech enhancement algorithms are dependent on having a correct estimate of the pitch period. One field of application where speech processing algorithms are widely used is in mobile telephones.

**[0003]** A well known way of estimating the pitch period is to use the autocorrelation function, or a similar conformity function, on the speech signal. An example of such a method is described in the article D. A. Krubsack, R. J. Niederjohn, "An Autocorrelation Pitch Detector and Voicing Decision with Confidence Measures Developed for Noise-Corrupted Speech", IEEE Transactions on Signal Processing, vol. 39, no. 2, pp. 319-329, Febr. 1991. The speech signal is divided into segments of 51.2 ms, and the standard short-time autocorrelation function is calculated for each successive speech segment. A peak picking algorithm is applied to the autocorrelation function of each segment. This algorithm starts by choosing the maximum peak (largest value) in the pitch range of 50 to 333 Hz. The period corresponding to this peak is selected as an estimate of the pitch period.

**[0004]** However, such a basic pitch estimation algorithm is not sufficient. In some cases pitch doubling can occur, i.e. the highest peak appears at twice the pitch period. The highest peak may also appear at another multiple of the true pitch period. In these cases a simple selection of the maximum peak will provide a wrong estimate of the pitch period.

**[0005]** The above-mentioned article also discloses a method of improving the algorithm in these situations. The algorithm checks for peaks at one-half, one-third, one-fourth, one-fifth, and one-sixth of the first estimate of the pitch period. If the half of the first estimate is within the pitch range, the maximum value of the autocorrelation within an interval around this half value is located. If this new peak is greater than one-half of the old peak, the new corresponding value replaces the old estimate, thus providing a new estimate which is presumably corrected for the possibility of the pitch period doubling error. This test is performed again to check for double doubling errors (fourfold errors). If this most recent test fails, a similar test is performed for tripling errors of this new estimate. This test checks for pitch period errors of sixfold. If the original test failed, the original estimate is tested (in a similar manner) for tripling errors and errors of

fivefold. The final value is used to calculate the pitch estimate.

**[0006]** However, this known algorithm is rather complex and requires a high number of calculations, and these drawbacks make it less usable in real time environments on small digital signal processors as they are used in mobile telephones and similar devices.

**[0007]** Thus, it is an object of the invention to provide a method of the above-mentioned type which is less complex than the prior art methods, such that the method is suitable for small digital signal processors.

**[0008]** According to the invention, this object is achieved in that the method further comprises the steps of providing an intermediate signal derived from the speech signal, converting the intermediate signal to a binary signal, which is set to logical "1" where the intermediate signal exceeds a pre-selected threshold and to logical "0" where the intermediate signal does not exceed the pre-selected threshold, calculating the autocorrelation of the binary signal, and using the distance between peaks in the autocorrelation of the binary signal as an estimate of the pitch.

**[0009]** The calculation of the autocorrelation of the binary signal takes only a fraction of the computational resources needed for the prior art algorithms. Since there are only values in some positions of the binary signal, the values of the resulting autocorrelation will occur around zero and around the pitch period of the speech signal, and there will only be a few values separated from zero. Thus, the pitch period can easily be estimated to the distance between the values at position zero and the values separated from zero. The large amount of operations needed in prior art algorithms where a specific value has to be found in a vector of numbers is thus avoided.

**[0010]** In one embodiment the intermediate signal may be provided by filtering the speech signal through a filter based on a set of filter parameters estimated by means of linear predictive analysis (LPA). In this way much of the smearing of the original speech signal is removed.

**[0011]** Alternatively, the intermediate signal may be provided by calculating the autocorrelation of a signal derived from the speech signal by filtering the speech signal through a filter based on a set of filter parameters estimated by means of linear predictive analysis (LPA). This solution also removes most of the smearing of the original speech signal, and further the possibility of clearer peaks in the intermediate signal is improved.

**[0012]** If the peak corresponding to the distance between the peaks is represented by a number of samples, the best estimate is achieved when the sample having the maximum amplitude of said conformity function is selected as the estimate of the pitch.

**[0013]** In an expedient embodiment of the invention the method is used in a mobile telephone, which is a typical example of a device having only limited computational resources.

**[0014]** As mentioned, the invention further relates to a device adapted to estimate the pitch of a speech signal. The device comprises means for sampling the speech signal to obtain a series of samples, means for dividing the series of samples into segments, each segment having a fixed number of consecutive samples, means for calculating for each segment a conformity function for the signal, and means for detecting peaks in the conformity function.

**[0015]** When the device further comprises means for providing an intermediate signal derived from the speech signal, means for converting said intermediate signal to a binary signal, said binary signal being set to logical "1" where the intermediate signal exceeds a pre-selected threshold and to logical "0" where the intermediate signal does not exceed the pre-selected threshold, means for calculating the autocorrelation of the binary signal, and means for using the distance between peaks in the autocorrelation of the binary signal as an estimate of the pitch, a device less complex than prior art devices is achieved, which also avoids the pitch halving situation.

**[0016]** In one embodiment the device may be adapted to provide the intermediate signal by filtering the speech signal through a filter based on a set of filter parameters estimated by means of linear predictive analysis (LPA). In this way much of the smearing of the original speech signal is removed.

**[0017]** Alternatively, the device may be adapted to provide the intermediate signal by calculating the autocorrelation of a signal derived from the speech signal by filtering the speech signal through a filter based on a set of filter parameters estimated by means of linear predictive analysis (LPA). This solution also removes most of the smearing of the original speech signal, and further the possibility of clearer peaks in the intermediate signal is improved.

**[0018]** If the peak corresponding to the distance between the peaks is represented by a number of samples, the best estimate is achieved when the device is adapted to select the sample having the maximum amplitude of said conformity function as the estimate of the pitch.

**[0019]** In an expedient embodiment of the invention, the device is a mobile telephone, which is a typical example of a device having only limited computational resources.

**[0020]** In another embodiment the device is an integrated circuit which can be used in different types of equipment.

**[0021]** The invention will now be described more fully below with reference to the drawing, in which

figure 1 shows a block diagram of a pitch detector according to the invention,

figure 2 shows the generation of a residual signal,

figure 3a shows a 20 ms segment of a voiced

speech signal,

figure 3b shows the autocorrelation function of a residual signal corresponding to the segment of figure 3a, and

figure 4 shows an example of an autocorrelation function where pitch doubling could arise.

**[0022]** Figure 1 shows a block diagram of an example of a pitch detector 1 according to the invention. A speech signal 2 is sampled with a sampling rate of 8 kHz in the sampling circuit 3 and the samples are divided into segments or frames of 160 consecutive samples. Thus, each segment corresponds to 20 ms of the speech signal. This is the sampling and segmentation normally used for the speech processing in a standard mobile telephone.

**[0023]** Each segment of 160 samples is then processed in a filter 4, which will be described in further detail below.

**[0024]** First, however, the nature of speech signals will be mentioned briefly. In a classical approach a speech signal is modelled as an output of a slowly time-varying linear filter. The filter is either excited by a quasi-periodic sequence of pulses or random noise depending on whether a voiced or an unvoiced sound is to be created. The pulse train which creates voiced sounds is produced by pressing air out of the lungs through the vibrating vocal cords. The period of time between the pulses is called the pitch period and is of great importance for the singularity of the speech. On the other hand, unvoiced sounds are generated by forming a constriction in the vocal tract and produce turbulence by forcing air through the constriction at a high velocity. This description deals with the detection of the pitch period of voiced sounds, and thus unvoiced sounds will not be further considered.

**[0025]** As speech is a varying signal also the filter has to be time-varying. However, the properties of a speech signal change relatively slowly with time. It is reasonable to believe that the general properties of speech remain fixed for periods of 10-20 ms. This has led to the basic principle that if short segments of the speech signal are considered, each segment can effectively be modelled as having been generated by exciting a linear time-invariant system during that period of time. The effect of the filter can be seen as caused by the vocal tract, the tongue, the mouth and the lips.

**[0026]** As mentioned, voiced speech can be interpreted as the output signal from a linear filter driven by an excitation signal. This is shown in the upper part of figure 2 in which the pulse train 21 is processed by the filter 22 to produce the voiced speech signal 23. A good signal for the detection of the pitch period is obtained if the excitation signal can be extracted from the speech. By estimating the filter parameters A in the block 24 and then filtering the speech through an inverse filter 25

based on the estimated filter parameters, a signal 26 similar to the excitation signal can be obtained. This signal is called the residual signal. This process is shown in the lower part of figure 2. The blocks 24 and 25 are included in the filter 4 in figure 1.

**[0027]** The estimation of the filter parameters is based on an all-pole modelling which is performed by means of the method called linear predictive analysis (LPA). The name comes from the fact that the method is equivalent with linear prediction. This method is well known in the art and will not be described in further detail here.

**[0028]** The estimation of the pitch is based on the autocorrelation of the residual signal, which is obtained as described above. Thus, the output signal from the filter 4 is taken to an autocorrelation calculation unit 5. Figure 3a shows an example of a 20 ms segment of a voiced speech signal and figure 3b the corresponding autocorrelation function of the residual signal. It will be seen from figure 3a that the actual pitch period is about 5.25 ms corresponding to 42 samples, and thus the pitch estimation should end up with this value.

**[0029]** The next step in the estimation of the pitch is to apply a peak picking algorithm to the autocorrelation function provided by the unit 5. This is done in the peak detector 6 which identifies the maximum peak (i.e. the largest value) in the autocorrelation function. The index value, i.e. the sample number or the lag, of the maximum peak is then used as a preliminary estimate of the pitch period. In the case shown in figure 3b it will be seen that the maximum peak is actually located at a lag of 42 samples. The search of the maximum peak is only performed in the range where a pitch period is likely to be located. In this case the range is set to 60-333 Hz.

**[0030]** However, this basic pitch estimation algorithm is not always sufficient. In some cases pitch doubling may occur, i.e. due to distortion the peak in the autocorrelation function corresponding to the true pitch period is not the highest peak, but instead the highest peak appears at twice the pitch period. The highest peak could also appear at other multiples of the actual pitch period (pitch tripling, etc.) although this occurs relatively rarely. A typical example where pitch doubling would arise is shown in figure 4 which again shows the autocorrelation function of the residual signal. Here, too, the correct pitch period would be around 42 samples, but the peak at twice the pitch period, i.e. around 84 samples, is actually higher than the one at 42 samples. The basic pitch estimation algorithm would therefore estimate the pitch period to 84 samples and pitch doubling would thus occur.

**[0031]** To avoid the problem of pitch doubling the pitch detection algorithm is therefore improved as described below.

**[0032]** After the preliminary pitch estimate has been determined, it is checked in the risk check unit 7 whether there is any risk of pitch doubling. All peaks with a peak value higher than 75% of the maximum peak are detected and the further processing depends on the result of

this detection. If only one peak is detected, i.e. the original maximum peak, there is no need to perform a process to avoid pitch doubling. In this situation the preliminary pitch estimate is used as the final pitch estimate.

If, however, more than one peak is detected, there is a risk of pitch doubling and a further algorithm must be performed to ensure that the correct peak is selected as the pitch estimate. This is performed in the unit 8.

**[0033]** To identify the peak corresponding to the actual pitch period a modified signal is provided based on the location of the peaks in the autocorrelation of the residual signal. This modified signal, referred to as binary signal, consists of only ones and zeros. The binary signal is set to one where the high peaks are found in the autocorrelation sequence. All other values are set to zero, and then the autocorrelation of the binary signal is calculated. Since there are only values in some positions in the binary signal, the resulting autocorrelation will only have a few values separated from zero, and these values will occur around the pitch period of the signal. The pitch period is estimated by observing the distance between the indexes of the values around zero and those separated from zero. If the group of values separated from zero contains only a single value, it is selected as the estimate of the pitch period. If there is more than one value in the group, the one with the highest amplitude in the autocorrelation of the residual signal is chosen.

**[0034]** Sometimes cases may arise where the peak at lag zero is the only peak present. This situation will occur when a peak has been split on two samples and there are no other high peaks in the autocorrelation of the residual signal. In this case the preliminary pitch estimate is chosen as the final pitch estimate.

**[0035]** This algorithm is very simple, and therefore it is well suited in e.g. mobile telephones in which the computational resources are severely limited, and a demand for a low-complexity algorithm is thus placed upon the system. The algorithm may also be implemented in an integrated circuit which may then be used in other types of equipment.

**[0036]** Although a preferred embodiment of the present invention has been described and shown, the invention is not restricted to it, but may also be embodied in other ways within the scope of the subject-matter defined in the following claims.

**[0037]** Thus, the autocorrelation function may be calculated directly of the speech signal instead of the residual signal, or other conformity functions may be used instead of the autocorrelation function. As an example, a cross correlation could be calculated between the speech signal and the residual signal.

**[0038]** Further, different sampling rates and sizes of the segments may be used.

## Claims

1. A method of estimating the pitch of a speech signal (2), said method comprising the steps of:

- sampling the speech signal to obtain a series of samples,
- dividing the series of samples into segments, each segment having a fixed number of consecutive samples,
- calculating for each segment a conformity function for the signal, and
- detecting peaks in the conformity function,

**characterized in that** the method further comprises the steps of:

- providing an intermediate signal derived from the speech signal,
- converting said intermediate signal to a binary signal, said binary signal being set to logical "1" where the intermediate signal exceeds a pre-selected threshold and to logical "0" where the intermediate signal does not exceed the pre-selected threshold,
- calculating the autocorrelation of the binary signal, and
- using the distance between peaks in the autocorrelation of the binary signal as an estimate of the pitch.

2. A method according to claim 1, **characterized in that** the intermediate signal is provided by filtering the speech signal through a filter (4) based on a set of filter parameters estimated by means of linear predictive analysis (LPA).

3. A method according to claim 1, **characterized in that** the intermediate signal is provided by calculating the autocorrelation of a signal derived from the speech signal by filtering the speech signal through a filter (4) based on a set of filter parameters estimated by means of linear predictive analysis (LPA).

4. A method according to any one of claims 1 to 3, **characterized in that** it further comprises the step of:

- selecting, if the peak corresponding to the distance between the peaks is represented by a number of samples, the sample having the maximum amplitude of said conformity function as the estimate of the pitch.

5. Use of the method according to any one of claims 1 to 4 in a mobile telephone.

6. A device adapted to estimate the pitch of a speech

signal, and comprising:

- means (3) for sampling the speech signal to obtain a series of samples,
- means for dividing the series of samples into segments, each segment having a fixed number of consecutive samples,
- means (5) for calculating for each segment a conformity function for the signal, and
- means (6) for detecting peaks in the conformity function,

**characterized in that** the device further comprises:

- means for providing an intermediate signal derived from the speech signal,
- means (8) for converting said intermediate signal to a binary signal, said binary signal being set to logical "1" where the intermediate signal exceeds a pre-selected threshold and to logical "0" where the intermediate signal does not exceed the pre-selected threshold,
- means (5) for calculating the autocorrelation of the binary signal, and
- means for using the distance between peaks in the autocorrelation of the binary signal as an estimate of the pitch.

7. A device according to claim 6, **characterized in that** the device is adapted to provide the intermediate signal by filtering the speech signal through a filter (4) based on a set of filter parameters estimated by means of linear predictive analysis (LPA).

8. A device according to claim 6, **characterized in that** the device is adapted to provide the intermediate signal by calculating the autocorrelation of a signal derived from the speech signal by filtering the speech signal through a filter (4) based on a set of filter parameters estimated by means of linear predictive analysis (LPA).

9. A device according to any one of claims 6 to 8, **characterized in that** it is further adapted to select, if the peak corresponding to the distance between the peaks is represented by a number of samples, the sample having the maximum amplitude of said conformity function as the estimate of the pitch.

10. A device according to any one of claims 6 to 9, **characterized in that** the device is a mobile telephone.

11. A device according to any one of claims 6 to 9, **characterized in that** the device is an integrated circuit.

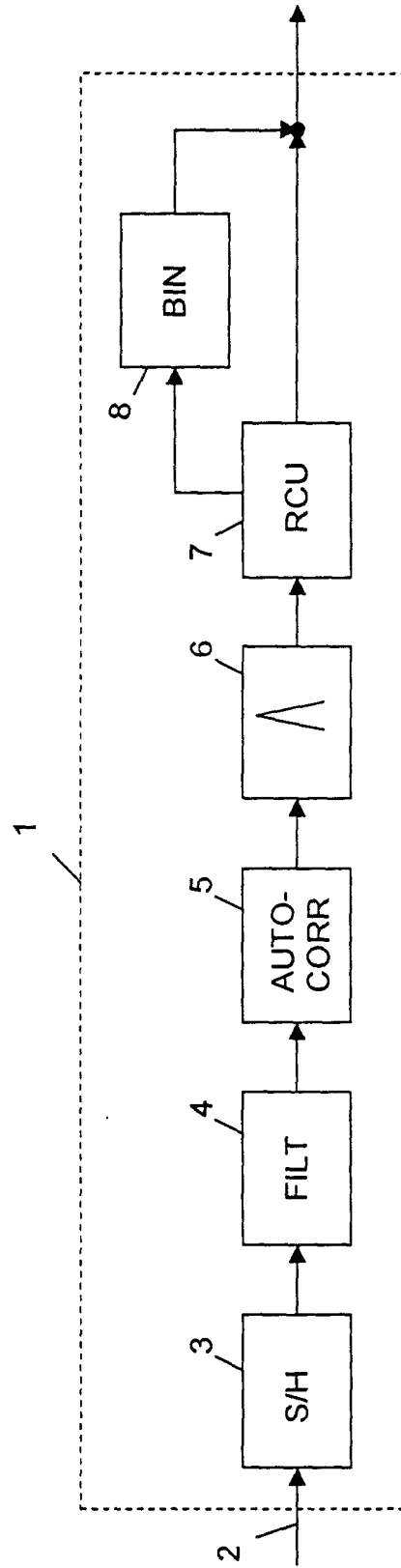


Fig. 1

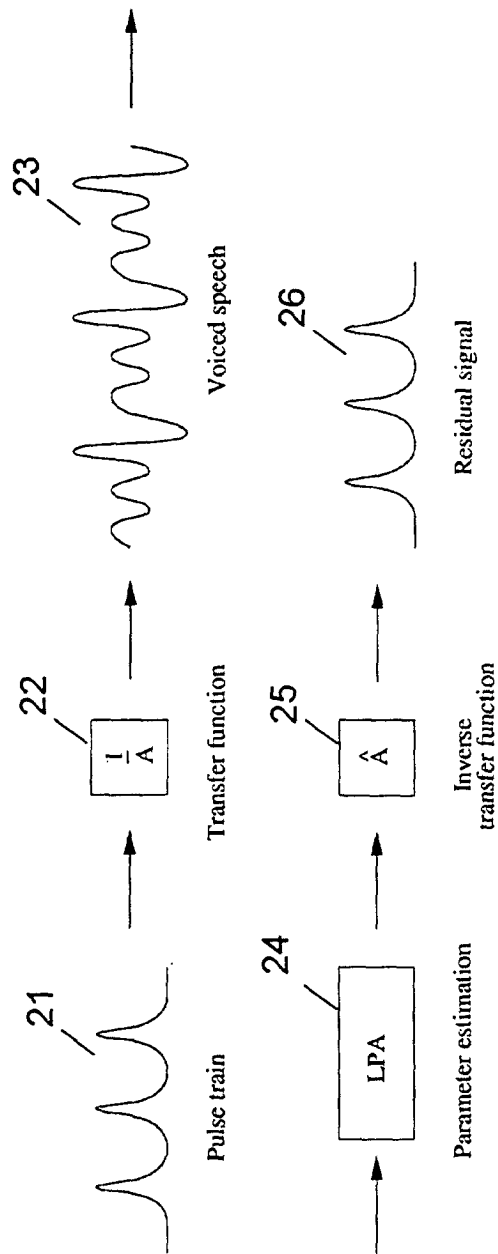


Fig. 2

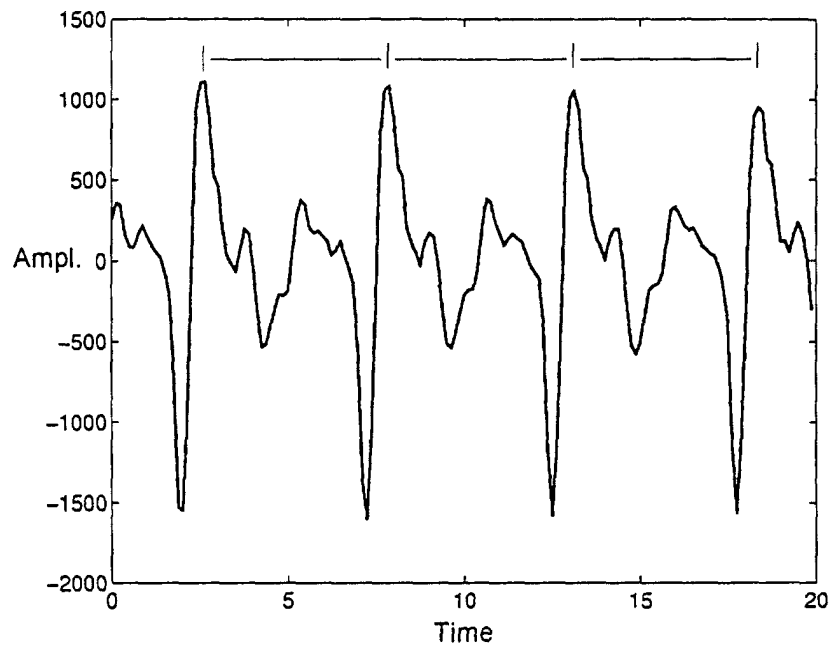


Fig. 3a

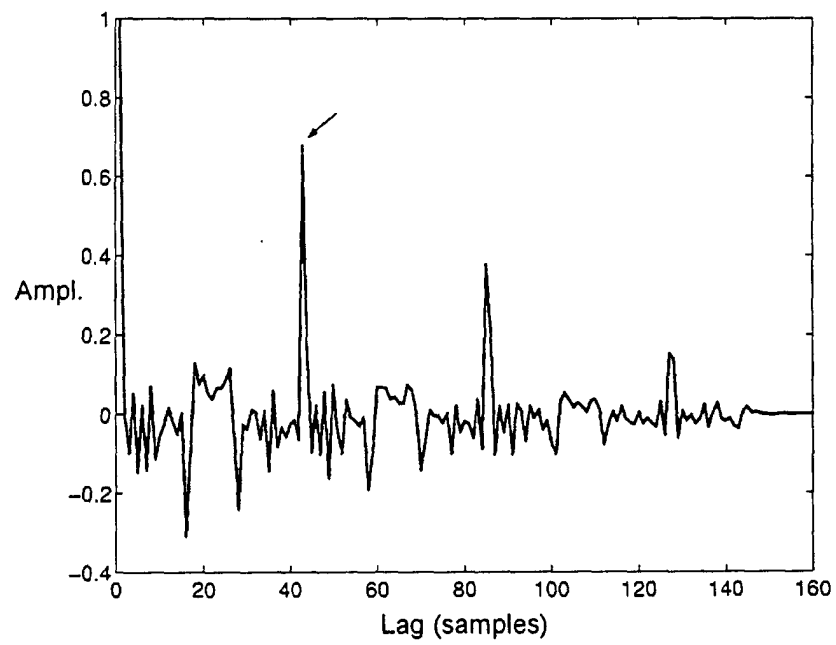


Fig. 3b



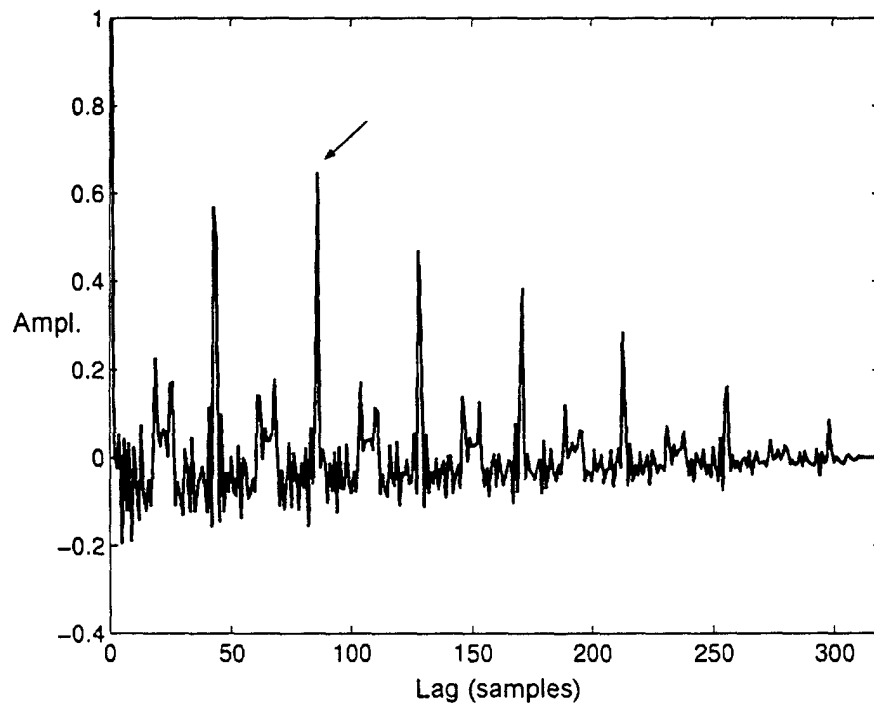


Fig. 4



European Patent  
Office

# EUROPEAN SEARCH REPORT

Application Number  
EP 00 61 0034

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int.Cl.7)
A	<p>ALKULAIBI A ET AL: "Fast 3-level binary higher order statistics for simultaneous voiced/unvoiced and pitch detection of a speech signal"</p> <p>SIGNAL PROCESSING. EUROPEAN JOURNAL DEVOTED TO THE METHODS AND APPLICATIONS OF SIGNAL PROCESSING, NL, ELSEVIER SCIENCE PUBLISHERS B.V. AMSTERDAM, vol. 63, no. 2, 1 December 1997 (1997-12-01), pages 133-140, XP004102257</p> <p>ISSN: 0165-1684</p> <p>* abstract *</p> <p>* page 134, right-hand column, last paragraph - page 136, right-hand column, paragraph 1 *</p>	1-11	G10L11/04
A	<p>US 5 970 441 A (MEKURIA)</p> <p>19 October 1999 (1999-10-19)</p> <p>* abstract *</p> <p>* column 2, line 66 - column 3, line 50 *</p>	1-3, 5-8, 10, 11	<p>TECHNICAL FIELDS SEARCHED (Int.Cl.7)</p> <p>G10L</p>
The present search report has been drawn up for all claims			
Place of search		Date of completion of the search	Examiner
THE HAGUE		4 September 2000	Quélavoine, R
<p>CATEGORY OF CITED DOCUMENTS</p> <p>X : particularly relevant if taken alone</p> <p>Y : particularly relevant if combined with another document of the same category</p> <p>A : technological background</p> <p>O : non-written disclosure</p> <p>P : intermediate document</p> <p>T : theory or principle underlying the invention</p> <p>E : earlier patent document, but published on, or after the filing date</p> <p>D : document cited in the application</p> <p>L : document cited for other reasons</p> <p>&amp; : member of the same patent family, corresponding document</p>			

EPO FORM 1503 03/82 (P04C01)

**ANNEX TO THE EUROPEAN SEARCH REPORT  
ON EUROPEAN PATENT APPLICATION NO.**

EP 00 61 0034

This annex lists the patent family members relating to the patent documents cited in the above-mentioned European search report.  
The members are as contained in the European Patent Office EDP file on  
The European Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

04-09-2000

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 5970441 A	19-10-1999	AU 8565998 A	16-03-1999
		EP 1008140 A	14-06-2000
		WO 9910879 A	04-03-1999
-----			