

(19)



(11)

EP 1 454 315 B1

(12)

EUROPEAN PATENT SPECIFICATION

(45) Date of publication and mention of the grant of the patent:
04.04.2007 Bulletin 2007/14

(51) Int Cl.:
G10L 19/08 (2006.01)

(21) Application number: **02784985.0**

(86) International application number:
PCT/CA2002/001948

(22) Date of filing: **13.12.2002**

(87) International publication number:
WO 2003/052744 (26.06.2003 Gazette 2003/26)

(54) SIGNAL MODIFICATION METHOD FOR EFFICIENT CODING OF SPEECH SIGNALS

**SIGNALÄNDERUNGSVERFAHREN ZUR EFFIZIENTEN KODIERUNG VON SPRACHSIGNALEN
PROCEDE DE MODIFICATION DU SIGNAL ASSURANT LE CODAGE EFFICACE DES SIGNAUX DE PAROLE**

(84) Designated Contracting States:
**AT BE BG CH CY CZ DE DK EE ES FI FR GB GR
IE IT LI LU MC NL PT SE SI SK TR**
Designated Extension States:
AL LT LV MK RO

(74) Representative: **Derry, Paul Stefan et al
Venner Shipley LLP
20 Little Britain
London EC1A 7DH (GB)**

(30) Priority: **14.12.2001 CA 2365203**

(56) References cited:
**US-A- 5 974 377 US-A1- 2001 023 395
US-B1- 6 449 590**

(43) Date of publication of application:
08.09.2004 Bulletin 2004/37

(60) Divisional application:
06125444.7 / 1 758 101

- **TAMMI M ET AL: "Signal modification for voiced wideband speech coding and its application for IS-95 system" 2002 IEEE PROCEEDINGS SPEECH CODING WORKSHOP., 6 - 9 October 2002, pages 35-37, XP002250153 IBARAKI, JAPAN, Piscataway, NJ, USA, IEEE, USA ISBN: 0-7803-7549-1**
- **CHUI S P ET AL: "Low delay CELP coding at 8 kbps using classified voiced and unvoiced excitation codebooks" PROCEEDINGS OF ICSIPNN '94. INTERNATIONAL CONFERENCE ON SPEECH, IMAGE PROCESSING AND NEURAL NETWORKS (CAT. NO.94TH0638-7), 13 - 16 April 1994, pages 472-475 vol.2, XP002250152 HONG KONG, New York, NY, USA, IEEE, USA ISBN: 0-7803-1865-X**

(73) Proprietor: **Nokia Corporation
02150 Espoo (FI)**

- (72) Inventors:
- **TAMMI, Mikko
33 720 Tampere (FI)**
 - **JELINEK, Milan
North Hatley, Quebec, J0B 2C0 (CA)**
 - **LAFLAMME, Claude
Orford, Quebec, J1X 6W1 (CA)**
 - **RUOPPILA, Vesa
Montreal, Quebec, H2L 3R7 (CA)**

EP 1 454 315 B1

Note: Within nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition fee has been paid. (Art. 99(1) European Patent Convention).

Description**FIELD OF THE INVENTION**

5 [0001] The present invention relates generally to the encoding and decoding of sound signals in communication systems. More specifically, the present invention is concerned with a signal modification technique applicable to, in particular but not exclusively, code-excited linear prediction (CELP) coding.

BACKGROUND OF THE INVENTION

10 [0002] Demand for efficient digital narrow- and wideband speech coding techniques with a good trade-off between the subjective quality and bit rate is increasing in various application areas such as teleconferencing, multimedia, and wireless communications. Until recently, the telephone bandwidth constrained into a range of 200-3400 Hz has mainly been used in speech coding applications. However, wideband speech applications provide increased intelligibility and naturalness in communication compared to the conventional telephone bandwidth. A bandwidth in the range 50-7000 Hz has been found sufficient for delivering a good quality giving an impression of face-to-face communication. For general audio signals, this bandwidth gives an acceptable subjective quality, but is still lower than the quality of FM radio or CD that operate in ranges of 20-16000 Hz and 20-20000 Hz, respectively.

15 [0003] A speech encoder converts a speech signal into a digital bit stream which is transmitted over a communication channel or stored in a storage medium. The speech signal is digitized, that is sampled and quantized with usually 16-bits per sample. The speech encoder has the role of representing these digital samples with a smaller number of bits while maintaining a good subjective speech quality. The speech decoder or synthesizer operates on the transmitted or stored bit stream and converts it back to a sound signal.

20 [0004] *Code-Excited Linear Prediction (CELP)* coding is one of the best techniques for achieving a good compromise between the subjective quality and bit rate. This coding technique is a basis of several speech coding standards both in wireless and wire line applications. In CELP coding, the sampled speech signal is processed in successive blocks of N samples usually called *frames*, where N is a predetermined number corresponding typically to 10-30 ms. A linear prediction (LP) filter is computed and transmitted every frame. The computation, of the LP filter typically needs a *look ahead*, i.e. a 5-10 ms speech segment from the subsequent frame. The N -sample frame is divided into smaller blocks called *subframes*. Usually the number of subframes is three or four resulting in 4-10 ms subframes. In each subframe, an excitation signal is usually obtained from two components: a past excitation and an innovative, fixed-codebook excitation. The component formed from the past excitation is often referred to as the adaptive codebook or pitch excitation. The parameters characterizing the excitation signal are coded and transmitted to the decoder, where the reconstructed excitation signal is used as the input of the LP filter.

25 [0005] In conventional CELP coding, long term prediction for mapping the past excitation to the present is usually performed on a subframe basis. Long term prediction is characterized by a delay parameter and a pitch gain that are usually computed, coded and transmitted to the decoder for every subframe. At low bit rates, these parameters consume a substantial proportion of the available bit budget. Signal modification techniques [1-7]

40 [1] W.B. Kleijn, P. Kroon, and D. Nahumi, "The RCELP speech-coding algorithm," *European Transactions on Telecommunications*, Vol. 4, No. 5, pp. 573-582, 1994.

[2] W.B. Kleijn, R.P. Ramachandran, and P. Kroon, "Interpolation of the pitch-predictor parameters in analysis-by-synthesis speech coders," *IEEE Transactions on Speech and Audio Processing*, Vol. 2, No. 1, pp. 42-54, 1994.

45 [3] Y. Gao, A. Benyassine, J. Thyssen, H. Su, and E. Shlomot, "EX-CELP: A speech coding paradigm," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Salt Lake City, Utah, U.S.A., pp. 689-692, 7-11 May 2001.

50 [4] US Patent 5,704,003, "RCELP coder," Lucent Technologies Inc., (W.B. Kleijn and D. Nahumi), Filing Date: 19 September 1995.

[5] European Patent Application 0 602 826 A2, "Time shifting for analysis-by-synthesis coding," AT&T Corp., (B. Kleijn), Filing Date: 1 December 1993.

55 [6] Patent Application WO 001 11653, "Speech encoder with continuous warping combined with long term prediction," Conexant Systems Inc., (Y. Gao), Filing Date: 24 August 1999.

[7] Patent Application WO 00/11654, "Speech encoder adaptively applying pitch preprocessing with continuous warping," Conexant Systems Inc., (H. Su and Y. Gao), Filing Date: 24 Aug. 1999.

5 [0006] improve the performance of long term prediction at low bit rates by adjusting the signal to be coded. This is done by adapting the evolution of the pitch cycles in the speech signal to fit the long term prediction delay, enabling to transmit only one delay parameter per frame. Signal modification is based on the premise that it is possible to render the difference between the modified speech signal and the original speech signal inaudible. The CELP coders utilizing signal modification are often referred to as generalized analysis-by-synthesis or *relaxed CELP* (RCELP) coders.

10 [0007] Signal modification techniques adjust the pitch of the signal to a predetermined delay contour. Long term prediction then maps the past excitation signal to the present subframe using this delay contour and scaling by a gain parameter. The delay contour is obtained straightforwardly by interpolating between two open-loop pitch estimates, the first obtained in the previous frame and the second in the current frame. Interpolation gives a delay value for every time instant of the frame. After the delay contour is available, the pitch in the subframe to be coded currently is adjusted to follow this artificial contour by warping, i.e. changing the time scale of the signal.

15 [0008] In *discontinuous warping* [1, 4 and 5]

[1] W.B. Kleijn, P. Kroon, and D. Nahumi, "The RCELP speech-coding algorithm," European Transactions on Telecommunications, Vol. 4, No. 5. pp. 573-582, 1994.

20 [4] US Patent 5,704,003, "RCELP coder," Lucent Technologies Inc., (W.B. Kleijn and D. Nahumi), Filing Date: 19 September 1995.

[5] European Patent Application 0 602 826 A2, "Time shifting for analysis-by-synthesis coding," AT&T Corp., (B. Kleijn), Filing Date: 1 December 1993.

25 a signal segment is shifted in time without altering the segment length. Discontinuous warping requires a procedure for handling the resulting overlapping or missing signal portions. *Continuous warping* [2, 3, 6, 7]

30 [2] W.B. Kleijn, R.P. Ramachandran, and P. Kroon, "Interpolation of the pitch-predictor parameters in analysis-by-synthesis speech coders," IEEE Transactions on Speech and Audio Processing, Vol. 2, No. 1, pp. 42-54, 1994.

[3] Y. Gao, A. Benyassine, J. Thyssen, H. Su, and E. Shlomot, "EX-CELP: A speech coding paradigm," IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Salt Lake City, Utah, U.S.A., pp. 689-692, 7-11 May 2001.

35 [6] Patent Application WO 00/11653, "Speech encoder with continuous warping combined with long term prediction," Conexant Systems Inc., (Y. Gao), Filing Date: 24 Aug. 1999.

40 [7] Patent Application WO 00/11654, "Speech encoder adaptively applying pitch preprocessing with continuous warping," Conexant Systems Inc., (H. Su and Y. Gao), Filing Date 24 Aug. 1999.

45 either contracts or expands a signal segment. This is done using a time continuous approximation for the signal segment and re-sampling it to a desired length with unequal sampling intervals determined based on the delay contour. For reducing artifacts in these operations, the tolerated change in the time scale is kept small. Moreover, warping is typically done using the LP residual signal or the weighted speech signal to reduce the resulting distortions. The use of these signals instead of the speech signal also facilitates detection of pitch pulses and low-power regions in between them, and thus the determination of the signal segments for warping. The actual modified speech signal is generated by inverse filtering.

50 [0009] After the signal modification is done for the current subframe, the coding can proceed in any conventional manner except the adaptive codebook excitation is generated using the predetermined delay contour. Essentially the same signal modification techniques can be used both in narrow- and wideband CELP coding.

[0010] Signal modification techniques can also be applied in other types of speech coding methods such as waveform interpolation coding and sinusoidal coding for instance in accordance with [8].

55 [8] US Patent 6,223,151, "Method and apparatus for pre-processing speech signals prior to coding by transform-based speech coders," Telefon Aktie Bolaget LM Ericsson, (W.B. Kleijn and T. Eriksson), Filing Date 10 Feb. 1999.

SUMMARY OF THE INVENTION

[0011] The invention is defined by the claims.

[0012] Follows a non-restrictive description of illustrative embodiments of the invention given by way of example only with reference to the accompanying drawings.

BRIEF DESCRIPTION OF THE DRAWINGS

[0013]

Figure 1 is an illustrative example of original and modified residual signals for one frame;

Figure 2 is a functional block diagram of an illustrative embodiment of a signal modification method according to the invention;

Figure 3 is a schematic block diagram of an illustrative example of speech communication system showing the use of speech encoder and decoder;

Figure 4 is a schematic block diagram of an illustrative embodiment of speech encoder that utilizes a signal modification method;

Figure 5 is a functional block diagram of an illustrative embodiment of pitch pulse search;

Figure 6 is an illustrative example of located pitch pulse positions and a corresponding pitch cycle segmentation for one frame;

Figure 7 is an illustrative example on determining a delay parameter when the number of pitch pulses is three ($c = 3$);

Figure 8 is an illustrative example of delay interpolation (thick line) over a speech frame compared to linear interpolation (thin line);

Figure 9 is an illustrative example of a delay contour over ten frames selected in accordance with the delay interpolation (thick line) of Figure 8 and linear interpolation (thin line) when the correct pitch value is 52 samples;

Figure 10 is a functional block diagram of the signal modification method that adjusts the speech frame to the selected delay contour in accordance with an illustrative embodiment of the present invention;

Figure 11 is an illustrative example on updating the target signal $\tilde{w}(t)$ using a determined optimal shift δ , and on replacing the signal segment $w_s(k)$ with interpolated values shown as gray dots;

Figure 12 is a functional block diagram of a rate determination logic in accordance with an illustrative embodiment of the present invention; and

Figure 13 is a schematic block diagram of an illustrative embodiment of speech decoder that utilizes the delay contour formed in accordance with an illustrative embodiment of the present invention.

DETAILED DESCRIPTION OF THE ILLUSTRATIVE EMBODIMENTS

[0014] Although the illustrative embodiments of the present invention will be described in relation to speech signals and the 3GPP AMR Wideband Speech Codec AMR-WB Standard (ITU-T G.722.2), it should be kept in mind that the concepts of the present invention may be applied to other types of sound signals as well as other speech and audio coders.

[0015] Figure 1 illustrates an example of modified residual signal 12 within one frame. As shown in Figure 1, the time shift in the modified residual signal 12 is constrained such that this modified residual signal is time synchronous with the original, unmodified residual signal 11 at frame boundaries occurring at time instants t_{n-1} and t_n . Here n refers to the index of the present frame.

[0016] More specifically, the time shift is controlled implicitly with a delay contour employed for interpolating the delay parameter over the current frame. The delay parameter and contour are determined considering the time alignment constraints at the above-mentioned frame boundaries. When linear interpolation is used to force the time alignment, the

resulting delay parameters tend to oscillate over several frames. This often causes annoying artifacts to the modified signal whose pitch follows the artificial oscillating delay contour. Use of a properly chosen nonlinear interpolation technique for the delay parameter will substantially reduce these oscillations.

[0017] A functional block diagram of the illustrative embodiment of the signal modification method according to the invention is presented in Figure 2.

[0018] The method starts, in "pitch cycle search" block 101, by locating individual pitch pulses and pitch cycles. The search of block 101 utilizes an open-loop pitch estimate interpolated over the frame. Based on the located pitch pulses, the frame is divided into pitch cycle segments, each containing one pitch pulse and restricted inside the frame boundaries t_{n-1} and t_n .

[0019] The function of the "delay curve selection" block 103 is to determine a delay parameter for the long term predictor and form a delay contour for interpolating this delay parameter over the frame. The delay parameter and contour are determined considering the time synchrony constrains at frame boundaries t_{n-1} and t_n . The delay parameter determined in block 103 is coded and transmitted to the decoder when signal modification is enabled for the current frame.

[0020] The actual signal modification procedure is conducted in the "pitch synchronous signal modification" block 105. Block 105 first forms a target signal based on the delay contour determined in block 103 for subsequently matching the individual pitch cycle segments into this target signal. The pitch cycle segments are then shifted one by one to maximize their correlation with this target signal. To keep the complexity at a low level, no continuous time warping is applied while searching the optimal shift and shifting the segments.

[0021] The illustrative embodiment of signal modification method as disclosed in the present specification is typically enabled only on purely voiced speech frames. For instance, transition frames such as voiced onsets are not modified because of a high risk of causing artifacts. In purely voiced frames, pitch cycles usually change relatively slowly and therefore small shifts suffice to adapt the signal to the long term prediction model. Because only small, cautious signal adjustments are made, the probability of causing artifacts is minimized.

[0022] The signal modification method constitutes an efficient classifier for purely voiced segments, and hence a rate determination mechanism to be used in a source-controlled coding of speech signals. Every block 101, 103 and 105 of Figure 2 provide several indicators on signal periodicity and the suitability of signal modification in the current frame. These indicators are analyzed in logic blocks 102, 104 and 106 in order to determine a proper coding mode and bit rate for the current frame. More specifically, these logic blocks 102, 104 and 106 monitor the success of the operations conducted in blocks 101, 103, and 105.

[0023] If block 102 detects that the operation performed in block 101 is successful, the signal modification method is continued in block 103. When this block 102 detects a failure in the operation performed in block 101, the signal modification procedure is terminated and the original speech frame is preserved intact for coding (see block 108 corresponding to normal mode (no signal modification)).

[0024] If block 104 detects that the operation performed in block 103 is successful, the signal modification method is continued in block 105. When, on the contrary, this block 104 detects a failure in the operation performed in block 103, the signal modification procedure is terminated and the original speech frame is preserved intact for coding (see block 108 corresponding to normal mode (no signal modification)).

[0025] If block 106 detects that the operation performed in block 105 is successful, a low bit rate mode with signal modification is used (see block 107). On the contrary, when this block 106 detects a failure in the operation performed in block 105 the signal modification procedure is terminated, and the original speech frame is preserved intact for coding (see block 108 corresponding to normal mode (no signal modification)). The operation of the blocks 101-108 will be described in detail later in the present specification.

[0026] Figure 3 is a schematic block diagram of an illustrative example of speech communication system depicting the use of speech encoder and decoder. The speech communication system of Figure 3 supports transmission and reproduction of a speech signal across a communication channel 205. Although it may comprise for example a wire, an optical link or a fiber link, the communication channel 205 typically comprises at least in part a radio frequency link. The radio frequency link often supports multiple, simultaneous speech communications requiring shared bandwidth resources such as may be found with cellular telephony. Although not shown, the communication channel 205 may be replaced by a storage device that records and stores the encoded speech signal for later playback.

[0027] On the transmitter side, a microphone 201 produces an analog speech signal 210 that is supplied to an analog-to-digital (A/D) converter 202.

The function of the A/D converter 202 is to convert the analog speech signal 210 into a digital speech signal 211. A speech encoder 203 encodes the digital speech signal 211 to produce a set of coding parameters 212 that are coded into binary form and delivered to a channel encoder 204. The channel encoder 204 adds redundancy to the binary representation of the coding parameters before transmitting them into a bitstream 213 over the communication channel 205.

[0028] On the receiver side, a channel decoder 206 is supplied with the above mentioned redundant binary representation of the coding parameters from the received bitstream 214 to detect and correct channel errors that occurred in

the transmission. A speech decoder 207 converts the channel-error-corrected bitstream 215 from the channel decoder 206 back to a set of coding parameters for creating a synthesized digital speech signal 216. The synthesized speech signal 216 reconstructed by the speech decoder 207 is converted to an analog speech signal 217 through a digital-to-analog (D/A) converter 208 and played back through a loudspeaker unit 209.

5 [0029] Figure 4 is a schematic block diagram showing the operations performed by the illustrative embodiment of speech encoder 203 (Figure 3) incorporating the signal modification functionality. The present specification presents a novel implementation of this signal modification functionality of block 603 in Figure 4. The other operations performed by the speech encoder 203 are well known to those of ordinary skill in the art and have been described, for example, in the publication [10]

10 [10] 3GPP TS 26.190, "AMR Wideband Speech Codec: Transcoding Functions," *3GPP Technical Specification*, which is incorporated herein by reference. When not stated otherwise, the implementation of the speech encoding and decoding operations in the illustrative embodiments and examples of the present invention will comply with the AMR Wideband Speech Codec (AMR-WB) Standard.

15 [0030] The speech encoder 203 as shown in Figure 4 encodes the digitized speech signal using one or a plurality of coding modes. When a plurality of coding modes are used and the signal modification functionality is disabled in one of these modes, this particular mode will operate in accordance with well established standards known to those of ordinary skill in the art.

20 [0031] Although not shown in Figure 4, the speech signal is sampled at a rate of 16 kHz and each speech signal sample is digitized. The digital speech signal is then divided into successive frames of given length, and each of these frames is divided into a given number of successive subframes. The digital speech signal is further subjected to pre-processing as taught by the AMR-WB standard. This preprocessing includes high-pass filtering, preemphasis filtering using a filter $P(z) = 1 - 0.68z^{-1}$ and down-sampling from the sampling rate of 16 kHz to 12.8 kHz. The subsequent operations of Figure 4 assume that the input speech signal $s(t)$ has been preprocessed and down-sampled to the sampling rate of 12.8 kHz.

25 [0032] The speech encoder 203 comprises an LP (Linear Prediction) analysis and quantization module 601 responsive to the input, preprocessed digital speech signal $s(t)$ 617 to compute and quantize the parameters $a_0, a_1, a_2, \dots, a_{n_A}$ of the LP filter $1/A(z)$, wherein n_A is the order of the filter and $A(z) = a_0 + a_1Z^{-1} + a_2Z^{-2} + \dots + a_{n_A}Z^{-n_A}$. The binary representation 616 of these quantized LP filter parameters is supplied to the multiplexer 614 and subsequently multiplexed into the bitstream 615. The non-quantized and quantized LP filter parameters can be interpolated for obtaining the corresponding LP filter parameters for every subframe.

30 [0033] The speech encoder 203 further comprises a pitch estimator 602 to compute open-loop pitch estimates 619 for the current frame in response to the LP filter parameters 618 from the LP analysis and quantization module 601. These open-loop pitch estimates 619 are interpolated over the frame to be used in a signal modification module 603.

35 [0034] The operations performed in the LP analysis and quantization module 601 and the pitch estimator 602 can be implemented in compliance with the above-mentioned AMR-WB Standard.

[0035] The signal modification module 603 of Figure 4 performs a signal modification operation prior to the closed-loop pitch search of the adaptive codebook excitation signal for adjusting the speech signal to the determined delay contour $d(t)$. In the illustrative embodiment, the delay contour $d(t)$ defines a long term prediction delay for every sample of the frame. By construction the delay contour is fully characterized over the frame $t \in (t_{n-1}, t_n]$ by a delay parameter 40 $d_n = d(t_n)$ and its previous value $d_{n-1} = d(t_{n-1})$ that are equal to the value of the delay contour at frame boundaries. The delay parameter 620 is determined as a part of the signal modification operation, and coded and then supplied to the multiplexer 614 where it is multiplexed into the bitstream 615.

45 [0036] The delay contour $d(t)$ defining a long term prediction delay parameter for every sample of the frame is supplied to an adaptive codebook 607. The adaptive codebook 607 is responsive to the delay contour $d(t)$ to form the adaptive codebook excitation $u_b(t)$ of the current subframe from the excitation $u(t)$ using the delay contour $d(t)$ as $u_b(t) = u(t - d(t))$. Thus the delay contour maps the past sample of the excitation signal $u(t - d(t))$ to the present sample in the adaptive codebook excitation $u_b(t)$.

50 [0037] The signal modification procedure produces also a modified residual signal $\tilde{r}(t)$ to be used for composing a modified target signal 621 for the closed-loop search of the fixed-codebook excitation $u_c(t)$. The modified residual signal $\tilde{r}(t)$ is obtained in the signal modification module 603 by warping the pitch cycle segments of the LP residual signal, and is supplied to the computation of the modified target signal in module 604. The LP synthesis filtering of the modified residual signal with the filter $1/A(z)$ yields then in module 604 the modified speech signal. The modified target signal 621 of the fixed-codebook excitation search is formed in module 604 in accordance with the operation of the AMR-WB Standard, but with the original speech signal replaced by its modified version.

55 [0038] After the adaptive codebook excitation $u_b(t)$ and the modified target signal 621 have been obtained for the current subframe, the encoding can further proceed using conventional means.

[0039] The function of the closed-loop fixed-codebook excitation search is to determine the fixed-codebook excitation signal $u_c(t)$ for the current subframe. To schematically illustrate the operation of the closed-loop fixed-codebook search,

the fixed-codebook excitation $u_c(t)$ is gain scaled through an amplifier 610. In the same manner, the adaptive-codebook excitation $u_b(t)$ is gain scaled through an amplifier 609. The gain scaled adaptive and fixed-codebook excitations $U_b(t)$ and $U_c(t)$ are summed together through an adder 611 to form a total excitation signal $u(t)$. This total excitation signal $u(t)$ is processed through an LP synthesis filter $1/A(z)$ 612 to produce a synthesis speech signal 625 which is subtracted from the modified target signal 621 through an adder 605 to produce an error signal 626. An error weighting and minimization module 606 is responsive to the error signal 626 to calculate, according to conventional methods, the gain parameters for the amplifiers 609 and 610 every subframe. The error weighting and minimization module 606 further calculates, in accordance with conventional methods and in response to the error signal 626, the input 627 to the fixed codebook 608. The quantized gain parameters 622 and 623 and the parameters 624 characterizing the fixed-codebook excitation signal $u_c(t)$ are supplied to the multiplexer 614 and multiplexed into the bitstream 615. The above procedure is done in the same manner both when signal modification is enabled or disabled.

[0040] It should be noted that, when the signal modification functionality is disabled, the adaptive excitation codebook 607 operates according to conventional methods. In this case, a separate delay parameter is searched for every subframe in the adaptive codebook 607 to refine the open-loop pitch estimates 619. These delay parameters are coded, supplied to the multiplexer 614 and multiplexed into the bitstream 615. Furthermore, the target signal 621 for the fixed-codebook search is formed in accordance with conventional methods.

[0041] The speech decoder as shown in Figure 13 operates according to conventional methods except when signal modification is enabled. Signal modification disabled and enabled operation differs essentially only in the way the adaptive codebook excitation signal $u_b(t)$ is formed. In both operational modes, the decoder decodes the received parameters from their binary representation. Typically the received parameters include excitation, gain, delay and LP parameters. The decoded excitation parameters are used in module 701 to form the fixed-codebook excitation signal $u_c(t)$ for every subframe. This signal is supplied through an amplifier 702 to an adder 703. Similarly, the adaptive codebook excitation signal $u_b(t)$ of the current subframe is supplied to the adder 703 through an amplifier 704. In the adder 703, the gain-scaled adaptive and fixed-codebook excitation signals $u_b(t)$ and $u_c(t)$ are summed together to form a total excitation signal $u(t)$ for the current subframe. This excitation signal $u(t)$ is processed through the LP synthesis filter $1/A(z)$ 708, that uses LP parameters interpolated in module 707 for the current subframe, to produce the synthesized speech signal $\hat{s}(t)$.

[0042] When signal modification is enabled, the speech decoder recovers the delay contour $d(t)$ in module 705 using the received delay parameter d_n and its previous received value d_{n-1} as in the encoder. This delay contour $d(t)$ defines a long term prediction delay parameter for every time instant of the current frame. The adaptive codebook excitation $u_b(t) = u(t - d(t))$ is formed from the past excitation for the current subframe as in the encoder using the delay contour $d(t)$.

[0043] The remaining description discloses the detailed operation of the signal modification procedure 603 as well as its use as a part of the mode determination mechanism.

Search of Pitch Pulses and Pitch Cycle Segments

[0044] The signal modification method operates pitch and frame synchronously, shifting each detected pitch cycle segment individually but constraining the shift at frame boundaries. This requires means for locating pitch pulses and corresponding pitch cycle segments for the current frame. In the illustrative embodiment of the signal modification method, pitch cycle segments are determined based on detected pitch pulses that are searched according to Figure 5.

[0045] Pitch pulse search can operate on the residual signal $r(t)$, the weighted speech signal $w(t)$ and/or the weighted synthesized speech signal $\hat{w}(t)$. The residual signal $r(t)$ is obtained by filtering the speech signal $s(t)$ with the LP filter $A(z)$, which has been interpolated for the subframes. In the illustrative embodiment, the order of the LP filter $A(z)$ is 16. The weighted speech signal $w(t)$ is obtained by processing the speech signal $s(t)$ through the weighting filter

$$W(z) = \frac{A(z/\gamma_1)}{1 - \gamma_2 z^{-1}}, \quad (1)$$

where the coefficients $\gamma_1 = 0.92$ and $\gamma_2 = 0.68$. The weighted speech signal $w(t)$ is often utilized in open-loop pitch estimation (module 602) since the weighting filter defined by Equation (1) attenuates the formant structure in the speech signal $s(t)$, and preserves the periodicity also on sinusoidal signal segments. That facilitates pitch pulse search because possible signal periodicity becomes clearly apparent in weighted signals. It should be noted that the weighted speech signal $w(t)$ is needed also for the look ahead in order to search the last pitch pulse in the current frame. This can be done by using the weighting filter of Equation (1) formed in the last subframe of the current frame over the look ahead portion.

[0046] The pitch pulse search procedure of Figure 5 starts in block 301 by locating the last pitch pulse of the previous

frame from the residual signal $r(t)$: A pitch pulse typically stands out clearly as the maximum absolute value of the low-pass filtered residual signal in a pitch cycle having a length of approximately $p(t_{n-1})$. A normalized Hamming window $H_5(z) = (0.08z^{-2} + 0.54z^{-1} + 1 + 0.54z + 0.08z^2)/2.24$ having a length of five (5) samples is used for the low-pass filtering in order to facilitate the locating of the last pitch pulse of the previous frame. This pitch pulse position is denoted by T_0 .

The illustrative embodiment of the signal modification method according to the invention does not require an accurate position for this pitch pulse, but rather a rough location estimate of the high-energy segment in the pitch cycle.

[0047] After locating the last pitch pulse at T_0 in the previous frame, a pitch pulse prototype of length $2l + 1$ samples is extracted in block 302 of Figure 5 around this rough position estimate as, for example:

$$m_n(k) = \hat{w}(T_0 - l + k) \quad \text{for } k = 0, 1, \dots, 2l. \quad (2)$$

This pitch pulse prototype is subsequently used in locating pitch pulses in the current frame.

[0048] The synthesized weighted speech signal $\hat{w}(t)$ (or the weighted speech signal $w(t)$) can be used for the pulse prototype instead of the residual signal $r(t)$. This facilitates pitch pulse search, because the periodic structure of the signal is better preserved in the weighted speech signal. The synthesized weighted speech signal $\hat{w}(t)$ is obtained by filtering the synthesized speech signal $\hat{s}(t)$ of the last subframe of the previous frame by the weighting filter $W(z)$ of Equation (1). If the pitch pulse prototype extends over the end of the previously synthesized frame, the weighted speech signal $w(t)$ of the current frame is used for this exceeding portion. The pitch pulse prototype has a high correlation with the pitch pulses of the weighted speech signal $w(t)$ if the previous synthesized speech frame contains already a well-developed pitch cycle. Thus the use of the synthesized speech in extracting the prototype provides additional information for monitoring the performance of coding and selecting an appropriate coding mode in the current frame as will be explained in more detail in the following description.

[0049] Selecting $l = 10$ samples provides a good compromise between the complexity and performance in the pitch pulse search. The value of l can also be determined proportionally to the open-loop pitch estimate.

[0050] Given the position T_0 of the last pulse in the previous frame, the first pitch pulse of the current frame can be predicted to occur approximately at instant $T_0 + p(T_0)$. Here $p(t)$ denotes the interpolated open-loop pitch estimate at instant (position) t . This prediction is performed in block 303.

[0051] In block 305, the predicted pitch pulse position $T_0 + p(T_0)$ is refined as

$$T_1 = T_0 + p(T_0) + \arg \max C(j), \quad (3)$$

where the weighted speech signal $w(t)$ in the neighborhood of the predicted position is correlated with the pulse prototype:

$$C(j) = \gamma(j) \sum_{k=0}^{2l} m_n(k) w(T_0 + p(T_0) + j - l + k), \quad j \in [-j_{\max}, j_{\max}]. \quad (4)$$

[0052] Thus the refinement is the argument j , limited into $[-j_{\max}, j_{\max}]$, that maximizes the weighted correlation $C(j)$ between the pulse prototype and one of the above mentioned residual signal, weighted speech signal or weighted synthesized speech signal. According to an illustrative example, the limit j_{\max} is proportional to the open-loop pitch estimate as $\min\{20, \langle p(0)/4 \rangle\}$, where the operator $\langle \cdot \rangle$ denotes rounding to the nearest integer. The weighting function

$$\gamma(j) = 1 - |j| / p(T_0 + p(T_0)) \quad (5)$$

in Equation (4) favors the pulse position predicted using the open-loop pitch estimate, since $\gamma(j)$ attains its maximum value 1 at $j = 0$. The denominator $p(T_0 + p(T_0))$ in Equation (5) is the open-loop pitch estimate for the predicted pitch pulse position.

[0053] After the first pitch pulse position T_1 has been found using Equation (3), the next pitch pulse can be predicted to be at instant $T_2 = T_1 + p(T_1)$ and refined as described above. This pitch pulse search comprising the prediction 303 and refinement 305 is repeated until either the prediction or refinement procedure yields a pitch pulse position outside the current frame. These conditions are checked in logic block 304 for the prediction of the position of the next pitch

pulse (block 303) and in logic block 306 for the refinement of this position of the pitch pulse (block 305). It should be noted that the logic block 304 terminates the search only if a predicted pulse position is so far in the subsequent frame that the refinement step cannot bring it back to the current frame. This procedure yields c pitch pulse positions inside the current frame, denoted by T_1, T_2, \dots, T_c .

[0054] According to an illustrative example, pitch pulses are located in the integer resolution except the last pitch pulse of the frame denoted by T_c . Since the exact distance between the last pulses of two successive frames is needed to determine the delay parameter to be transmitted, the last pulse is located using a fractional resolution of 1/4 sample in Equation (4) for j . The fractional resolution is obtained by upsampling $w(t)$ in the neighborhood of the last predicted pitch pulse before evaluating the correlation of Equation (4). According to an illustrative example, Hamming-windowed sinc interpolation of length 33 is used for upsampling. The fractional resolution of the last pitch pulse position helps to maintain the good performance of long term prediction despite the time synchrony constrain set to the frame end. This is obtained with a cost of the additional bit rate needed for transmitting the delay parameter in a higher accuracy.

[0055] After completing pitch cycle segmentation in the current frame, an optimal shift for each segment is determined. This operation is done using the weighted speech signal $w(t)$ as will be explained in the following description. For reducing the distortion caused by warping, the shifts of individual pitch cycle segments are implemented using the LP residual signal $r(t)$. Since shifting distorts the signal particularly around segment boundaries, it is essential to place the boundaries in low power sections of the residual signal $r(t)$. In an illustrative example, the segment boundaries are placed approximately in the middle of two consecutive pitch pulses, but constrained inside the current frame. Segment boundaries are always selected inside the current frame such that each segment contains exactly one pitch pulse. Segments with more than one pitch pulse or "empty" segments without any pitch pulses hamper subsequent correlation-based matching with the target signal and should be prevented in pitch cycle segmentation. The s^{th} extracted segment of l_s samples is denoted as $w_s(k)$ for $k = 0, 1, \dots, l_s - 1$. The starting instant of this segment is t_s , selected such that $w_s(0) = w(t_s)$. The number of segments in the present frame is denoted by c .

[0056] While selecting the segment boundary between two successive pitch pulses T_s and T_{s+1} inside the current frame, the following procedure is used. First the central instant between two pulses is computed as $\Lambda = (T_s + T_{s+1})/2$. The candidate positions for the segment boundary are located in the region $[\Lambda - \varepsilon_{\max}, \Lambda + \varepsilon_{\max}]$, where ε_{\max} corresponds to five samples. The energy of each candidate boundary position is computed as

$$Q(\varepsilon') = r^2(\Lambda + \varepsilon' - 1) + r^2(\Lambda + \varepsilon'), \quad \varepsilon' \in [-\varepsilon_{\max}, \varepsilon_{\max}]. \quad (6)$$

The position giving the smallest energy is selected because this choice typically results in the smallest distortion in the modified speech signal. The instant that minimizes Equation (6) is denoted as ε . The starting instant of the new segment is selected as $t_s = \Lambda + \varepsilon$. This defines also the length of the previous segment, since the previous segment ends at instant $\Lambda + \varepsilon - 1$.

[0057] Figure 6 shows an illustrative example of pitch cycle segmentation. Note particularly the first and the last segment $w_1(k)$ and $w_4(k)$, respectively, extracted such that no empty segments result and the frame boundaries are not exceeded.

Determination of the Delay Parameter

[0058] Generally the main advantage of signal modification is that only one delay parameter per frame has to be coded and transmitted to the decoder (not shown). However, special attention has to be paid to the determination of this single parameter. The delay parameter not only defines together with its previous value the evolution of the pitch cycle length over the frame, but also affects time asynchrony in the resulting modified signal.

[0059] In the methods described in [1, 4-7]

[1] W.B. Kleijn, P. Kroon, and D. Nahumi, "The RCELP speech-coding algorithm," European Transactions on Telecommunications, Vol. 4, No. 5, pp. 573-582, 1994.

[4] US Patent 5,704,003, "RCELP coder," Lucent Technologies Inc., (W.B. Kleijn and D. Nahumi), Filing Date 19 Sep. 1995.

[5] European Patent Application 0 602 826 A2, "Time shifting for analysis-by-synthesis coding," AT&T Corp., (B. Kleijn), Filing Date 1 Dec. 1993.

[6] Patent Application WO 00/11653, "Speech encoder with continuous warping combined with long term prediction,"

Conexant Systems Inc., (Y. Gao), Filing Date 24 Aug. 1999.

[7] Patent Application WO 00/11654, "Speech encoder adaptively applying pitch preprocessing with continuous warping," Conexant Systems Inc., (H. Su and Y. Gao), Filing Date 24 Aug. 1999.

no time synchrony is required at frame boundaries, and thus the delay parameter to be transmitted can be determined straightforwardly using an open-loop pitch estimate. This selection usually results in a time asynchrony at the frame boundary, and translates to an accumulating time shift in the subsequent frame because the signal continuity has to be preserved. Although human hearing is insensitive to changes in the time scale of the synthesized speech signal, increasing time asynchrony complicates the encoder implementation. Indeed, long signal buffers are required to accommodate the signals whose time scale may have been expanded, and a control logic has to be implemented for limiting the accumulated shift during encoding. Also, time asynchrony of several samples typical in RCELP coding may cause mismatch between the LP parameters and the modified residual signal. This mismatch may result in perceptual artifacts to the modified speech signal that is synthesized by LP filtering the modified residual signal.

[0060] On the contrary, the illustrative embodiment of the signal modification method according to the present invention preserves the time synchrony at frame boundaries. Thus, a strictly constrained shift occurs at the frame ends and every new frame starts in perfect time match with the original speech frame.

[0061] To ensure time synchrony at the frame end, the delay contour $d(t)$ maps, with the long term prediction, the last pitch pulse at the end of the previous synthesized speech frame to the pitch pulses of the current frame. The delay contour defines an interpolated long-term prediction delay parameter over the current n^{th} frame for every sample from instant $t_{n-1} + 1$ through t_n . Only the delay parameter $d_n = d(t_n)$ at the frame end is transmitted to the decoder implying that $d(t)$ must have a form fully specified by the transmitted values. The long-term prediction delay parameter has to be selected such that the resulting delay contour fulfils the pulse mapping. In a mathematical form this mapping can be presented as follows: Let κ_c be a temporary time variable and T_o and T_c the last pitch pulse positions in the previous and current frames, respectively. Now, the delay parameter d_n has to be selected such that, after executing the pseudo-code presented in Table 1, the variable κ_c has a value very close to T_o minimizing the error $|\kappa_c - T_o|$. The pseudo-code starts from the value $\kappa_0 = T_c$ and iterates backwards c times by updating $\kappa_i := \kappa_{i-1} - d(\kappa_{i-1})$. If κ_c then equals to T_o , long term prediction can be utilized with maximum efficiency without time asynchrony at the frame end.

Table 1. Loop for searching the optimal delay parameter.

```

% initialization
 $\kappa_0 := T_c$ ;

% loop
for  $i = 1$  to  $c$ 
 $\kappa_i := \kappa_{i-1} - d(\kappa_{i-1})$ ;
end;
    
```

[0062] An example of the operation of the delay selection loop in the case $c = 3$ is illustrated in Figure 7. The loop starts from the value $\kappa_0 = T_c$ and takes the first iteration backwards as $\kappa_1 = \kappa_0 - d(\kappa_0)$. Iterations are continued twice more resulting in $\kappa_2 = \kappa_1 - d(\kappa_1)$ and $\kappa_3 = \kappa_2 - d(\kappa_2)$. The final value κ_3 is then compared against T_o in terms of the error $e_n = |\kappa_3 - T_o|$. The resulting error is a function of the delay contour that is adjusted in the delay selection algorithm as will be taught later in this specification.

[0063] Signal modification methods [1, 4, 6, 7] such as described in the following documents:

[1] W.B. Kleijn, P. Kroon, and D. Nahumi, "The RCELP speech-coding algorithm," European Transactions on Telecommunications, Vol. 4, No. 5, pp. 573-582, 1994.

[4] US Patent 5,704,003, "RCELP coder," Lucent Technologies Inc., (W.B. Kleijn and D. Nahumi), Filing Date 19 Sep. 1995.

[6] Patent Application WO 00/11653, "Speech encoder with continuous warping combined with long term prediction," Conexant Systems Inc., (Y. Gao), Filing Date 24 Aug. 1999.

[7] Patent Application WO 00/11654, "Speech encoder adaptively applying pitch preprocessing with continuous warping," Conexant Systems Inc., (H. Su and Y. Gao), Filing Date 24 Aug. 1999.

interpolate the delay parameters linearly over the frame between d_{n-1} and d_n . However, when time synchrony is required at the frame end, linear interpolation tends to result in an oscillating delay contour. Thus pitch cycles in the modified speech signal contract and expand periodically causing easily annoying artifacts. The evolution and amplitude of the oscillations are related to the last pitch position. The further the last pitch pulse is from the frame end in relation to the pitch period, the more likely the oscillations are amplified. Since the time synchrony at the frame end is an essential requirement of the illustrative embodiment of the signal modification method according to the present invention, linear interpolation familiar from the prior methods cannot be used without degrading the speech quality. Instead, the illustrative embodiment of the signal modification method according to the present invention discloses a piecewise linear delay contour

$$d(t) = \begin{cases} (1 - \alpha(t))d_{n-1} + \alpha(t)d_n & t_{n-1} < t < t_{n-1} + \sigma_n \\ d_n & t_{n-1} + \sigma_n \leq t \leq t_n \end{cases} \quad (7)$$

where

$$\alpha(t) = (t - t_{n-1}) / \sigma_n \quad (8)$$

Oscillations are significantly reduced by using this delay contour. Here t_n and t_{n-1} are the end instants of the current and previous frames, respectively, and d_n and d_{n-1} are the corresponding delay parameter values. Note that $t_{n-1} + \sigma_n$ is the instant after which the delay contour remains constant.

[0064] In an illustrative example, the parameter σ_n varies as a function of d_{n-1} as

$$\sigma_n = \begin{cases} 172 \text{ samples, } & d_{n-1} \leq 90 \text{ samples} \\ 128 \text{ samples, } & d_{n-1} > 90 \text{ samples} \end{cases} \quad (9)$$

and the frame length N is 256 samples. To avoid oscillations, it is beneficial to decrease the value of σ_n as the length of the pitch cycle increases. On the other hand, to avoid rapid changes in the delay contour $d(t)$ in the beginning of the frame as $t_{n-1} < t < t_{n-1} + \sigma_n$, the parameter σ_n has to be always at least a half of the frame length. Rapid changes in $d(t)$ degrade easily the quality of the modified speech signal.

[0065] Note that depending on the coding mode of the previous frame, d_{n-1} can be either the delay value at the frame end (signal modification enabled) or the delay value of the last subframe (signal modification disabled). Since the past value d_{n-1} of the delay parameter is known at the decoder, the delay contour is unambiguously defined by d_n and the decoder is able to form the delay contour using Equation (7).

[0066] The only parameter which can be varied while searching the optimal delay contour is d_n , the delay parameter value at the end of the frame constrained into $[34, 231]$. There is no simple explicit method for solving the optimal d_n in a general case. Instead, several values have to be tested to find the best solution. However, the search is straightforward. The value of d_n can be first predicted as

$$d_n^{(0)} = 2 \frac{T_c - T_0}{c} - d_{n-1} \quad (10)$$

In the illustrative embodiment, the search is done in three phases by increasing the resolution and focusing the search range to be examined inside $[34, 231]$ in every phase. The delay parameters giving the smallest error $e_n = |k_c - T_0|$ in the procedure of Table 1 in these three phases are denoted by $d_n^{(1)}$, $d_n^{(2)}$, and $d_n = d_n^{(3)}$, respectively.

In the first phase, the search is done around the value $d_n^{(0)}$ predicted using Equation (10) with a resolution of four samples in the range $[d_n^{(0)} - 11, d_n^{(0)} + 12]$ when $d_n^{(0)} < 60$, in the range $[d_n^{(0)} - 15, d_n^{(0)} + 16]$ otherwise.

The second phase constrains the range into $[d_n^{(1)} - 3, d_n^{(1)} + 3]$ and uses the integer resolution. The last, third phase examines the range $[d_n^{(2)} - 3/4, d_n^{(2)} + 3/4]$ with a resolution of 1/4 sample for $d_n^{(2)} < 92^{1/2}$. Above that range $[d_n^{(2)} - 1/2, d_n^{(2)} + 1/2]$ and a resolution of 1/2 sample is used. This third phase yields the optimal delay parameter d_n to be transmitted to the decoder. This procedure is a compromise between the search accuracy and complexity. Of course, those of ordinary skill in the art can readily implement the search of the delay parameter under the time synchrony constrains using alternative means without departing from the nature of the present invention.

[0067] The delay parameter $d_n \in [34, 231]$ can be coded using nine bits per frame using a resolution of 1/4 sample for $d_n < 92^{1/2}$ and 1/2 sample for $d_n > 92^{1/2}$.

[0068] Figure 8 illustrates delay interpolation when $d_{n-1} = 50$, $d_n = 53$, $\sigma_n = 172$, and the frame length $N = 256$. The interpolation method used in the illustrative embodiment of the signal modification method is shown in thick line whereas the linear interpolation corresponding to prior methods is shown in thin line. Both interpolated contours perform approximately in a similar manner in the delay selection loop of Table 1, but the disclosed piecewise linear interpolation results in a smaller absolute change $|d_{n-1} - d_n|$. This feature reduces potential oscillations in the delay contour $d(t)$ and annoying artifacts in the modified speech signal whose pitch will follow this delay contour.

[0069] To further clarify the performance of the piecewise linear interpolation method, Figure 9 shows an example on the resulting delay contour $d(t)$ over ten frames with thick line. The corresponding delay contour $d(t)$ obtained with conventional linear interpolation is indicated with thin line. The example has been composed using an artificial speech signal having a constant delay parameter of 52 samples as an input of the speech modification procedure. A delay parameter $d_0 = 54$ samples was intentionally used as an initial value for the first frame to illustrate the effect of pitch estimation errors typical in speech coding. Then, the delay parameters d_n both for the linear interpolation and the herein disclosed piecewise linear interpolation method were searched using the procedure of Table 1. All the parameters needed were selected in accordance with the illustrative embodiment of the signal modification method according to the present invention. The resulting delay contours $d(t)$ show that piecewise linear interpolation yields a rapidly converging delay contour $d(t)$ whereas the conventional linear interpolation cannot reach the correct value within the ten frame period. These prolonged oscillations in the delay contour $d(t)$ often cause annoying artifacts to the modified speech signal degrading the overall perceptual quality.

Modification of the Signal

[0070] After the delay parameter d_n and the pitch cycle segmentation have been determined, the signal modification procedure itself can be initiated. In the illustrative embodiment of the signal modification method, the speech signal is modified by shifting individual pitch cycle segments one by one adjusting them to the delay contour $d(t)$. A segment shift is determined by correlating the segment in the weighted speech domain with the target signal. The target signal is composed using the synthesized weighted speech signal $\hat{w}(t)$ of the previous frame and the preceding, already shifted segments in the current frame. The actual shift is done on the residual signal $r(t)$.

[0071] Signal modification has to be done carefully to both maximize the performance of long term prediction and simultaneously to preserve the perceptual quality of the modified speech signal. The required time synchrony at frame boundaries has to be taken into account also during modification.

[0072] A block diagram of the illustrative embodiment of the signal modification method is shown in Figure 10. Modification starts by extracting a new segment $w_s(k)$ of l_s samples from the weighted speech signal $w(t)$ in block 401. This segment is defined by the segment length l_s and starting instant t_s giving $w_s(k) = w(t_s + k)$ for $k = 0, 1, \dots, l_s - 1$. The segmentation procedure is carried out in accordance with the teachings of the foregoing description.

[0073] If no more segments can be selected or extracted (block 402), the signal modification operation is completed (block 403). Otherwise, the signal modification operation continues with block 404.

[0074] For finding the optimal shift of the current segment $w_s(k)$, a target signal $\tilde{w}(t)$ is created in block 405. For the first segment $w_s(k)$ in the current frame, this target signal is obtained by the recursion

$$\begin{aligned} \tilde{w}(t) &= \hat{w}(t), & t \leq t_{n-1} \\ \tilde{w}(t) &= \tilde{w}(t - d(t)), & t_{n-1} < t \leq t_{n-1} + l_1 + \delta_1. \end{aligned} \quad (11)$$

Here $\hat{w}(t)$ is the weighted synthesized speech signal available in the previous frame for $t \leq t_{n-1}$. The parameter δ_1 is the maximum shift allowed for the first segment of length l_1 . Equation (11) can be interpreted as simulation of long term prediction using the delay contour over the signal portion in which the current shifted segment may potentially be situated.

The computation of the target signal for the subsequent segments follows the same principle and will be presented later in this section.

[0075] The search procedure for finding the optimal shift of the current segment can be initiated after forming the target signal. This procedure is based on the correlation $c_s(\delta')$ computed in block 404 between the segment $w_s(k)$ that starts at instant t_s and the target signal $\tilde{w}(t)$ as

$$c_s(\delta') = \sum_{k=0}^{l_s-1} w_s(k) \tilde{w}(k + t_s + \delta'), \quad \delta' \in [-\lceil \delta_s \rceil, \lceil \delta_s \rceil], \quad (12)$$

where δ_s determines the maximum shift allowed for the current segment $w_s(k)$ and $\lceil \cdot \rceil$ denotes rounding towards plus infinity. Normalized correlation can be well used instead of Equation (12), although with increased complexity. In the illustrative embodiment, the following values are used for δ_s :

$$\delta_s = \begin{cases} 4 \frac{1}{2} \text{ samples,} & d_n < 90 \text{ samples} \\ 5 \text{ samples,} & d_n \geq 90 \text{ samples} \end{cases} \quad (13)$$

As will be described later in this section, the value of δ_s is more limited for the first and the last segment in the frame.

[0076] Correlation (12) is evaluated with an integer resolution, but higher accuracy improves the performance of long term prediction. For keeping the complexity low it is not reasonable to upsample directly the signal $w_s(k)$ or $\tilde{w}(t)$ in Equation (12). Instead, a fractional resolution is obtained in a computationally efficient manner by determining the optimal shift using the upsampled correlation $C_s(\delta')$.

[0077] The shift δ maximizing the correlation $c_s(\delta')$ is searched first in the integer resolution in block 404. Now, in a fractional resolution the maximum value must be located in the open interval $(\delta - 1, \delta + 1)$, and bounded into $[-\delta_s, \delta_s]$. In block 406, the correlation $c_s(\delta')$ is upsampled in this interval to a resolution of 1/8 sample using Hamming-windowed sinc interpolation of a length equal to 65 samples. The shift δ corresponding to the maximum value of the upsampled correlation is then the optimal shift in a fractional resolution. After finding this optimal shift, the weighted speech segment $w_s(k)$ is recalculated in the solved fractional resolution in block 407. That is, the precise new starting instant of the segment is updated as $t_s := t_s - \delta + \delta_f$, where $\delta_f = \lceil \delta \rceil$. Further, the residual segment $r_s(k)$ corresponding to the weighted speech segment $w_s(k)$ in fractional resolution is computed from the residual signal $r(t)$ at this point using again the sinc interpolation as described before (block 407). Since the fractional part of the optimal shift is incorporated into the residual and weighted speech segments, all subsequent computations can be implemented with the upward-rounded shift $\delta_f = \lceil \delta \rceil$.

[0078] Figure 11 illustrates recalculation of the segment $w_s(k)$ in accordance with block 407 of Figure 10. In this illustrative example, the optimal shift is searched with a resolution of 1/8 sample by maximizing the correlation giving the value $\delta = -1\frac{3}{8}$. Thus the integer part δ_f becomes $\lceil -1\frac{3}{8} \rceil = -1$ and the fractional part. Consequently, the starting instant of the segment is updated as $t_s = t_s + 3/8$. In Figure 11, the new samples of $w_s(k)$ are indicated with gray dots.

[0079] If the logic block 106, which will be disclosed later, permits to continue signal modification, the final task is to update the modified residual signal $\tilde{r}(t)$ by copying the current residual signal segment $r_s(k)$ into it (block 411):

$$\tilde{r}(t_s + \delta_f + k) = r_s(k), \quad k = 0, 1, \dots, l_s - 1. \quad (14)$$

Since shifts in successive segments are independent from each others, the segments positioned to $F(t)$ either overlap or have a gap in between them. Straightforward weighted averaging can be used for overlapping segments. Gaps are filled by copying neighboring samples from the adjacent segments. Since the number of overlapping or missing samples is usually small and the segment boundaries occur at low-energy regions of the residual signal, usually no perceptual artifacts are caused. It should be noted that no continuous signal warping as described in [2], [6], [7],

[2] W.B. Kleijn, R.P. Ramachandran, and P. Kroon, "Interpolation of the pitch-predictor parameters in analysis-by-synthesis speech coders," IEEE Transactions on Speech and Audio Processing, Vol. 2, No. 1, pp. 42-54, 1994.

[6] Patent Application WO 00/11653, "Speech encoder with continuous warping combined with long term prediction," Conexant Systems Inc., (Y. Gao), Filing Date 24 Aug. 1999.

[7] Patent Application WO 00/11654, "Speech encoder adaptively applying pitch preprocessing with continuous warping," Conexant Systems Inc., (H. Su and Y. Gao), Filing Date 24 Aug. 1999.

is employed, but modification is done discontinuously by shifting pitch cycle segments in order to reduce the complexity.
 5 **[0080]** Processing of the subsequent pitch cycle segments follows the above-disclosed procedure, except the target signal $\tilde{w}(t)$ in block 405 is formed differently than for the first segment. The samples of $\tilde{w}(t)$ are first replaced with the modified weighted speech samples as

$$10 \quad \tilde{w}(t_s + \delta_s + k) = w_s(k), \quad k = 0, 1, \dots, l_s - 1. \quad (15)$$

This procedure is illustrated in Figure 11. Then the samples following the updated segment are also updated,

$$15 \quad \tilde{w}(k) = \tilde{w}(k - d(k)), \quad k = t_s + \delta_s + l_s, \dots, t_s + \delta_s + l_s + l_{s+1} + \delta_{s+1} - 2. \quad (16)$$

The update of target signal $\tilde{w}(t)$ ensures higher correlation between successive pitch cycle segments in the modified speech signal considering the delay contour $d(t)$ and thus more accurate long term prediction. While processing the last segment of the frame, the target signal $\tilde{w}(t)$ does not need to be updated.

20 **[0081]** The shifts of the first and the last segments in the frame are special cases which have to be performed particularly carefully. Before shifting the first segment, it should be ensured that no high power regions exist in the residual signal $r(t)$ close to the frame boundary t_{n-1} , because shifting such a segment may cause artifacts. The high power region is
 25 searched by squaring the residual signal $r(t)$ as

$$30 \quad E_0(k) = r^2(k), \quad k \in [t_{n-1} - \zeta_0, t_{n-1} + \zeta_0], \quad (17)$$

where $\zeta_0 = \langle p(t_{n-1})/2 \rangle$. If the maximum of $E_0(k)$ is detected close to the frame boundary in the range $[t_{n-1} - 2, t_{n-1} + 2]$, the allowed shift is limited to 1/4 samples. If the proposed shift $|\delta|$ for the first segment is smaller than this limit, the signal modification procedure is enabled in the current frame, but the first segment is kept intact.

35 **[0082]** The last segment in the frame is processed in a similar manner. As was described in the foregoing description, the delay contour $d(t)$ is selected such that in principle no shifts are required for the last segment. However, because the target signal is repeatedly updated during signal modification considering correlations between successive segments in Equations (16) and (17), it is possible the last segment has to be shifted slightly. In the illustrative embodiment, this shift is always constrained to be smaller than 3/2 samples. If there is a high power region at the frame end, no shift is allowed. This condition is verified by using the squared residual signal

$$40 \quad E_1(k) = r^2(k), \quad k \in [t_n - \zeta_1 + 1, t_n + 1], \quad (18)$$

45 where $\zeta_1 = p(t_n)$. If the maximum of $E_1(k)$ is attained for k larger than or equal to $t_n - 4$, no shift is allowed for the last segment. Similarly as for the first segment, when the proposed shift $|\delta| < 1/4$, the present frame is still accepted for modification, but the last segment is kept intact.

50 **[0083]** It should be noted that, contrary to the known signal modification methods, the shift does not translate to the next frame, and every new frame starts perfectly synchronized with the original input signal. As another fundamental difference particularly to RCELP coding, the illustrative embodiment of signal modification method processes a complete speech frame before the subframes are coded. Admittedly, subframe-wise modification enables to compose the target signal for every subframe using the previously coded subframe potential improving the performance. This approach cannot be used in the context of the illustrative embodiment of the signal modification method since the allowed time asynchrony at the frame end is strictly constrained. Nevertheless, the update of the target signal with Equations (15) and (16) gives practically speaking equal performance with the subframe-wise processing, because modification is
 55 enabled only on smoothly evolving voiced frames.

Mode Determination Logic Incorporated into the Signal Modification Procedure

[0084] The illustrative embodiment of signal modification method according to the present invention incorporates an efficient classification and mode determination mechanism as depicted in Figure 2. Every operation performed in blocks 101, 103 and 105 yields several indicators quantifying the attainable performance of long term prediction in the current frame. If any of these indicators is outside its allowed limits, the signal modification procedure is terminated by one of the logic blocks 102, 104, or 106. In this case, the original signal is preserved intact.

[0085] The pitch pulse search procedure 101 produces several indicators on the periodicity of the present frame. Hence the logic block 102 analyzing these indicators is the most important component of the classification logic. The logic block 102 compares the difference between the detected pitch pulse positions and the interpolated open-loop pitch estimate using the condition

$$|T_k - T_{k-1} - p(T_k)| < 0.2 p(T_k), \quad k = 1, 2, \dots, c, \quad (19)$$

and terminates the signal modification procedure if this condition is not met.

[0086] The selection of the delay contour $d(t)$ in block 103 gives also additional information on the evolution of the pitch cycles and the periodicity of the current speech frame. This information is examined in the logic block 104. The signal modification procedure is continued from this block 104 only if the condition $|d_n - d_{n-1}| < 0.2 d_n$ is fulfilled. This condition means that only a small delay change is tolerated for classifying the current frame as purely voiced frame. The logic block 104 also evaluates the success of the delay selection loop of Table 1 by examining the difference $|k_c - T_0|$ for the selected delay parameter value d_n . If this difference is greater than one, sample, the signal modification procedure is terminated.

[0087] For guaranteeing a good quality for the modified speech signal, it is advantageous to constrain shifts done for successive pitch cycle segments in block 105. This is achieved in the logic block 106 by imposing the criteria

$$|\delta^{(s)} - \delta^{(s-1)}| \leq \begin{cases} 4.0 \text{ samples, } & d_n < 90 \text{ samples} \\ 4.8 \text{ samples, } & d_n \geq 90 \text{ samples} \end{cases} \quad (20)$$

to all segments of the frame. Here $\delta^{(s)}$ and $\delta^{(s-1)}$ are the shifts done for the s^{th} and $(s-1)^{th}$ pitch cycle segments, respectively. If the thresholds are exceeded, the signal modification procedure is interrupted and the original signal is maintained.

[0088] When the frames subjected to signal modification are coded at a low bit rate, it is essential that the shape of pitch cycle segments remains similar over the frame. This allows faithful signal modeling by long term prediction and thus coding at a low bit rate without degrading the subjective quality. The similarity of successive segments can be quantified simply by the normalized correlation

$$g_s = \frac{\sum_{k=0}^{l_s-1} w_s(k) \tilde{w}(k + t_s + \delta_s)}{\sqrt{\sum_{k=0}^{l_s-1} w_s^2(k) \sum_{k=0}^{l_s-1} \tilde{w}^2(k + t_s + \delta_s)}} \quad (21)$$

between the current segment and the target signal at the optimal shift after the update of $w_s(k)$ in block 407 of Figure 10. The normalized correlation g_s is also referred to as pitch gain.

[0089] Shifting of the pitch cycle segments in block 105 maximizing their correlation with the target signal enhances the periodicity and yields a high pitch prediction gain if the signal modification is useful in the current frame. The success of the procedure is examined in the logic block 106 using the criteria

$$g_s \geq 0.84.$$

If this condition is not fulfilled for all segments, the signal modification procedure is terminated (block 409) and the original signal is kept intact. When this condition is met (block 106), the signal modification continues in block 411. The pitch gain g_s is computed in block 408 between the recalculated segment $w_s(k)$ from block 407 and the target signal $\tilde{w}(t)$ from block 405. In general, a slightly lower gain threshold can be allowed on male voices with equal coding performance. The gain thresholds can be changed in different operation modes of the encoder for adjusting the usage percentage of the signal modification mode and thus the resulting average bit rate.

[0090] Mode Determination Logic for a Source-controlled Variable Bit Rate Speech Codec .

[0091] This section discloses the use of the signal modification procedure as a part of the general rate determination mechanism in a source-controlled variable bit rate speech codec. This functionality is immersed into the illustrative embodiment of the signal modification method, since it provides several indicators on signal periodicity and the expected coding performance of long term prediction in the present frame. These indicators include the evolution of pitch period, the fitness of the selected delay contour for describing this evolution, and the pitch prediction gain attainable with signal modification. If the logic blocks 102, 104 and 106 shown in Figure 2 enable signal modification, long term prediction is able to model the modified speech frame efficiently facilitating its coding at a low bit rate without degrading subjective quality. In this case, the adaptive codebook excitation has a dominant contribution in describing the excitation signal, and thus the bit rate allocated for the fixed-codebook excitation can be reduced. When a logic block 102, 104 or 106 disables signal modification, the frame is likely to contain a non-stationary speech segment such as a voiced onset or rapidly evolving voiced speech signal. These frames typically require a high bit rate for sustaining good subjective quality.

[0092] Figure 12 depicts the signal modification procedure 603 as a part of the rate determination logic that controls four coding modes. In this illustrative embodiment, the mode set comprises a dedicated mode for non-active speech frames (block 508), unvoiced speech frames (block 507), stable voiced frames (block 506), and other types of frames (block 505). It should be noted that all these modes except the mode for stable voiced frames 506 are implemented in accordance with techniques well known to those of ordinary skill in the art.

[0093] The rate determination logic is based on signal classification done in three steps in logic blocks 501, 502, and 504, from which the operation of blocks 501 and 502 is well known to those of ordinary skill in the art.

[0094] First, a voice activity detector (VAD) 501 discriminates between active and inactive speech frames. If an inactive speech frame is detected, the speech signal is processed according to mode 508.

[0095] If an active speech frame is detected in block 501, the frame is subjected to a second classifier 502 dedicated to making a voicing decision. If the classifier 502 rates the current frame as unvoiced speech signal, the classification chain ends and the speech signal is processed in accordance with mode 507. Otherwise, the speech frame is passed through to the signal modification module 603.

[0096] The signal modification module then provides itself a decision on enabling or disabling the signal modification of the current frame in a logic block 504. This decision is in practice made as an integral part of the signal modification procedure in the logic blocks 102, 104 and 106 as explained earlier with reference to Figure 2. When signal modification is enabled, the frame is deemed as a stable voiced, or purely voiced speech segment.

[0097] When the rate determination mechanism selects mode 506, the signal modification mode is enabled and the speech frame is encoded in accordance with the teachings of the previous sections. Table 2 discloses the bit allocation used in the illustrative embodiment for the mode 506. Since the frames to be coded in this mode are characteristically very periodic, a substantially lower bit rate suffices for sustaining good subjective quality compared for instance to transition frames. Signal modification allows also efficient coding of the delay information using only nine bits per 20-ms frame saving a considerable proportion of the bit budget for other parameters. Good performance of long term prediction allows to use only 13 bits per 5-ms subframe for the fixed-codebook excitation without sacrificing the subjective speech quality. The fixed-codebook comprises one track with two pulses, both having 64 possible positions.

Table 2. Bit allocation in the voiced 6.2-kbps mode for a 20-ms frame comprising four subframes.

Parameter	Bits/Frame
LP Parameters	34
Pitch Delay	9
Pitch Filtering	4 = 1 + 1 + 1 + 1
Gains	24 = 6+ 6+ 6+ 6
Algebraic Codebook	52 =13+13+13+13
Mode Bit	1
Total	124 bits = 6.2 kbps

Table 3. Bit allocation in the 12.65-kbps mode in accordance with the AMR-WB standard.

Parameter	Bits/Frame
LP Parameters	46
Pitch Delay	$30 = 9 + 6 + 9 + 6$
Pitch Filtering	$4 = 1 + 1 + 1 + 1$
Gains	$24 = 7 + 7 + 7 + 7$
Algebraic Codebook	$144 = 36 + 36 + 36 + 36$
Mode Bit	1
Total	253 bits = 12,65 kbps

[0098] The other coding modes 505, 507 and 508 are implemented following known techniques. Signal modification is disabled in all these modes. Table 3 shows the bit allocation of the mode 505 adopted from the AMR-WB standard.

[0099] The technical specifications [11] and [12] related to the AMR-WB standard are enclosed here as references on the comfort noise and VAD functionalities in 501 and 508, respectively.

[11] 3GPP TS 26.192, "AMR Wideband Speech Codec: Comfort Noise Aspects," 3GPP Technical Specification.

[12] 3GPP TS 26.193, "AMR Wideband Speech Codec: Voice Activity Detector (VAD)," 3GPP Technical Specification.

[0100] In summary, the present specification has described a frame synchronous signal modification method for purely voiced speech frames, a classification mechanism for detecting frames to be modified, and to use these methods in a source-controlled CELP speech codec in order to enable high-quality coding at a low bit rate.

[0101] The signal modification method incorporates a classification mechanism for determining the frames to be modified. This differs from prior signal modification and preprocessing means in operation and in the properties of the modified signal. The classification functionality embedded into the signal modification procedure is used as a part of the rate determination mechanism in a source-controlled CELP speech codec.

[0102] Signal modification is done pitch and frame synchronously, that is, adapting one pitch cycle segment at a time in the current frame such that a subsequent speech frame starts in perfect time alignment with the original signal. The pitch cycle segments are limited by frame boundaries. This feature prevents time shift translation over frame boundaries simplifying encoder implementation and reducing a risk of artifacts in the modified speech signal. Since time shift does not accumulate over successive frames, the signal modification method disclosed does not need long buffers for accommodating expanded signals nor a complicated logic for controlling the accumulated time shift. In source-controlled speech coding, it simplifies multi-mode operation between signal modification enabled and disabled modes, since every new frame starts in time alignment with the original signal.

[0103] Of course, many other modifications and variations are possible. In view of the above detailed illustrative description of the present invention and associated drawings, such other modifications and variations will now become apparent to those of ordinary skill in the art. It should also be apparent that such other variations may be effected without departing from the scope of the present invention.

Claims

1. A method of forming a delay contour characterising a long term prediction in a technique using signal modification for digitally encoding a speech signal, the method comprising:

dividing the speech signal into a series of successive frames;
 locating a pitch pulse of the speech signal in a previous frame; and
 locating a corresponding pitch pulse of the speech signal in a current frame;

characterised by forming a delay contour by selecting a long term prediction delay parameter for the current frame by iterating backwards a function of a temporary time variable, from the location of the pitch pulse of the speech signal in the current frame towards the location of the corresponding pitch pulse of the speech signal in the previous frame.

2. A method as claimed in claim 1, comprising:

forming the delay contour as a function of distances of successive pitch pulses between a last pitch pulse of the previous frame and a last pitch pulse of the current frame.

3. A method as claimed in claim 1 or claim 2, further comprising:

fully characterising the delay contour with a long-term-prediction delay parameter of the previous frame and the long-term-prediction delay parameter of the current frame.

4. A method as claimed in claim 3, wherein forming the delay contour comprises:

nonlinearly interpolating the delay contour between the long-term-prediction delay parameter of the previous frame and the long-term-prediction delay parameter of the current frame.

5. A method as claimed in claim 3, wherein forming the delay contour comprises:

determining a piecewise linear delay contour between the long-term-prediction delay parameter of the previous frame and the long-term-prediction delay parameter of the current frame.

6. A method as claimed in any preceding claim, wherein locating a pitch pulse comprises deriving a linear prediction residual signal from the speech signal.

7. A method as claimed in any of claims 1 to 5, wherein locating a pitch pulse comprises deriving a weighted speech signal from the speech signal.

8. A method as claimed in any of claims 1 to 5, wherein locating a pitch pulse comprises deriving a synthesised weighted speech signal from the speech signal

9. A method as claimed in any preceding claim, wherein the backwards iteration comprises searching for a long term prediction delay parameter value in plural phases and beginning with a long term prediction delay parameter value predicted for the end of the current frame, each successive phase having increased resolution and a more focused search range.

10. A method as claimed in claim 9, comprising predicting the long term prediction delay parameter value as being equal to the difference between the long term prediction delay parameter value at the end of the previous frame and twice the difference between the locations of the pitch pulses of the speech signal in the previous and current frames divided by the number of iterations of the function.

11. A method as claimed in any preceding claim, comprising modifying the speech signal by shifting pitch cycle segments one by one to adjust them to the delay contour.

12. A method as claimed in claim 11, comprising determining a segment shift by correlating a segment in the weighted speech domain with a target signal.

13. A method as claimed in claim 12, comprising composing the target signal using the synthesised weighted speech signal of the previous frame and any preceding, shifted segments in the current frame.

14. A device (603) for forming a delay contour characterising a long term prediction in a technique using signal modification for digitally encoding a speech signal, the device comprising:

a divider of the speech signal into a series of successive frames;
 a detector of a location of a pitch pulse of the speech signal in a previous frame; and
 a detector of a location of a corresponding pitch pulse of the speech signal in a current frame,

characterised by a former of a delay contour for selecting a long term prediction delay parameter for the current frame by backwards iteration of a function of a temporary time variable, from the location of the pitch pulse of the speech signal in the current frame towards the location of the corresponding pitch pulse of the speech signal in the

previous frame.

15. A device as claimed in claim 14, wherein the former is:

5 a calculator of the long-term-prediction delay parameter as a function of distances of successive pitch pulses between the last pitch pulse of the previous frame and the last pitch pulse of the current frame.

16. A device as claimed in claim 14 or claim 15, further incorporating:

10 a function fully characterising the delay contour with a long-term- prediction delay parameter of the previous frame and the long-term-prediction delay parameter of the current frame.

17. A device as claimed in claim 16, wherein the former is:

15 a selector of a nonlinearly interpolated delay contour between the long-term-prediction delay parameter of the previous frame and the long-term-prediction delay parameter of the current frame.

18. A device as claimed in claim 16, wherein the former is:

20 a selector of a piecewise linear delay contour determined from the long-term-prediction delay parameter of the previous frame and the long-term-prediction delay parameter of the current frame.

19. A device as claimed in any of claims 14 to 18, wherein the former is a searcher of a long term prediction delay parameter value by backwards iteration in plural phases and beginning with a long term prediction delay parameter value predicted for the end of the current frame, each successive phase having increased resolution and a more focused search range.

25

20. A device as claimed in claim 19, comprising a predictor of the long term prediction delay parameter value as being equal to the difference between the long term prediction delay parameter value at the end of the previous frame and twice the difference between the locations of the pitch pulses of the speech signal in the previous and current frames divided by the number of iterations of the function.

30

21. A device as claimed in any of claims 14 to 20, comprising a modifier of the speech signal by shifting pitch cycle segments one by one to adjust them to the delay contour.

35

22. A device as claimed in claim 21, comprising a determiner of a segment shift by correlating a segment in the weighted speech domain with a target signal.

23. A device as claimed in claim 22, comprising a composer of the target signal using a synthesised weighted speech signal of the previous frame and any preceding, shifted segments in the current frame.

40

Patentansprüche

45 1. Verfahren zum Bilden einer Verzögerungskontur, die eine Langzeitvorhersage in einer Methode charakterisiert, die Signalmodifikation zur digitalen Codierung eines Sprachsignals verwendet, wobei das Verfahren umfasst:

Aufteilen des Sprachsignals in eine Reihe aufeinanderfolgender Rahmen;
 Lokalisieren eines Tonhöhenpulses des Sprachsignals in einem vorhergehenden Rahmen; und
 50 Lokalisieren eines entsprechenden Tonhöhenpulses des Sprachsignals in einem derzeitigen Rahmen;

55

gekennzeichnet durch das Bilden einer Verzögerungskontur, indem ein Langzeitvorhersage-Verzögerungsparameter für den derzeitigen Rahmen gewählt wird, indem eine Funktion einer temporären Zeitvariable rückwärts iteriert wird, von der Stelle des Tonhöhenpulses des Sprachsignals in dem derzeitigen Rahmen in Richtung der Stelle des entsprechenden Tonhöhenpulses des Sprachsignals im vorhergehenden Rahmen.

55

2. Verfahren nach Anspruch 1, umfassend:

Bilden der Verzögerungskontur als eine Funktion von Abständen aufeinanderfolgender Tonhöhenpulse zwischen einem letzten Tonhöhenpuls des vorhergehenden Rahmens und einem letzten Tonhöhenpuls des derzeitigen Rahmens.

- 5 3. Verfahren nach Anspruch 1 oder 2, weiter umfassend:
- vollständiges Charakterisieren der Verzögerungskontur mit einem Langzeitvorhersage-Verzögerungsparameter des vorhergehenden Rahmens und dem Langzeitvorhersage-Verzögerungsparameter des derzeitigen Rahmens.
- 10 4. Verfahren nach Anspruch 3, wobei das Bilden der Verzögerungskontur umfasst:
- nichtlineares Interpolieren der Verzögerungskontur zwischen dem Langzeitvorhersage-Verzögerungsparameter des vorhergehenden Rahmens und dem Langzeitvorhersage-Verzögerungsparameter des derzeitigen Rahmens.
- 15 5. Verfahren nach Anspruch 3, wobei das Bilden der Verzögerungskontur umfasst:
- Bestimmen einer stückweise linearen Verzögerungskontur zwischen dem Langzeitvorhersage-Verzögerungsparameter des vorhergehenden Rahmens und dem Langzeitvorhersage-Verzögerungsparameter des derzeitigen Rahmens.
- 20 6. Verfahren nach einem der vorhergehenden Ansprüche, wobei das Lokalisieren eines Tonhöhenpulses das Ableiten eines Linear-Vorhersage-Restsignals aus dem Sprachsignal umfasst.
- 25 7. Verfahren nach einem der Ansprüche 1 bis 5, wobei das Lokalisieren eines Tonhöhenpulses ein Ableiten eines gewichteten Sprachsignals aus dem Sprachsignal umfasst.
- 30 8. Verfahren nach einem der Ansprüche 1 bis 5, wobei das Lokalisieren eines Tonhöhenpulses ein Ableiten eines synthetisierten gewichteten Sprachsignals aus dem Sprachsignal umfasst.
- 35 9. Verfahren nach einem der vorhergehenden Ansprüche, wobei die Rückwärts-Iteration ein Suchen nach einem Langzeitvorhersage-Verzögerungsparameterwert in mehreren Phasen und ein Beginnen mit einem Langzeitvorhersage-Verzögerungsparameterwert, der für das Ende des derzeitigen Rahmens vorhergesagt wird, umfasst, wobei jede aufeinander folgende Phase eine gesteigerte Auflösung und einen stärker fokussierten Suchbereich aufweist.
- 40 10. Verfahren nach Anspruch 9, umfassend ein Vorhersagen des Langzeitvorhersage-Verzögerungsparameterwerts als gleich der Differenz zwischen dem Langzeitvorhersage-Verzögerungsparameterwerte am Ende des vorhergehenden Rahmens und zweimal der Differenz zwischen den Stellen der Tonhöhenpulse des Sprachsignals in dem vorhergehenden und derzeitigen Rahmen, geteilt durch die Anzahl von Iterationen der Funktion.
- 45 11. Verfahren nach einem der vorhergehenden Ansprüche, umfassend ein Modifizieren des Sprachsignals durch Verschieben von Tonhöhen-Zyklus-Segmenten, eins nach dem anderen, um sie an die Verzögerungskontur anzupassen.
- 50 12. Verfahren nach Anspruch 11, umfassend ein Bestimmen einer Segmentverschiebung durch Korrelieren eines Segments in der gewichteten Sprachdomäne mit einem Zielsignal.
13. Verfahren nach Anspruch 12, umfassend ein Zusammensetzen des Zielsignals unter Verwendung des synthetisierten gewichteten Sprachsignals des vorhergehenden Rahmens und aller vorhergehenden verschobenen Segmente im derzeitigen Rahmen.
- 55 14. Vorrichtung (603) zum Bilden einer Verzögerungskontur, die eine Langzeitvorhersage charakterisiert, in einer Methode, welche Signalmodifikation zur digitalen Codierung eines Sprachsignals verwendet, wobei die Vorrichtung umfasst:
- eine Aufteilungseinrichtung für das Sprachsignal in eine Reihe aufeinander folgender Rahmen;
- einen Detektor für eine Stelle eines Tonhöhenpulses des Sprachsignals in einem vorhergehenden Rahmen; und
- einen Detektor für eine Stelle eines entsprechenden Tonhöhenpulses des Sprachsignals in einem derzeitigen

Rahmen,

gekennzeichnet durch eine Bildungseinrichtung einer Verzögerungskontur zum Wählen eines Langzeitvorhersage-Verzögerungsparameters für den derzeitigen Rahmen **durch** Rückwärts-Iteration einer Funktion einer temporären Zeitvariablen, von der Stelle des Tonhöhenpulses des Sprachsignals in dem derzeitigen Rahmen in Richtung des entsprechenden Tonhöhenpulses des Sprachsignals in dem vorhergehenden Rahmen.

5

15. Vorrichtung nach Anspruch 14, wobei die Bildungseinrichtung eine Berechnungseinrichtung des Langzeitvorhersage-Verzögerungsparameters als eine Funktion von Abständen aufeinander folgender Tonhöhenpulse zwischen dem letzten Tonhöhenpuls des vorhergehenden Rahmens und dem letzten Tonhöhenpuls des derzeitigen Rahmens ist.

10

16. Vorrichtung nach Anspruch 14 oder 15, weiter einschließend:

eine Funktion, die die Verzögerungskontur vollständig mit einem Langzeitvorhersage-Verzögerungsparameter des vorhergehenden Rahmens und dem Langzeitvorhersage-Verzögerungsparameters des derzeitigen Rahmens charakterisiert.

15

17. Vorrichtung nach Anspruch 16, wobei die Bildungseinrichtung ist:

eine Auswahleinrichtung einer nichtlinear interpolierten Verzögerungskontur zwischen dem Langzeitvorhersage-Verzögerungsparameter des vorhergehenden Rahmens und dem Langzeitvorhersage-Verzögerungsparameter des derzeitigen Rahmens.

20

18. Vorrichtung nach Anspruch 16, wobei die Bildungseinrichtung ist:

eine Auswahleinrichtung einer stückweise linearen Verzögerungskontur, die aus dem Langzeitvorhersage-Verzögerungsparameter des vorhergehenden Rahmens und dem Langzeitvorhersage-Verzögerungsparameter des derzeitigen Rahmens bestimmt wird.

25

19. Vorrichtung nach einem der Ansprüche 14 bis 18, wobei die Bildungseinrichtung eine Sucheinrichtung eines Langzeitvorhersage-Verzögerungsparameterwerts durch Rückwärtsiteration in mehreren Phasen ist, und wobei begonnen wird mit einem Langzeitvorhersage-Verzögerungsparameterwert, der für das Ende des derzeitigen Rahmens vorhergesagt wird, wobei jede aufeinanderfolgende Phase eine gesteigerte Auflösung und einen stärker fokussierten Suchbereich aufweist.

30

20. Vorrichtung nach Anspruch 19, umfassend eine Vorhersageeinrichtung des Langzeitvorhersage-Verzögerungsparameterwerts als gleich der Differenz zwischen dem Langzeitvorhersage-Verzögerungsparameterwert am Ende des vorhergehenden Rahmens und zweimal der Differenz zwischen den Stellen der Tonhöhenpulse des Sprachsignals in dem vorhergehenden und dem derzeitigen Rahmen, geteilt durch die Anzahl von Iterationen der Funktion.

35

21. Vorrichtung nach einem der Ansprüche 14 bis 20, umfassend eine Modifizierungseinrichtung des Sprachsignals durch Verschieben von Tonhöhen-Zyklus-Segmenten, eins nach dem anderen, um sie an die Verzögerungskontur anzupassen.

40

22. Vorrichtung nach Anspruch 21, umfassend eine Bestimmungseinrichtung einer Segmentverschiebung durch Korrelieren eines Segments in der gewichteten Sprachdomäne mit einem Zielsignal.

45

23. Vorrichtung nach Anspruch 22, umfassend eine Zusammensetzungseinrichtung des Zielsignals unter Verwendung eines synthetisierten gewichteten Sprachsignals des vorhergehenden Rahmens und aller vorhergehenden verschobenen Segmente im derzeitigen Rahmen.

50

Revendications

1. Procédé de formation d'un contour de délai caractérisant une prédiction à long terme dans une technique utilisant une modification du signal pour coder numériquement un signal de parole, le procédé comprenant les étapes consistant à :

55

diviser le signal de parole en une série de trames successives ;

localiser une impulsion de ton du signal de parole dans une trame précédente ; et
localiser une impulsion de ton correspondante du signal de parole dans une trame courante ;

caractérisé par la formation d'un contour de délai en sélectionnant un paramètre de délai de prédiction à long terme pour la trame courante en itérant en sens inverse une fonction d'une variable de temps temporaire, depuis l'emplacement de l'impulsion de ton du signal de parole dans la trame courante et l'emplacement de l'impulsion de ton correspondante du signal de parole dans la trame précédente.

2. Procédé tel que revendiqué dans la revendication 1, comprenant l'étape consistant à :

former le contour de délai comme une fonction de distances d'impulsions de ton successives entre au moins une dernière impulsion de ton de la trame précédente et une dernière impulsion de ton de la trame courante.

3. Procédé tel que revendiqué dans la revendication 1 ou la revendication 2, comprenant en outre l'étape consistant à :

intégralement caractériser le contour de délai avec un paramètre de délai de prédiction à long terme de la trame précédente et le paramètre de délai de prédiction à long terme de la trame courante.

4. Procédé tel que revendiqué dans la revendication 3, dans lequel la formation du contour de délai comprend l'étape consistant à :

interpoler non linéairement le contour de délai entre le paramètre de délai de prédiction à long terme de la trame précédente et le paramètre de délai de prédiction à long terme de la trame courante.

5. Procédé tel que revendiqué dans la revendication 3, dans lequel la formation du contour de délai comprend l'étape consistant à :

déterminer un contour de délai linéaire pièce par pièce entre le paramètre de délai de prédiction à long terme de la trame précédente et le paramètre de délai de prédiction à long terme de la trame courante.

6. Procédé tel que revendiqué dans l'une quelconque des revendications précédentes, dans lequel la localisation d'une impulsion de ton comprend de dériver un signal résiduel de prédiction linéaire à partir du signal de parole.

7. Procédé tel que revendiqué dans l'une quelconque des revendications 1 à 5, dans lequel la localisation d'une impulsion de ton comprend de dériver un signal de parole pondéré à partir du signal de parole.

8. Procédé tel que revendiqué dans l'une quelconque des revendications 1 à 5, dans lequel la localisation d'une impulsion de ton comprend de dériver un signal de parole pondéré synthétisé à partir du signal de parole.

9. Procédé tel que revendiqué dans l'une quelconque des revendications précédentes, dans lequel l'itération en sens inverse comprend de rechercher une valeur de paramètre de délai de prédiction à long terme dans plusieurs phases et de commencer avec une valeur de paramètre de délai de prédiction à long terme prédite pour la fin de la trame courante, chaque phase successive ayant une résolution accrue et une plage de recherche plus concentrée.

10. Procédé tel que revendiqué dans la revendication 9, comprenant de prédire la valeur de paramètre de délai de prédiction à long terme comme étant égale à la différence entre la valeur de paramètre de délai de prédiction à long terme à la fin de la trame précédente et deux fois la différence entre les emplacements des impulsions de ton du signal de parole dans la trame précédente et la trame courante divisée par le nombre d'itérations de la fonction.

11. Procédé tel que revendiqué dans l'une quelconque des revendications précédentes, comprenant de modifier le signal de parole en décalant des segments de cycle de ton un par un pour les ajuster au contour de délai.

12. Procédé tel que revendiqué dans la revendication 11, comprenant de déterminer un décalage de segment en corrélant un segment dans le domaine de parole pondéré avec un signal cible.

13. Procédé tel que revendiqué dans la revendication 12, comprenant de composer le signal cible en utilisant le signal de parole pondéré synthétisé de la trame précédente et n'importe quels segments décalés précédents dans la trame courante.

14. Dispositif (603) pour former un contour de délai caractérisant une prédiction à long terme dans une technique utilisant une modification de signal pour coder numériquement un signal de parole, le dispositif comprenant :

un diviseur du signal de parole en une série de trames successives ;
un détecteur d'un emplacement d'une impulsion de ton du signal de parole dans une trame précédente ; et
un détecteur d'un emplacement d'une impulsion de ton correspondante du signal de parole dans une trame courante,

caractérisé par un précédent d'un contour de délai pour sélectionner un paramètre de délai de prédiction à long terme pour la trame courante par l'intermédiaire d'une itération en sens inverse d'une fonction d'une variable de temps temporaire, depuis l'emplacement de l'impulsion de ton du signal de parole dans la trame courante et l'emplacement de l'impulsion de ton correspondante du signal de parole dans la trame précédente.

15. Dispositif tel que revendiqué dans la revendication 14, dans lequel le précédent est :

un calculateur du paramètre de délai de prédiction à long terme comme une fonction des distances d'impulsions de ton successives entre la dernière impulsion de ton de la trame précédente et la dernière impulsion de ton de la trame courante.

16. Dispositif tel que revendiqué dans la revendication 14 ou la revendication 15, incorporant en outre :

une fonction caractérisant intégralement le contour de délai avec un paramètre de délai de prédiction à long terme de la trame précédente et le paramètre de délai de prédiction à long terme de la trame courante.

17. Dispositif tel que revendiqué dans la revendication 16, dans lequel le précédent est :

un sélecteur d'un contour de délai interpolé non linéairement entre le paramètre de délai de prédiction à long terme de la trame précédente et le paramètre de délai de prédiction à long terme de la trame courante.

18. Dispositif tel que revendiqué dans la revendication 16, dans lequel le précédent est :

un sélecteur d'un contour de délai linéaire pièce par pièce déterminé à partir du paramètre de délai de prédiction à long terme de la trame précédente et du paramètre de délai de prédiction à long terme de la trame courante.

19. Dispositif tel que revendiqué dans l'une quelconque des revendications 14 à 18, dans lequel le précédent est :

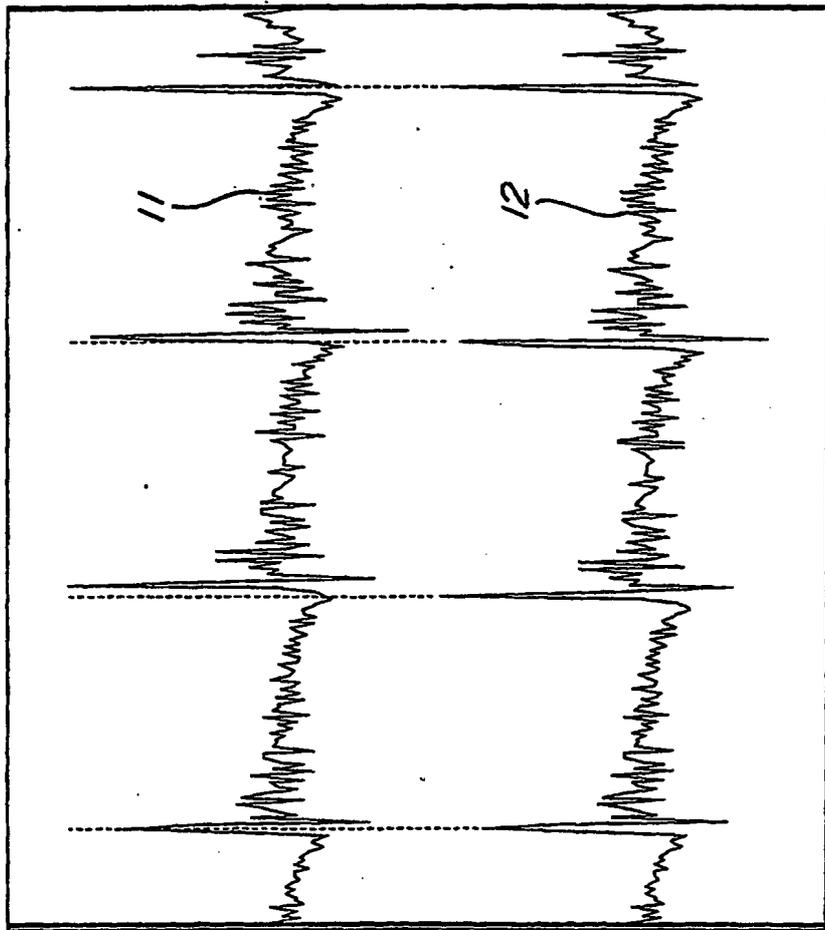
un chercheur de valeur de paramètre de délai de prédiction à long terme par itération en sens inverse dans plusieurs phases et commençant avec une valeur de paramètre de délai de prédiction à long terme prédite pour la fin de la trame courante, chaque phase successive ayant une résolution accrue et une plage de recherche plus concentrée.

20. Dispositif tel que revendiqué dans la revendication 19, comprenant un prédicteur de la valeur de paramètre de délai de prédiction à long terme comme étant égale à la différence entre la valeur de paramètre de délai de prédiction à long terme à la fin de la trame précédente et deux fois la différence entre les emplacements des impulsions de ton du signal de parole dans la trame précédente et la trame courante divisée par le nombre d'itérations de la fonction.

21. Dispositif tel que revendiqué dans l'une quelconque des revendications 14 à 20, comprenant un modificateur du signal de parole en décalant des segments de cycle de ton un par un pour les ajuster au contour de délai.

22. Dispositif tel que revendiqué dans la revendication 21, comprenant un déterminateur d'un décalage de segment en corrélant un segment dans le domaine de parole pondéré avec un signal cible.

23. Dispositif tel que revendiqué dans la revendication 22, comprenant un compositeur du signal cible utilisant un signal de parole pondéré synthétisé de la trame précédente et n'importe quels segments décalés précédents dans la trame courante.



Original residual signal

Modified residual signal

t_{n-1} time t_n

FIG. 1

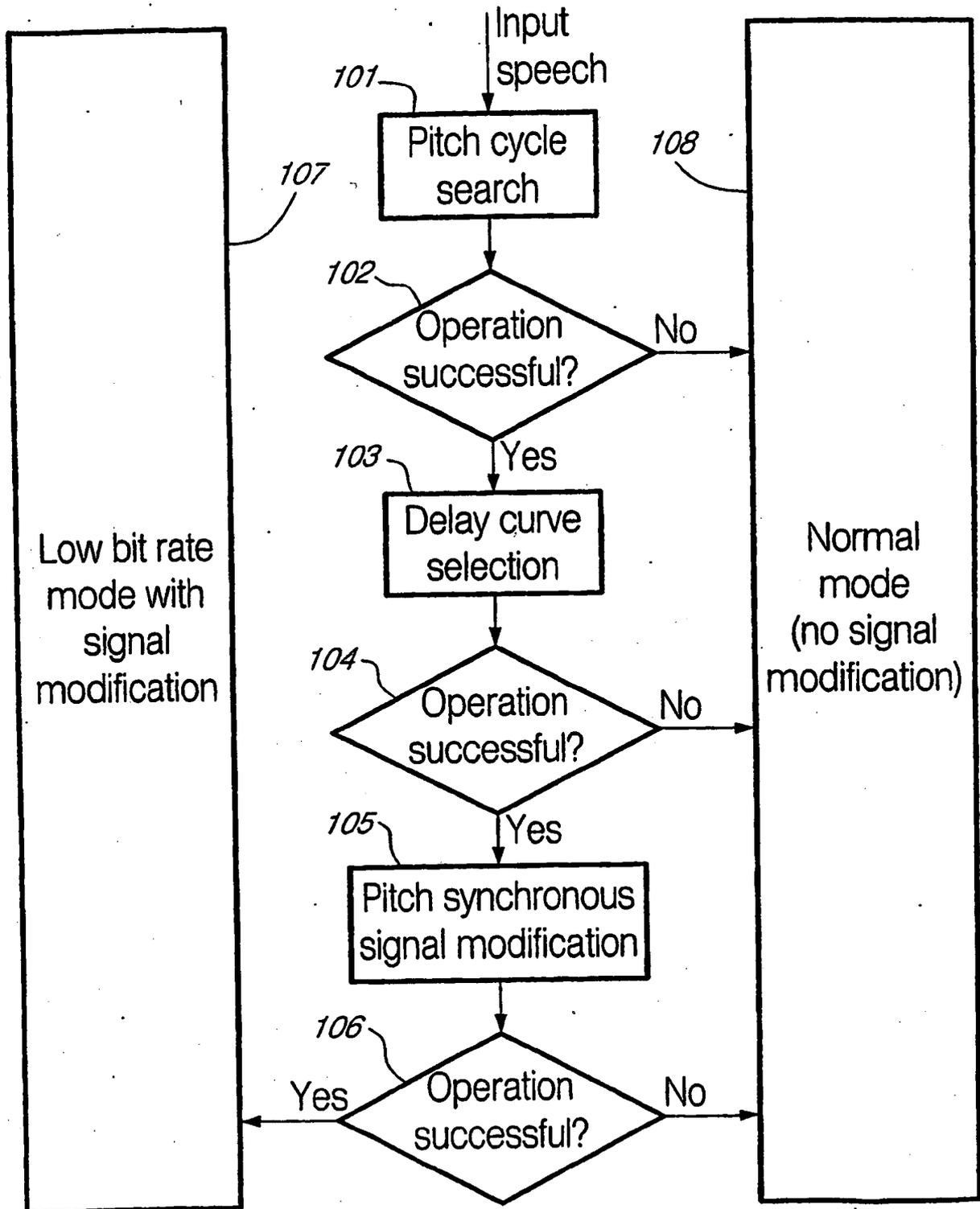


FIG. 2

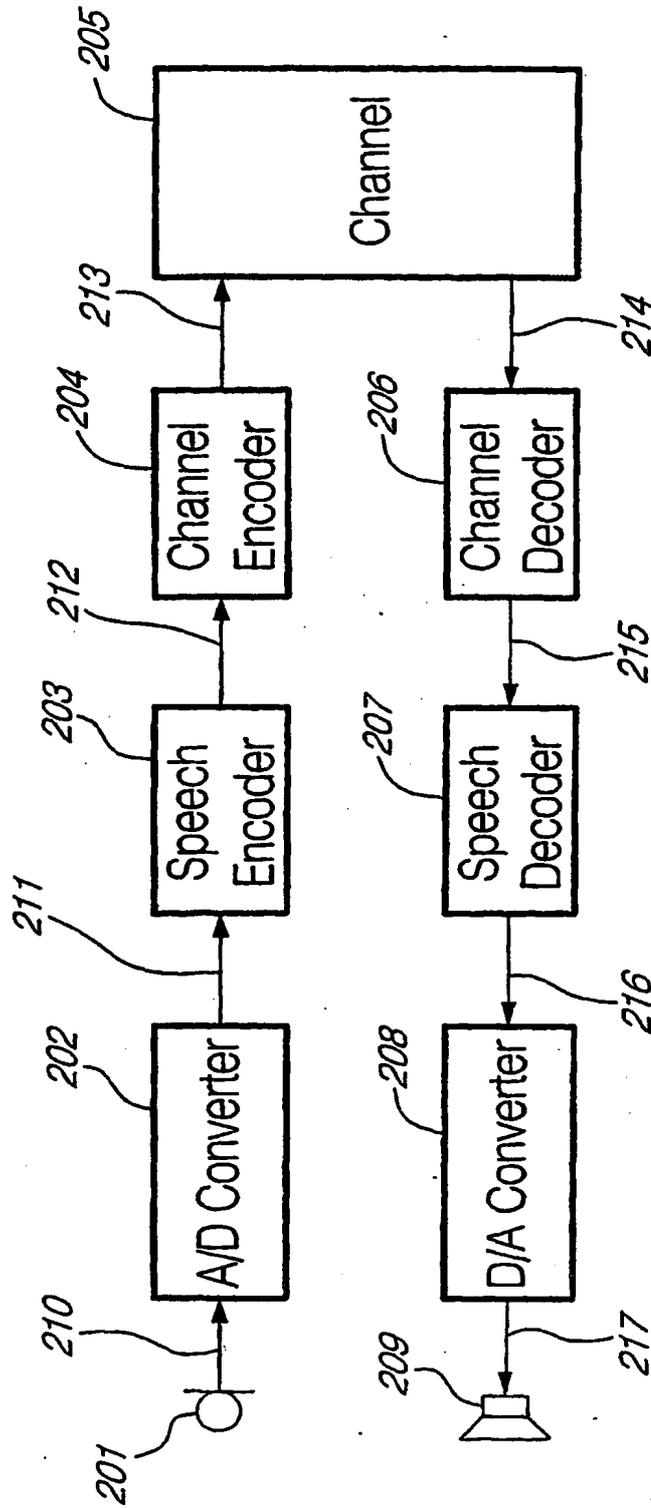


FIG. 3

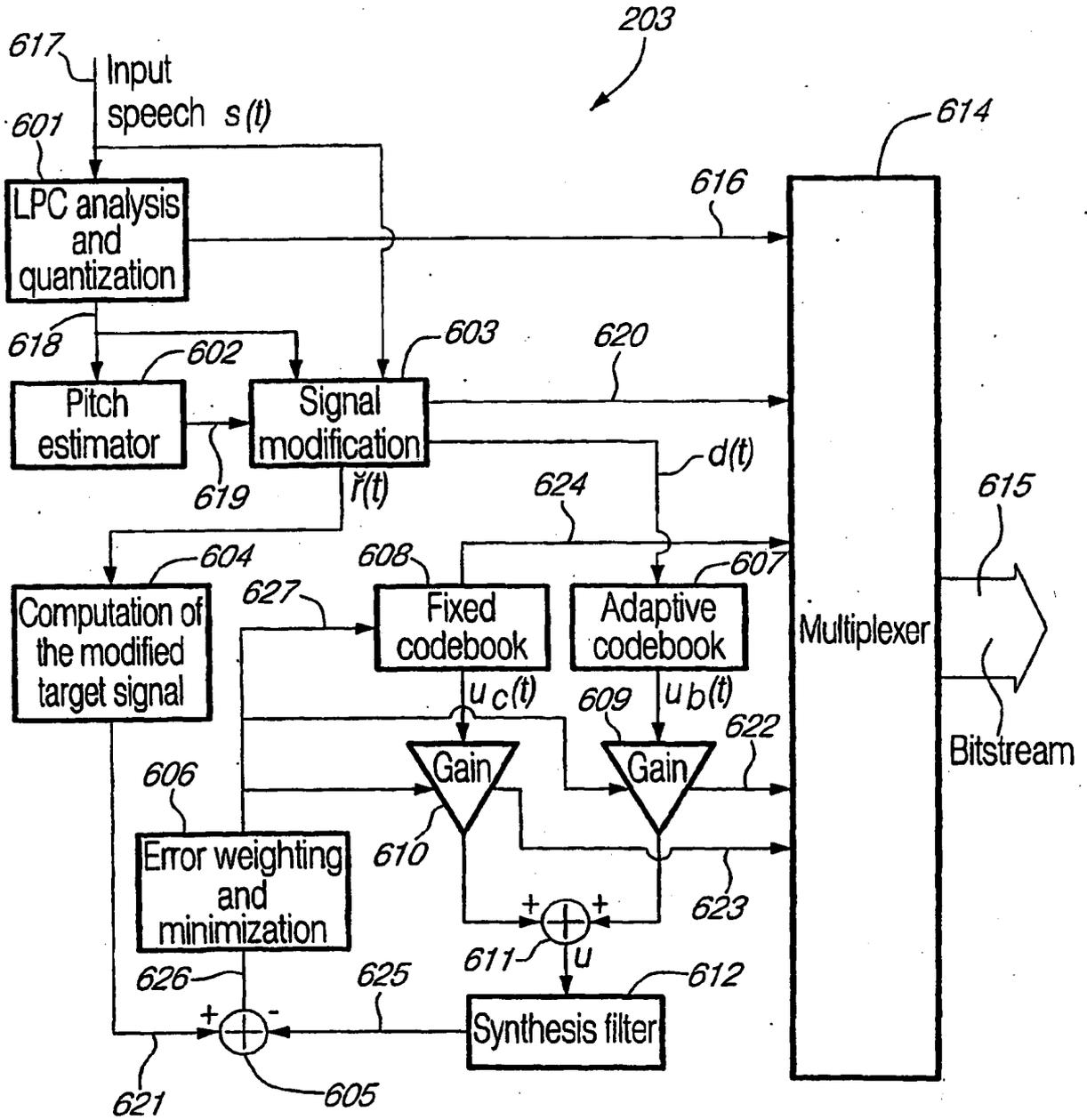


FIG. 4

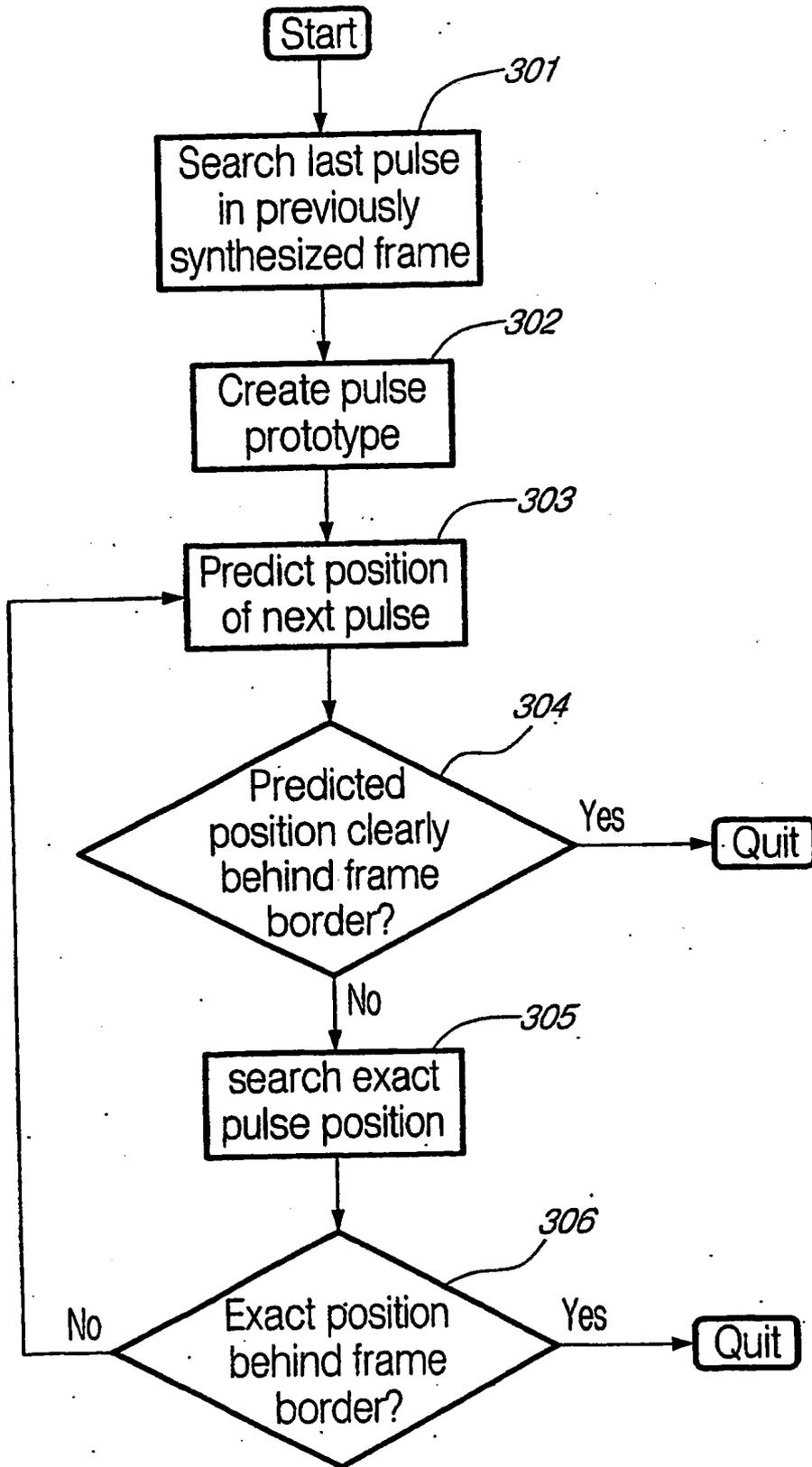


FIG. 5

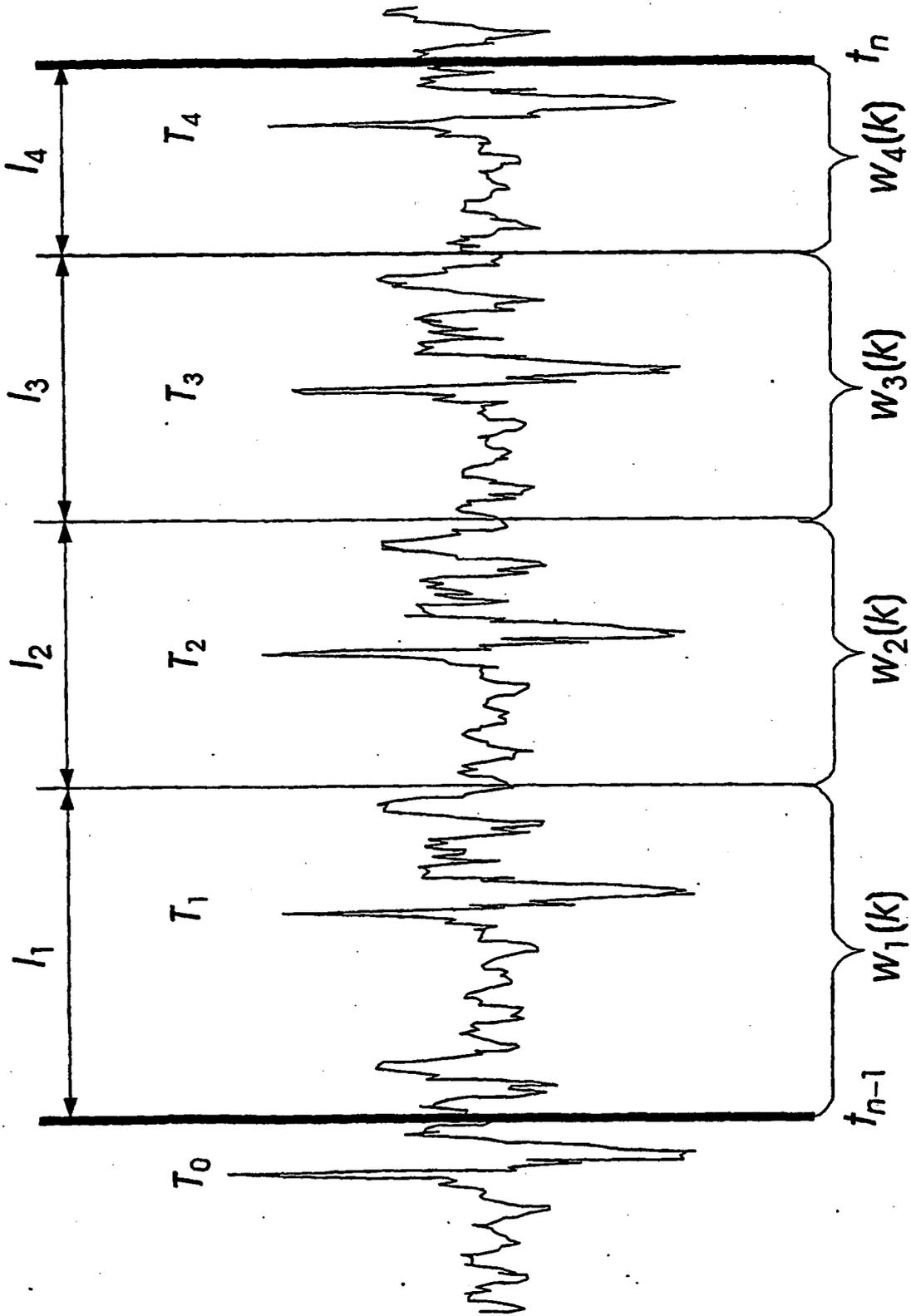


FIG. 6

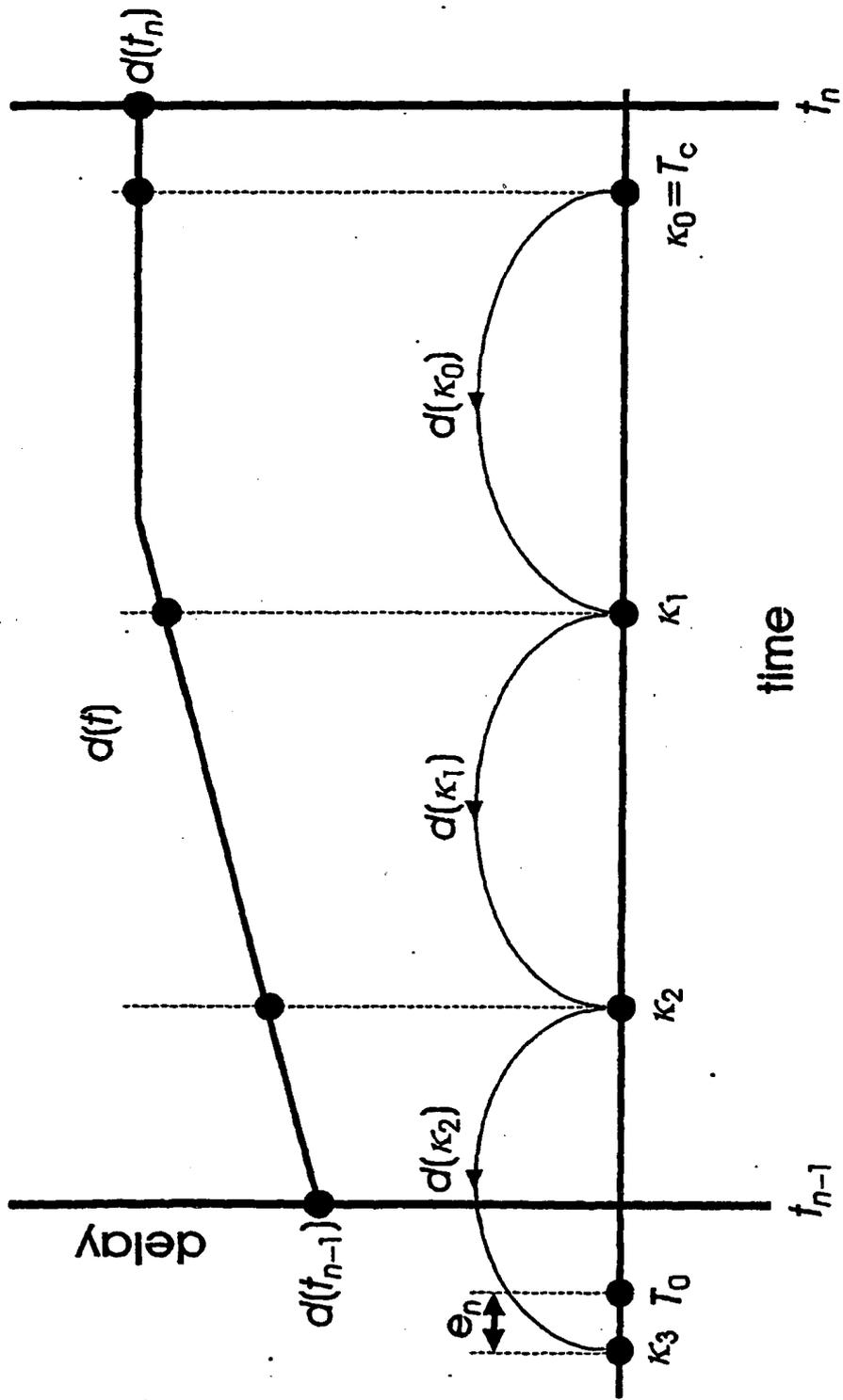


FIG. 7

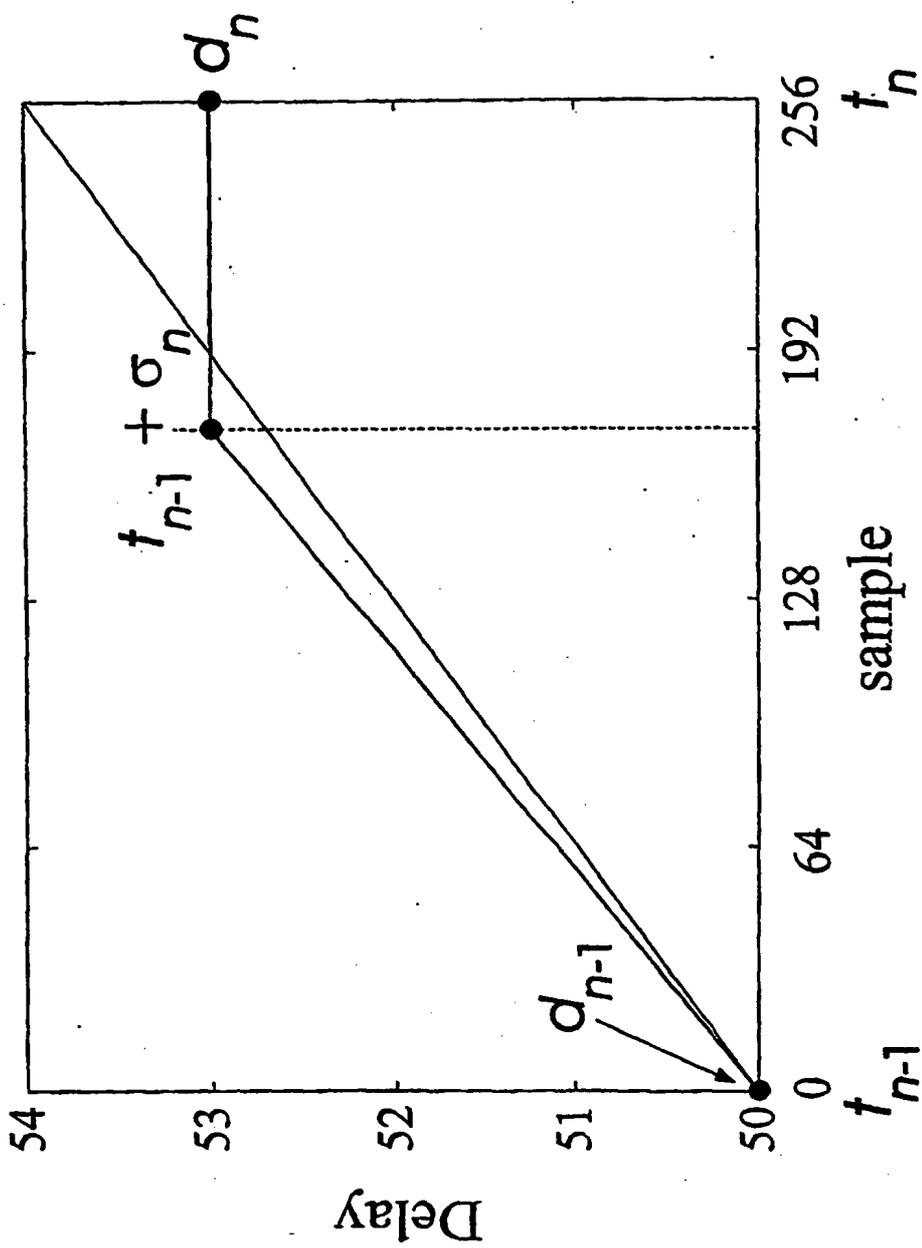


FIG. 8

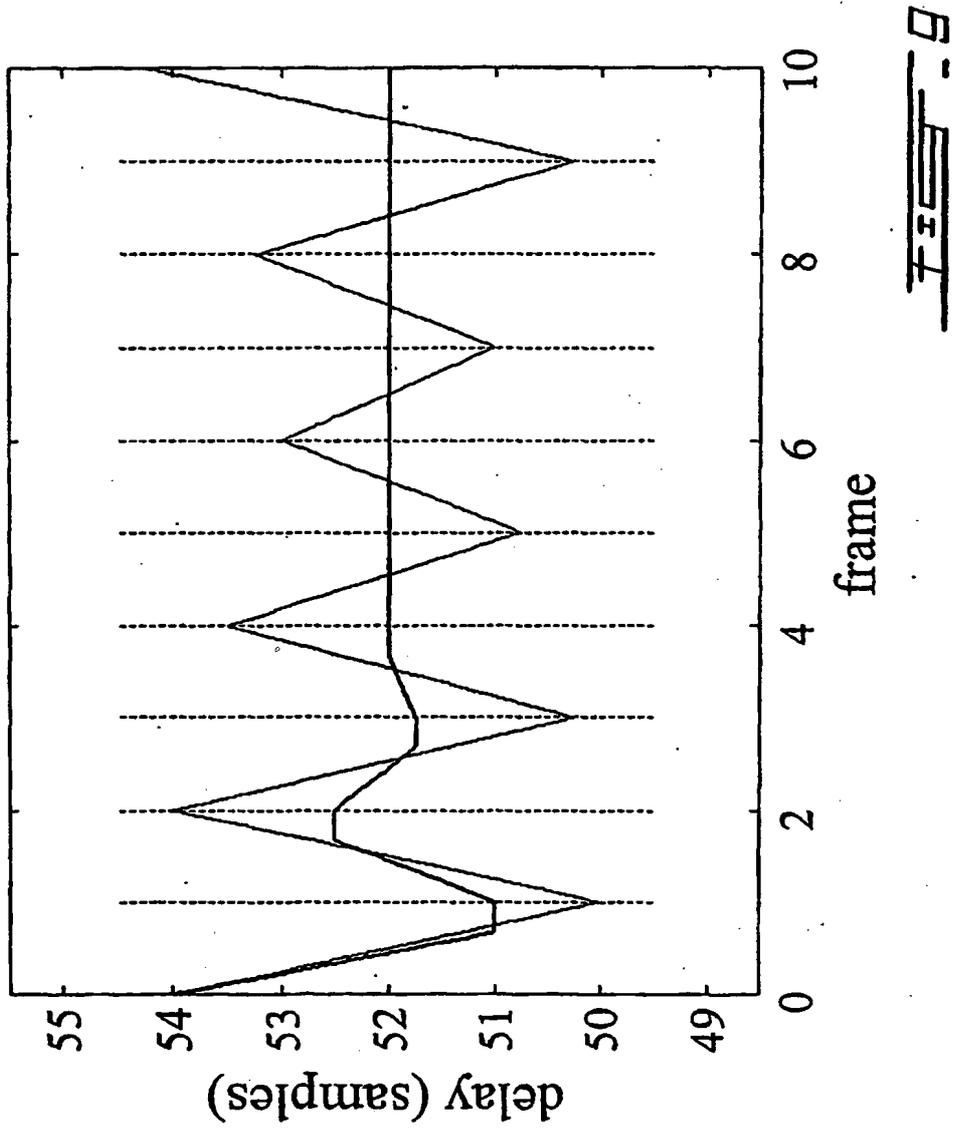


FIG. 9

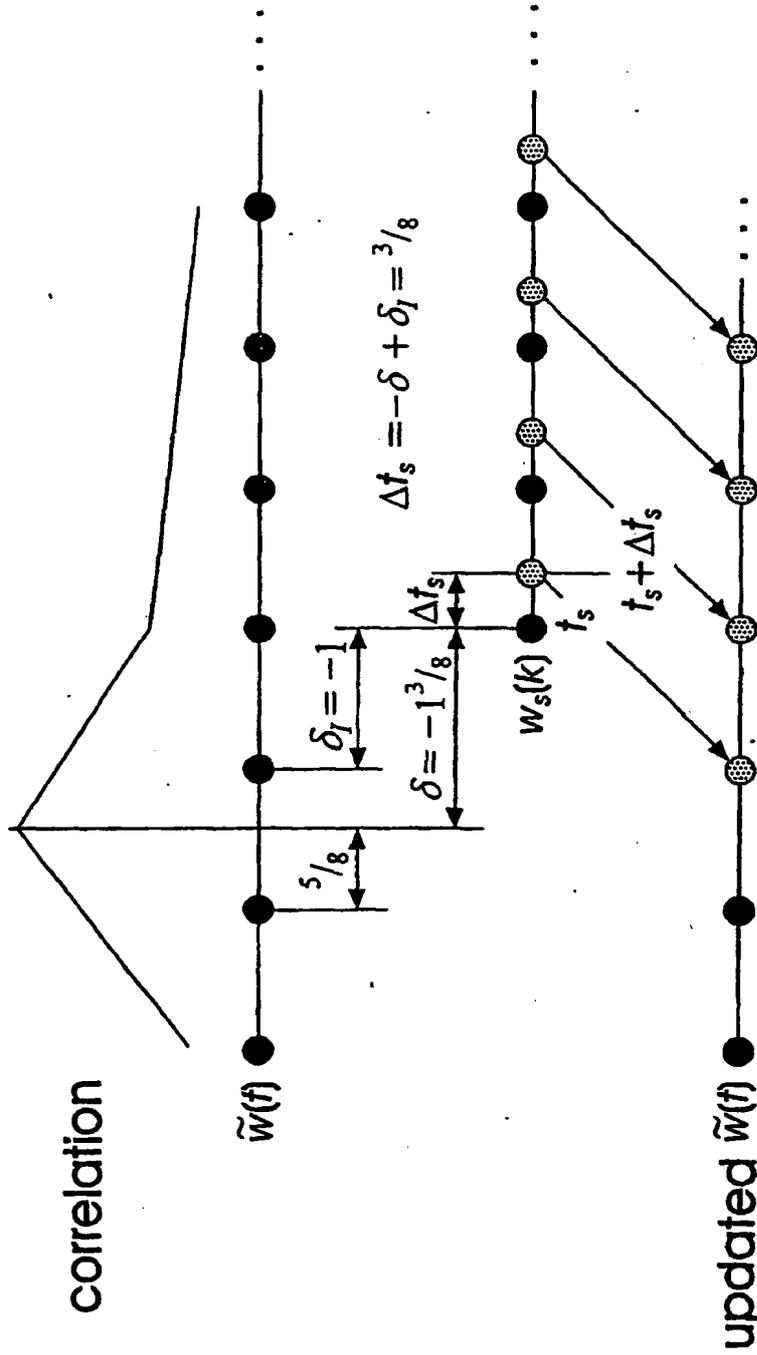
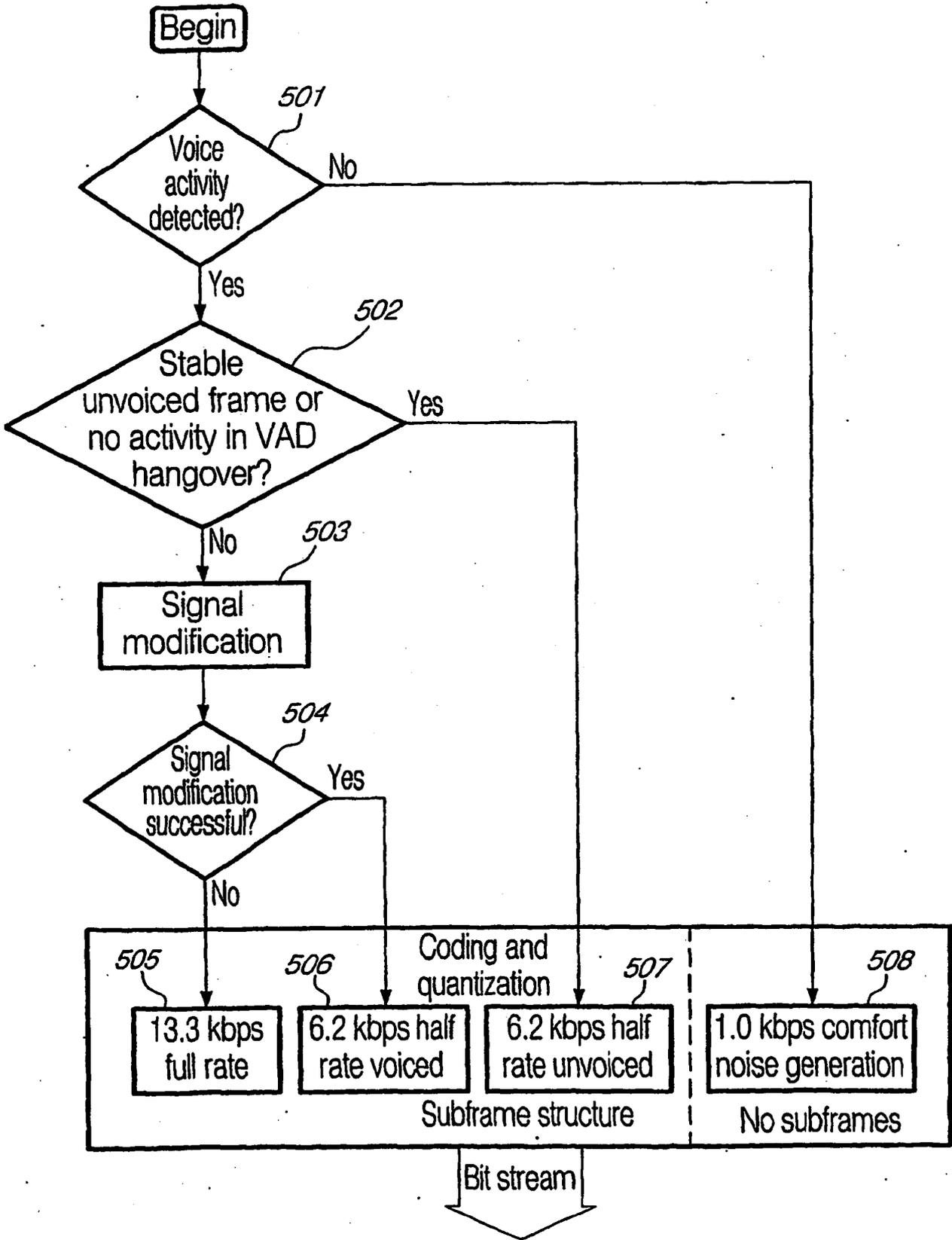


FIG. 11



~~FIG. 12~~

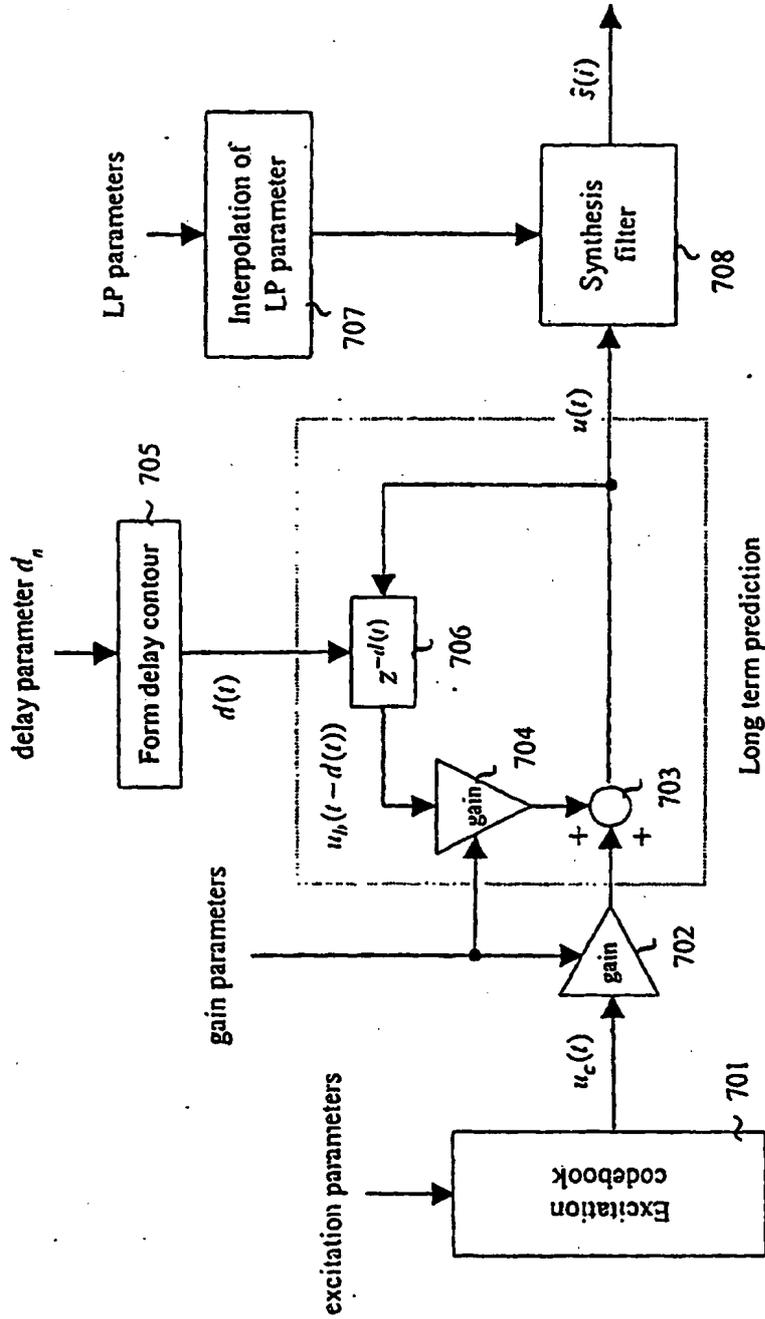


Figure 13

REGISTER ENTRY FOR EP1454315

European Application No EP02784985.0 filing date 13.12.2002

Priority claimed:

14.12.2001 in Canada - doc: 2365203

PCT EUROPEAN PHASE

PCT Application PCT/CA2002/001948 Publication No WO2003/052744 on
26.06.2003

Designated States AT BE BG CH CY CZ DE DK EE ES FI FR GB GR IE IT LI LU MC NL
PT SE SI SK TR

Title SIGNAL MODIFICATION METHOD FOR EFFICIENT CODING OF SPEECH SIGNALS

Applicant/Proprietor

NOKIA CORPORATION, Keilalahdentie 4, 02150 Espoo, Finland

[ADP No. 69578813001]

Inventors

MIKKO TAMMI, Kemiaankatu 9 E 51, 33 720 Tampere, Finland

[ADP No. 73596025001]

MILAN JELINEK, 245 Merrill Park, North Hatley, Quebec, J0B 2C0, Canada

[ADP No. 73596033001]

CLAUDE LAFLAMME, 294 chemin Dépôt, Orford, Quebec, J1X 6W1, Canada

[ADP No. 73596041001]

VESA RUOPPILA, 3913 Mentana Street, Montreal, Quebec, H2L 3R7, Canada

[ADP No. 73596058001]

Classified to

G10L

Address for Service

VENNER SHIPLEY LLP, 20 Little Britain, LONDON, EC1A 7DH, United Kingdom

[ADP No. 08897431001]

EPO Representative

PAUL STEFAN DERRY, Venner Shipley LLP, 20 Little Britain, London EC1A 7DH,
United Kingdom

[ADP No. 71351928001]

Publication No EP1454315 dated 08.09.2004

Publication in English

Examination requested 05.07.2004

Patent Granted with effect from 04.04.2007 (Section 25(1)) with title SIGNAL
MODIFICATION METHOD FOR EFFICIENT CODING OF SPEECH SIGNALS

24.11.2006 VENNER SHIPLEY LLP, 20 Little Britain, LONDON, EC1A 7DH, United
Kingdom

[ADP No. 08897431001]

registered as address for service

Entry Type 8.11 Staff ID. CRED Auth ID. A1

TIMED: 03/07/07 10:25:37

REGISTER ENTRY FOR EP1454315 (Cont.)

PAGE: 2

09.03.2007 EPO: Search report published on 05.02.2004

Entry Type 25.11 Staff ID. RD06 Auth ID. EPT

**** END OF REGISTER ENTRY ****

OA80-01
PA

OPTICS - PATENTS

03/07/07 10:25:48
PAGE: 1

RENEWAL DETAILS

PUBLICATION NUMBER

EP1454315

PROPRIETOR(S)

Nokia Corporation, Keilalahdentie 4, 02150 Espoo, Finland

DATE FILED

13.12.2002

DATE GRANTED

04.04.2007

DATE NEXT RENEWAL DUE

13.12.2007

DATE NOT IN FORCE

DATE OF LAST RENEWAL

YEAR OF LAST RENEWAL

00

STATUS

PATENT IN FORCE

**** END OF REPORT ****