



Europäisches Patentamt  
European Patent Office  
Office européen des brevets



(11) **EP 0 865 026 B1**

(12) **EUROPÄISCHE PATENTSCHRIFT**

(45) Veröffentlichungstag und Bekanntmachung des Hinweises auf die Patenterteilung:  
**03.12.2003 Patentblatt 2003/49**

(51) Int Cl.7: **G10L 21/04**

(21) Anmeldenummer: **98104455.5**

(22) Anmeldetag: **12.03.1998**

(54) **Effizientes Verfahren zur Geschwindigkeitsmodifikation von Sprachsignalen**

Method for modifying speech speed

Méthode pour la modification du débit de parole

(84) Benannte Vertragsstaaten:  
**AT DE FR GB NL**

(73) Patentinhaber: **GRUNDIG Aktiengesellschaft  
90471 Nürnberg (DE)**

(30) Priorität: **14.03.1997 DE 19710545**

(72) Erfinder: **Carl, Holger, Dr.  
90762 Fürth (DE)**

(43) Veröffentlichungstag der Anmeldung:  
**16.09.1998 Patentblatt 1998/38**

(56) Entgegenhaltungen:  
**EP-A- 0 427 953                      EP-A- 0 608 833  
EP-A- 0 726 560**

**EP 0 865 026 B1**

Anmerkung: Innerhalb von neun Monaten nach der Bekanntmachung des Hinweises auf die Erteilung des europäischen Patents kann jedermann beim Europäischen Patentamt gegen das erteilte europäische Patent Einspruch einlegen. Der Einspruch ist schriftlich einzureichen und zu begründen. Er gilt erst als eingelegt, wenn die Einspruchsgebühr entrichtet worden ist. (Art. 99(1) Europäisches Patentübereinkommen).

## Beschreibung

**[0001]** Gegenstand der Erfindung ist ein Verfahren zur Geschwindigkeitsmodifikation von Sprachsignalen im Zeitbereich, insbesondere eine effiziente Overlap-Add-Methode.

**[0002]** In verschiedenen Bereichen der Verarbeitung von Sprach- und Audiosignalen ist eine Veränderung der Wiedergabegeschwindigkeit dieser Signale erwünscht, möglichst ohne daß damit eine Beeinträchtigung ihrer Natürlichkeit und - im Fall von Sprache - ihrer Verständlichkeit verbunden wäre. Dieses Ziel, den Klangcharakter zu erhalten, kann man aus technischer Sicht folgendermaßen formulieren: Trotz einer Modifikation der Zeitskala dieser Signale sollen ihre Kurzzeitspektraleigenschaften unverändert bleiben. Insbesondere bedeutet das für Sprachsignale, daß Grundfrequenz und Formanten bei der Geschwindigkeitsmodifikation erhalten bleiben müssen.

**[0003]** Die Zeitstauchung oder Zeitdehnung von Audiosignalen wird in Studios eingesetzt, zum Beispiel mit dem Ziel, Werbesendungen auf die vorgesehene Länge zu trimmen. Auch in der Diktiertechnik ist die Anpassung der Wiedergabegeschwindigkeit an die Bedürfnisse bzw. Fähigkeiten der Schreibkraft von Bedeutung. Eine weitere Anwendung besteht bei der Echtzeitübertragung von Sprachsignalen, bei der Datenpakete mit variabler Verzögerung beim Empfänger eintreffen. Durch Anwendung der Geschwindigkeitsmodifikation kann man hier die Über-Alles-Verzögerung im Mittel geringer halten als das Worst-Case Delay der Übertragungstrecke, ohne daß ein zu spät eintreffendes Datenpaket zu Aussetzern oder anderen, ähnlich störenden Effekten führen würde.

Für viele Anwendungen ergeben sich neben dem Wunsch nach möglichst hoher Klangqualität die folgenden zusätzlichen Anforderungen an das Verfahren:

**[0004]** Eine kostengünstige Echtzeitrealisierung muß erzielbar sein, und es muß zur Laufzeit eine nach Möglichkeit stufenlose Änderung des Geschwindigkeitsmodifikationsfaktors möglich sein. Von Vorteil ist ohne Zweifel auch, wenn der Algorithmus ohne eine stets fehlerbehaftete Pitch-Schätzung auskommt.

**[0005]** Aus "Method for Time or Frequency Compression-Expansion of Speed", von G. Fairbaks und R. P. Jaeger, Inst. of Radio Engineers Trans. on Audio, Vol. AU-2, No. 1, pp. 7-12, Jan. 1954, sind erste Untersuchungen zur Sprachsignalstauchung bzw. Sprachsignaldehnung bekannt. Häufig wurden seitdem Frequenzbereichsverfahren eingesetzt - naheliegend, da, wie eingangs erwähnt, die Kurzzeitspektraleigenschaften des Sprachsignals erhalten bleiben sollen. Seit Mitte der achtziger Jahre sind vergleichsweise einfache im Zeitbereich arbeitende Overlap-Add-Verfahren bekannt, mit denen sehr gut klingende zeitskalierte Sprachsignale erzeugt werden können.

**[0006]** In "Signal Estimation from Modified Short-Time Fourier Transform", von D. W. Griffin, in IEEE Trans. Acoust., Speech, Signal Processing, Vol. ASSP-32, No. 2, pp. 236-242, Apr. 1984, berichten Griffin und Lim von Experimenten mit einer sehr aufwendigen iterativ arbeitenden Phasenbestimmung. Auf diesen Ansatz nimmt wiederum die Veröffentlichung von S. Roucos und A. M. Wilgus "High Quality Time-Scale Modification for Speech", IEEE Proc. Int. Conf. Acoust., Speech, Signal Processing, pp. 493-496, 1985, Bezug, die eine Zeitbereichsmethode vorschlagen, die mittels eines Overlap-Add-Ansatzes zeitskalierte Sprachsignale erzeugt. Bei diesem sogenannten SOLA-Verfahren (SOLA = Synchronized Overlap-Add) erfolgt eine Synchronisation der in regelmäßigen Abständen dem Originalsignal entnommenen Abschnitte durch Verschiebung vor der jeweils entsprechenden Fensterung und Addition im Zielsignal. Dies entspricht im weiteren Sinne der Phasenoptimierung, wie sie in den Frequenzbereichsverfahren durchgeführt wird. Eng mit dem SOLA-Algorithmus verwandt ist das sogenannte WSOLA-Verfahren (WSOLA = Waveform Similarity Overlap-Add), das W. Verhelst und M. Roelands in "An Overlap-Add Technique Based on Waveform Similarity (WSOLA) for High Quality Time-Scale Modification of Speed", IEEE Proc. Int. Conf. Acoust., Speech, Signal Processing, pp. 554-567, 1993, und "Waveform Similarity Based Overlap-Add (WSOLA) for Time-Scale Modification of Speech: Structures and Evaluation", Int. Conf. on Speech Communication and Technology, pp. 337-340, 1993, vorstellen. Der Hauptunterschied zwischen diesen beiden Ansätzen besteht in der Synchronisation, die im WSOLA-Verfahren durch versetztes Entnehmen von Segmenten aus dem Originalsignal durchgeführt wird, was sich gegenüber dem SOLA-Prinzip vor allem aufwandsmindernd auswirkt.

**[0007]** Aufgabe der Erfindung ist es, ein Verfahren zur Geschwindigkeitsmodifikation von Sprachsignalen im Zeitbereich anzugeben, das besonders effizient arbeitet.

**[0008]** Diese Aufgabe wird durch die Merkmale der Ansprüche 1 und 3 gelöst. Vorteilhafte Ausgestaltungen der Erfindung sind in der nachfolgenden Beschreibung angegeben.

**[0009]** Die Erzeugung der mit dem Faktor  $\alpha$  zeitskalierten Version  $y(k)$  eines Sprachsignals  $x(k)$  erfolgt gemäß der Synthese

$$y(k) = \sum_{\lambda = -\infty}^{\infty} x(k + \lambda(\alpha - 1)L + \Delta_{\lambda})w(k - \lambda L)$$

mit einer Fensterfunktion

$$w(k) = \begin{cases} v(k) & \text{für } 0 < k < N \\ 1 & \text{für } N \leq k < L \\ 1 - v(k - L) & \text{für } L \leq k < L + N \\ 0 & \text{sonst} \end{cases}$$

**[0010]** Die hierin vorkommende für  $k=0, \dots, N-1$  definierte Funktion  $v(k)$  ist dabei sinnvollerweise zwischen ihren Extrema  $v(0)=\varepsilon_0$  mit  $0 < \varepsilon_0 < 1$  und  $v(N-1)=1-\varepsilon_1$  mit  $0 < \varepsilon_1 < 1$  monoton wachsend.

**[0011]** Die angegebene  $w(k)$ -Definition stellt sicher, daß die für sinnvolles Overlap-Add notwendige Bedingung

$$\sum_{\lambda=-\infty}^{\infty} w(k - \lambda L) \equiv 1 \quad \forall k \in \{-\infty, \dots, \infty\}$$

erfüllt ist.

**[0012]** Die in obiger Synthesegleichung enthaltene Verschiebevariable  $\Delta_\lambda$  ist zwecks der erwähnten Synchronisation aus einem "Toleranzbereich"  $-\Delta_{\max}, \dots, \Delta_{\max}$  zu bestimmen.

**[0013]** Die prinzipielle Vorgehensweise ist wie folgt:

Aus dem Originalsignal  $x(k)$  werden in - abgesehen von einem synchronisationsbedingtem "Jitter" - regelmäßigen  $\alpha L$  Werte betragenden Abständen Segmente der Länge  $L+N$  entnommen und nach Gewichtung mit  $w(k)$  jeweils um  $L$  Abtastwerte versetzt aufaddiert. Das auf diese Weise erhaltene Signal  $y(k)$  ist gegenüber  $x(k)$  um den Faktor  $\alpha$  beschleunigt, das heißt, daß eine im Originalsignal  $x(k)$  enthaltene Äußerung von  $K$  Abtastwerten Länge durch dieses Vorgehen auf einen  $y(k)$  -Abschnitt der Länge  $K/\alpha$  abgebildet, also verkürzt und damit in der Wiedergabe beschleunigt für  $\alpha > 1$ , bzw. verlängert, das heißt verlangsamt, wird, wenn  $\alpha < 1$  ist.

**[0014]** Die Synchronisation der zu überlappenden Abschnitte ist für die resultierende Klangqualität von großer Bedeutung. Hierzu wird der folgende Ansatz verwendet: Während der Abarbeitung des Verfahrens kann zu jedem dem Signal  $x(k)$  entnommenen Segment für den nächsten Schritt als "Idealsegment" der um  $L$  Abtastwerte versetzte Abschnitt von  $x(k)$  angesehen werden, da durch diese Wahl die Overlap-Add-Operation wieder das Originalsignal  $x(k)$  reproduzieren würde. Die erwünschte Zeitskalierung erfordert nun aber, daß für die Overlap-Add-Synthese i. a. ein anderer, gegenüber dem "Idealsegment" versetzter Abschnitt von  $x(k)$  ausgewählt wird. Die bestmögliche Synchronisation ist gegeben, wenn der für die Overlap-Add-Operation benutzte Abschnitt größtmögliche Ähnlichkeit ("Waveform Similarity") mit dem "Idealsegment" aufweist.

**[0015]** Als Kriterium für die Ähnlichkeit der genannten Segmente bieten sich verschiedene Maße an. Naheliegender ist beispielsweise die Benutzung des Korrelationskoeffizienten. Während W. Verhelst und M. Roelands in "An Overlap-Add Technique Based on Waveform Similarity (WSOLA) for High Quality Time-Scale Modification of Speed", in IEEE Proc. Int. Conf. Acoust., Speech, Signal Processing, pp. 554-557, 1993, und "Waveform Similarity Based Overlap-Add (WSOLA) for Time-Scale Modification of Speech: Structures and Evaluation" in Int. Conf. on Speech Communication and Technology, pp. 337-340, 1993, für die Auswertung des Ähnlichkeitsmaßes das komplette Segment der Länge  $L+N$  herangezogen haben, erscheint es als vollkommen ausreichend, die Berechnung auf den Bereich der  $N$  Abtastwerte zu beschränken, in dem die Segmente tatsächlich überlappen.

**[0016]** Für die weiteren Darstellungen ist es hilfreich, die folgende Vektormodifikation einzuführen:

Der  $N$  Werte lange Abschnitt des "Idealsegments", in dem die Überlappung mit dem neu zu bestimmenden Segment stattfinden wird, sei mit  $x$  bezeichnet, die ersten  $N$  Werte des verschobenen Segments mit  $x_q$ . Die Gewichtung dieses Abschnitts mit der steigenden Flanke des Fensters wird durch Multiplikation dieses Vektors mit einer Diagonalmatrix  $V$  repräsentiert, die mit den Werten

$v(0), \dots, v(N-1)$  besetzt ist. Entsprechend wird die Gewichtung des Idealsegmentabschnitts  $x$  mit der fallenden Flanke des Fensters durch Multiplikation mit  $\mathbf{1} - V$  dargestellt, wobei  $\mathbf{1}$  die  $N \times N$ -Einheitsmatrix bezeichnet. Der im kritischen Überlappungsbereich aus der Overlap-Add-Synthese resultierende  $y(k)$ -Abschnitt lautet damit

$$y = (\mathbf{1} - V)x + Vx_q$$

**[0017]** Beispielsweise läßt sich nun als Maß für die Ähnlichkeit der hierbei beteiligten Komponenten eine Kreuzkorreliertenberechnung gemäß

5

$$C_{\delta} = x^T (1-V)^T V x_q$$

angeben. Die Maximierung dieses Ausdrucks bezüglich der sich in  $x_q$  wiederfindenden Verschiebung  $\delta \in \{-\Delta_{\max}, \dots, \Delta_{\max}\}$  liefert die für das betrachtete Segment im Sinne des angesetzten Ähnlichkeitsmaßes optimale Verschiebung  $\Delta_{\lambda}$ .

10

**[0018]** Die Berechnung der  $C_{\delta}$  erfordert alle  $L$  Abtastwerte  $2N$  Multiplikationen für die Vorabberechnung des Ausdrucks  $x^T(1-V)^T V$  sowie anschließend  $(2\Delta_{\max}+1)N$  Multiplikationen und Additionen.

15

**[0019]** Dies stellt gegenüber W. Verhelst und M. Roelands in "An Overlap-Add Technique Based on Waveform Similarity (WSOLA) for High Quality Time-Scale Modification of Speed", in IEEE Proc. Int. Conf. Acoust., Speech, Signal Processing, pp. 554-557, 1993, und "Waveform Similarity Based Overlap-Add (WSOLA) for Time-Scale Modification of Speech: Structures and Evaluation" in Int. Conf. on Speech Communication and Technology, pp. 337-340, 1993, eine Aufwandsreduktion um den Faktor zwei dar, der sich für  $L > N$  sogar noch erhöht. Die Beschränkung der Ähnlichkeitsberechnung auf den Bereich der Überlappung hat keinerlei negative Auswirkungen auf die Qualität der zeitskalierten Sprachproben.

20

**[0020]** Ein anderer Ansatz für die Synchronisation ist, anstelle der Maximierung der "Waveform Similarity" den Fehler zwischen dem synthetisierten Signal  $y$  und dem Originalsignal  $x$  zu minimieren. Eine einfache willkürliche Wahl ist, für diesen Fehler den quadratischen Ausdruck

$$E_{\delta} = \|x - y\|^2$$

25

anzusetzen.

**[0021]** Bei Vernachlässigung der Vorabberechnungen beläuft sich der für die Auswertung von  $E_{\delta}$  anfallende Aufwand auf  $(2\Delta_{\max}+1)4N$  DSP-Operationen alle  $L$  Abtastwerte. Hierunter werden solche Operationen verstanden, die ein Signalprozessor mit gängiger Architektur in einem Schritt abarbeiten kann.

30

**[0022]** Ein weiterer Ansatz besteht darin, anstelle des absoluten Fehlers den relativen Fehler

$$R_{\delta} = \frac{\|x - y\|^2}{\|y\|^2}$$

35

zu minimieren, was als SNR-Maximierung interpretiert werden kann.  $(2\Delta_{\max}+1)5N$  Operationen sind hier vor jeder Overlap-Add-Operation erforderlich.

## Patentansprüche

40

1. Verfahren zur Geschwindigkeitsmodifikation von Sprachsignalen, insbesondere digitalisierten Sprachsignalen, bei dem

45

- ein analoges Sprachsignal digitalisiert wird, wodurch ein digitalisiertes Sprachsignal entsteht, welches in einem Speicher gespeichert wird,

- ein Faktor  $\alpha$  definiert wird, um welchen das Sprachsignal verlängert oder verkürzt wird,

50

- eine Fensterfunktion mit einem ersten steigenden Abschnitt der Länge  $N$ , einem zweiten, sich direkt an den ersten Abschnitt anschließenden, konstanten Abschnitt der Länge  $L-N$  und einem dritten, sich direkt an den zweiten Abschnitt anschließenden, fallenden Abschnitt definiert wird, wobei bei einer Überlagerung des ersten steigenden Abschnittes eines Fensters mit dem dritten fallenden Abschnitt eines anderen Fensters und einer Addition beider Abschnitte im Überlappungsbereich sich das Ergebnis eines ergibt, was dem Wert des zweiten Abschnittes der Fensterfunktion entspricht,

55

- aus dem digitalisierten, gespeicherten Sprachsignal in unregelmäßigen Abständen einer mittleren Länge  $\alpha L$  Segmente einer definierten Länge  $L+N$  entnommen werden,

- diese, aus dem digitalisierten, gespeicherten Sprachsignal entnommenen Segmente mit der Fensterfunktion im Zeitbereich gewichtet werden,

- die gewichteten Segmente jeweils um eine definierte Anzahl von Abtastwerten  $L$  versetzt aufaddiert werden, wodurch das so entstehende Sprachsignal für  $\alpha > 1$  verkürzt und für  $\alpha < 1$  verlängert wird,

- nacheinander an den Stellen der Entnahme der Segmente aus dem digitalisierten Sprachsignal das dort entnommene, mit der Fensterfunktion gewichtete Segment mit dem nachfolgend entnommenen, ebenfalls mit der Fensterfunktion gewichteten Segment unter Ähnlichkeitsaspekten verglichen wird, **dadurch gekennzeichnet,**
  - **dass** zum schnellen Vergleich der Ähnlichkeit der Segmente lediglich der N Werte lange dritte, mit dem fallenden Fensterabschnitt gewichtete Abschnitt des Segmentes mit dem jeweils ersten, mit dem steigenden N Werte langen Fensterabschnitt gewichteten Abschnitt des nachfolgend entnommenen Segmentes verglichen wird,
  - **dass** diese Segmente so zueinander versetzt aufaddiert werden daß die Ähnlichkeit der beiden Segmentabschnitte maximal wird, und
  - **dass** zur Berechnung der Ähnlichkeit, als deren Maß, eine Korrelation verwendet wird.
2. Verfahren nach Anspruch 1,  
**dadurch gekennzeichnet, dass**
- die Ähnlichkeit beider vergleichener Segmentabschnitte maximal wird, wenn eine Maximierung des Ähnlichkeitsmaßes in Bezug zur Verschiebung zueinander durchgeführt wird.
3. Verfahren zur Geschwindigkeitsmodifikation von Sprachsignalen, insbesondere digitalisierten Sprachsignalen, bei dem
- ein analoges Sprachsignal digitalisiert wird, wodurch ein digitalisiertes Sprachsignal entsteht, welches in einem Speicher gespeichert wird,
  - ein Faktor  $\alpha$  definiert wird, um welchen das Sprachsignal verlängert oder verkürzt wird,
  - eine Fensterfunktion mit einem ersten steigenden Abschnitt der Länge N, einem zweiten, sich direkt an den ersten Abschnitt anschließenden, konstanten Abschnitt der Länge L-N und einem dritten, sich direkt an den zweiten Abschnitt anschließenden, fallenden Abschnitt definiert wird, wobei bei einer Überlagerung des ersten steigenden Abschnittes eines Fensters mit dem dritten fallenden Abschnitt eines anderen Fensters und einer Addition beider Abschnitte im Überlappungsbereich sich das Ergebnis eins ergibt, was dem Wert des zweiten Abschnittes der Fensterfunktion entspricht,
  - aus dem digitalisierten, gespeicherten Sprachsignal in unregelmäßigen Abständen einer mittleren Länge  $\alpha L$  Segmente einer Länge L+N entnommen werden,
  - diese, aus dem digitalisierten, gespeicherten Sprachsignal entnommenen Segmente mit der Fensterfunktion im Zeitbereich gewichtet werden,
  - die gewichteten Segmente jeweils um eine definierte Anzahl von Abtastwerten L versetzt aufaddiert werden, wodurch das so entstehende Sprachsignal für  $\alpha > 1$  verkürzt und für  $\alpha < 1$  verlängert wird,
  - nacheinander an den Stellen der Entnahme der Segmente aus dem digitalisierten Sprachsignal, jeweils das dort entnommene Segment mit dem Segment des verlängerten oder verkürzten Sprachsignals, welches dieses entnommene Segment repräsentiert, verglichen wird, **dadurch gekennzeichnet,**
  - **dass** zum schnellen Vergleich der Abweichung des verlängerten oder verkürzten Sprachsignals vom digitalisierten Sprachsignal lediglich der N Werte lange dritte Abschnitt des zuletzt entnommenen Segmentes als Referenz herangezogen wird,
  - **dass** die entnommenen Segmente so zueinander versetzt aufaddiert werden, daß die ermittelte Abweichung minimal ist und
  - **dass** als Maß für die Abweichung der relative Fehler oder der absolute quadratische Fehler herangezogen wird.

## Claims

1. Method of modifying the speed of voice signals, in particular digitized voice signals, in which
- an analog voice signal is digitized, thereby producing a digitized voice signal that is stored in a memory,
  - a factor  $\alpha$  is defined by which the voice signal is lengthened or shortened,
  - a window function having a first, rising section of length N, a second, constant section of length L+N directly adjoining the first section and a third, falling section directly adjoining the second section is defined, wherein, if the first, rising section of a window overlaps the third, falling section of another window and the two sections are added in the overlap region, the result amounts to one, which corresponds to the value of the second section of the window function,

- segments of a defined length  $L+N$  are extracted from the digitized, stored voice signal at irregular intervals of mean length  $\alpha L$ ,
- said segments extracted from the digitized, stored voice signal are weighted with the window function in the time domain,
- the weighted segments, each offset by a defined number of sample values  $L$ , are added, which shortens the voice signal thus produced for  $\alpha > 1$  and lengthens it for  $\alpha < 1$ ,
- the segment extracted successively at the points of extraction of the segments from the digitized voice signal is compared there with the subsequently extracted segment, likewise weighted with the window function, for similarity,

**characterized**

- **in that**, for the purpose of rapidly comparing the similarity of the segments, only the  $N$ -value-long, third section, weighted with the falling window section, of the segment is compared with each first section, weighted with the rising  $N$ -value-long window section, of the subsequently extracted segment,
- **in that** said segments, are added in an offset manner with respect to one another in such a way that the similarity of the two segment sections becomes maximal,
- **in that**, to calculate the similarity, a correlation is used as a measure thereof.

2. Method according to Claim 1, **characterized in that** the similarity of the two compared segment sections is a maximum if a maximization of the similarity measure is performed in relation to the displacement with respect to one another.

3. Method of modifying the speed of voice signals, in particular digitized voice signals in which

- an analog voice signal is digitized, thereby producing a digitized voice signal that is stored in a memory,
- a factor  $\alpha$  is defined by which the voice signal is lengthened or shortened,
- a window function having a first, rising section of length  $N$ , a second, constant section of length  $L+N$  directly adjoining the first section and a third, falling section directly adjoining the second section is defined, wherein, if the first, rising section of a window overlaps the third, falling section of another window and the two sections are added in the overlap region, the result amounts to one, which corresponds to the value of the second section of the window function,
- segments of a length  $L+N$  are extracted from the digitized, stored voice signal at irregular intervals of mean length  $\alpha L$ ,
- said segments extracted from the digitized, stored voice signal are weighted with the window function in the time domain,
- the weighted segments, each offset by a defined number of sample values  $L$ , are added, which shortens the voice signal thus produced for  $\alpha > 1$  and lengthens it for  $\alpha < 1$ ,
- the segment extracted successively at the points of extraction of the segments from the digitized voice signal is compared there in each case with the segment of the lengthened or shortened voice signal that represents said extracted segment,

**characterized**

- **in that**, for the purpose of rapidly comparing the deviation of the lengthened or shortened voice signal from the digitized voice signal, only the  $N$ -value-long, third section of the segment extracted last is used as reference,
- **in that** the segments extracted are added in an offset manner to one another in such a way that the deviation determined is a minimum, and
- the relative error or the absolute square error is used as a measure of the deviation.

**Revendications**

1. Procédé pour modifier la vitesse de signaux vocaux, notamment de signaux vocaux numérisés, selon lequel

- un signal vocal analogique est numérisé, ce qui fait apparaître un signal vocal numérisé qui est mémorisé dans une mémoire,
- un facteur  $\alpha$  est défini, facteur avec lequel le signal vocal est allongé ou raccourci,

- une fonction fenêtre comportant une première section montante de longueur N, une seconde section constante de longueur L-N, qui se raccorde directement à la première section, une troisième section retombante, qui se raccorde à la seconde section, est définie, auquel cas lors d'une superposition de la première section montante d'une fenêtre avec la troisième section retombante d'une autre fenêtre et lors d'une addition des deux sections dans la zone de chevauchement, on obtient le résultat un, qui correspond à la valeur de la seconde section de la fonction fenêtre,
- $\alpha L$  segments ayant une longueur définie L+N sont prélevés du signal vocal numérisé et mémorisé, à des intervalles irréguliers ayant une longueur moyenne,
- ces segments prélevés du signal numérisé et mémorisé sont pondérés avec la fonction fenêtre dans le domaine temporel,
- les segments pondérés sont additionnés en étant décalés respectivement d'un nombre défini de valeurs d'échantillonnage L, ce qui a pour effet que le signal vocal ainsi obtenu est raccourci pour  $\alpha > 1$  et est allongé pour  $\alpha < 1$ ,
- le segment prélevé dans le signal vocal numérisé, et pondéré avec la fonction fenêtre est comparé, et ce successivement aux emplacements du prélèvement des segments à partir du signal vocal numérisé, au segment prélevé ensuite, également pondéré avec la fonction fenêtre, selon des aspects de similitude,

**caractérisé en ce**

- **que** pour la comparaison rapide de la similitude des segments, seule la troisième section du segment, qui possède une longueur de N valeurs et est pondérée par la section fenêtre retombante du segment est comparée à la section montante d'une longueur de N valeurs, qui est pondérée par la section fenêtre, du segment prélevé ensuite,
- **que** ces segments sont additionnés en étant décalés les uns par rapport aux autres de telle sorte que la similitude des sections de segments devienne maximale, et
- **que** pour le calcul de la similitude, on utilise une corrélation en tant que mesure de cette similitude.

**2. Procédé selon la revendication 1, caractérisé en ce que**

- la similitude des deux sections comparées de segment devient maximale lorsqu'on rend maximum le degré de similitude par rapport au décalage réciproque des segments.

**3. Procédé pour modifier la vitesse de signaux vocaux, notamment de signaux vocaux numérisés, selon lequel un signal vocal analogique est numérisé, ce qui fait apparaître un signal vocal numérisé qui est mémorisé dans une mémoire,**

- un facteur  $\alpha$  est défini, facteur avec lequel le signal vocal est allongé ou raccourci,
- une fonction fenêtre comportant une première section montante de longueur N, une seconde section constante de longueur L-N, qui se raccorde directement à la première section, et une troisième section retombante, qui se raccorde à la seconde section, est définie, auquel cas lors d'une superposition de la première section montante d'une fenêtre avec la troisième section retombante d'une autre fenêtre et lors d'une addition des deux sections dans la zone de chevauchement, on obtient le résultat un, qui correspond à la valeur de la seconde section de la fonction fenêtre,
- $\alpha L$  segments ayant une longueur définie L+N sont prélevés du signal vocal numérisé et mémorisé, à des intervalles irréguliers ayant une longueur moyenne,
- ces segments prélevés du signal numérisé et mémorisé sont pondérés avec la fonction fenêtre dans le domaine temporel,
- les segments pondérés sont additionnés en étant décalés respectivement d'un nombre défini de valeurs d'échantillonnage L, ce qui a pour effet que le signal vocal ainsi obtenu est raccourci pour  $\alpha > 1$  et est allongé pour  $\alpha < 1$ ,
- le segment prélevé dans le signal vocal numérisé, et pondéré avec la fonction fenêtre est comparé, et ce successivement aux emplacements de prélèvement des segments à partir du signal vocal numérisé, au segment du signal vocal allongé ou raccourci, qui représente ce segment prélevé,

**caractérisé en ce**

- **que** pour la comparaison rapide de l'écart entre le signal vocal allongé ou le signal vocal raccourci par rapport au signal vocal numérisé, on utilise comme référence uniquement la troisième section, d'une grandeur de N

## EP 0 865 026 B1

valeurs, du signal prélevé en dernier,

- **qu'**on additionne les segments prélevés d'une manière décalée entre eux de telle sorte que l'écart déterminé est minimum, et
- **qu'**on utilise comme mesure de l'écart l'erreur relative ou l'erreur quadratique absolue.

5

10

15

20

25

30

35

40

45

50

55