

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4963345号
(P4963345)

(45) 発行日 平成24年6月27日(2012.6.27)

(24) 登録日 平成24年4月6日(2012.4.6)

(51) Int.Cl. F I
G 1 O L 13/06 (2006.01) G 1 O L 13/06 2 4 O D

請求項の数 3 (全 13 頁)

<p>(21) 出願番号 特願2004-270307 (P2004-270307)</p> <p>(22) 出願日 平成16年9月16日 (2004.9.16)</p> <p>(65) 公開番号 特開2006-84859 (P2006-84859A)</p> <p>(43) 公開日 平成18年3月30日 (2006.3.30)</p> <p>審査請求日 平成19年8月10日 (2007.8.10)</p> <p>(出願人による申告) 平成16年度独立行政法人情報通信研究機構、研究テーマ「大規模コーパスベース音声対話翻訳技術の研究開発」に関する委託研究、産業活力再生特別措置法第30条の適用を受ける特許出願</p>	<p>(73) 特許権者 393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2</p> <p>(74) 代理人 100099933 弁理士 清水 敏</p> <p>(72) 発明者 津崎 実 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内</p> <p>(72) 発明者 小坂 直敏 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内</p> <p>(72) 発明者 河井 恒 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内</p> <p style="text-align: right;">最終頁に続く</p>
---	---

(54) 【発明の名称】 音声合成方法及び音声合成プログラム

(57) 【特許請求の範囲】

【請求項1】

末尾に第1の音素を有する第1の音声サンプルと、先頭に、前記第1の音素と同じ音素である第2の音素を有する第2の音声サンプルとを、所定の時間期間内において接続合成する音声合成方法であって、

前記第1の音素の駆動波形と前記第2の音素の駆動波形とを接続して得られる音声波形の継続時間を決定するステップを含み、前記継続時間の先頭は、前記第1の音素の駆動波形の開始時刻であり、前記継続時間の末尾は、前記第2の音素の駆動波形の終了時刻であり、前記所定の時間期間は、前記継続時間内に含まれ、

前記音声合成方法はさらに、

前記所定の時間期間内の第1の時刻、および前記所定の時間期間内で当該第1の時刻より遅い第2の時刻により画定される移行区間を決定するステップと、

前記移行期間内の前記第1の時刻から前記第2の時刻までの間の時点であって、前記第1の音声サンプルと前記第2の音声サンプルとの混合割合が所定の関係を充足する時点を決するステップと、

前記時点を決するステップにおいて決定された時点での前記第1の音声サンプルと前記第2の音声サンプルとの駆動波形の位相を整合させるステップと、

前記所定の時間期間の先頭時刻から前記第1の時刻までの区間の合成音声を前記第1の音声サンプルから生成するステップと、

前記第1の音声サンプルから前記第2の音声サンプルへと、前記第1の時刻から前記第

2の時刻までの間の、時間に対する所定の滑らかな関数にしたがって両者の混合割合を変化させて混合することにより、前記移行区間における合成音声を生成するステップとを含み、

前記時点は、前記第1の音素の瞬時音圧値と前記第2の音素の瞬時音圧値との混合割合が実質的に等しくなる時間位置である、音声合成方法。

【請求項2】

前記合成音声を生成するステップは、前記第1の時刻の前記第1の音声サンプルの駆動波形から、前記第2の時刻の前記第2の音声サンプルへと、前記第1の音素の駆動波形の瞬時音圧値と前記第2の音素の駆動波形の瞬時音圧値とを、前記移行期間に対する前記第1の時刻からの経過時刻の重みで加重平均することによって合成音声を生成するステップを含む、請求項1に記載の音声合成方法。

10

【請求項3】

コンピュータにより実行されると、請求項1又は請求項2に記載の音声合成方法を実行するよう当該コンピュータを制御する、コンピュータで実行可能な音声合成プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、例えばTTS(Text-To-Speech)システムなどに用いられ、音声コーパスから抽出された音声サンプルを互いに接続して音声合成を行なう音声合成方法及び音声合成プログラムに関する。

20

【背景技術】

【0002】

音声合成技術の中に、素片接続型(または波形接続型)音声合成と呼ばれるものがある。この技術では、実際の音声(特定の話者であることが多い。)を収集して音声コーパスを作成する。音声コーパス中の音声データを音素、ダイフオンなどの所定の単位(波形セグメント)に分ける。各波形セグメントには、対応する音声データの音響・韻律に関する特徴量を示す情報が付されている。

【0003】

音声合成時には、テキストが与えられると、このテキストに対し形態素解析などのテキスト処理を実行し、テキストを音声化した際の各音素などの所定単位ごとに合成目標の音韻・韻律情報および言語情報を生成する。この合成目標にできるだけ合致する音響的な特徴量を有し、かつ互いに接続したときに自然な音声となるような音声サンプルを音声コーパスから抽出する。そして、抽出された音声サンプルを互いに接続することで合成音声波形を生成する。

30

【0004】

素片接続型音声合成の基本的な考え方は、次の二つに大別できる。第1は発話単位での接続であり、第2はダイフオン単位での接続である。

【0005】

発話単位での接続は、発話の基本単位(例えば日本語の場合は、子音と母音との連鎖であるいわゆるCV単位)を互いに接続する方法である。この方法は、自然な発話でも基本単位の間には音響的な不連続性が生じているのであるから、別々の音声サンプルを発話単位でつないで人工的な音響特性の不連続性が生じても、人間の聴覚にとっては許容可能なものであると想定している。

40

【0006】

ただし、有声の子音、半母音などは音響的に前後の音と切離されている度合いが低くなるし、また母音の連鎖ということも音声合成上では生じ得る。すなわち、実際には前後の音と連続している音声であるにもかかわらず、前後と切離して得た音声素片を音声合成時に用いることがある。その結果、自然音声でも発話単位の間連続的な遷移部が存在すると、素片接続により音声合成した結果が不自然になる可能性が高い。

【0007】

50

一方、ダイフオン単位での接続合成方法は、母音を代表とする変化の比較的緩やかな区間の中間で接続する手法である。現在信頼性の高い音響的特徴のほとんどが静的な特性を捉えたものであり、変化が少ない部分では静的な特性でマッチした区間同士を中間でつなげば、物理的に非常に少ない不連続の範囲でつなぎやすいという利点がある。

【非特許文献1】H.カワハラ他、「加重平均群遅延に対する定点法に基づく正確な音声事象検出法」、ICSLP-2000予稿集、北京、pp.664-667、2000年 (Hideki Kawahara et al., "Accurate vocal event detection method based on a fixed-point to weighted average group delay", ICSLP-2000, pp. 64-667, Beijing, 2000)

【発明の開示】

10

【発明が解決しようとする課題】

【0008】

しかし、このダイフオン単位接続法のように定常部で接続する方法では、パワースペクトル、基本周波数F0、波形としての連続性、などのすべての物理的側面でまったく等価な素片を常に保守することは不可能であり、通常はごく微小ではあるが物理的不連続の発生が不可避だという問題がある。しかも、人間の聴覚系は定常的状态に生じた僅かな差分に対しては非常に感度が鋭いのが現実であり、ごく微小な物理的不連続であっても、これを敏感に知覚してしまうという問題がある。

【0009】

したがって、本発明の一つの目的は、接続部における物理的不連続の発生を抑制して知覚的な自然性を高めることができる音声合成方法を提供することである。

20

【0010】

本発明の他の目的は、接続部において音響特徴が連続的に変化するように波形セグメントを接続し、知覚的な自然性を高めることができる音声合成方法を提供することである。

【課題を解決するための手段】

【0011】

本発明の第1の局面に係る音声合成方法は、末尾に第1の音素を有する第1の音声サンプルと、先頭に第2の音素を有する第2の音声サンプルとを、所定の時間期間内において接続合成する音声合成方法であって、所定の時間期間内の第1の時刻、および所定の時間期間内で当該第1の時刻より遅い第2の時刻により画定される移行区間を決定するステップと、所定の時間期間の先頭時刻から第1の時刻までの区間の合成音声を第1の音声サンプルから生成するステップと、第1の音声サンプルから第2の音声サンプルへと、時間に対する所定の滑らかな関数にしたがって両者の混合割合を変化させて混合することにより、移行区間における合成音声を生成するステップとを含む。

30

【0012】

移行区間において、第1の音声サンプルから第2の音声サンプルへと、両者の混合割合を滑らかな関数にしたがって変化させて混合させる。音声の不連続が発生することが避けられ、より自然な合成音声を生成できる。合成のための所定の時間期間の先頭から第1の時刻までの間は第1の音声サンプルを合成音声とする。したがってその直前の音素との連続が保たれ、自然な音声が合成できる。

40

【0013】

さらに好ましくは、合成音声を生成するステップは、第1の音声サンプルから第2の音声サンプルへと、第1の音素の瞬時音圧値と第2の音素の瞬時音圧値とを、移行期間に対する第1の時刻からの経過時刻の重みで加重平均することによって合成音声を生成するステップを含む。

【0014】

このような加重平均により第1の音声サンプルと第2の音声サンプルとを混合することで、両者の間の変化が直線的になり、簡単な処理で合成音声の接続部分を自然なものにすることができる。

【0015】

50

より好ましくは、合成音声を生成するステップは、移行期間内であって、第1の音声サンプルと第2の音声サンプルとの混合割合が所定の関係を充足する時点を決するステップと、このステップにおいて決定された時点での第1の音声と第2の音声との駆動波形の位相を整合させて、第1の音素の瞬時音圧値と第2の音素の瞬時音圧値とを重みで加重平均することによって合成音声を生成するステップを含むようにしてもよい。

【0016】

両者の駆動波形の位相を整合させることにより、合成後の音声波形は合成前の音声の特徴とよく一致する特徴を示し、接続部分がより自然なものとなる。

【0017】

好ましくは、合成音声を生成するステップは、移行期間内における第1の音声サンプルと第2の音声サンプルとの駆動波形の位相のずれの重み付きの和が最小化するように第1の音声サンプルと第2の音声サンプルとの駆動波形の位相を整合させて、第1の音素の瞬時音圧値と第2の音素の瞬時音圧値とを加重平均することによって合成音声を生成するステップを含む。ずれの和を計算する際の重みは、第1の音素の瞬時音圧値と第2の音素の瞬時音圧値との混合割合が実質的に等しくなる時間位置で最大となり、混合割合が前記時間位置から遠ざかるにつれて減少するように選択される。

【0018】

第1の音声サンプルと第2の音声サンプルとは、互いに異なる音声データから得られたものであることが通常である。したがって両者の駆動波形の周期が多少異なる場合があり、両者の駆動波形を移行期間の全体にわたって一致させることはできない。そこで、駆動波形のずれの和を最小化するように両者の位相を整合させるのが合理的である。ただし、この場合に各時点でのずれを対等に扱うのではなく、両者の混合割合が実質的に等しくなる時間位置で最大となり、そこから遠ざかるにつれて減少するような重みを用いると、音声に対する悪影響を小さくとどめることができ、自然な合成音声を生成できる。

【0019】

本発明の第2の局面に係るコンピュータプログラムは、コンピュータにより実行されると、上記したいずれかの音声合成方法を実行するよう当該コンピュータを制御するものである。

【発明を実施するための最良の形態】

【0020】

以下に述べる本発明の実施の形態は、コンピュータおよびコンピュータ上で動作するソフトウェアにより実現される。もちろん、以下に述べる機能の一部又は全部を、ソフトウェアでなくハードウェアで実現することも可能である。

【0021】

図1に、本発明の実施の形態で利用されるコンピュータシステム20の外観図を、図2にコンピュータシステム20のブロック図を、それぞれ示す。なおここに示すコンピュータシステム20はあくまで一例であり、この他にも種々の構成が可能である。

【0022】

図1を参照して、コンピュータシステム20は、コンピュータ40と、いずれもこのコンピュータ40に接続されたモニタ42、キーボード46、およびマウス48を含む。コンピュータ40にはさらに、CD-ROM(Compact Disk Read-Only Memory)ドライブ50と、FD(Flexible Disk)ドライブ52とが内蔵されている。

【0023】

図2を参照して、コンピュータシステム20はさらに、コンピュータ40に接続されるプリンタ44を含むが、これは図1には示していない。またコンピュータ40はさらに、CD-ROMドライブ50およびFDドライブ52に接続されたバス66と、いずれもバス66に接続された中央演算装置(Central Processing Unit: CPU)56、コンピュータ40のブートアッププログラムなどを記憶したROM(Read-Only Memory)58、CPU56が使用する作業エリアおよびCPU5

10

20

30

40

50

6により実行されるプログラムの格納エリアを提供するRAM(Random Access Memory)60、および後述する音声データベースを格納したハードディスク54を含む。

【0024】

以下に述べる実施の形態のシステムを実現するソフトウェアは、たとえば、CD-ROM62のような記録媒体上に記録されて流通し、CD-ROMドライブ50のような読取装置を介してコンピュータ40に読込まれ、ハードディスク54に格納される。CPU56がこのプログラムを実行する際には、ハードディスク54からこのプログラムを読み出してRAM60に格納し、図示しないプログラムカウンタによって指定されるアドレスから命令を読み出して実行する。CPU56は、処理対象のデータをハードディスク54から読出し、処理結果を同じくハードディスク54に格納する。

10

【0025】

コンピュータシステム20の動作自体は周知であるので、ここではその詳細については繰返さない。

【0026】

なお、ソフトウェアの流通形態は上記したように記憶媒体に固定された形には限定されない。たとえば、ネットワークを通じて接続された他のコンピュータからデータを受取る形で流通することもあり得る。また、ソフトウェアの一部が予めハードディスク54中に格納されており、ソフトウェアの残りの部分をネットワーク経由でハードディスク54に取込んで実行時に統合するような形の流通形態もあり得る。

20

【0027】

一般的に、現代のプログラムはコンピュータのオペレーティングシステム(OS)によって提供される汎用の機能を利用し、それらを所望の目的にしたがって組織化した形態で実行することにより前記した所望の目的を達成する。したがって、以下に述べる本実施の形態の各機能のうち、OSまたはサードパーティが提供する汎用的な機能を含まず、それら汎用的な機能の実行順序の組合せだけを指定するプログラム(群)であっても、それらを利用して全体的として所望の目的を達成する制御構造を有するプログラム(群)である限り、それらが本発明の技術的範囲に含まれることは明らかである。

【0028】

[原理]

本実施の形態のプログラムによってコンピュータ40が実行する音声合成方法の原理を説明する。なお、以下の説明では、音素P1、音素V2、音素P3、音素V2a、音素V2b、音素Pa、音素Pbをそれぞれ単にP1、V2、P3、V2a、V2b、Pa、Pbという。またV2aとV2bは、同じ音素V2に対応する、互いに異なる実サンプル中の音素(波形)であるものとする。

30

【0029】

例として3つの母音からなる音素の連鎖を接続合成によって実現する場合を考える。図3(a)に示すように、実現したい音素の連鎖がP1-V2-P3であるとする。また、音声コーパス中に存在する、P1とV2の連鎖P-V2aからなる実サンプル(第1の音声サンプル)101(図3(b)に示す)と、音素V2とP3の連鎖V2b-P3からなる実サンプル(第2の音声サンプル)102(図3(c)に示す)を接続する場合を想定する。なお、図3(b)のPaは、音声コーパス中で第1の音声サンプル101のV2aの直後に存在している音素であり、図3(c)に示すPbは、音声コーパス中で第2の音声サンプル102の音素V2bの直前に存在している音素であるものとする。

40

【0030】

従来法によるダイフオン接続の場合は、図3(d)に示すように、V2aとV2bの中間点110で接続をする。V2aもV2bも合成目標であるV2に合致する音響特徴量を有しており、実質的にその継続長もほぼ同じである。したがってこのように互いの中間点でV2aとV2bとを接続することで、V2と一致する継続長を持つ音素V2a/V2bが得られる。

50

【 0 0 3 1 】

しかしこの場合、接続点である中間点 1 1 0 を境に音響的特徴が微妙に変化し、聴覚的な不自然さを感じる。両者の差が偶然小さければ聴覚的にも滑らかな接続が達成できることになるが、それが必ず保証されているわけではない。

【 0 0 3 2 】

そこで、この実施の形態では、V 2 a と V 2 b を中間点 1 1 0 でいきなり繋ぐのではなく、以下のようにする。図 3 (e) を参照して、V 2 の継続時間 1 4 0 内で第 1 の時点 1 4 2 と、V 2 の継続時間 1 4 0 内でかつ第 1 の時点 1 4 2 より後の第 2 の時点 1 4 4 とを決定する。この二つの時点 1 4 2 および 1 4 4 により移行期間 1 0 3 が画定される。V 2 の継続時間 1 4 0 内で移行期間 1 0 3 の前の期間 1 0 4 を V 2 a の保存期間、後の期間 1 0 5 を V 2 b の保存期間 1 0 5 とする。すなわち、V 2 の継続時間 1 4 0 を V 2 a の保存期間 1 0 4、V 2 b の保存期間 1 0 5、および V 2 a から V 2 b への移行期間 1 0 3 に分ける。そして、この移行期間 1 0 3 で V 2 a 的な音から V 2 b 的な音へと音響的特徴を徐々に移行させることによって知覚的に有害な急激な変化を排除する。なお、図 3 (e) における記号「V 2 a b」は、V 2 の継続期間において V 2 a と V 2 b が混合されて接続されていることを表している。

10

【 0 0 3 3 】

本実施の形態では、第 1 の時点 1 4 2 および第 2 の時点 1 4 4 については、継続時間 1 4 0 の全体に対しどのような割合の時点とするかを予め決めておき、合成目標が与えられ継続時間 1 4 0 が定まったところで決定する。

20

【 0 0 3 4 】

V 2 a から V 2 b への音声の変化の様子を時間の関数として示したのが図 3 (f) である。図 3 を参照して、実線 1 3 0 で示したように、本実施の形態では V 2 a の保存期間 1 0 4 では V 2 a が 1 0 0 %、V 2 b の保存期間 1 0 5 では V 2 b が 1 0 0 %、移行期間 1 0 3 では V 2 a から V 2 b まで直線的に変化する割合で両者を混合する。

【 0 0 3 5 】

これに対し図 3 (d) に示した従来のダイフオン接続に対応して時間の関数で表せば、そのグラフは図 3 (f) において破線 1 2 0 で示すようになる。破線 1 2 0 で示すように、接続点である中間点 1 1 0 を境に音素が V 2 a から V 2 b に完全に変化し、その結果、音響的特徴もこの時点で変化する。そのため生じる音響的特徴の不連続により、聴覚的な不自然さを感じることもある。

30

【 0 0 3 6 】

本実施の形態において目標とする効果を上げるためには、図 3 (e) 及び (f) に示す移行期間 1 0 3 を、実現したい V 2 の継続時間 1 4 0 を上限として、その中で十分に長く取る必要がある。最大では V 2 の継続時間 1 4 0 の全体を移行期間とし、V 2 a、V 2 b の保存期間 1 0 4、1 0 5 の長さをゼロとすることもできる。

【 0 0 3 7 】

しかし、V 2 a の開始時点は P 1 からの調音結合による影響が残っている可能性が高い。この部分で V 2 b を混入させ始めると、その部分には V 2 b の抽出前の環境である P b からの影響も混入することとなり、P 1 から V 2 a への連続性に悪影響を及ぼすとともに、V 2 としての音韻性が低下する可能性が考えられる。したがって一般的には、移行期間 1 0 3 の直前に V 2 a の保存期間 1 0 4 をある程度の長さ設けた方が望ましい。

40

【 0 0 3 8 】

一方、V 2 b の後半は P 3 に対する遷移部と見なすことができ、自然音声であっても遷移部として音韻性が曖昧になる部分である。この部分に多少他の環境からの混入が生じたとしても、V 2 としての知覚は既に確立されており、P 3 への滑らかな移行が生じている限り、知覚的にはそれほど害とされないと考えられる。したがって、本実施の形態では、V 2 a の保存期間 1 0 4 を V 2 の継続時間 1 4 0 の半分の長さとし、残りの半分の移行期間 1 0 3 とする。この場合、V 2 b の保存期間 1 0 5 の長さはゼロとなる。

【 0 0 3 9 】

50

但し、本方法はTTSシステムの中で使われる可能性が高いので、P1、V2、P3、Pa、Pbに関する情報を事前に持っている想定できる。V2a、V2bの保存期間、移行期間の最適な割合についてはそれらの環境の組合せに依存する可能性が高い。その場合は、抽出環境と接続の組合せ毎に適切な形で混合することも可能である。

【0040】

以下、接続合成方法の具体的な内容について説明する。

【0041】

[第1の実施の形態(駆動周期同期型のクロスフェード法)]

この実施の形態では、移行期間103において、V2aとV2bとの瞬時音圧値を移行期間103の先頭からの経過時間に応じた重みで加重平均することによってクロスフェードする。

10

【0042】

この方法においてさらに品質向上が必要な場合は、次のような方策が可能となる。一般的にクロスフェードの難点として途中で二つの音が混ざる場合に、二つの音の中間的なひとつの音が聞こえて欲しいにもかかわらず、単純に二つの音が混ざって同時に聞こえてしまうということがある。この欠点を最小化するために、カワハラ(非特許文献1)によって提案された事象検出アルゴリズムにより駆動時点のマーキングを取り、2音間で位相のずれを最小化した形でクロスフェードをかける。または、音声コーパスの各サンプルに、駆動波形のピークを示すピッチマークを付しておき、そうしたマークを利用して位相のずれを最小化するようにしてもよい。

20

【0043】

図4は、位相のずれをまったく考慮しない場合に生じうる問題を示したものである。図4(a)の波形201が母音V2aの波形であり、図4(b)の波形202が母音V2bの波形であるものとする。符号211、符号212で示す矢印が、それぞれV2a、V2bのピッチマークである。V2a、V2b間の位相がずれたままで混合した結果が図4(c)に示されている。波形203上から駆動周期を観察することが困難になることが分かる。このような場合、聴覚系は二つの母音が平行して存在しているという知覚像を持ちやすい。

【0044】

図5はこれに対してV2a、V2b間で位相を整合した場合を示す。この処理では、ピッチマーク211、212を参照して二つの波形の位相を整合させる。図5(a)(b)はそれぞれ混合前のV2a、V2bの波形201、202であり、図5(c)は混合後の波形204を示す。混合前のそれぞれの母音の波形とよく似た周期構造が保存されていることが分かる。したがって、このように二つの波形の位相を整合させることにより、二つの母音が同時に聞こえる不都合を回避できる。

30

【0045】

現実的には2母音間で位相(駆動周期)が完全に一致することは期待できないため、ある時点でのずれをなくすと別の時点ではずれが生じざるを得ない。但し、本方法では混合の期間は短時間であるので、混合の割合が釣合う時点でのマーカのずれが生じた場合のペナルティーを重く評価することが合理的な位相整合の取り方となる。したがって本実施の形態では、混合の割合がV2a、V2bそれぞれ50%となる時点で両者の位相を整合させるようにする。

40

【0046】

第1の実施の形態の方法を図6にまとめる。これはコンピュータにより実現するときのプログラムの制御構造を示すものでもある。まず、移行期間103の開始時点および終了時点を計算により設定し、第1の音声サンプル101と第2の音声サンプル102の接続対象音素V2aとV2bの駆動波形の位相を整合させる(ステップ300)。次に、V2aとV2bの瞬時音圧値を移行期間103に対する移行期間の先頭からの経過時間を重みとして加重平均して移行期間103における混合波形を生成し(ステップ301)、第1の音声サンプル101と第2の音声サンプル102を接続する(ステップ302)。

50

【 0 0 4 7 】

この実施の形態によれば、音声波形の瞬時音圧値の加重平均という比較的簡単な処理で二つの音素をその中間で結合し、一つの音素（特にダイフオン接続における母音部）を生成できる。その結合個所では、合成波形は第1の実サンプルの波形から第2の実サンプルの波形に滑らかに変化する。したがって人間が音響的な特徴の不連続を知覚する可能性が低くなる。さらに、二つの波形の位相を整合させることにより、接続後の音声がより自然なものとなるという効果がある。

【 0 0 4 8 】

上記した実施の形態では、移行期間では、単純に加重平均をとることで第1の波形から第2の波形に滑らかに、かつ直線的に波形をクロスフェードしている。しかし本発明はそのような実施の形態には限定されない。例えば、時間に関して2次以上の関数で、かつ移行期間の両端でそれぞれ $V2a$ および $V2b$ となるという境界条件を満足するような関数によって、移行期間における両者の混合割合を決定するようにしてもよい。この場合、この関数の値がある時間における $V2a$ の混合割合を表すものとするれば、関数の値が時間に対して単調減少となることが好ましい。

【 0 0 4 9 】

[第2の実施の形態（正弦波モデルによるモルフィング法）]

第2の実施の形態にかかる音声合成方法は、二つの音素 $V2a$ と $V2b$ の混合を単純な時間波形レベルではなく、正弦波成分に分解した後に行なう方法である。母音音声に代表されるような音響信号は、振幅・周波数の異なる複数の正弦振動の加算として表現可能である。この方法では、混合する基となる二つの音声信号 $V2a$ と $V2b$ をフーリエ変換によりそれぞれ複数の正弦波成分に分解し、その間の対応付けを取った後、 $V2a$ と $V2b$ の間で各成分の周波数、振幅項が連続的な変化を生じるように変化させることによって、中間的な音を実現する。

【 0 0 5 0 】

例えば $V2a$ 、 $V2b$ の間に基本周波数のずれがあるような場合、前述の第1の実施の形態では、中間部に2種類の基本周波数成分が出現することが避けられない。聴覚系はそのような場合には二つの音を知覚しがちである。しかしこの第2の実施の形態に係る方法では、基本周波数成分は基本周波数成分として連続的に変化を起こす。したがって混合部が2音に分離して聞こえる印象を回避することが可能となる。

【 0 0 5 1 】

第2の実施の形態の方法を図7で説明する。この方法をコンピュータプログラムで実現する場合、そのためのプログラムの制御構造もこの様な形となる。第1の音声サンプル101の接続対象音素 $V2a$ を、 n 個の正弦波の成分に分解する（ステップ400）。同様に、第2の音声サンプル102の接続対象音素 $V2b$ を、 n 個の正弦波の成分に分解する（ステップ401）。なお n は予め定められた整数であるものとする。

【 0 0 5 2 】

次に、分解した正弦波成分の対応するもの同士を混合するに際して、両波形がデジタルデータであるため、一方の波形のある成分と他方の波形のどの成分とを組合せて加重平均するかを決定する（ステップ402）。

【 0 0 5 3 】

次に、移行期間を設定し、対応する n 個の正弦波同士の、対応を決定した点により表される振幅を、移行期間103の重みで加重平均することにより混合した後（ステップ403）、 n 個の波形をフーリエ逆変換により合成して移行期間103における混合波形を生成し（ステップ404）、次いで第1の音声サンプル101と第2の音声サンプル102を接続する（ステップ405）。

【 0 0 5 4 】

なおこの実施の形態では、処理を簡単にするために n を固定するものとしている。しかし本発明はそのような実施の形態には限定されず、フレームごとに n が変化するような方

10

20

30

40

50

法も可能である。

【 0 0 5 5 】

[第 3 の実施の形態 (音響ボコーダ・モデルによるモルフィング法)]

第 2 の実施の形態が音声信号を正弦波成分に分解するのに対して、この実施例はソースフィルター原理に基づいて音声信号を駆動源情報成分と共振特性成分とに分離し、それぞれの次元での連続的な変化を実現した後、それらからボコーダにより音声を合成する。

【 0 0 5 6 】

母音信号に代表されるような音は、駆動音源情報の成分とその伝達系の共振特性の成分へと分解可能である。前者は声帯振動の周期性によってフーリエスペクトル上には基本周波数とその高調波の位置として主に反映される。後者はそれらの高調波成分の包絡を決定するもので、主に声道の形状の変化によって人間の発声の場合は実現される。したがって、声道形状を変化させずに声の高さだけを変えたり、反対に声の高さを変えずに声道形状だけを変えたりといった独立の制御が原理的には可能である。

【 0 0 5 7 】

第 3 の実施の形態では、このように分解した音源情報成分と共振特性成分という独立の次元の変数上で V 2 a、V 2 b の混合を実施する。第 2 の実施の形態に比べてより人間の音声生成器官で生じていることに近いたため、不自然な混合を起こす可能性の低減が見込まれる。

【 0 0 5 8 】

第 3 の実施の形態の方法を図 8 で説明する。コンピュータにより本実施の形態にかかる方法を実現する場合、そのためのプログラムの制御構造は図 8 に示されるようなものとなる。図 8 を参照して、第 1 の音声サンプル 1 0 1 の接続対象音素 V 2 a を、音源情報成分と共振特性成分とに分解する (ステップ 5 0 0)。同様に、第 2 の音声サンプル 1 0 2 の接続対象音素 V 2 b を、音源情報成分と共振特性成分とに分解する (ステップ 5 0 1)。

【 0 0 5 9 】

次に、移行期間を設定し、音源情報成分同士の振幅を、移行期間 1 0 3 (図 3 (e) および (f) 参照) に対する移行期間 1 0 3 の先頭からの経過時間の重みで加重平均することにより混合し、1 個の音源情報成分を生成する (ステップ 5 0 4)。具体的には、標本周波数値のパラメータ上での加重平均をとる。同様に、共振特性成分同士の振幅を、移行期間 1 0 3 内での経過時間の重みで加重平均することにより混合し、1 個の共振特性成分を生成する (ステップ 5 0 5)。具体的には、共振特性のスペクトルのパラメータ上での加重平均をとる。ただしこの共振特性の加重平均の場合、フォルマント間の対応をとって周波数軸の不均等圧縮・伸長を行なう。

【 0 0 6 0 】

次に、生成した音源情報成分と共振特性成分とを音響ボコーダにより合成して、移行期間 1 0 3 における混合波形を生成し (ステップ 5 0 6)、次いで第 1 の音声サンプル 1 0 1 と第 2 の音声サンプル 1 0 2 とを接続する (ステップ 5 0 7)。

【 0 0 6 1 】

[変形例]

上記した実施の形態では、移行期間では、単純に加重平均をとることで第 1 の波形から第 2 の波形に滑らかに直線的に波形をクロスフェードしている。しかし本発明はそのような実施の形態には限定されない。例えば、時間に関して 2 次以上の関数で、かつ移行期間の両端でそれぞれ V 2 a および V 2 b となるという境界条件を満足するような関数によって、移行期間における両者の混合割合を決定するようにしてもよい。この場合、この関数の値がある時間における V 2 a の混合割合を表すものとすれば、関数の値が時間に対して単調減少となることが好ましい。

【 0 0 6 2 】

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味および範囲内

10

20

30

40

50

でのすべての変更を含む。

【図面の簡単な説明】

【0063】

【図1】この発明の一実施の形態の音声合成プログラムを実行するコンピュータシステムの外観図である。

【図2】図1のコンピュータシステムのブロック図である。

【図3】この発明の一実施の形態の音声合成方法の原理を説明するための模式図である。

【図4】位相を整合させることなく音素を混合する場合の波形図である。

【図5】位相を整合させて音素を混合する場合の波形図である。

【図6】第1の実施の形態の音声合成方法を説明するためのフローチャートである。

【図7】第2の実施の形態の音声合成方法を説明するためのフローチャートである。

【図8】第3の実施の形態の音声合成方法を説明するためのフローチャートである。

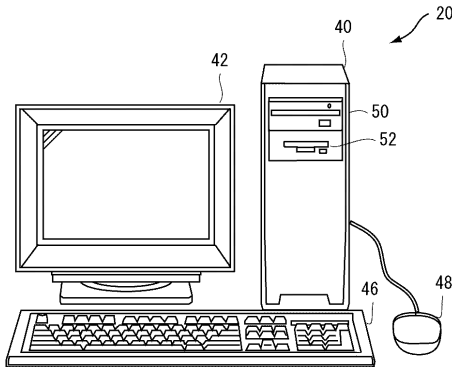
【符号の説明】

【0064】

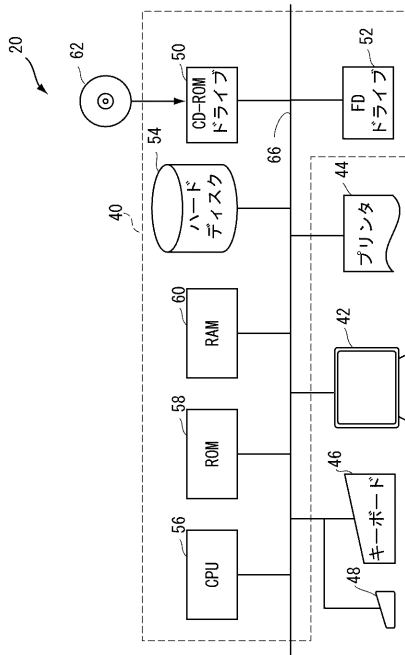
40 コンピュータ、101 第1の音声サンプル、102 第2の音声サンプル、103 移行期間、104 第1の音素の保存期間、105 第2の音素の保存期間、V2a 第1の音素、V2b 第2の音素

10

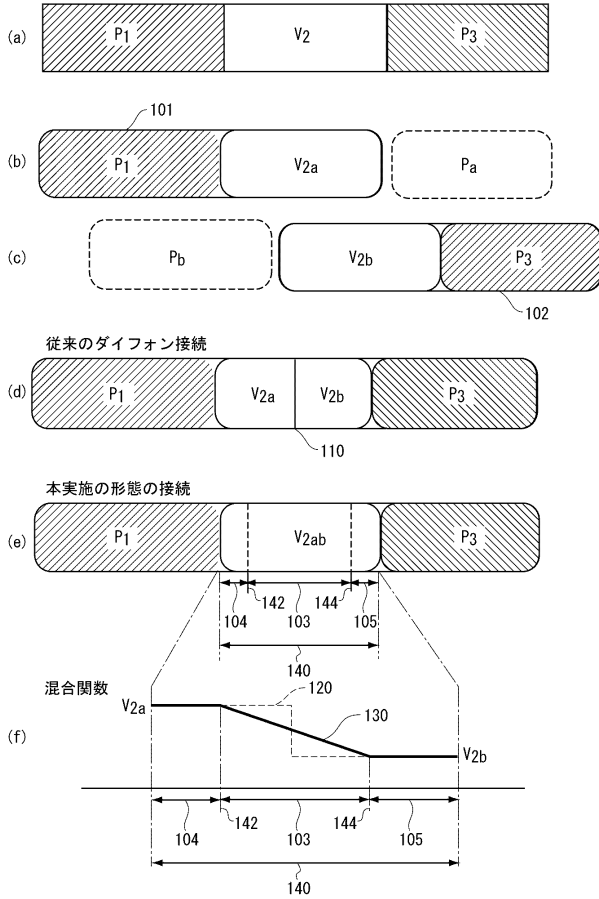
【図1】



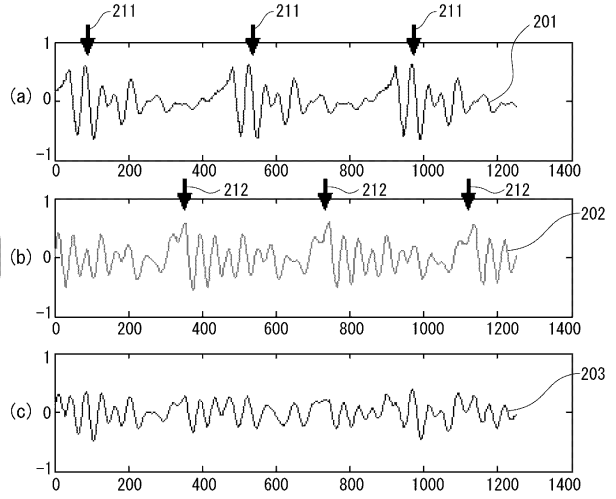
【図2】



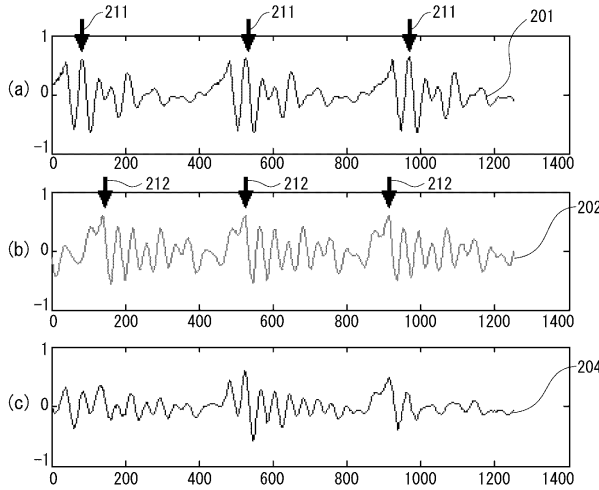
【図3】



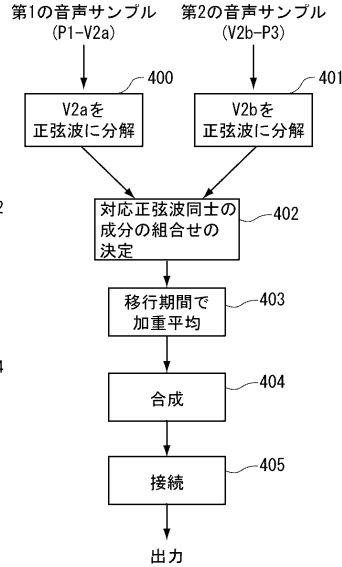
【図4】



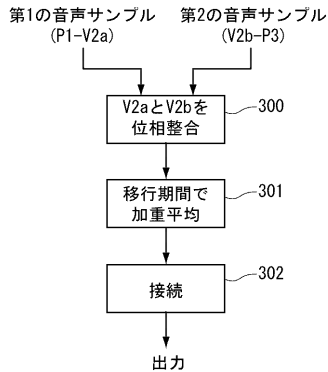
【図5】



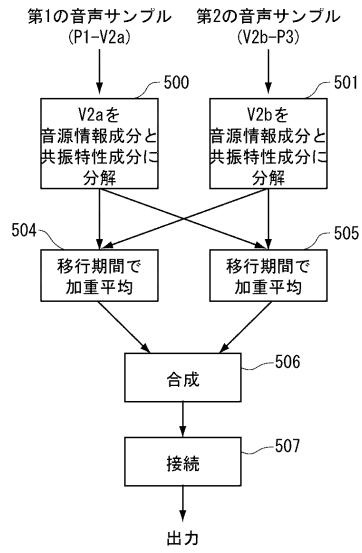
【図7】



【図6】



【図8】



フロントページの続き

審査官 毛利 太郎

- (56)参考文献 特開平11-038989(JP,A)
特開平10-124082(JP,A)
特開平11-194796(JP,A)
特開平05-297891(JP,A)
特開2004-102118(JP,A)
特開平11-224096(JP,A)
特表2005-523478(JP,A)
特開2004-258660(JP,A)
Eric MOULINES and Francis CHARPENTIER, Pitch-Synchronous Waveform Processing Techniques for Text-to-Speech Synthesis using Diphones, Speech Communication, EP, ISCA, 1990年, Vol.9, p.453-467
- (58)調査した分野(Int.Cl., DB名)
G10L 13/00-13/08