

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3588302号
(P3588302)

(45) 発行日 平成16年11月10日(2004.11.10)

(24) 登録日 平成16年8月20日(2004.8.20)

(51) Int. Cl.⁷

F I

G 1 0 L 13/06

G 1 0 L 5/04

F

G 1 0 L 13/08

G 1 0 L 3/00

H

G 1 0 L 5/04

D

請求項の数 13 (全 11 頁)

(21) 出願番号	特願2000-65106 (P2000-65106)	(73) 特許権者	000005821
(22) 出願日	平成12年3月9日(2000.3.9)		松下電器産業株式会社
(65) 公開番号	特開2000-310997 (P2000-310997A)		大阪府門真市大字門真1006番地
(43) 公開日	平成12年11月7日(2000.11.7)	(74) 代理人	100062144
審査請求日	平成12年6月7日(2000.6.7)		弁理士 青山 稜
(31) 優先権主張番号	09/264981	(74) 代理人	100086405
(32) 優先日	平成11年3月9日(1999.3.9)		弁理士 河宮 治
(33) 優先権主張国	米国 (US)	(74) 代理人	100098280
			弁理士 石野 正弘
		(72) 発明者	ニコラス・キブレ
			アメリカ合衆国93436カリフォルニア
			州ロンボック、ウエスト・パイン401番
			、ナンバー31

最終頁に続く

(54) 【発明の名称】 連結型音声合成のための単位重複領域の識別方法および連結型音声合成方法

(57) 【特許請求の範囲】

【請求項1】

音声の時変特性を表す統計モデルを画定するステップと、
 同じ母音を含む異なる音声単位に対応する複数の時系列データを提供するステップと、
 前記時系列データから音声信号パラメータを抽出し、前記音声信号パラメータを用いて前記統計モデルを学習するステップと、
 学習させた前記統計モデルを用いて前記時系列データ内の繰り返しシーケンスを識別し、
 前記繰り返しシーケンスを前記母音の中心の核をなす状態遷移部と関連付けるステップと、
 前記繰り返しシーケンスを用いて、前記音声単位の少なくとも1つに対する連結型音声合成のための単位重複領域を定めるステップとを含み、前記単位重複領域は、前記繰り返しシーケンスの直前の時系列データ又は直後の時系列データであることを特徴とする、連結型音声合成のための単位重複領域の識別方法。

10

【請求項2】

前記統計モデルは隠れマルコフモデルである、請求項1に記載の方法。

【請求項3】

前記統計モデルはリカレントニューラルネットワークである、請求項1に記載の方法。

【請求項4】

前記音声信号パラメータは音声フォルマントを含む、請求項1に記載の方法。

【請求項5】

20

前記統計モデルは、前記母音の中心の核をなす状態遷移部と、前記中心の核をなす状態遷移部の周囲の遷移部とを別々にモデル化するデータ構造を有する、請求項 1 に記載の方法。

【請求項 6】

前記統計モデルは、前記母音の中心の核をなす状態遷移部と、前記中心の核をなす状態遷移部に先行する第 1 の遷移部と、前記中心の核をなす状態遷移部に後続する第 2 の遷移部とを別々にモデル化するデータ構造を有し、

前記データ構造を用いて、前記第 1 の遷移部および前記第 2 の遷移部の 1 つに対応する前記時系列データの 1 部分を破棄するステップを含む、請求項 1 に記載の方法。

【請求項 7】

音声の時変特性を表す統計モデルを画定するステップと、
同じ母音を含む異なる音声単位に対応する複数の時系列データを提供するステップと、
前記時系列データから音声信号パラメータを抽出し、前記音声信号パラメータを用いて前記統計モデルを学習するステップと、

学習させた前記統計モデルを用いて前記時系列データ内の繰り返しシーケンスを識別し、
前記繰り返しシーケンスを前記母音の中心の核をなす状態遷移部と関連付けるステップと、

前記繰り返しシーケンスを用いて、連結型音声合成のための単位重複領域を定めるステップとを含み、
前記単位重複領域は、前記繰り返しシーケンスの直前の時系列データ又は直後の時系列データであり、

前記音声単位の各単位重複領域に基づいて、2つの異なる前記音声単位からの前記時系列データを重複させ、
マージすることにより、新たな音声単位を連結して合成するステップとを含むことを特徴とする、
連結型音声合成方法。

【請求項 8】

前記合成するステップを行う前に、前記単位重複領域の少なくとも 1 つの継続時間を選択的に変化させて、
前記単位重複領域の他方の継続時間に一致させるステップをさらに含む、
請求項 7 に記載の方法。

【請求項 9】

前記統計モデルは隠れマルコフモデルである、請求項 7 に記載の方法。

【請求項 10】

前記統計モデルはリカレントニューラルネットワークである、請求項 7 に記載の方法。

【請求項 11】

前記音声信号パラメータは音声フォルマントを含む、請求項 7 に記載の方法。

【請求項 12】

前記統計モデルは、前記母音の中心の核をなす状態遷移部と、前記中心の核をなす状態遷移部の周囲の遷移部とを別々にモデル化するデータ構造を有する、
請求項 7 に記載の方法。

【請求項 13】

前記統計モデルは、前記母音の中心の核をなす状態遷移部と、前記中心の核をなす状態遷移部に先行する第 1 の遷移部と、
前記中心の核をなす状態遷移部に後続する第 2 の遷移部とを別々にモデル化するデータ構造を有し、

前記データ構造を用いて、前記第 1 の遷移部および前記第 2 の遷移部の 1 つに対応する前記時系列データの 1 部分を破棄するステップを含む、
請求項 7 に記載の方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、本発明は連結型 (concatenative) 音声合成システムに関する。より詳しくは、本発明は、連結した音声単位 (音声ユニット: speech unit) について適切なエッジ境界領域を識別するためのシステムおよび方法に関する。システムは、音声単位モデルを用いて設けられた音声単位データベースを利用する。

10

20

30

40

50

【 0 0 0 2 】

【 従来 の 技 術 】

連結型音声の合成は、今日、数多くの様々な形態で世の中に存在しており、それは、どのように連結音声単位が格納され、処理されるかに依存している。これらの形態は、時間領域波形表現や、（例えば、フォルマント線形予測コーディング L P C 表現などの）周波数領域表現、またはこれらの組み合わせを含む。

【 0 0 0 3 】

音声単位の形態にかかわらず、連結型音声の合成は、各単位（ユニット：u n i t）のエッジで適切な境界領域を識別することにより行われる。ここで、単位は滑らかに重複され、それにより語や句を含む新たな音声単位に合成される。連結型音声合成システムにおける音声単位は、典型的には2音（d i p h o n e s）または半音節（d e m i s y l l a b l e s）である。この場合には、境界重複領域は音素内にある（p h o n e m e - m e d i a l）。したがって、例えば、「t o o l」という語は、「t o o t h」および「f o o l」という語から導き出された単位「t u」および「u l」により組み立てられる。決定すべきは、どの程度の量のソース語が音声単位にセーブされるかであり、また一緒に置かれたときにどの程度重複するべきかである。

【 0 0 0 4 】

連結型テキスト - 音声（t e x t - t o - s p e e c h : T T S）システムに関する従来の研究では、重複領域を判定するのに多くの方法が利用されてきた。このようなシステムを設計するに際しては、3つの因子が考慮される。すなわち、

- ・シームレスな連結：音声単位の重複により、ある単位とテキストとの間は十分滑らかに遷移し、急激な変化は聞こえないようにすべきである。リスナーには、音声片から組み立てられた音声を聞いているとはわからないようする必要がある。

【 0 0 0 5 】

- ・歪みのない遷移：音声単位の重複により、それ自身の歪みを生じてはならない。単位は、非重複音声との識別ができないように混在する必要がある。

【 0 0 0 6 】

- ・最小のシステム負荷：音声合成部における計算に必要な要件および/または記憶容量の要件は、できるだけ小さくする必要がある。

【 0 0 0 7 】

【 発明 が 解 決 し よ う と す る 課 題 】

現在のシステムではこれらの3つの目標の間にはトレードオフが存在し、3つのすべてに関して最適なシステムは存在していない。現在のアプローチは、一般的に3つの目標のバランスをとった、2つの選択に基づいてグループ化できる。第1の選択は、短い重複領域を用いるか、長い重複領域を用いるかである。短い重複領域を用いると、単一の声門パルスと同じ程度に早くできる。一方、長い重複領域を用いると、全音素の大部分を含むことができる。第2の選択は、重複領域は前後関係が整合しているか、または変化してもよいかである。前者の場合には、各音声単位の対応する部分は、先行する単位および後続の単位がどのような単位であるかにかかわらず重複している。後者の場合には、その単位が用いられる度に、隣接する単位に依存して、用いられる部分が変化する。

【 0 0 0 8 】

重複が長いと、単位間の遷移がよりシームレスになるという利点がある。その理由は、それらの間の微妙な相違が取り除かれる機会が多いからである。しかし、重複が長いと歪みを生じやすい。信号と異なり、混合すると歪みが生じる。

【 0 0 0 9 】

重複が短いと、歪みを最小にできるという利点がある。重複を短くすると、重複部分を十分に一致させることが簡単かつ確実にできる。短い重複領域は、（動的変化状態とは異なり）ほぼその瞬間の状態の特徴を表すと考えられる。しかし重複を短くすると、重複が長いシステムで実現できるシームレスな連結が犠牲になる。

【 0 0 1 0 】

10

20

30

40

50

重複が長い場合でシームレスが実現できることが望ましく、重複が短い場合に歪みを少なくできることが望ましいが、現在までのところ、これを達成できるシステムは存在しない。最新のシステムの中には、重複が長い場合の利点を保持しながら歪みを最小にするという目的で、可変重複領域を用いる実験が行われているものがある。しかし、このようなシステムは、計算負荷が高い処理に非常に大きく頼っているために、多くの用途には非実用的である。

【 0 0 1 1 】

本発明の目的は、シームレスで、かつ歪みのない重複を与える音声単位の領域を識別する方法、および連結型音声合成する方法を提供することである。

【 0 0 1 2 】**【 課題を解決するための手段 】**

本発明の連結型音声合成のための単位重複領域の識別方法は、音声の時変特性を表す統計モデルを画定するステップと、同じ母音を含む異なる音声単位に対応する複数の時系列データを提供するステップと、前記時系列データから音声信号パラメータを抽出し、前記音声信号パラメータを用いて前記統計モデルを学習するステップと、学習させた前記統計モデルを用いて前記時系列データ内の繰り返しシーケンスを識別し、前記繰り返しシーケンスを前記母音の中心の核をなす状態遷移部と関連付けるステップと、前記繰り返しシーケンスを用いて、連結型音声合成のための単位重複領域を定めるステップとからなり、それにより上記目的が達成される。

【 0 0 1 3 】

前記統計モデルは隠れマルコフモデルであってもよい。

【 0 0 1 4 】

前記統計モデルはリカレントニューラルネットワークであってもよい。

【 0 0 1 5 】

前記音声信号パラメータは音声フォルマントを含んでいてもよい。

【 0 0 1 6 】

前記統計モデルは、前記母音の中心の核をなす状態遷移部と、前記中心の核をなす状態遷移部の周囲の遷移部とを別々にモデル化するデータ構造を有していてもよい。

【 0 0 1 7 】

統計モデルを学習する前記ステップは、埋め込み再評価により行われ、前記時系列データによって表される全データセットにわたって整列のために収束したモデルを生成してもよい。

【 0 0 1 8 】

前記統計モデルは、前記母音の中心の核をなす状態遷移部と、前記中心の核をなす状態遷移部に先行する第1の遷移部と、前記中心軌線領域に後続する第2の遷移部とを別々にモデル化するデータ構造を有し、前記データ構造を用いて、前記第1の遷移部および前記第2の遷移部の1つに対応する前記時系列データの1部分を破棄するステップを含んでいてもよい。

【 0 0 1 9 】

本発明による連結型音声合成方法は、音声の時変特性を表す統計モデルを画定するステップと、同じ母音を含む異なる音声単位に対応する複数の時系列データを提供するステップと、前記時系列データから音声信号パラメータを抽出し、前記音声信号パラメータを用いて前記統計モデルを学習するステップと、学習させた前記統計モデルを用いて前記時系列データ内の繰り返しシーケンスを識別し、前記繰り返しシーケンスを前記母音の中心の核をなす状態遷移部と関連付けるステップと、前記繰り返しシーケンスを用いて、連結型音声合成のための単位重複領域を定めるステップと、前記音声単位の各単位重複領域に基づいて、2つの異なる前記音声単位からの前記時系列データを重複させ、マージすることにより、新たな音声単位を連結して合成するステップとからなり、それにより上記目的が達成される。

【 0 0 2 0 】

前記合成するステップを行う前に、前記単位重複領域の少なくとも1つの継続時間を選択的に変化させて、前記単位重複領域の他方の継続時間に一致させるステップをさらに含んでいてもよい。

【0021】

前記統計モデルは隠れマルコフモデルであってもよい。

【0022】

前記統計モデルはリカレントニューラルネットワークであってもよい。

【0023】

前記音声信号パラメータは音声フォルマントを含んでいてもよい。

【0024】

前記統計モデルは、前記母音の中心の核をなす状態遷移部と、前記中心の核をなす状態遷移部の周囲の遷移部とを別々にモデル化するデータ構造を有していてもよい。

【0025】

統計モデルを学習する前記ステップは、埋め込み再評価により行われ、前記時系列データによって表される全データセットにわたって整列のために収束したモデルを生成してもよい。

【0026】

前記統計モデルは、前記母音の中心の核をなす状態遷移部と、前記中心の核をなす状態遷移部に先行する第1の遷移部と、前記中心の核をなす状態遷移部に後続する第2の遷移部とを別々にモデル化するデータ構造を有し、前記データ構造を用いて、前記第1の遷移部および前記第2の遷移部の1つに対応する前記時系列データの1部分を破棄するステップを含んでいてもよい。

【0027】

本発明は統計的モデル化技術を利用することにより、音声単位内で中心軌跡領域を識別する。これらの領域は最適な重複境界を識別するのに用いられる。好ましい本実施の形態では、時系列データが、隠れマルコフモデルを用いて統計的にモデル化される。隠れマルコフモデルは、各音声単位の音素領域上に構築され、学習または埋め込み(embedd)再評価を経て整列(align)される。

【0028】

好ましい実施の形態では、各音声単位の最初と最後の音素は3要素からなると考えられる。すなわち中心の核をなす状態遷移部(中心軌跡:nuclear trajectory)、中心の核をなす状態遷移部に先行する遷移部および中心の核をなす状態遷移部に後続する遷移部である。モデル化プロセスはこれらの3要素を最適に識別し、それにより中心の核をなす状態遷移部は問題となる音素のすべてのインスタンスに対して、相対的な整合を維持する。

【0029】

識別された中心の核をなす状態遷移部を用いると、中心の核をなす状態遷移部の先頭境界および終端境界は重複領域を画定する。重複領域はその後、連結合成に用いられる。

【0030】

好ましい本実施の形態では、母音の中心の核をなす状態遷移部、中心の核をなす状態遷移部に先行する第1の遷移部、および中心の核をなす状態遷移部に後続する第2の遷移部を別個にモデル化するためのデータ構造を有する統計的モデルを利用する。データ構造は、音声単位データの一部を破棄するのに用いられる。音声単位データの一部のデータは、連結プロセスの間には用いられない音声単位の部分に対応する。

【0031】

本発明には多数の利点および使用法が存在するが、本発明は、連結型音声合成システムに用いられる音声単位データベースの自動構築の基礎として用いることができる。自動化技術は、導き出された合成音声の品質を向上し、データベース収集プロセスにおける労力を大幅に削減することができる。

【0032】

10

20

30

40

50

音声信号パラメータは、同じ母音を含む、異なる音声単位に対応する時系列データから抽出される。抽出されたパラメータは、隠れマルコフモデルといった統計的モデルを学習するのに用いられる。統計的モデルは、母音の中心の核をなす状態遷移部と、その周りの遷移部とを別々にモデル化するデータ構造を有する。このモデルは、埋め込み再評価を経て学習され、中心の核をなす状態遷移部を識別する最適に整列されたモデルを決定する。中心の核をなす状態遷移部の境界は、後の音声単位との連結のために重複領域を定めるよう機能する。

【0033】

【発明の実施の形態】

本発明は、以下の添付の図面を参照して説明される。

10

【0034】

本発明により利用される技術をもっともよく理解するためには、連結合成の基本的な理解が必要である。図1は、例を通した連結合成プロセスを示す。この例では、異なる2つの語からの音声単位(この場合は音節)が連結され、第3の語を形成する。より具体的には、「suffice」および「tight」という語からの音声単位が組み合わせられ、新たな「fight」という語が合成される。

【0035】

図1を参照して、「suffice」および「tight」という語からの時系列データが、好ましくは音節の境界で抽出され、音声単位10、12を規定する。この場合、音声単位10は14においてさらに細分割され、連結に必要な関連部分を分離する。

20

【0036】

その後、音声単位は16で整列され、それにより各部分18および20により規定される重複領域が作られる。整列後、時系列データがマージされ、新たな語22が合成される。

【0037】

本発明は特に、重複領域16と最適部分18、20に関連し、ある音声単位から別の音声単位までの遷移をシームレスで、かつ歪みがないようにする。

【0038】

本発明は、自動化された手順を経てこの最適な重複を実現する。この手順では、母音内で中心の核をなす(中心軌跡: nuclear trajectory)領域が探し出される(なお、「中心軌跡」の「軌跡」とは、本明細書において、目標周波数に向かって変化する概念を表すのに用いられる)。ここで母音内で「中心の核をなす」領域とは、母音の中心にある、安定した領域をいう。音声波形は、それを構成するフォーマット周波数によって表すことができる。これらの周波数は、ある音節が次の音節に融和して発音されると一定の変化を生じる。伝統的には、発声は、安定した目標周波数に向かって変化するこれらのフォーマット周波数を利用して、典型的には母音を利用してなされている。このとき周波数の波形は、直ちにより安定した波形になる。本明細書で母音内で「中心の核をなす」とは、母音によって占められる、中心にある安定した領域をいう。音声信号は、動的ではあるが同じ音素の異なる例に対しては相対的に変化がない動的パターンに続く。母音の境界領域は、隣接する子音によって影響を受けるが、中心にある安定した領域は強く影響を受けない。

30

40

【0039】

これらの最適な重複領域を改良するための手順が、図2に示される。まず、音声単位30が提供されている。データベース30は時系列データを含んでおり、時系列データは、連結合成システムを構成する異なる音声単位に対応する。好ましい本実施の形態では、音声単位は発声された語の例の中から抽出される。発声された語の例は、後に音節境界でさらに分割される。図2では、図解的に音声単位32, 34が描かれている。音声単位32は「tight」という語から抽出され、音声単位34は「suffice」という語から抽出されている。

【0040】

データベース30に格納されている時系列データはまず、36においてパラメータ化され

50

る。概して、音声単位は任意の方法論を用いてパラメータ化できる。好ましい本実施の形態では、各音声単位内で音素領域をフォルマント解析してパラメータ化を行う。フォルマント解析は、必然的に音声フォルマント周波数の抽出を伴う。本実施の形態ではフォルマント周波数 F 1、F 2 および F 3 が抽出される。必要であれば、R M S 信号レベルもまたパラメータ化できる。

【0041】

現在のところはフォルマント解析が好ましいが、パラメータ化の他の形態もまた利用できる。例えば、音声の特徴抽出は線形予測コーディング (Linear Predictive Coding: LPC) などの手順を用いて行い、適切な特徴パラメータを識別し、抽出できる。

10

【0042】

適切なパラメータが抽出され、各音声単位の音素領域が表されると、38で示されるようにモデルが構築され、各単位の音素領域が表される。好ましい本実施の形態はこの目的のために隠れマルコフモデルを用いる。しかし、概して時変または動的挙動を表す、適切な任意の統計的モデルを用いることができる。例えば、リカレントニューラルネットワークモデルを利用できる。

【0043】

好ましい本実施の形態は、音素領域を3つの異なる中間領域に分割してモデル化する。これらの領域は40で示されており、中心の核をなす状態遷移部(中心の核をなす領域)42と、中心の核をなす状態遷移部42に先行する状態遷移部(先行状態遷移領域)44と、中心の核をなす状態遷移部42に後続する状態遷移部(後続状態遷移領域)46とを含む。好ましい実施の形態では、これらの3領域の各々について別々の隠れマルコフモデルを用いる。先行および後続の状態遷移部44、46には、3状態モデルが用いられる。一方、中心の核をなす状態遷移部42には4または5状態モデルが用いられる。図2には5状態モデルが示されている。より大きな状態数を中心の核をなす状態遷移部42に用いると、後の手順は、整合のある非ヌル中心軌線に収束する。

20

【0044】

まず、音声モデル40が平均的な初期値で設けられる。その後、48で示されたこれらのモデルに関して、埋め込み(embedded)再評価が行われる。再評価とは、実質的には学習プロセスを継続することである。学習プロセスによりモデルは最適化されて、時系列データ内でもっともよい繰り返しシーケンスを表す。繰り返しシーケンスとは、母音内で中心にある安定した領域に関連する時系列データが呈する、より規則的な反復パターンのシーケンスをいう。これは、音声データが時系列データとして表されたときに、子音に対応する音声部分が規則性をもって反復しない非常に無秩序なパターンを呈しやすいこととは対照的である。したがって、母音が発生される度に繰り返して生じやすい時系列データ内のパターンは、母音領域内で識別できる。時系列データの繰り返しシーケンスは、識別されて所与の母音に対応する発声部分の識別手段として用いられる。例えば、音節「ya」の終端における母音音声は、音節「a」の統計的パターンと非常に関連のある統計的パターンを呈する。同じ統計的パターンは、例えば、音節「ka」、「ma」、「ha」内の安定領域において見出すことができる。対照的に、安定的な母音領域に先行する音節部分では、統計的な関連がない場合が多く、したがって識別可能な繰り返しパターンも存在しない。さらなる例示のために、時系列データが統計モデルを学習するのに用いられ、各モデルがパラメータの組を規定すると仮定する。モデルを学習させた後、母音音声「a」はパラメータ番号のシーケンス: 4-5-3, 1-6に対応する。母音が存在するたびに同一の番号のパターンが発生しているとすると、そのパターンは、その母音が存在することを示すのに信頼性高く利用できる繰り返しシーケンスを構成する。本発明では、子音、または安定的な母音に融和する音声などの他の音声は、非常に繰り返しのあるシーケンスを生成することが統計的に存在しないと判断する。したがって、発せられた音声内に安定した母音領域があることを検出する手段として、非常によく反復するシーケンス(繰り返しシーケンス)を見つけ出す。

30

40

50

【 0 0 4 5 】

中心の核をなす状態遷移部 4 2、先行および後続の状態遷移部 4 4、4 6 は、データベース 3 0 を介して供給される現実のデータに基づいて、学習プロセスにより各音素領域に整合するモデルが構築されるよう設計される。この点に関して、中心の核をなす部分 4 2 は母音の核心を表し、先行および後続の状態遷移部 4 4、4 6 は、現在の音素および現在の音素に先行するおよび後続する音声に固有の母音の相を表す。例えば、「t i g h t」という語から抽出された音声単位 3 2 では、先行する遷移部は、前にある子音字「t」により母音「a y」の音声に与えられた音調(c o l o r a t i o n)を表す。

【 0 0 4 6 】

整合プロセスは本来、最適な整列モデルに収束する。どのようにしてそのようになるのかを理解するために、音声単位 3 0 のデータベースが、少なくとも 2 つ、好ましくは多数の各母音の音声の例を含むとする。例えば図 2 には、「t i g h t」および「s u f f i c e」の双方に見受けられる母音の音声「a y」が、音声単位 3 2、3 4 により表されている。埋め込み再評価プロセスまたは学習プロセスは、音声「a y」のこのような複数のインスタンスを用いて初期音声モデル 4 0 の学習を行い、それにより最適に整列された音声モデル 5 0 を生成する。音声「a y」の例のすべてにわたって整合のある時系列データの部分は、中核、または中心の核をなす領域を表す。5 0 で図示されるように、システムは、先行および後続の状態遷移部を別々に学習する。これらは、母音に先行するおよび後続する音声に依存して当然に異なっている。

10

【 0 0 4 7 】

一旦モデルが学習され、最適に整列されたモデルを生成すると、中心の核をなす領域 4 2 の両側の境界が確定し、連結合成のための重複領域の位置が決定される。そのため、ステップ 5 2 では最適に整列されたモデルが重複境界を決定するのに用いられる。図 2 は、重複境界 A および B を示す。重複境界 A および B は、「s u f f i c e」および「t i g h t」という語から導かれた音声単位に対するフォルマント周波数データに重ね合わされている。

20

【 0 0 4 8 】

パラメータデータ(この場合はフォルマント周波数データ)で識別された重複境界により、システムはステップ 5 4 において時系列データを分類して時系列データ内の重複境界を定める。必要であれば、分類されたデータは連結型音声合成について後に使用するために、データベース 3 0 に格納してもよい。

30

【 0 0 4 9 】

図示の関係上、オーバーレイテンプレート 5 6 として模式的に示されている重複境界領域が、「s u f f i c e」という語の時系列データの模式的表現に重ね合わされて示されている。具体的には、テンプレート 5 6 は、後半の音節「. . . f i c e」内で括弧 5 8 によって示すように整列されている。この音声単位が連結音声に用いられると、先行領域 6 2 は破棄され、境界 A および B により定められている中心の核をなす領域 6 4 は、クロスフェード領域または連結領域として働く。

【 0 0 5 0 】

ある実施形態では、連結合成を行うために、重複領域の継続時間を調整する必要がある。このプロセスが図 3 に示される。入力テキスト 7 0 が解析され、ステップ 7 2 に示されるようにデータベース 3 0 から適切な音声単位が選択される。例えば、「f i g h t」という語が入力テキストとして与えられると、システムは「t i g h t」および「s u f f i c e」という語から抽出した、あらかじめ格納してある音声単位を選択する。

40

【 0 0 5 1 】

各音声単位を中心の核をなす領域は必ずしも同じ時間にわたっている必要はない。そのためステップ 7 4 では、各中心の核をなす領域の継続時間が伸張または短縮され、それにより継続時間を一致させる。図 3 では、中心の核をなす領域 6 4 a が領域 6 4 b に伸張される。音声単位 B も同様に変更される。図 3 は中心の核をなす領域 6 4 c が領域 6 4 d に圧縮され、それにより 2 つの単位の各領域が同じ継続時間を持つことになる。

50

【 0 0 5 2 】

一旦継続時間が調整されて一致すると、ステップ76において、音声単位からのデータがマージされて、78で示される新しく連結された単語を形成する。

【 0 0 5 3 】

【 発明の効果 】

これまでの説明によれば、本発明は連結型音声合成システムに用いられる音声単位データベースを構築する自動化手段を提供することが理解される。中心の核をなす領域を分離することによって、このシステムは、シームレスで、かつ歪みのない重複を与える。有利なのは、重複領域は共通の固定サイズに伸張または圧縮され、連結プロセスを簡単化できることである。統計的モデル化プロセスを用いることで、中心の核をなす領域は音声信号の1部分を表すことができる。ここでは、音響学上の音声特性は、同じ音素の異なる例に対しては相対的に変化がない動的パターンを生じる結果となる。変化がないことにより、シームレスで、かつ歪みのない遷移が可能になる。

10

【 0 0 5 4 】

本発明の原理により生成された音声単位は、コンピュータ処理システムにかかる負担を最小にして、後の抽出および連結に用いるデータベースに容易に格納できる。したがって、このシステムは、処理能力が制限されている合成音声に関する製品および応用の開発には理想的といえる。さらに、音声単位を生成する自動化プロセスは、目的が特化された音声単位データベースを構築するのに必要な時間と労力を大幅に減少させる。例えば音声単位を生成する自動化プロセスは、専門的なボキャブラリに対して、または多言語音声合成システムの開発に対して必要とされるであろう。

20

【 0 0 5 5 】

現時点での好ましい形態で本発明を説明してきたが、当業者であれば、特許請求の範囲に記載された本発明の精神から逸脱することなく本システムを修正できる。

【 図面の簡単な説明 】

【 図 1 】 連結型音声を合成する技術の理解に有用なブロック図である。

【 図 2 】 本発明による、音声単位が構築される手順を示すフローチャートである。

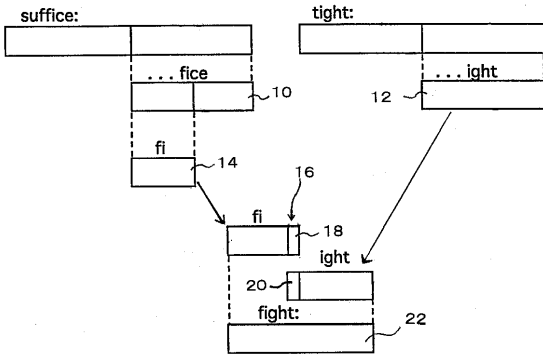
【 図 3 】 本発明の音声単位データベースを用いた、連結型音声を合成するプロセスを示すブロック図である。

【 符号の説明 】

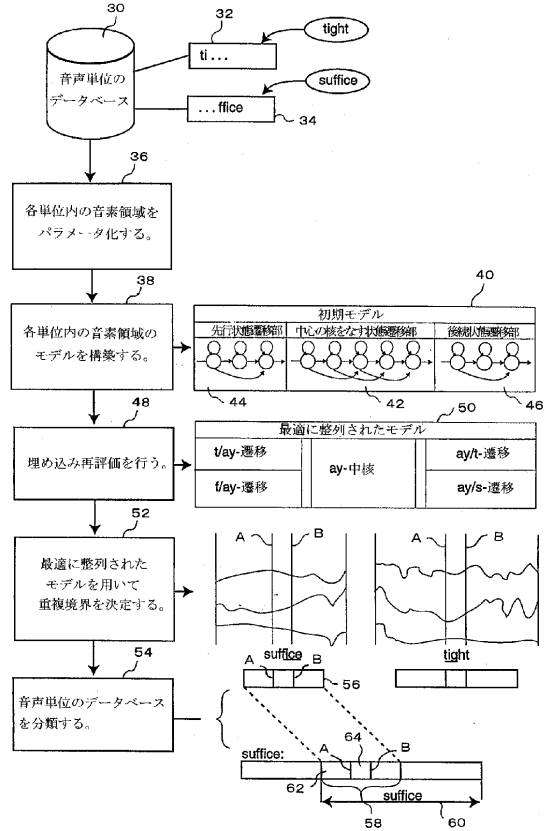
- 4 0 音声モデル
- 4 2 中心の核をなす状態遷移部
- 4 4 先行状態遷移部
- 4 6 後続状態遷移部
- 5 0 音声モデル
- 5 6 オーバレイテンプレート
- 6 2 先行領域
- 6 4 中心の核をなす領域

30

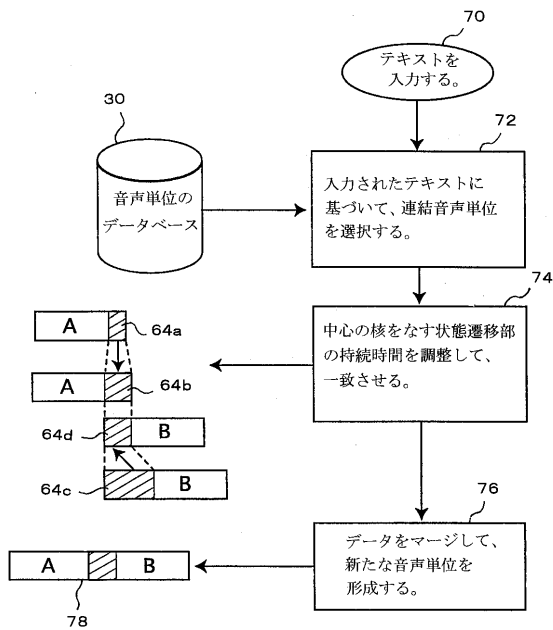
【図1】



【図2】



【図3】



フロントページの続き

(72)発明者 スティーブ・ピアソン

アメリカ合衆国93110カリフォルニア州サンタ・バーバラ、カリエ・シタ3909番

審査官 渡邊 聡

(56)参考文献 特開平10-247097(JP,A)

特開平10-049193(JP,A)

特開平07-244497(JP,A)

吉村 他, HMMに基づく音声合成のための状態継続長モデルの構築, 電子情報通信学会技術研究報告, 電子情報通信学会, 1998年 9月, SP98-64, 45-50

益子 他, HMMを用いた音声合成における音素モデルの検討, 日本音響学会講演論文集, 日本音響学会, 1996年 3月, 平成8年春季I, 273-274

(58)調査した分野(Int.Cl.⁷, DB名)

G10L 13/08

G10L 13/06