

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4586386号
(P4586386)

(45) 発行日 平成22年11月24日(2010.11.24)

(24) 登録日 平成22年9月17日(2010.9.17)

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

請求項の数 12 (全 12 頁)

(21) 出願番号	特願2004-73977 (P2004-73977)	(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
(22) 出願日	平成16年3月16日(2004.3.16)	(74) 代理人	100099933 弁理士 清水 敏
(65) 公開番号	特開2005-265895 (P2005-265895A)	(72) 発明者	西澤 信行 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
(43) 公開日	平成17年9月29日(2005.9.29)	(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
審査請求日	平成19年2月26日(2007.2.26)	審査官	井上 健一

最終頁に続く

(54) 【発明の名称】 素片接続型音声合成装置及び方法

(57) 【特許請求の範囲】

【請求項1】

合成音声の目標と音声素片候補との間で、複数のサブコストを含むコストを算出し、当該コストに基づいて音声素片データベースから音声素片を選択し接続することにより音声合成を行なう素片接続型音声合成装置であって、

前記音声素片データベースに含まれる音声素片候補から、前記複数のサブコストのうちの一部のみを用い、2以上の段階に分けて複数の音声素片候補列を選択するための多段予備選択手段と、

前記合成音声の目標との間で、前記複数のサブコストを全て含んで算出されるコストが所定の条件を充足する一つの音声素片候補列を、前記多段予備選択手段により予備的に選択された前記複数の音声素片候補列から選択するための選択手段と、

前記選択手段により選択された音声素片候補列の音声波形を前記合成器指令に従って接続し合成音声波形を出力するための接続手段とを含む、素片接続型音声合成装置。

【請求項2】

前記多段予備選択手段は、前記音声素片データベースに含まれる音声素片候補から、前記複数のサブコストのうちの一部のみを用い、2以上の段階に分けて、かつ後段の予備選択では前段の予備選択で用いられたサブコストより多数種類のサブコストを用いて予備選択を行なって、複数の音声素片候補列を選択するための手段を含む、請求項1に記載の素片接続型音声合成装置。

【請求項3】

前記多段予備選択手段は、前記音声素片データベースに含まれる音声素片候補から、前記複数のサブコストのうちの一部のみを用い、2以上の段階に分けて、かつ後段の予備選択では、前段の予備選択で用いられたサブコストより多数種類のサブコストであって、かつ前段の予備選択で用いられたサブコストを含むサブコストを用いて予備選択を行なって、複数の音声素片候補列を選択するための手段を含む、請求項2に記載の素片接続型音声合成装置。

【請求項4】

前記多段予備選択手段は、

前記合成音声の目標と前記音声素片データベース中の各音声素片候補との間で、第1のサブコストを算出し、算出された第1のサブコストを用いて複数の音声素片候補列を選択するための第1段の予備選択手段と、

10

前記第1のサブコストと、前記第1のサブコストと異なる第2のサブコストとの双方を用いて、前記第1段の予備選択手段により選択された複数の音声素片候補列の中から複数の音声素片候補列を選択するための第2段の予備選択手段とを含む、請求項1に記載の素片接続型音声合成装置。

【請求項5】

前記第1段の予備選択手段は、

前記合成音声の目標と前記音声素片データベース中の音声素片候補からなる各音声素片候補列との間で、第1のサブコストを算出するための第1のサブコスト算出手段と、

前記第1のサブコスト算出手段により算出された第1のサブコストを記憶するための第1のサブコスト記憶手段と、

20

前記第1のサブコスト算出手段により算出された第1のサブコストが所定のしきい値よりも小さな音声素片候補列を選択するための手段とを含む、請求項4に記載の素片接続型音声合成装置。

【請求項6】

前記多段予備選択手段は、

前記合成音声の目標と前記音声素片データベース中の音声素片候補からなる各音声素片候補列との間で、ターゲットコストのみからなる第1のサブコストを算出し、算出された第1のサブコストを用いて複数の音声素片候補列を選択するための第1段の予備選択手段と、

30

前記第1のサブコストと、接続コストを含む第2のサブコストとの双方を用いて、前記第1段の予備選択手段により選択された複数の音声素片候補列の中から複数の音声素片候補列を選択するための第2段の予備選択手段とを含む、請求項1に記載の素片接続型音声合成装置。

【請求項7】

合成音声の目標と音声素片候補との間で、複数のサブコストを含むコストを算出し、当該コストに基づいて音声素片データベースから音声素片を選択し接続することにより音声合成を行なう素片接続型音声合成方法であって、

前記音声素片データベースに含まれる音声素片候補から、前記複数のサブコストのうちの一部のみを用い、2以上の段階に分けて複数の音声素片候補列を選択する多段予備選択ステップと、

40

前記合成音声の目標との間で、前記複数のサブコストを全て含んで算出されるコストが所定の条件を充足する一つの音声素片候補列を、前記多段予備選択ステップにおいて予備的に選択された前記複数の音声素片候補列から選択する選択ステップと、

前記選択ステップにおいて選択された音声素片候補列の音声波形を前記合成器指令に従って接続し合成音声波形を出力する接続ステップとを含む、素片接続型音声合成方法。

【請求項8】

前記多段予備選択ステップは、前記音声素片データベースに含まれる音声素片候補から、前記複数のサブコストのうちの一部のみを用い、2以上の段階に分けて、かつ後段の予備選択では前段の予備選択で用いられたサブコストより多数種類のサブコストを用いて予備

50

選択を行なって、複数の音声素片候補列を選択するステップを含む、請求項 7 に記載の素片接続型音声合成方法。

【請求項 9】

前記多段予備選択ステップは、前記音声素片データベースに含まれる音声素片候補から、前記複数のサブコストのうちの一部のみを用い、2 以上の段階に分けて、かつ後段の予備選択では、前段の予備選択で用いられたサブコストより多数種類のサブコストであって、かつ前段の予備選択で用いられたサブコストを含むサブコストを用いて予備選択を行なって、複数の音声素片候補列を選択するステップを含む、請求項 8 に記載の素片接続型音声合成方法。

【請求項 10】

前記多段予備選択ステップは、

前記合成音声の目標と前記音声素片データベース中の各音声素片候補との間で、第 1 のサブコストを算出し、算出された第 1 のサブコストを用いて複数の音声素片候補列を選択する第 1 段の予備選択ステップと、

前記第 1 のサブコストと、前記第 1 のサブコストと異なる第 2 のサブコストとの双方を用いて、前記第 1 段の予備選択ステップにおいて選択された複数の音声素片候補列の中から複数の音声素片候補列を選択する第 2 段の予備選択ステップとを含む、請求項 7 に記載の素片接続型音声合成方法。

【請求項 11】

前記第 1 段の予備選択ステップは、

前記合成音声の目標と前記音声素片データベース中の音声素片からなる各音声素片候補列との間で、第 1 のサブコストを算出する第 1 のサブコスト算出ステップと、

前記第 1 のサブコスト算出ステップにおいて算出された第 1 のサブコストを、第 1 のサブコスト記憶手段に記憶させるステップと、

前記第 1 のサブコスト算出ステップにおいて算出された第 1 のサブコストが所定のしきい値よりも小さな音声素片候補列を選択するステップとを含む、請求項 10 に記載の素片接続型音声合成方法。

【請求項 12】

前記多段予備選択ステップは、

前記合成音声の目標と前記音声素片データベース中の音声素片からなる各音声素片候補列との間で、ターゲットコストのみからなる第 1 のサブコストを算出し、算出された第 1 のサブコストを用いて複数の音声素片候補列を選択する第 1 段の予備選択ステップと、

前記第 1 のサブコストと、接続コストを含む第 2 のサブコストとの双方を用いて、前記第 1 段の予備選択ステップにおいて選択された複数の音声素片候補列の中から複数の音声素片候補列を選択する第 2 段の予備選択ステップとを含む、請求項 7 に記載の素片接続型音声合成方法。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は音声合成装置に関し、特に、所定のコスト関数に基づいて音声素片を選択し接続することにより合成器指令に合致した音声合成を行なう音声合成装置に関する。

【背景技術】

【0002】

音声認識、音声合成は、人間とコンピュータを用いた諸システムとのインターフェースを実現する技術として重要である。これらと人工知能技術とを併用することにより、利用者は相手がコンピュータシステムであることを意識せずに様々なサービスを利用することができる。

【0003】

中でも音声合成は、人間に対するシステム出力のためのインターフェースとしてその重要性は大きい。人間は、合成された音声の不自然さを敏感に感じ取る。合成された音声

10

20

30

40

50

不自然であると利用者が感じると、発話にも影響を及ぼし、その結果、人間とシステムとの間の対話がうまく行かなくなるおそれもある。

【0004】

最近の音声合成技術としては、予め人間の発話を多数集めて語・音節・音素等を単位とする音声素片を音素ラベルと関連付けてデータベース化しておき、合成時には、指定された語・音節・音素等に対応する音声素片の中から、最も適切と思われるものを選択して接続するものが知られている。これを素片接続型音声合成と呼ぶ。なお、音素ラベルとは、通常は各音素の音素記号とその開始・終了時刻を記述したものをいう。これに加えて、その区間におけるMFCC (Mel-Frequency Cepstrum Coefficient)、基本周波数(F0)等の音響特徴量、さらに前後の素片の音素記号を含む場合もある。

10

【0005】

素片接続型音声合成では、与えられた合成目標を基準として、いかにして適切な音声素片をデータベース中から取出すかが問題となる。

【0006】

合成目標を構成するデータは、典型的には音素と、F0、持続時間、MFCC、及びパワー等の音声特徴量とを含む。これらを以下「合成器指令」と呼ぶ。

【0007】

素片接続型音声合成では、合成器指令と音声素片のF0、持続時間、MFCC、パワー等とのずれ、及び接続に伴う自然劣化を表現するための「コスト」と呼ばれる評価関数を定義し、コストを最小とする音声素片を求めることにより、最適な音声素片系列を決定する。

20

【0008】

本件出願の出願人は、上記した「コスト」を、それぞれある音声の特徴に対応するような「サブコスト」に分解し、それらを結合したもの(例えば線形和)により定義した素片接続型音声合成を提案している。例えば特許文献1を参照されたい。

【0009】

サブコストには、物理量から計算されるものと、シンボリックな情報から事前に作成した規則から基づき得られるものとがある。前者は、複数のサンプル値に対する非線形演算であることも多く、その計算量は相対的に大きい。後者は、単純なテーブル参照の形であることが多く、テーブル参照で実現される場合にはサブコスト計算に必要な計算量は非常に少ない。

30

【0010】

以上はあくまで一例であるが、この例に限らず、各サブコストの計算量はその種類により大きなばらつきがある場合が多い。

【0011】

一方、上記とは別に、サブコストは、ターゲットコストに関係するものと接続コストに属するものとの二つに大別することもできる。ターゲットコストは、合成目標と素片候補との間の誤差を表す。接続コストは、合成音声において隣接する素片間の誤差(不連続性)を表す。

40

【0012】

素片接続型音声合成では、コストの最小化に基づく素片選択が行なわれるが、特に素片候補数が多い場合にはコストの計算に要する計算量が問題となる。

【0013】

最小コストとなる素片候補系列の推定において、可能な組合せのコストを全て調べることは、組合せ爆発により記憶容量・計算時間の双方において非現実的である。そこで、各時刻の素片候補を予備選択により絞り込む方法が考えられる。この際、計算量を考えて、前後の素片関係の影響を受けないターゲットコストに属するサブコストのみで予備選択を行なう方法が有力である。

【0014】

50

【特許文献1】特開2003-208188号公報(段落0014~0047)

【発明の開示】

【発明が解決しようとする課題】

【0015】

しかし、合成音声に対する接続コストの影響が比較的大きいことから、予備選択段階でターゲットコストのみに基づいて絞り込むことができる候補数には限界がある。ターゲットコストのみに基づいてあまり候補を絞り込むと、高品質な音声合成を行なうことが可能な候補が捨てられてしまうおそれがあるためである。その結果、本選択時の計算量の削減にも限界がある。

【0016】

それゆえに、本発明の目的は、高品質な音声合成が可能で、かつ選択のための計算量を削減できる素片接続型音声合成装置及び方法を提供することである。

【課題を解決するための手段】

【0017】

本発明の第1の局面に係る素片接続型音声合成装置は、合成音声の目標と音声素片候補との間で、複数のサブコストを含むコストを算出し、当該コストに基づいて音声素片データベースから音声素片を選択し接続することにより音声合成を行なう素片接続型音声合成装置であって、音声素片データベースに含まれる音声素片候補から、複数のサブコストのうちの一部のみを用い、2以上の段階に分けて複数の音声素片候補を選択するための多段予備選択手段と、合成音声の目標との間で、複数のサブコストを全て含んで算出されるコストが所定の条件を充足する一つの音声素片候補を、予備選択手段により予備的に選択された複数の音声素片候補から選択するための選択手段と、選択手段により選択された音声素片候補の音声波形を合成器指令に従って接続し合成音声波形を出力するための接続手段とを含む。

【0018】

好ましくは、多段予備選択手段は、音声素片データベースに含まれる音声素片候補から、複数のサブコストのうちの一部のみを用い、2以上の段階に分けて、かつ後段の予備選択では前段の予備選択で用いられたサブコストより多数種類のサブコストを用いて予備選択を行なって、複数の音声素片候補を選択するための手段を含む。

【0019】

より好ましくは、多段予備選択手段は、音声素片データベースに含まれる音声素片候補から、複数のサブコストのうちの一部のみを用い、2以上の段階に分けて、かつ後段の予備選択では、前段の予備選択で用いられたサブコストより多数種類のサブコストであって、かつ前段の予備選択で用いられたサブコストを含むサブコストを用いて予備選択を行なって、複数の音声素片候補を選択するための手段を含む。

【0020】

多段予備選択手段は、合成音声の目標と音声素片データベース中の各音声素片候補との間で、第1のサブコストを算出し、算出された第1のサブコストを用いて複数の音声素片を選択するための第1段の予備選択手段と、第1のサブコストと、第1のサブコストと異なる第2のサブコストとの双方を用いて、第1段の予備選択手段により選択された複数の音声素片の中から複数の音声素片を選択するための第2段の予備選択手段とを含んでもよい。

【0021】

さらに好ましくは、第1段の予備選択手段は、合成音声の目標と音声素片データベース中の各音声素片候補との間で、第1のサブコストを算出するための第1のサブコスト算出手段と、第1のサブコスト算出手段により算出された第1のサブコストを記憶するための第1のサブコスト記憶手段と、第1のサブコスト算出手段により算出された第1のサブコストが所定のしきい値よりも小さな音声素片候補を選択するための手段とを含む。

【0022】

好ましくは、多段予備選択手段は、合成音声の目標と音声素片データベース中の各音声

10

20

30

40

50

素片候補との間で、ターゲットコストのみからなる第1のサブコストを算出し、算出された第1のサブコストを用いて複数の音声素片を選択するための第1段の予備選択手段と、第1のサブコストと、接続コストを含む第2のサブコストとの双方を用いて、第1段の予備選択手段により選択された複数の音声素片からなる音声素片系列の中から複数の音声素片系列を選択するための第2段の予備選択手段とを含む。

【0023】

本発明の第2の局面に係る素片接続型音声合成方法は、合成音声の目標と音声素片候補との間で、複数のサブコストを含むコストを算出し、当該コストに基づいて音声素片データベースから音声素片を選択し接続することにより音声合成を行なう素片接続型音声合成方法であって、音声素片データベースに含まれる音声素片候補から、複数のサブコストのうちの一部のみを用い、2以上の段階に分けて複数の音声素片候補を選択する多段予備選択ステップと、合成音声の目標との間で、複数のサブコストを全て含んで算出されるコストが所定の条件を充足する一つの音声素片候補を、予備選択ステップにおいて予備的に選択された複数の音声素片候補から選択する選択ステップと、選択ステップにおいて選択された音声素片候補の音声波形を合成器指令に従って接続し合成音声波形を出力する接続ステップとを含む。

10

【0024】

好ましくは、多段予備選択ステップは、音声素片データベースに含まれる音声素片候補から、複数のサブコストのうちの一部のみを用い、2以上の段階に分けて、かつ後段の予備選択では前段の予備選択で用いられたサブコストより多数種類のサブコストを用いて予備選択を行なって、複数の音声素片候補を選択するステップを含む。

20

【0025】

さらに好ましくは、多段予備選択ステップは、音声素片データベースに含まれる音声素片候補から、複数のサブコストのうちの一部のみを用い、2以上の段階に分けて、かつ後段の予備選択では、前段の予備選択で用いられたサブコストより多数種類のサブコストであって、かつ前段の予備選択で用いられたサブコストを含むサブコストを用いて予備選択を行なって、複数の音声素片候補を選択するステップを含む。

【0026】

多段予備選択ステップは、合成音声の目標と音声素片データベース中の各音声素片候補との間で、第1のサブコストを算出し、算出された第1のサブコストを用いて複数の音声素片を選択する第1段の予備選択ステップと、第1のサブコストと、第1のサブコストと異なる第2のサブコストとの双方を用いて、第1段の予備選択ステップにおいて選択された複数の音声素片の中から複数の音声素片を選択する第2段の予備選択ステップとを含んでもよい。

30

【0027】

さらに好ましくは、第1段の予備選択ステップは、合成音声の目標と音声素片データベース中の各音声素片候補との間で、第1のサブコストを算出する第1のサブコスト算出ステップと、第1のサブコスト算出ステップにおいて算出された第1のサブコストを、第1のサブコスト記憶手段に記憶させるステップと、第1のサブコスト算出ステップにおいて算出された第1のサブコストが所定のしきい値よりも小さな音声素片候補を選択するステップとを含む。

40

【0028】

好ましくは、多段予備選択ステップは、合成音声の目標と音声素片データベース中の各音声素片候補との間で、ターゲットコストのみからなる第1のサブコストを算出し、算出された第1のサブコストを用いて複数の音声素片を選択する第1段の予備選択ステップと、第1のサブコストと、接続コストを含む第2のサブコストとの双方を用いて、第1段の予備選択ステップにおいて選択された複数の音声素片からなる音声素片系列の中から複数の音声素片系列を選択する第2段の予備選択ステップとを含む。

【発明を実施するための最良の形態】

【0029】

50

[第 1 の実施の形態]

図 1 に、本発明の第 1 の実施の形態に係る音声合成システム 20 のブロック図を示す。図 1 を参照して、この音声合成システム 20 は、従来と同様の音声素片 DB 34 と、合成目標となるテキストを分析した結果得られる合成器指令 36 を入力として受け、音声素片 DB 34 に含まれる拡張された音声素片から適切な音声素片を選択し接続して合成音声波形 40 を出力するための音声合成装置 38 とを含む。

【 0030 】

音声合成装置 38 は、合成器指令 36 を受け、合成器指令 36 により指定された音声素片のうちで、後述するように多段の予備選択を行なって予備選択候補群 62 を選択するための多段予備選択部 60 と、合成器指令 36 を受け、予備選択候補群 62 から全サブコストを用いて計算したコストの最も小さな素片を選択するための素片選択部 64 と、素片選択部 64 により選択された音声素片を接続して合成音声波形 40 を出力するための接続部 66 とを含む。なお、予備選択候補群 62 は素片の選択のみに用いられるので、コスト計算に必要な特徴量のみを含み、音声素片データそのものは含まない。接続部 66 は、素片選択部 64 により選択された素片の音声素片データを音声素片 DB 34 を参照して得ることになる。

【 0031 】

本実施の形態で使用されるサブコストは、基本周波数 (F 0) 誤差、継続時間長誤差、MFCC 誤差、F 0 不連続誤差、MFCC 不連続誤差、音素環境誤差にそれぞれ対応する 6 種類のサブコストを含む。これらのうち、前 3 者はターゲットコストに属し、後 3 者は接続コストに属する。

【 0032 】

本実施の形態に係る素片選択部 64 によるコスト計算では、コスト C_0 は以下のようにしてサブコストから計算される。

【 0033 】

【 数 1 】

$$C_o = \left(\sum_{i1=1}^{N_1} (w_{i1} C_{i1})^{p_1} \right)^{\frac{1}{p_1}} + \left(\sum_{i2=1}^{N_2} (w_{i2} C_{i2})^{p_2} \right)^{\frac{1}{p_2}} \quad (1)$$

ただし、 C_{i1} ($i1 = 1 \sim 3$) はターゲットサブコスト、 C_{i2} ($i2 = 1 \sim 3$) は接続コスト、 w_{i1} ($i1 = 1 \sim 3$) はターゲットサブコスト間に定義された重み、 w_{i2} ($i2 = 1 \sim 3$) は接続サブコスト間に定義された重み、 p_1 及び p_2 はそれぞれ、ターゲットコストと接続コスト間に定義された重みである。ただし、本実施の形態では後述するように多段予備選択における計算量を削減するため、 p_1 及び p_2 はいずれも 1 とする。

【 0034 】

一般的に、音素環境誤差のサブコストは比較的単純なテーブル参照である。したがってその計算量は非常に小さい。それ以外については、サブコストの計算量は比較的大きい。例えば MFCC は多次元量であるため、そのサブコストの計算に要する時間は他のサブコストより大きくなる。

【 0035 】

図 1 を参照して、多段予備選択部 60 は、4 つの予備選択部 70、80、90 及び 100 を含む。予備選択をどのような順番でどのサブコストに基づいて行なうかは、アプリケーション、より具体的には各サブコストに対し予想される計算量の相違により異なる。F 0 誤差、継続時間長誤差に関するサブコスト計算が比較的小さい場合には、図 1 に示すような構成が考えられる。4 つの予備選択部 70、80、90 及び 100、並びに素片選択部 64 の機能は以下のとおりである。なお、接続コストに関するサブコストが予備選択コストに含まれる場合、コストは前後の素片の影響を受ける。したがって、その時点での予備選択コスト関数を最小化する解について、各時刻において独立に素片候補を予備選択するのではなく、素片候補の選択系列の N - ベスト解を得ておく必要がある。その後は、そ

10

20

30

40

50

のN - ベスト解について後段の予備選択関数で再度コスト計算を行なってその結果のN - ベスト解を得る、という処理を繰返す必要がある。

【0036】

第1の予備選択部70：合成器指令36を受け、音声素片DB34中の素片候補から各時刻におけるF0誤差、継続時間長誤差による予備選択をして第1の候補群72を出力する。

【0037】

第2の予備選択部80：第1の候補群72中の素片から、各時刻におけるF0誤差、継続時間長誤差、MFCC誤差による予備選択をして第2の候補群82を出力する。

【0038】

第3の予備選択部90：第2の候補群82中の素片から、各時刻におけるF0誤差、継続時間長誤差、MFCC誤差、及び音素環境誤差を考慮したN - ベスト探索を行ない第3の候補群（選択系列群）92を出力する。

【0039】

第4の予備選択部100：第3の候補群92中の素片候補に、F0誤差、継続時間長誤差、MFCC誤差、音素環境誤差、及びF0不連続誤差を考慮したN - ベスト探索を行ない、予備選択候補群（選択系列群）62を出力する。

【0040】

素片選択部64：予備選択候補群62に含まれるN - ベスト選択系列に対して、全てのサブコストを考慮して行なう1 - ベスト探索を行ない、素片を一つ選択し接続部66に与える。

【0041】

なお、N - ベスト解はビームサーチ又はN - ベストDP (Dynamic Programming) サーチにより行なうことができる。(ここでN - ベストDPサーチとは、DP探索における各ノードでN - ベスト解を保持する方法のことをいう。通常のDPサーチは各ノードで1 - ベストの解のみを保持している。)

【0042】

ここで、ビームサーチについては、N - ベスト解で選択される候補系列の数Nに対して、ビーム幅が小さいほど最適解に近い解が得られる可能性が小さくなる。一方、N - ベストDPサーチでは、各ノードが保持するN - ベスト解の数が少ないほど、最適解が得られる可能性が低くなる。(ここで、各ノードにおけるN - ベスト解の数が、最終的に必要となるN - ベスト解の数と同数以上であれば、解が真のN - ベスト解であることは保証される。しかし、多段選択の途中におけるN - ベスト解の中に真の最適解が含まれている保証はなく、計算途中で真のN - ベスト解を得ること自体にはそれほど意味はない。)ただし、素片候補が大量に存在する場合には、仮に最終的に最適解でない解が得られたとしても実用上十分な品質が得られる可能性が高い。

【0043】

この実施の形態では、前段の予備選択部で算出されたサブコストは、後段の予備選択部でも素片選択に使用される。したがって、サブコストが式(1)で表され、かつ p_1 及び p_2 がいずれも1として設計した場合(すなわちコストがサブコストの線形和で表される場合)、前段の予備選択部で算出したサブコストをそのまま次の予備選択部でのコスト計算に用いることができる。そのために多段予備選択部60は、それぞれ予備選択部70、80、90及び100で行なわれたサブコスト計算の結果を記憶するための第1～第4のコスト記憶部74、84、94及び104をさらに含む。これら第1～第4のコスト記憶部74、84、94及び104に記憶されたサブコストは、それぞれ予備選択部80、90、及び100並びに素片選択部64に与えられ、コスト計算に用いられる。

【0044】

この音声合成システム20は以下のように動作する。まず合成器指令36が音声合成装置38に与えられる。多段予備選択部60の第1の予備選択部70は、合成器指令36に基づいて、合成器指令36により指定された音素に対応する音声素片であってかつF0誤

10

20

30

40

50

差及び継続時間長誤差により算出されたサブコストの線形和が所定のしきい値以下であるものを音声素片 D B 3 4 から抽出し、第 1 の候補群 7 2 として出力する。このときのサブコストの計算結果は第 1 のコスト記憶部 7 4 に記憶される。

【 0 0 4 5 】

第 2 の予備選択部 8 0 は、合成器指令 3 6 に基づいて、F 0 誤差、継続時間長誤差及び M F C C 誤差により算出されたサブコストの線形和が所定のしきい値以下であるものを第 1 の候補群 7 2 から抽出し、第 2 の候補群 8 2 として出力する。このとき、第 2 の予備選択部 8 0 は、第 1 のコスト記憶部 7 4 に記憶された F 0 誤差及び継続時間長誤差により算出されたサブコストをサブコスト計算に用いる。したがって実質的には M F C C によるサブコスト計算のみが行なわれる。第 2 の予備選択部 8 0 によるサブコストの計算結果は第 2 のコスト記憶部 8 4 に記憶される。

10

【 0 0 4 6 】

第 3 の予備選択部 9 0 は、合成器指令 3 6 に基づいて、F 0 誤差、継続時間長誤差、M F C C 誤差、及び音素環境誤差に基づいて算出されたサブコストの線形和に基づき、素片候補の N - ベスト解を第 2 の候補群 8 2 から抽出し、第 3 の候補群 9 2 として出力する。このとき、第 3 の予備選択部 9 0 は、第 2 のコスト記憶部 8 4 に記憶された F 0 誤差、継続時間長誤差、及び M F C C 誤差により算出されたサブコストをサブコスト計算に用いる。したがって、実質的には第 3 の予備選択部 9 0 では音素環境誤差のみに基づくサブコスト計算が行なわれる。第 3 の予備選択部 9 0 によるサブコストの計算結果は第 3 のコスト記憶部 9 4 に記憶される。

20

【 0 0 4 7 】

第 4 の予備選択部 1 0 0 は、合成器指令 3 6 に基づいて、F 0 誤差、継続時間長誤差、M F C C 誤差、音素環境誤差、及び F 0 不連続誤差に基づいて算出されたサブコストの線形和に基づき、素片候補の N - ベスト解を第 3 の候補群 9 2 から抽出し、予備選択候補群 6 2 として出力する。このとき、第 4 の予備選択部 1 0 0 は、第 3 のコスト記憶部 9 4 に記憶された F 0 誤差、継続時間長誤差、M F C C 誤差、及び音素環境誤差により算出されたサブコストをサブコスト計算に用いる。したがって、実質的には第 4 の予備選択部 1 0 0 では F 0 不連続誤差のみに基づくサブコスト計算が行なわれる。第 4 の予備選択部 1 0 0 によるサブコストの計算結果は第 4 のコスト記憶部 1 0 4 に記憶される。

【 0 0 4 8 】

素片選択部 6 4 は、合成器指令 3 6 を受け、予備選択候補群 6 2 に含まれる音声素片のうち、式 (1) により算出されるコストが最も小さなものを選択して接続部 6 6 に与える。

30

【 0 0 4 9 】

接続部 6 6 は、素片選択部 6 4 により選択された音声素片に対応する音声素片データを音声素片 D B 3 4 から読出し、音声が滑らかに接続されるように変形して接続し、合成音声波形 4 0 として出力する。

【 0 0 5 0 】

多段予備選択部 6 0 により合成器指令 3 6 に対し計算されるサブコストが小さなものを予備的に選択しておくため、素片選択部 6 4 が式 (1) にしたがって素片選択を行なう際のコスト計算の計算量は少なく済む。多段予備選択部 6 0 内の各予備選択部 7 0、8 0、9 0 及び 1 0 0 によるサブコストの算出では、それぞれ前段でのサブコスト計算の結果を用いる。したがって各予備選択部 7 0、8 0、9 0 及び 1 0 0 における計算量は実質的には少なく済む。

40

【 0 0 5 1 】

また、予備選択部 7 0、8 0、9 0 及び 1 0 0 による予備選択では、徐々に選択の基準が細かくなっていくため、素片候補の限定は徐々に行なわれる。その結果、予備選択の段階で適切な素片候補が捨てられる危険性も低くなる。多段予備選択で得られた予備選択候補群 6 2 の中からコスト最小の音声素片を選択して接続した場合、接続時の変形による品質低下はほとんどない。その結果、最終的に得られる合成音声波形 4 0 にも、音声素片の

50

接続による品質低下はほとんどない。

【0052】

なお、上記した実施の形態では、予備選択部70、80、90及び100による4段階の多段予備選択を行なっているが、予備選択の各段階でのサブコスト計算及び段数がこの実施の形態に限定されないことはもちろんである。アプリケーションにより、種々の形で多段予備選択を行なうことができる。

【0053】

例えば、F0誤差、継続時間長誤差による予備選択での計算量が比較的大きいと思われる場合には、MFCC不連続誤差による音質への影響が比較的小さいことを考慮し、次のような多段予備選択を行なうことも考えられる。

【0054】

(1) 音素環境誤差を考慮したN-ベスト探索。

【0055】

(2) (1)により得られたN-ベスト解に対して、F0誤差、継続時間長誤差、及び音素環境誤差を考慮したN-ベスト探索。

【0056】

(3) (2)により得られたN-ベスト解に対して、F0誤差、継続時間長誤差、F0不連続誤差、及び音素環境誤差を考慮したN-ベスト探索。

【0057】

(4) (3)により得られたN-ベスト解に対して、F0誤差、継続時間長誤差、F0不連続誤差、MFCC不連続誤差、及び音素環境誤差を考慮したN-ベスト探索。

【0058】

(5) (4)により得られたN-ベスト解に対して、F0誤差、継続時間長誤差、F0不連続誤差、MFCC不連続誤差、MFCC誤差、及び音素環境誤差を考慮したN-ベスト探索。

【0059】

(6) (5)により得られたN-ベスト解に対して、全てのサブコストを考慮した1-ベスト探索による素片選択。

【0060】

上記した実施の形態の説明では、式(1)における重み p_1 及び p_2 の値をいずれも1として説明した。しかし、本発明はそのような実施の形態には限定されず、重み p_1 及び p_2 の値のいずれか、又は双方を1以外の値としてもよい。

【0061】

また、サブコスト関数も式(1)に示すものには限定されず、設計思想により様々なサブコスト関数を考えることができる。その場合も、多段予備選択でのサブコスト計算を考慮して、後段のサブコスト計算では前段のサブコスト計算の結果を使用できるようにすると効率がよい。

【0062】

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でのすべての変更を含む。

【図面の簡単な説明】

【0063】

【図1】本発明の第1の実施の形態にかかる音声合成システム20のブロック図である。

【符号の説明】

【0064】

20 音声合成システム、34 音声素片DB、36 合成器指令、38 音声合成装置、40 合成音声波形、60 多段予備選択部、62 予備選択候補群、64 素片選択部、66 接続部、70、80、90、100 予備選択部、72、82、92 候補

10

20

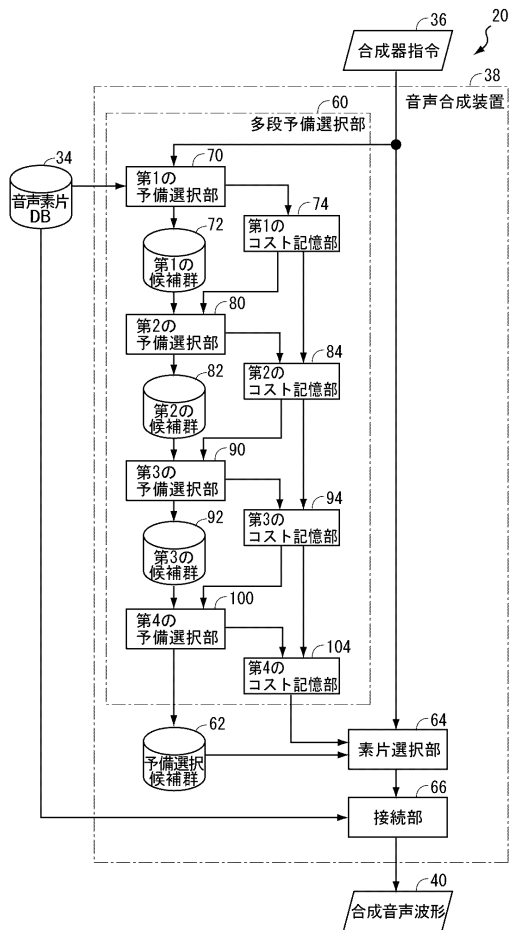
30

40

50

群、74, 84, 94, 104 コスト記憶部

【図1】



フロントページの続き

(56)参考文献 特開平08-263095(JP,A)

特開平08-248972(JP,A)

特開2003-208188(JP,A)

戸田 智基, 波形接続型音声合成における知覚的評価に基づく素片選択サブコスト関数の最適化
 , 電子情報通信学会技術研究報告, 日本, 社団法人電子情報通信学会, 2003年 8月15日
 , Vol.103 No.264

(58)調査した分野(Int.Cl., DB名)

G10L 13/06