

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4424024号
(P4424024)

(45) 発行日 平成22年3月3日(2010.3.3)

(24) 登録日 平成21年12月18日(2009.12.18)

(51) Int.Cl. F I
G 1 0 L 13/06 (2006.01) G 1 0 L 13/06 2 2 0 C

請求項の数 14 (全 16 頁)

<p>(21) 出願番号 特願2004-75185 (P2004-75185) (22) 出願日 平成16年3月16日 (2004.3.16) (65) 公開番号 特開2005-266010 (P2005-266010A) (43) 公開日 平成17年9月29日 (2005.9.29) 審査請求日 平成19年2月26日 (2007.2.26)</p> <p>特許権者において、実施許諾の用意がある。</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 (74) 代理人 100099933 弁理士 清水 敏 (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 (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 審査官 井上 健一</p> <p style="text-align: right;">最終頁に続く</p>
--	--

(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、複数の音声素片候補を含む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、各音声素片候補は、各音声素片の音響特徴量データと、各音声素片の波形データとを含み、

前記音声素片データベースを記憶するための第1の記憶装置と、

前記音声素片データベースに記憶された複数の音声素片候補の音響特徴量データと、前記音声素片データベースに記憶された複数の音声素片候補の中で、所定の基準で選択された音声素片候補の波形データとを記憶するための、前記第1の記憶装置より高速アクセス可能な第2の記憶装置とを含み、

前記複数のサブコストは、各音声素片候補の波形データが記憶されている記憶装置へのアクセス速度に関するアクセス速度コストを含み、

前記複数の音声素片候補の各々の音響特徴量データは、当該音声素片候補の波形データが前記第1及び第2の記憶装置のいずれに記憶されているかを示す第1のフラグを含み、前記音声合成装置はさらに、

前記合成音声の目標との間で、前記複数のサブコストを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、前記第2の記憶装置に記憶された音響特徴量に基づいて前記複数の音声素片候補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の音声波形を、当該選択された音声波形に

対応する第 1 のフラグに基づいて、前記第 1 の記憶装置又は前記第 2 の記憶装置のいずれかから読出して前記合成音声の目標に従って接続し、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】

前記第 1 の記憶装置に記憶された前記音声素片データベースは、当該音声素片データベースに含まれる複数の音声素片候補のうち、対応する波形データを前記第 2 の記憶装置に記憶すべき音声素片候補を選択するための基準となる選択基準情報が付されており、

前記素片接続型音声合成装置はさらに、前記第 1 の記憶装置に記憶された前記音声素片データベースに含まれる前記複数の音声素片候補のうち、前記選択基準情報により選択された音声素片候補の波形データを前記第 2 の記憶装置にロードするためのロード手段を含む、請求項 1 に記載の素片接続型音声合成装置。

10

【請求項 3】

テストデータに基づいて、前記音声素片データベースの前記選択基準情報を生成するための選択基準情報生成手段をさらに含む、請求項 2 に記載の素片接続型音声合成装置。

【請求項 4】

前記選択基準情報生成手段は、

前記テストデータに基づき、前記音声素片データベースに含まれる音声素片候補を使用して音声合成をシミュレートするための音声合成シミュレート手段と、

前記音声素片データベースに含まれる前記複数の音声素片候補の各々について、前記音声合成シミュレート手段による音声合成の際に選択された頻度を記録するための頻度記録手段とを含み、

20

前記選択基準情報は、前記頻度記録手段により記録された、前記音声素片データベースに含まれる前記複数の音声素片候補の各々の頻度情報である、請求項 3 に記載の素片接続型音声合成装置。

【請求項 5】

前記ロード手段は、前記頻度情報に基づき、前記音声素片データベースに含まれる音声素片のうち、前記音声合成シミュレート手段による音声合成で選択された頻度の高いものを上位から所定個数選択し、選択された音声素片の波形データを前記第 2 の記憶装置にロードするための手段を含む、請求項 4 に記載の素片接続型音声合成装置。

【請求項 6】

30

前記音声合成シミュレート手段は、前記テストデータから得られる合成音声の目標と、前記音声素片データベースに含まれる前記複数の音声素片候補との間で、前記複数のサブコストのうち、前記アクセス速度コストを除くアクセスコストと、前記頻度記録手段により記録された各音声素片候補が選択された頻度に基づいて算出される選択頻度コストとを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、前記音声素片データベースから選択するための手段を含む、請求項 4 又は請求項 5 に記載の素片接続型音声合成装置。

【請求項 7】

前記素片接続型音声合成装置は、

キャッシュメモリと、

40

前記キャッシュメモリに対して設けられたキャッシュメモリ管理機構と、

前記複数の音声素片候補に対応して設けられ、対応の音声素片候補が前記選択手段により選択されたか否かを記録するための第 2 のフラグを記憶するための手段と、

前記第 2 のフラグの値を所定の第 1 の値に初期化するための手段と、

前記複数の音声素片候補のいずれかが前記選択手段により選択されるたびに、選択された音声素片候補に対応する第 2 のフラグを前記第 1 の値と異なる所定の第 2 の値に更新するためのフラグ更新手段とをさらに含み、

前記アクセス速度コストは、前記第 1 及び第 2 のフラグに基づいて算出され、

前記選択手段は、前記合成音声の目標との間で、前記アクセス速度コストを含む前記複数のサブコストを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

50

、前記第2の記憶装置に記憶された音響特徴量に基づいて前記複数の音声素片候補から選択するための手段を含む、請求項1～請求項6のいずれかに記載の素片接続型音声合成装置。

【請求項8】

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、複数の音声素片候補を含む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、各音声素片候補は、各音声素片の音響特徴量データと、各音声素片の波形データとを含み、

前記音声素片データベースを第1の記憶ステップに記憶させるステップと、

前記音声素片データベースに記憶された複数の音声素片候補の音響特徴量データと、前記音声素片データベースに記憶された複数の音声素片候補の中で、所定の基準で選択された音声素片候補の波形データとを、前記第1の記憶装置より高速アクセス可能な第2の記憶装置に記憶させるステップとを含み、

前記複数のサブコストは、各音声素片候補の波形データが記憶されている記憶装置へのアクセス速度に関するアクセス速度コストを含み、

前記複数の音声素片候補の各々の音響特徴量データは、当該音声素片候補の波形データが前記第1及び第2の記憶装置のいずれに記憶されているかを示す第1のフラグを含み、

前記音声合成方法はさらに、

前記合成音声の目標との間で、前記複数のサブコストを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、前記第2の記憶装置に記憶された音響特徴量に基づいて前記複数の音声素片候補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の音声波形を、当該選択された音声波形に対応する第1のフラグに基づいて、前記第1の記憶装置又は前記第2の記憶装置のいずれかから読出して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9】

前記第1の記憶装置に記憶された前記音声素片データベースは、当該音声素片データベースに含まれる複数の音声素片候補のうち、対応する波形データを前記第2の記憶装置に記憶すべき音声素片候補を選択する基準となる選択基準情報が付されており、

前記素片接続型音声合成方法はさらに、前記第1の記憶装置に記憶された前記音声素片データベースに含まれる前記複数の音声素片候補のうち、前記選択基準情報により選択された音声素片候補の波形データを前記第2の記憶装置にロードするロードステップを含む、請求項8に記載の素片接続型音声合成方法。

【請求項10】

テストデータに基づいて、前記音声素片データベースの前記選択基準情報を生成する選択基準情報生成ステップをさらに含む、請求項9に記載の素片接続型音声合成方法。

【請求項11】

前記選択基準情報生成ステップは、

前記テストデータに基づき、前記音声素片データベースに含まれる音声素片候補を使用して音声合成をシミュレートする音声合成シミュレートステップと、

前記音声素片データベースに含まれる前記複数の音声素片候補の各々について、前記音声合成シミュレートステップによる音声合成の際に選択された頻度を記録する頻度記録ステップとを含み、

前記選択基準情報は、前記頻度記録ステップにおいて記録された、前記音声素片データベースに含まれる前記複数の音声素片候補の各々の頻度情報である、請求項10に記載の素片接続型音声合成方法。

【請求項12】

前記ロードステップは、前記頻度情報に基づき、前記音声素片データベースに含まれる音声素片のうち、前記音声合成シミュレートステップによる音声合成で選択された頻度の高いものを上位から所定個数選択し、選択された音声素片の波形データを前記第2の記憶装

置にロードするステップを含む、請求項 1 1 に記載の素片接続型音声合成方法。

【請求項 1 3】

前記音声合成シミュレートステップは、前記テストデータから得られる合成音声の目標と、前記音声素片データベースに含まれる前記複数の音声素片候補との間で、前記複数のサブコストのうち、前記アクセス速度コストを除くアクセスコストと、前記頻度記録ステップにおいて記録された各音声素片候補が選択された頻度に基づいて算出される選択頻度コストとを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、前記音声素片データベースから選択するステップを含む、請求項 1 1 又は請求項 1 2 に記載の素片接続型音声合成方法。

【請求項 1 4】

前記素片接続型音声合成方法は、キャッシュメモリと、前記キャッシュメモリに対して設けられたキャッシュメモリ管理機構とを有するコンピュータ上で実行され、さらに、

前記複数の音声素片候補に対応して設けられ、対応の音声素片候補が前記選択ステップにおいて選択されたか否かを記録する第 2 のフラグを所定の記憶装置に記憶するステップと、

前記第 2 のフラグの値を所定の第 1 の値に初期化するステップと、

前記複数の音声素片候補のいずれかが前記選択ステップにおいて選択されるたびに、選択された音声素片候補に対応する第 2 のフラグを前記第 1 の値と異なる所定の第 2 の値に更新するステップとをさらに含み、

前記アクセス速度コストは、前記第 1 及び第 2 のフラグに基づいて算出され、

前記選択ステップは、前記合成音声の目標との間で、前記アクセス速度コストを含む前記複数のサブコストを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、前記第 2 の記憶装置に記憶された音響特徴量に基づいて前記複数の音声素片候補から選択するステップを含む、請求項 8 ~ 請求項 1 3 のいずれか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。合成された音声の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

10

20

30

40

50

む場合もある。

【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。

10

【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。

20

【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。

30

【0012】

このような音声合成をリアルタイムで行なおうとする場合、いかにして素片選択と合成とを高速に行なうかが問題となる。この処理を高速化するためには、素片選択のコスト計算を高速に行なうとともに、選択された音声素片を接続する処理も高速にすることが望ましい。

【0013】

実際に音声素片を接続する際には音声素片の波形データが必要となるが、個々の波形データのデータ量が比較的大きく、また、合成音声の品質を高くするためには、音声素片データベースに格納される音声素片の数を大きくする必要があり、その結果、音声素片データベース全体の容量は大きくなる。従って従来は、素片データのうち音響特徴量のみをメモリに格納して素片選択のコスト計算を行ない、音声素片データベースは、固定ハードディスク等比較的容量が大きな記憶装置に格納しておき、素片が選択された後、素片の接続時に波形データを読み出すようにしている。

40

【0014】

【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、容量が大きな記憶装置のアクセス速度は比較的低速である。そのため、素片選択後、音声素片の波形データの取得に要する時間が大きく、音声合成に要する時間も長く

50

かかるという問題があった。

【0016】

それゆえに、本発明の目的は、利用可能な音声素片の数を多く保ったまま、高速に音声合成を行なうことが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、複数の音声素片候補を含む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、各音声素片候補は、各音声素片の音響特徴量データと、各音声素片の波形データとを含み、音声素片データベースを記憶するための第1の記憶装置と、音声素片データベースに記憶された複数の音声素片候補の音響特徴量データと、音声素片データベースに記憶された複数の音声素片候補の中で、所定の基準で選択された音声素片候補の波形データとを記憶するための、第1の記憶装置より高速アクセス可能な第2の記憶装置とを含み、複数のサブコストは、各音声素片候補の波形データが記憶されている記憶装置へのアクセス速度に関するアクセス速度コストを含み、複数の音声素片候補の各々の音響特徴量データは、当該音声素片候補の波形データが第1及び第2の記憶装置のいずれに記憶されているかを示す第1のフラグを含み、音声合成装置はさらに、合成音声の目標との間で、複数のサブコストを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、第2の記憶装置に記憶された音響特徴量に基づいて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、当該選択された音声波形に対応する第1のフラグに基づいて、第1の記憶装置又は第2の記憶装置のいずれかから読出して合成音声の目標に従って接続し、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、第1の記憶装置に記憶された音声素片データベースは、当該音声素片データベースに含まれる複数の音声素片候補のうち、対応する波形データを第2の記憶装置に記憶すべき音声素片候補を選択するための基準となる選択基準情報が付されており、素片接続型音声合成装置はさらに、第1の記憶装置に記憶された音声素片データベースに含まれる複数の音声素片候補のうち、選択基準情報により選択された音声素片候補の波形データを第2の記憶装置にロードするためのロード手段を含む。

【0019】

さらに好ましくは、テストデータに基づいて、音声素片データベースの選択基準情報を生成するための選択基準情報生成手段をさらに含む。

【0020】

選択基準情報生成手段は、テストデータに基づき、音声素片データベースに含まれる音声素片候補を使用して音声合成をシミュレートするための音声合成シミュレート手段と、音声素片データベースに含まれる複数の音声素片候補の各々について、音声合成シミュレート手段による音声合成の際に選択された頻度を記録するための頻度記録手段とを含んでもよい。選択基準情報は、頻度記録手段により記録された、音声素片データベースに含まれる複数の音声素片候補の各々の頻度情報でもよい。

【0021】

好ましくは、ロード手段は、頻度情報に基づき、音声素片データベースに含まれる音声素片のうち、音声合成シミュレート手段による音声合成で選択された頻度の高いものを上位から所定個数選択し、選択された音声素片の波形データを第2の記憶装置にロードするための手段を含む。

【0022】

さらに好ましくは、音声合成シミュレート手段は、テストデータから得られる合成音声の目標と、音声素片データベースに含まれる複数の音声素片候補との間で、複数のサブコストのうち、アクセス速度コストを除くアクセスコストと、頻度記録手段により記録され

10

20

30

40

50

た各音声素片候補が選択された頻度に基づいて算出される選択頻度コストとを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、音声素片データベースから選択するための手段を含む。

【0023】

素片接続型音声合成装置は、キャッシュメモリと、キャッシュメモリに対して設けられたキャッシュメモリ管理機構と、複数の音声素片候補に対応して設けられ、対応の音声素片候補が選択手段により選択されたか否かを記録するための第2のフラグを記憶するための手段と、第2のフラグの値を所定の第1の値に初期化するための手段と、複数の音声素片候補のいずれかが選択手段により選択されるたびに、選択された音声素片候補に対応する第2のフラグを第1の値と異なる所定の第2の値に更新するためのフラグ更新手段とをさらに含んでもよい。アクセス速度コストは、第1及び第2のフラグに基づいて算出され、選択手段は、合成音声の目標との間で、アクセス速度コストを含む複数のサブコストを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、第2の記憶装置に記憶された音響特徴量に基づいて複数の音声素片候補から選択するための手段を含んでもよい。

10

【0024】

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、複数の音声素片候補を含む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、各音声素片候補は、各音声素片の音響特徴量データと、各音声素片の波形データとを含み、音声素片データベースを第1の記憶装置に記憶させるステップと、音声素片データベースに記憶された複数の音声素片候補の音響特徴量データと、音声素片データベースに記憶された複数の音声素片候補の中で、所定の基準で選択された音声素片候補の波形データとを、第1の記憶装置より高速アクセス可能な第2の記憶装置に記憶させるステップとを含み、複数のサブコストは、各音声素片候補の波形データが記憶されている記憶装置へのアクセス速度に関するアクセス速度コストを含み、複数の音声素片候補の各々の音響特徴量データは、当該音声素片候補の波形データが第1及び第2の記憶装置のいずれに記憶されているかを示す第1のフラグを含み、音声合成方法はさらに、合成音声の目標との間で、複数のサブコストを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、第2の記憶装置に記憶された音響特徴量に基づいて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、当該選択された音声波形に対応する第1のフラグに基づいて、第1の記憶装置又は第2の記憶装置のいずれかから読出して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。

20

30

【0025】

好ましくは、第1の記憶装置に記憶された音声素片データベースは、当該音声素片データベースに含まれる複数の音声素片候補のうち、対応する波形データを第2の記憶装置に記憶すべき音声素片候補を選択する基準となる選択基準情報を有し、素片接続型音声合成方法はさらに、第1の記憶装置に記憶された音声素片データベースに含まれる複数の音声素片候補のうち、選択基準情報により選択された音声素片候補の波形データを第2の記憶装置にロードするロードステップを含む。

40

【0026】

さらに好ましくは、素片接続型音声合成方法はテストデータに基づいて、音声素片データベースの選択基準情報を生成する選択基準情報生成ステップをさらに含む。

【0027】

選択基準情報生成ステップは、テストデータに基づき、音声素片データベースに含まれる音声素片候補を使用して音声合成をシミュレートする音声合成シミュレートステップと、音声素片データベースに含まれる複数の音声素片候補の各々について、音声合成シミュレートステップによる音声合成の際に選択された頻度を記録する頻度記録ステップとを含んでもよい。選択基準情報は、頻度記録ステップにおいて記録された、音声素片データベ

50

ースに含まれる複数の音声素片候補の各々の頻度情報でもよい。

【0028】

好ましくは、ロードステップは、頻度情報に基づき、音声素片データベースに含まれる音声素片のうち、音声合成シミュレートステップによる音声合成で選択された頻度の高いものを上位から所定個数選択し、選択された音声素片の波形データを第2の記憶装置にロードするステップを含む。

【0029】

さらに好ましくは、音声合成シミュレートステップは、テストデータから得られる合成音声の目標と、音声素片データベースに含まれる複数の音声素片候補との間で、複数のサブコストのうち、アクセス速度コストを除くアクセスコストと、頻度記録ステップにおいて記録された各音声素片候補が選択された頻度に基づいて算出される選択頻度コストとを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、音声素片データベースから選択するステップを含む。

【0030】

素片接続型音声合成方法は、キャッシュメモリと、キャッシュメモリに対して設けられたキャッシュメモリ管理機構とを有するコンピュータ上で実行される方法でもよい。素片接続型音声合成方法はさらに、複数の音声素片候補に対応して設けられ、対応の音声素片候補が選択ステップにおいて選択されたか否かを記録する第2のフラグを所定の記憶装置に記憶するステップと、第2のフラグの値を所定の第1の値に初期化するステップと、複数の音声素片候補のいずれかが選択ステップにおいて選択されるたびに、選択された音声素片候補に対応する第2のフラグを第1の値と異なる所定の第2の値に更新するステップとをさらに含んでもよい。アクセス速度コストは、第1及び第2のフラグに基づいて算出される。選択ステップは、合成音声の目標との間で、アクセス速度コストを含む複数のサブコストを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、第2の記憶装置に記憶された音響特徴量に基づいて複数の音声素片候補から選択するステップを含んでもよい。

【発明を実施するための最良の形態】

【0031】

素片選択型音声合成において、実際に選ばれる音声素片には偏りが生じる。従って、選ばれやすい音声素片は音声素片データベースより高速にアクセス可能な記憶装置（例えばメモリ）に格納しておくことで、音声合成の速度を全体として上げることができる。さらに、素片選択の際に、各音声素片の波形データが記憶されている記憶装置のアクセス速度をコストに加える。記憶装置のアクセス速度に対応するコストを、本明細書では「アクセス速度コスト」と呼ぶ。アクセス速度コストは、記憶装置のアクセス速度が高いほど小さく（0に近く）、低いほど大きくなるように、予め算出式を設計する。高速な記憶装置に波形データが記憶されている音声素片ほど、実際に音声合成で選択される可能性が高くなり、音声合成の速度を全体として高くすることができる。

【0032】

図1に、本発明の一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DB30と、多数の音声合成のためのテキストからなるテストデータ42を使用して、音声素片DB30を用いた音声合成をシミュレートし、音声素片DB30に含まれる音声素片ごとに音声合成で使用された頻度を算出して、頻度情報44を有する頻度情報付き音声素片DB34を生成するための頻度情報生成装置32とを含む。

【0033】

音声合成システム20はさらに、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、頻度情報付き音声素片DB34に含まれる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を含む。音声合成装置38は、素片選択のための、頻度情報付き音声素片DB34内の各音声素片の音響情報と、頻度情報付き音声素片DB34内の音声素片のうち、頻度情報44により

10

20

30

40

50

表される出現頻度が高いものを予め記憶しておくためのメモリ 48 を含む。

【0034】

音声合成システム 20 はさらに、音声合成装置 38 の起動時に、頻度情報付き音声素片 DB 34 内の各音声素片の音響情報と、頻度情報 44 を参照することにより、頻度情報付き音声素片 DB 34の中から選択した出現頻度の上位の所定個数の音声素片の波形データとをメモリ 48 にロードするための音声素片データロード装置 46 とを含む。

【0035】

図 2 を参照して、頻度情報生成装置 32 は、音声素片 DB 30 を使用し、テストデータ 42 に含まれる各入力に対して実際に音声合成と同様の処理をして、各音声素片の使用頻度を算出する機能を持つ。頻度情報生成装置 32 は、テストデータ 42 の各テキスト文を受け、合成目標となる合成器指令 62 を作成するための合成器指令作成部 60 と、この合成器指令 62 と音声素片 DB 30 中の各音声素片の音響特徴量とのターゲットコストを算出するためのターゲットコスト算出部 68 と、合成器指令 62 と音声素片 DB 30 中の各音声素片の音響特徴量との接続コストを算出するための接続コスト算出部 70 と、音声素片 DB 30 中の各音声素片が音声合成で選択された頻度を反映する選択頻度コストを算出するための選択頻度コスト算出部 72 とを含む。

【0036】

選択頻度コストとは、頻度情報生成装置 32 において選択される音声素片に対し、意図的に偏りを生じさせるために導入されたコストである。選択頻度コストは、音声素片が選択された頻度が高くなるほど小さく、低くなるほど大きくなるような算出式で算出される。本実施の形態では、 i 番目の音声素片の選択頻度コスト Cs_i を次の式により算出する。

【0037】

【数 1】

$$Cs_i = a(1 - (\frac{n_i}{N_0 + N})^r) \quad (1)$$

ただし n_i は i 番目の音声素片の選択回数、 N は全ての素片の出現回数の和、 a 、 N_0 、及び r は適当な定数である（ただし $r > 0$ 、 $a > 0$ 、及び $N > 0$ ）。この式によれば、 i 番目の音声素片の選択頻度コスト Cs_i の値は、選択頻度 n_i が 0 であれば a であり、選択頻度 n_i が多くなるにつれて 0 に近づいて行く。選択頻度コスト算出部 72 はこのため、頻度情報 44 に記憶されている各音声素片の頻度情報を使用する。

【0038】

頻度情報生成装置 32 はさらに、ターゲットコスト算出部 68 により算出されたターゲットコストと、接続コスト算出部 70 により算出された接続コストと、選択頻度コスト算出部 72 により算出された選択頻度コストとに基づいて総コストを算出し、総コストの最も小さな音声素片を選択することにより、実際の音声合成時の素片選択をシミュレート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に関し、頻度情報 44 を更新するための頻度情報更新部 66 とを含む。

【0039】

図 3 を参照して、図 1 に示す音声素片データロード装置 46 がメモリ 48 にロードするデータについて説明する。図 3 に示すように、メモリ 48 は、頻度情報付き音声素片 DB 34 の全ての音声素片の音響特徴量 130、... を格納するための音響特徴量格納領域 120 と、頻度情報付き音声素片 DB 34 内の音声素片のうち、頻度情報 44 に記録された出現頻度が所定の値以上のものの波形データを記憶するための波形データ格納領域 122 とを含む。これらはいずれも音声素片データロード装置 46 により、音声合成装置 38 の起動時にメモリ 48 にロードされる。

【0040】

音響特徴量 130 の各々は、前述したとおり、音素ラベル、基本周波数 (F_0)、MFCC、パワー（図示せず）、持続時間（図示せず）を含むが、これらに加えて、音声素片

10

20

30

40

50

の波形データが頻度情報付き音声素片DB34に格納されているか、メモリ48の波形データ格納領域122に格納されているかをあらわす第1のフラグ(F₁)140と、音声素片の波形データが最近読出されたか否かを示す第2のフラグ(F₂)142とを含む。フラグ140は、音声素片データロード装置46が音響特徴量格納領域120に音響特徴量をロードした後、高頻度の音声素片の波形データを波形データ格納領域122にロードする際に、音声素片データロード装置46によって設定される。

【0041】

本実施の形態では、第1のフラグ140が0の場合には波形データは頻度情報付き音声素片DB34に格納されていることを表し、フラグが1の場合には波形データがメモリ48の波形データ格納領域122に格納されていることを示す。従ってこの第1のフラグ140の値を見ることで波形データをどこから読出せばよいかを判定できる。

10

【0042】

フラグ142は最初に0に初期化され、実際に音声合成を行ないながら、波形データが読出されたときに「1」に更新される。これは、装置をコンピュータで実現する場合、ハードディスク又はメモリから読出されたデータはメモリよりもさらに高速アクセス可能なキャッシュに格納されることがあることを考慮したものである。この第2のフラグ142の値はアクセス速度コストに反映される。

【0043】

図4を参照して、音声合成装置38は、前述のメモリ48に加え、それぞれ合成目標を定める合成器指令36を受け、メモリ48に記憶されている音声素片であって、かつ合成器指令36により指定された音素ラベルを持つ音声素片の音響特徴量と合成器指令36との間のターゲットコストを算出するためのターゲットコスト算出部82と、同じくメモリ48内の音声素片と合成器指令36との間の接続コストを算出するための接続コスト算出部84と、メモリ48に記憶されている音声素片に対応する第1及び第2のフラグ140、142に基づいてアクセス速度コストを算出するためのアクセス速度コスト算出部86とを含む。

20

【0044】

音声合成装置38はさらに、ターゲットコスト算出部82により算出されたターゲットコスト、接続コスト算出部84により算出された接続コスト、及びアクセス速度コスト算出部86により算出されたアクセス速度コストに基づいて総コストを算出し、総コストの最小の音声素片を選択するための素片選択部80と、素片選択部80により選択された音声素片の波形データを接続して合成音声波形40を出力するための接続部88とを含む。

30

【0045】

接続部88は、素片選択部80により指定された音声素片を頻度情報付き音声素片DB34又はメモリ48から読出すため、次のような信号を出力する。すなわち、接続部88は、素片選択部80により指定された音声素片の第1のフラグ140に対応するレベルをとり、波形データをメモリ48と頻度情報付き音声素片DB34とのいずれから読出すかを指定するための選択信号100と、波形データを読出すアドレスを指定するアドレス信号102とを出力する機能を持つ。選択信号100は、指定された音声素片のフラグが1のときにはHレベルをとり、それ以外のときにはLレベルをとる。

40

【0046】

音声合成装置38は、接続部88による波形データの読出を行なうための機能ブロックとして、第1及び第2の入力を持ち、選択信号100のレベルに応じて第1及び第2の入力の信号のいずれかを選択して出力するための選択回路90と、選択信号100のレベルを反転して反転選択信号104を出力するための反転回路92と、Hレベルの反転選択信号104を受けると、アドレス信号102により指定されるアドレスの波形データを頻度情報付き音声素片DB34から読出して選択回路90の第1の入力に与えるためのアクセス部94とを含む。一方、メモリ48は、Hレベルの選択信号100を受けると、アドレス信号102により指定されるアドレスの波形データを選択回路90の第2の入力に与える。選択回路90は、選択信号100がHレベルのときは第2の入力の信号を、Lレベル

50

のときには第 1 の入力 of 信号を選択して出力する。

【 0 0 4 7 】

音声合成装置 3 8 はさらに、接続部 8 8 により読出が指示された波形データについて、メモリ 4 8 のうち、対応する音声素片の第 2 のフラグ 1 4 2 (F_2) を「 1 」に更新するためのフラグ更新部 9 6 を含む。頻度情報付き音声素片 D B 3 4 又はメモリ 4 8 から読出されたデータは、いずれの場合も、コンピュータのキャッシュメモリに格納されることが通常である。キャッシュメモリはメモリ 4 8 と比較してもさらに高速にアクセス可能である。一度でも読出された波形データはキャッシュメモリに格納されている可能性が高いので、このように第 2 のフラグを更新し、次のアクセス速度算出部でのコスト計算に反映させ、より選択されやすくする。

10

【 0 0 4 8 】

なお、キャッシュメモリの容量には限りがあるため、何らかのアルゴリズムによってキャッシュに格納されているデータを選択して削除し、そこに新しいデータを格納する。本来であればキャッシュメモリにどの波形データが格納されているかを把握できればよいが、キャッシュメモリはコンピュータハードウェアにより管理されており、ソフトウェアでキャッシュの内容について知ることはできない。従って本実施の形態では、キャッシュメモリに実際にどのようなデータが格納されているかとは別に、一度でも読出されたことのあるデータについてはキャッシュメモリに格納されているものと想定した設計としている。もちろん、キャッシュメモリに記憶されている波形データがどの音声素片に対応するものであるかを容易に知ることができれば、それをアクセス速度コストの計算に反映させる

20

【 0 0 4 9 】

本実施の形態で使用されるサブコストは、前述したアクセス速度コスト以外に、基本周波数 (F_0) 誤差、継続長誤差、M F C C 誤差、 F_0 不連続誤差、M F C C 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。

【 0 0 5 0 】

本実施の形態に係る素片選択部 6 4 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。

【 0 0 5 1 】

【 数 2 】

$$C_o = \left(\sum_{i1=1}^{N_1} (w_{i1} C_{i1})^{p_1} \right)^{\frac{1}{p_1}} + \left(\sum_{i2=1}^{N_2} (w_{i2} C_{i2})^{p_2} \right)^{\frac{1}{p_2}} + w_s (2 - F_1 - F_2) \quad (2)$$

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、 w_s はアクセス速度コストの重み、 F_1 及び F_2 はそれぞれ第 1 及び第 2 のフラグの値である。また重み w_s の大きさは、他のサブコストの値の範囲を考慮して適当に決定し、品質に影響を及ぼす極端な偏りが生じないようにする。

40

【 0 0 5 2 】

上記した第 1 の実施の形態に係る音声合成システム 2 0 は以下のように動作する。大きく分けてこの音声合成システム 2 0 の動作には二つの局面がある。第 1 の局面は頻度情報付き音声素片 D B 3 4 の作成であり、第 2 の局面は第 1 の局面で作成された頻度情報付き音声素片 D B 3 4 を用いた音声合成である。以下、順に説明する。

【 0 0 5 3 】

まず第 1 の局面では頻度情報生成装置 3 2 が以下のようにして頻度情報付き音声素片 D B 3 4 を作成する。図 2 を参照して、まず音声素片 D B 3 0 を用意する。音声素片 D B 3 0 には頻度情報は付されていない。また、頻度情報 4 4 を記憶すべき領域をメモリ上に確

50

保しておく。さらに、頻度情報 4 4 中の、各音声素片の選択頻度回数を全て 0 に初期化する。

【 0 0 5 4 】

テストデータ 4 2 の第 1 のテキストを頻度情報生成装置 3 2 に与えると、合成器指令作成部 6 0 がそのテキストに基づいて合成目標の音素ごとに合成器指令 6 2 を作成し、素片選択部 6 4 に与える。素片選択部 6 4 は、この合成器指令 6 2 により指定された音響特徴量をターゲットコスト算出部 6 8 及び接続コスト算出部 7 0 に与える。ターゲットコスト算出部 6 8 及び接続コスト算出部 7 0 は、与えられた音響特徴量に基づき、音声素片 DB 3 0 に含まれる各音声素片との間でターゲットコスト及び接続コストを算出し素片選択部 6 4 に与える。選択頻度コスト算出部 7 2 は、式 (1) に従って各音声素片候補の選択頻度コストを算出し素片選択部 6 4 に与える。第 1 回目の処理では選択頻度はいずれも 0 であるから、選択頻度コストはいずれも式 (1) より「 a 」となる。

10

【 0 0 5 5 】

素片選択部 6 4 は、ターゲットコスト算出部 6 8 から与えられたターゲットコスト、接続コスト算出部 7 0 から与えられた接続コスト、及び選択頻度コスト算出部 7 2 から与えられた選択頻度コストから式 (2) によって総コストを算出する。素片選択部 6 4 はこの総コストが最小の音声素片を選択し、選択された音声素片を示す情報を頻度情報更新部 6 6 に与える。

【 0 0 5 6 】

頻度情報更新部 6 6 は、頻度情報 4 4 中の頻度情報のうち、素片選択部 6 4 により選択された音素の頻度に 1 を加算する。以上で第 1 番目の音素に対する処理を終了する。

20

【 0 0 5 7 】

同様の処理を、最初のテキストの各音素に対して繰り返す。この繰り返しにより、頻度情報 4 4 は徐々に更新されていく。選択頻度コスト算出部 7 2 による選択頻度コストの算出においては、頻度情報 4 4 の内容が反映される。すなわち、選択された回数が増えるほど選択頻度コストは小さくなる。従って、選択されたことのある素片候補についてはその後の素片選択で選択される可能性が高くなる。その結果、こうした処理を繰り返すと、互いによく似た音響特徴量をもつ音声素片同士であって、選択されたことのある素片候補はさらに選択されやすく、選択されたことのない候補はさらに選択されにくくなる。

【 0 0 5 8 】

テストデータ 4 2 の全てのテキストについて上記した処理を繰り返すことにより、頻度情報付き音声素片 DB 3 4 が完成する。頻度情報付き音声素片 DB 3 4 が完成すると、音声合成装置 3 8 による音声合成が可能となる。

30

【 0 0 5 9 】

音声素片の出現頻度を意図的に偏らせることにより、特徴空間において素片の密度が高い領域において、一部の素片のみがよく選択されるようになり、それ以外の素片の出現頻度は下がる。その分、素片の密度が低い部分でそれ以外の素片の頻度の順位が相対的に上がる。これにより、実際の合成時に高速アクセス可能な記憶装置に波形データが格納される音声素片の分布が広がり、高速な記憶装置に格納された波形データに対応する音声素片がより頻繁に選択されるようになる。

40

【 0 0 6 0 】

実際の音声合成は以下のようにして行なわれる。図 1 を参照して、最初に音声素片データロード装置 4 6 により、頻度情報付き音声素片 DB 3 4 中の音声素片の音響特徴量がメモリ 4 8 に格納される。音響特徴量に付随する第 1 及び第 2 のフラグの値は 0 で初期化される。次に、頻度情報 4 4 を基準とし、上位の所定個数の音声素片の波形データがメモリ 4 8 に格納される。波形データをメモリ 4 8 にロードした音声素片については、メモリ 4 8 に格納された音響特徴量に付随する第 1 のフラグ 1 4 0 (図 3 参照) の値を「 1 」に設定する。

【 0 0 6 1 】

音声素片データロード装置 4 6 によるメモリ 4 8 へのデータのロードが終わると、実際

50

の音声合成が開始される。図4を参照して、合成器指令36が与えられると、素片選択部80はターゲットコスト算出部82、接続コスト算出部84及びアクセス速度コスト算出部86に合成器指令62により指定された音響特徴量を与える。ターゲットコスト算出部82及び接続コスト算出部84は、与えられた音響特徴量を用い、メモリ48に格納されている音声素片候補のうち、指定された音素ラベルの音声素片候補の音響特徴量との間でターゲットコスト及び接続コストを算出し、素片選択部80に与える。アクセス速度コスト算出部86は、各音声素片候補の第1のフラグ140及び第2のフラグ142の値に基づきアクセス速度コストを算出し素片選択部80に与える。

【0062】

素片選択部80は、ターゲットコスト算出部82、接続コスト算出部84、及びアクセス速度コスト算出部86から与えられたターゲットコスト、接続コスト、及びアクセスコストに基づき、式(2)に従って総コストを算出する。素片選択部80はさらに、そのようにして算出された総コストが最小の音声素片候補を選択し、その音声素片候補を示す情報と、音響特徴量とを接続部88に与える。接続部88は、与えられた音響特徴量の中の、波形データアドレスをアドレス信号102に、第1のフラグ140を選択信号100に、それぞれ出力する。

10

【0063】

例えば選択された音声素片の波形データがメモリ48に格納されている場合、その音声素片の第1のフラグの値は1であり、選択信号100はHレベルとなる。アドレス信号102はメモリ48の、選択された音声素片の波形データのアドレスとなる。アドレス信号102はメモリ48に与えられる。Hレベルの選択信号100がメモリ48に与えられるので、メモリ48はアドレス信号102により指定されるアドレスの波形データを読み出し、選択回路90の第2の入力に与える。反転選択信号104はLレベルなのでアクセス部94は何もしない。

20

【0064】

選択信号100がHレベルなので、選択回路90は第2の入力を選択する。すなわち、選択回路90はメモリ48からの出力を接続部88に与える。接続部88はこの波形データを用いて波形接続を行なう。

【0065】

また、選択された音声素片の波形データが頻度情報付き音声素片DB34に格納されている場合、その音声素片の第1のフラグの値は0であり、選択信号100はLレベルとなる。アドレス信号102は頻度情報付き音声素片DB34の、選択された音声素片の波形データのアドレスとなる。アドレス信号102はアクセス部94に与えられる。Hレベルの反転選択信号104がアクセス部94に与えられるので、アクセス部94はアドレス信号102により指定されるアドレスの波形データを頻度情報付き音声素片DB34から読み出し、選択回路90の第1の入力に与える。選択信号100はLレベルなのでメモリ48からは何も出力されない。

30

【0066】

選択信号100がLレベルなので選択回路90は第1の入力の信号を選択して接続部88に与える。すなわちこの場合、頻度情報付き音声素片DB34から読み出された波形データが接続部88に与えられ、波形接続に用いられる。

40

【0067】

メモリ48の波形データにせよ、頻度情報付き音声素片DB34の波形データにせよ、最近接続部88により読み出されたものは図示しないキャッシュメモリに格納される可能性が高い。従ってフラグ更新部96は、メモリ48中の、接続部88により波形データが読み出された音声素片の音響特徴量に付随する第2のフラグ142の値を「1」に更新する。この結果、同じ音声素片が次に選択された場合、アクセス速度コストはより小さくなり、同じ音声素片が選択される可能性が高くなる。この音声素片に対応する波形データはキャッシュに格納されていて高速アクセス可能である可能性が高い。従って全体として波形データの読み出しが高速化される可能性が高い。

50

【 0 0 6 8 】

以上のように本実施の形態に係る音声合成システム 20 では、音声素片 DB 30 を使用してテストデータ 42 による音声合成実験を行なって、各音声頻度が選択された頻度を調べる。音声合成実験の途中では、音声素片が選択された回数に応じ、選択された素片はさらに選択されやすく、そうでない素片はさらに選択されにくくなるように、選択頻度コストというサブコストを導入し、音声素片の選択に人為的な偏りが生じるようにする。

【 0 0 6 9 】

音声合成時には、この頻度が上位の所定個数の音声素片候補の波形データをメモリに記憶しておく。さらに、サブコストとして、波形データが記憶されている記憶装置のアクセス速度を反映したアクセス速度コストを定義し、高速アクセス可能な記憶装置に波形データが記憶されている音声素片候補が選択されやすくする。その結果、メモリなどの高速アクセス可能な記憶媒体に記憶された音声素片が選択されやすくなり、全体として音声合成処理が高速化される。

【 0 0 7 0 】

さらに、キャッシュメモリなどが使用されることを考慮し、最近選択された音声波形についてはアクセス速度コストが少なく算出されるようにアクセス速度コストを設計する。これにより、コンピュータによるキャッシュ制御を利用した処理速度の向上を図ることができる。

【 0 0 7 1 】

また、上記した実施の形態の装置では、一旦頻度情報 44 が作成された後は、頻度情報 44 を更新しないことを前提としている。しかし本発明はそのような実施の形態には限定されない。例えば音声合成時、頻度情報 44 のをメモリ 48 に転記し、素片候補が選択されるたびに頻度情報 44 を更新し、全体の処理が終了したとき、または所定回数の素片候補の選択が行なわれるたびごとに、もとの頻度情報 44 に書き戻すようにしてもよい。こうすることで、実験だけでなく実際のデータに基づく音声合成での選択頻度に基づいて、波形データの格納場所を決定することができる。

【 0 0 7 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でのすべての変更を含む。

【 図面の簡単な説明 】

【 0 0 7 3 】

【 図 1 】 本発明の一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。

【 図 2 】 図 1 に示す頻度情報生成装置 32 のブロック図である。

【 図 3 】 図 1 に示すメモリ 48 の一部の記憶領域の構成を模式的に示す図である。

【 図 4 】 図 1 に示す音声合成装置 38 のブロック図である。

【 符号の説明 】

【 0 0 7 4 】

20 音声合成システム、30 音声素片 DB、32 頻度情報生成装置、34 頻度情報付き音声素片 DB、36, 62 合成器指令、38 音声合成装置、40 合成音声波形、42 テストデータ、44 頻度情報、46 音声素片データロード装置、48 メモリ、60 合成器指令作成部、64, 80 素片選択部、66 頻度情報更新部、68, 82 ターゲットコスト算出部、70, 84 接続コスト算出部、72 選択頻度コスト算出部、86 アクセス速度コスト算出部 88 接続部、90 選択回路、100 選択信号、102 アドレス信号、104 反転選択信号

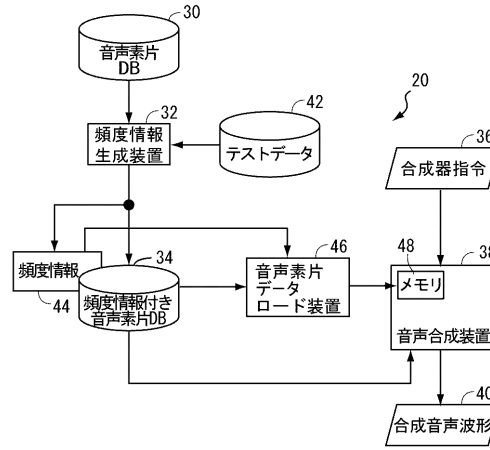
10

20

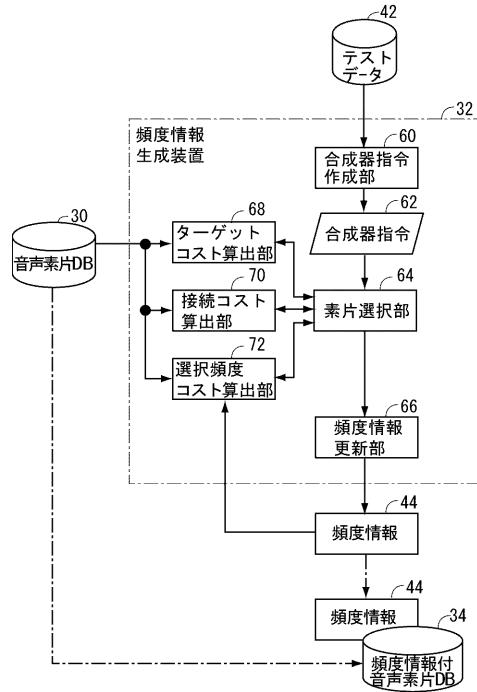
30

40

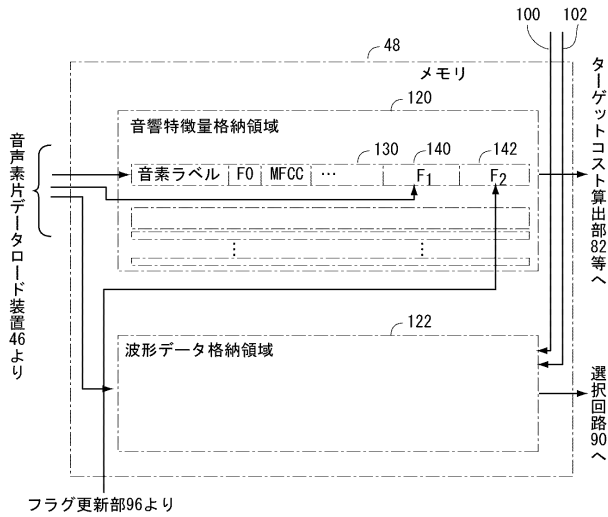
【図1】



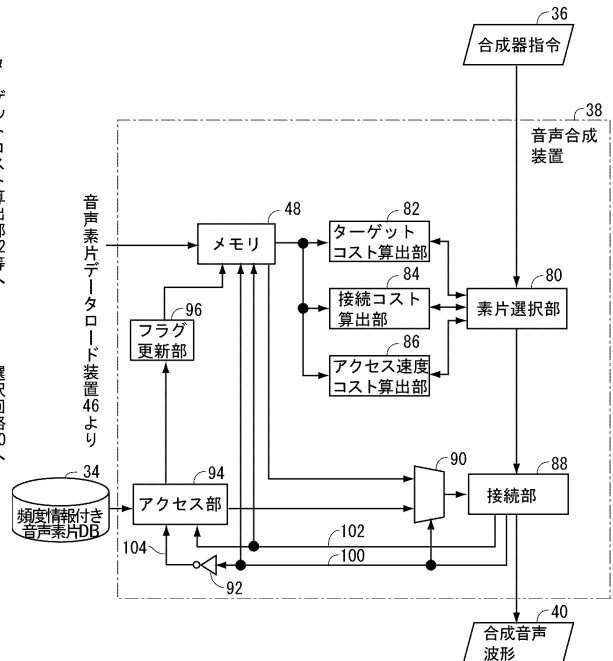
【図2】



【図3】



【図4】



フロントページの続き

- (56)参考文献 特開昭63-052199(JP,A)
特開2000-075877(JP,A)
特開平07-210194(JP,A)
特開2001-282273(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G10L 13/06