

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4292191号
(P4292191)

(45) 発行日 平成21年7月8日(2009.7.8)

(24) 登録日 平成21年4月10日(2009.4.10)

(51) Int.Cl. F I
G 1 0 L 13/06 (2006.01) G 1 0 L 13/06 2 4 0 C

請求項の数 4 (全 18 頁)

<p>(21) 出願番号 特願2006-57304 (P2006-57304) (22) 出願日 平成18年3月3日(2006.3.3) (65) 公開番号 特開2007-233216 (P2007-233216A) (43) 公開日 平成19年9月13日(2007.9.13) 審査請求日 平成21年3月2日(2009.3.2)</p> <p>(出願人による申告)平成17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 (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</p> <p>審査官 井上 健一</p>
---	---

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コンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

多数の音声素片データを格納した音声素片データベースとともに用いられる素片接続型音声合成装置であって、

合成ターゲットが与えられると、当該合成ターゲットを構成する各ターゲット音素のコンテキストに基づいて、音声合成において各ターゲット音素の合成に用いられるべき候補として予備選択されるべき音声素片データの数を予測するための素片候補数予測手段と、

合成ターゲットが与えられると、当該合成ターゲットを構成する各ターゲット音素について、当該ターゲット音素と前記音声素片データベース中の音声素片データの各々との間に算出されるターゲットコストに基づいて、前記音声素片データベース中から、前記素片候補数予測手段により予測された数と所定の関係にある数の音声素片データを、前記各ターゲットの音声合成のために予備的に選択するための素片候補予備選択手段と、

合成ターゲットを構成する各ターゲット音素について、前記素片候補予備選択手段により予備的に選択された音声素片データの各々との間に算出されるターゲットコスト及び接続コストに基づいて、音声合成に用いるべき音声素片データを選択するための素片選択手段と、

前記素片選択手段により選択された音声素片データの音声波形を前記合成ターゲットに従って接続するための波形接続手段とを含む、音声合成装置。

【請求項2】

前記素片候補数予測手段は、

各ターゲット音素のコンテキストに基づいて、音声合成において各ターゲット音素の合成に用いられるべき候補として予備選択されるべき音声素片データの数を、予め準備された回帰木を用いて予測するための回帰木による予測手段を含み、

当該回帰木は、一つのルートノードと、複数の葉ノードと、前記ルートノードと前記葉ノードとの間に存在する複数の中間ノードとを含み、

前記ルートノードと前記複数の中間ノードとの各々には、ターゲット音素のコンテキストに関する所定の条件が割当てられており、かつ当該所定の条件が充足されるか否かによって、前記ルートノードと前記複数の中間ノードとの各々から枝分かれする枝のいずれをたどるべきかが予め定められており、

前記複数の葉ノードの各々には、音声素片データの予備選択幅の予測値が割当てられており、

前記回帰木による予測手段は、

あるターゲット音素のコンテキストが与えられると、前記ルートノードから始めて、当該コンテキストが、各ノードでの条件を充足するか否かを判定し、判定結果に従って前記回帰木をたどっていくための判定手段と、

前記判定手段による判定結果に従って前記回帰木をたどって到達した葉ノードに割当てられた予備選択幅の予測値を前記予備選択されるべき音声素片データの数として出力するための手段とを含む、請求項 1 に記載の音声合成装置。

【請求項 3】

コンピュータにより実行されると、当該コンピュータを、請求項 1 又は請求項 2 に記載の音声合成装置として動作させる、コンピュータプログラム。

【請求項 4】

多数の音声素片データを格納した音声素片データベースとともに用いられ、合成ターゲットが与えられると、当該合成ターゲットを構成する各ターゲット音素のコンテキストに基づいて、前記音声素片データベースから当該ターゲット音素の音声合成に用いるべき音声素片データの候補を予備選択した後、予備選択された素片候補中から音声合成のための音声素片データを決定する、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から音声素片データの候補を予備選択するにあたり、予備選択される候補の数を、各ターゲット音素のコンテキストに基づいて動的に決定する事を特徴とする、素片接続型音声合成装置。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は音声合成技術に関し、特に、素片接続型音声合成技術における音声素片の予備選択技術に関する。

【背景技術】

【0002】

音声合成技術の一つとして、素片接続型音声合成がある。素片接続型音声合成では、多数の音声素片を予め準備しておく。各音声素片には、音素ラベル、音響パラメータ、音声コーパス内での出現環境などの情報が付されている。音声合成のターゲットが与えられると、これら多数の音声素片から、ターゲットとして与えられた音素で、与えられたパラメータに近く、かつ前後の音声との接続関係も良好なものを選択する。この選択動作を素片選択と呼ぶ。選択された音声素片の波形を接続して連続波形を生成する事により、目標となる音声を合成する。

【0003】

一般に波形接続型音声合成における素片選択は、コストと呼ばれるひずみ尺度の最小化に基づき行なわれる。コストは、通常、合成ターゲットと素片との間の誤差として定義されるターゲットコスト、及び素片間の不連続として定義される接続コストから構成される。

【0004】

10

20

30

40

50

こうした音声合成において最も重要なのは、適切な音声素片をいかにして選択するか、という問題である。

【0005】

波形接続型音声合成では、より自然性の高い合成音声を得るために、大規模な音声素片データベース（以下「音声素片DB」と呼ぶ。）が用いられる事が多い。この結果、考慮すべき音声素片の組合せの数が増え、合成に適した素片選択が困難となる。

【0006】

そこで、実際のシステムでは、素片の組合せを考える前に各時刻の素片の候補を、別の高速な手法により絞り込む処理（予備選択）が行なわれる事が多い。

【0007】

こうした予備選択を行なう手法の一例として、後掲の特許文献1には、次のような手法が開示されている。この手法では、第1段階の選択（予備選択）で、所定の条件に従って所定個数の音声素片を候補として選ぶ。第2段階では、それら候補の各々について、適切な比較をするための変形を行なった後に、変形後の音声素片と他の音声素片との間の変形ひずみの平均を算出する。それらを比較して変形ひずみの平均が最も小さな音声素片を、最適な音声素片として選択する。

【0008】

第1段階での予備選択には、対象となる音声素片のうちで、他の音声素片との間のピッチ長（又は継続時間長、音素環境、ピッチパターン等）の差分の絶対値の総和が小さな上位の所定個数を選ぶ方法、予め設定されたピッチ長又は継続時間長と音声素片のピッチ長又は継続時間長との差分の小さなものから所定個数を選ぶ方法などが挙げられている。

【特許文献1】特開2005-300919号公報（図2、図4、図6、図11、段落0013～0014、0019～0020）

【発明の開示】

【発明が解決しようとする課題】

【0009】

上記した予備選択手法では、予備選択のしかたによって最終的な素片選択結果が影響を受ける事が分かっている。最終的な素片選択結果をできるだけ適切なものとするためには、予備選択で残す音声素片候補の数をできるだけ多くする事が望ましい。しかし、予備選択で残す音声素片候補の数が増えれば、結果として素片選択に必要な処理が増加する事になり、予備選択の意義が薄れてしまう。処理量の抑制を主目的として音声素片候補の数を少なくすれば、最終的に得られる合成音声の音質が低下してしまう。

【0010】

一方で、自然な音声合成を行なうために音声素片DBはますます大規模化する傾向がある。従って、予備選択での処理量を抑えながら、最終的に適切な素片候補を選択できるような音声合成装置が求められている。

【0011】

それゆえに本発明の目的は、素片選択型の音声合成装置において、高速で、かつ高品質の音声信号を合成できる素片接続型音声合成装置を提供する事である。

【課題を解決するための手段】

【0012】

本発明の第1の局面に係る素片接続型音声合成装置は、多数の音声素片データを格納した音声素片データベースとともに用いられる素片接続型音声合成装置である。この素片接続型音声合成装置は、合成ターゲットが与えられると、当該合成ターゲットを構成する各ターゲット音素のコンテキストに基づいて、音声合成において各ターゲット音素の合成に用いられるべき候補として予備選択されるべき音声素片データの数を予測するための素片候補数予測手段と、合成ターゲットが与えられると、当該合成ターゲットを構成する各ターゲット音素について、当該ターゲット音素と音声素片データベース中の音声素片データの各々との間に算出されるターゲットコストに基づいて、音声素片データベース中から、素片候補数予測手段により予測された数と所定の関係にある数の音声素片データを、各々

10

20

30

40

50

ターゲットの音声合成のために予備的に選択するための素片候補予備選択手段と、合成ターゲットを構成する各ターゲット音素について、素片候補予備選択手段により選択された音声素片データの候補の各々との間に算出されるターゲットコスト及び接続コストに基づいて、音声合成に用いるべき音声素片データを選択するための素片選択手段と、素片選択手段により選択された音声素片データの音声波形を合成ターゲットに従って接続するための波形接続手段とを含む。

【0013】

素片候補数予測手段は、合成ターゲットが与えられると、そのターゲット音素の各々に対し、そのターゲット音素のコンテキストに基づいて、予備選択されるべき音声素片データの候補の数を予測する。素片候補予備選択手段は、ターゲットコストに基づいて、予測された数の音声素片データの候補を音声素片データベースから予備的に選択する。素片選択手段は、こうして予備的に選択された音声素片データの候補に対し、ターゲットコストと接続コストとの双方を用いて、音声合成に用いるべき音声素片データを選択する。波形接続手段は、こうして選択された音声素片データの音声波形を接続する事により音声合成を行なう。音声合成の処理で最も負荷の高いのは、接続コストの算出である。素片候補予備選択手段により音声素片データの候補が予め絞られているため、接続コストの算出の処理の負荷が小さくなる。素片候補予備選択手段では、ターゲットコストのみを用いているため、予備選択のための負荷は小さい。また、ターゲット音素のコンテキストに基づき、予備選択すべき音声素片データの数が素片候補数予測手段により予測される。予備選択において不必要に大きな数の音声素片データが選択されたために後の処理の負荷が高くなったり、予備選択において必要な数だけの音声素片データが選択されなかったために、最終的に得られる音声信号の品質が大きく損なわれたりするおそれが少なく、高品質の音声合成を、少ない負荷で高速に行なう事ができる。

【0014】

好ましくは、素片候補数予測手段は、各ターゲット音素のコンテキストに基づいて、音声合成において各ターゲット音素の合成に用いられるべき候補として予備選択されるべき音声素片データの数を、予め準備された回帰木を用いて予測するための回帰木による予測手段を含む。当該回帰木は、一つのルートノードと、複数の葉ノードと、ルートノードと葉ノードとの間に存在する複数の中間ノードとを含む。ルートノードと複数の中間ノードとの各々には、ターゲット音素のコンテキストに関する所定の条件が割当てられており、かつ当該所定の条件が充足されるか否かによって、ルートノードと複数の中間ノードとの各々から枝分かれする枝のいずれをたどるべきかが予め定められている。複数の葉ノードの各々には、音声素片の予備選択幅の予測値が割当てられている。回帰木による予測手段は、あるターゲット音素のコンテキストが与えられると、ルートノードから始めて、当該コンテキストが、各ノードでの条件を充足するか否かを判定し、判定結果に従って回帰木をたどっていくための判定手段と、判定手段による判定結果に従って回帰木をたどって到達した葉ノードに割当てられた予備選択幅の予測値を予備選択されるべき素片の数として出力するための手段とを含む。

【0015】

回帰木という簡単な判定手段によって予備選択すべき音声素片データの数を予測する事ができる。この回帰木を作成するためには予め学習が必要となるが、一旦学習をしておけば、同じ音声素片データベースを用いる限りは繰返して使用できる。

【0016】

さらに好ましくは、コンテキストは、音素情報からなる音素コンテキスト情報を含む。

【0017】

音素コンテキスト情報は、合成ターゲットには必ず含まれる。これ以外の情報が利用可能でないときにも、音素コンテキスト情報を使用する事により、予備選択すべき音声素片データの数を確実に予測できる。

【0018】

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

10

20

30

40

50

と、当該コンピュータを、上記したいずれかの音声合成装置として動作させる。

【0019】

本発明の第3の局面に係る素片接続型音声合成装置は、多数の音声素片データを格納した音声素片データベースとともに用いられ、合成ターゲットが与えられると、当該合成ターゲットを構成する各ターゲット音素のコンテキストに基づいて、音声素片データベースから当該ターゲット音素の音声合成に用いるべき音声素片データの候補を予備選択した後、予備選択された音声素片データの候補中から音声合成のための音声素片データを決定する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から音声素片データ候補を予備選択するにあたり、予備選択される音声素片データの候補の数を、各ターゲット音素のコンテキストに基づいて動的に決定する事を特徴とする。

10

【0020】

予備選択されるべき音声素片データの数が動的に決定される。その決定には、ターゲット音素のコンテキストが使用される。予備選択数をこの様に動的に決定する事により、予備選択される音声素片データの数が過大になって選択処理の負荷が過度に高くなったり、予備選択される音声素片データの数が過少になって最終的に得られる音声信号の音質が下がったりする事が防止できる。その結果、音声合成の音質を維持しながら、大量の音声合成を短時間で行なう事ができる。

【発明を実施するための最良の形態】

【0021】

以下に説明する、本発明の一実施の形態に係る音声合成装置は、予備選択において、ターゲット音素のコンテキストが与えられると、どの程度の数の素片候補を選択すれば最終的に適切な素片候補が得られるかを、予め行なった学習の結果によって予測する。この予測によって、予備選択で選択される素片候補の数は、コンテキストごとに動的に変化する。

20

【0022】

なお、本実施の形態において、ターゲット音素のコンテキストとは、ターゲット音素を含む、その前後の所定数の音素とからなる音素列をさすものとする。本実施の形態では、コンテキストとして、ターゲット音素と、その前後の二つずつの音素とからなる音素列を用いる。

【0023】

また、本実施の形態でも、コスト最小化によって音声素片を選択する。ターゲットコストとしては、継続時間、基本周波数 F_0 、及び平均MFCC(Mel Frequency Cepstrum Coefficient)が考慮される。接続コストとしては、素片境界における F_0 不連続、MFCC不連続、音素環境及びその代替に応じた一般的な接続の困難さが考慮される。ただし、環境代替の影響については、後述する予備選択時にその影響が考慮される様に、その影響の一部を素片選択時にターゲットコストとして取り扱っている。

30

【0024】

<構成>

図1に、本実施の形態に係る音声合成装置30のブロック図を示す。図1を参照して、音声合成装置30は、入力テキスト32が与えられると、当該テキストの出力音声波形34という形で音声合成を行なうためのものである。

40

【0025】

図1を参照して、音声合成装置30は、入力テキスト32に対してテキスト処理を行ない、形態素解析、構文解析、単語辞書の参照などによって音声合成の単位である音素単位に分割し、解析によって得られた韻律情報を付して合成ターゲットとして出力するためのテキスト処理部40と、テキスト処理部40の出力する合成ターゲットを構成するターゲット音素列に対し、合成すべき音声の韻律に対応する音響特徴量パラメータ(ターゲットパラメータ)を生成し、各音素に付するターゲット音素からなるターゲット音素列として出力するための合成パラメータ生成部42とを含む。

50

【 0 0 2 6 】

音声合成装置 3 0 はさらに、多数の音声素片をそれらの音響特徴量パラメータとともに格納するための素片 D B 5 2 と、ある音素を中心音素とするコンテキストが与えられると、当該中心音素に対応する素片候補として予備選択すべき素片候補の数（以下これを「予備選択幅」と呼ぶ。）を予測し出力するための素片候補数予測部 4 8 と、ある音素を中心音素とするコンテキスト及びターゲットコスト計算用データが与えられると、当該中心音素に対応する素片候補を、素片候補数予測部 4 8 により予測される数だけ、素片 D B 5 2 中の各素片に対して算出されるターゲットコストに基づいて素片 D B 5 2 から予備的に選択し出力するための素片候補予備選択部 5 0 とを含む。

【 0 0 2 7 】

音声合成装置 3 0 はさらに、合成パラメータ生成部 4 2 からターゲット音素列が与えられると、各音素のコンテキストを素片候補数予測部 4 8 に与え、それに応答して素片候補予備選択部 5 0 から与えられる素片候補の各々に対し、前述したターゲットコストと接続コストとの双方を用いた最適素片の選択を行なうための素片選択部 4 4 と、素片選択部 4 4 により選択された音声素片の波形を合成ターゲットに従って互いに接続し、出力音声波形 3 4 を出力するための波形接続部 4 6 とを含む。

【 0 0 2 8 】

素片候補数予測部 4 8 は、本実施の形態では、コンテキスト情報を用いて予め作成した回帰木により実現される。図 2 に、素片候補数予測部 4 8 で使用する回帰木 6 0 のルートノード付近の構成を示す。本実施の形態では、回帰木 6 0 の作成のために、予め所定数の質問（一実施例では 3 1 8 問を用いた。）を準備しておく。これら所定数の質問を回帰木 6 0 の各ノードに割当てて、各ノードには、予備選択すべき素片候補数が予め付されている。この素片候補数予測部 4 8 は、与えられたコンテキスト情報に対し、この回帰木 6 0 の各ノードの質問に答える形で回帰木 6 0 を順番にたどり、最終的に到達した葉ノードに付された数を予備選択幅として素片候補予備選択部 5 0 に返す機能を持つ。

【 0 0 2 9 】

なお、ここでは「質問」と述べたが、これはコンテキスト情報が充足すべき条件であると考えられる。コンテキスト情報がこの条件を充足する場合、及び充足しない場合に、そのノードから枝分かれしている枝のいずれに進むかは、回帰木 6 0 の作成過程で各枝に割当てられる。従って、回帰木 6 0 は本実施の形態では二分木となっている。もちろん、回帰木 6 0 を二分木とする必然性はなく、条件によって枝分かれが 3 つ以上になってもよい。

【 0 0 3 0 】

どの様にして回帰木 6 0 を作成するかについては図 4 を参照して後述する事にし、図 2 に示す回帰木 6 0 について具体的に説明する。回帰木 6 0 はルートノード 7 0 と、ルートノード 7 0 から分岐するノード 7 2 及び 7 4 と、これらノード 7 2 及び 7 4 からそれぞれ分岐するノード 7 6 及び 7 8、並びにノード 8 0 及び 8 2 とを含む。回帰木 6 0 はノード 7 6 及び 7 8 より下、並びにノード 8 0 及び 8 2 より下にさらに多数のノードを含むが、図 2 では図示を簡略化するためにそれらは示していない。

【 0 0 3 1 】

ルートノード 7 0 には、例えば「（ターゲット音素が）半音素前半か」という質問が割当てられたとする。与えられたターゲット音素が半音素の前半であればノード 7 4 に進み、それ以外であればノード 7 2 に進む。図 2 の回帰木 6 0 において、分岐の枝に付された「Y」及び「N」という記述は、それぞれ質問に対する答えが「イエス」の場合及び「ノー」の場合に進むべき枝を示す。

【 0 0 3 2 】

図 3 に、図 1 の素片候補予備選択部 5 0 の構成を示す。図 3 を参照して、素片候補予備選択部 5 0 は、素片選択部 4 4 からコンテキストが与えられると、素片 D B 5 2 から当該コンテキストの中心音素と一致する音素の音声素片を全て抽出するための素片抽出部 1 0 0 と、素片抽出部 1 0 0 により抽出された音声素片の各々に対して、素片選択部 4 4 から

10

20

30

40

50

与えられたターゲットコスト算出用データを用いて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102と、ターゲットコスト算出部102により算出されたターゲットコストが少ないものの上位から、素片候補数予測部48により予測された予備選択幅の数だけを素片選択部44に返すための順位比較部104とを含む。

【0033】

図4に、回帰木60を作成するための回帰木作成システム120の構成を示す。要するに、回帰木作成システム120は、実際に音声合成のための素片選択を多数回行ない、その際に最適なものとして最終的に選択された素片を、ターゲットコストのみによる予備選択で捨てない様にするためには、どの程度の予備選択幅としたらよいかをコンテキスト別に推定するためのものである。

10

【0034】

図4を参照して、回帰木作成システム120は、多数の音声合成用テキストからなる学習用データ140と、図1に示すものと同じ素片DB52と、学習用データ140から音声合成用テキストを読み出し、各音素に対しターゲットコスト及び接続コストの重み付け合計により得られるコストに基づいて、素片DB52から最適素片を選択する事により、素片選択データ144を作成するための素片選択データ作成部142とを含む。

【0035】

素片選択データ作成部142が作成する素片選択データ144は、ターゲット音素のコンテキストと、このコンテキストに対して最終的に得られた素片データについて、ターゲットコストが全体の中で何番目に小さかったかを示す順位データとの組からなる。

20

【0036】

なお、予備選択を行わずに素片選択を行なう事は容易ではないため、本実施の形態では素片選択データ144を作成する際の素片選択は、固定した予備選択幅及びビーム幅(例えば予備選択幅2000、ビーム幅500)の探索で行なう。素片DB52上において連続する音声素片を優先して探索する連続素片優先探索により仮説展開された素片候補が最終的に選択された場合、その選択された素片候補のターゲットコスト上での順位は0とする。

【0037】

回帰木作成システム120はさらに、予め準備された所定数の質問をコンピュータ読取可能な形式で格納する質問データ格納部152と、質問データ格納部152に格納された質問と、素片選択データ144とに基づき、回帰木60を作成するための予測回帰木作成部150とを含む。

30

【0038】

本実施の形態では、回帰木60の作成には以下の考え方をを用いている。すなわち、本実施の形態では、コンテキスト情報を用いて予備選択幅を削減する。そのために、必要な予備選択幅を基準にコンテキストクラスタリングを行なう事で、予備選択幅を予測する回帰木60を作成する。

【0039】

あるコンテキストが、あるクラスタに属しているとき、そのコンテキストにおいて必要な予備選択幅は、クラスタに属するサンプル中の予備選択順位の最悪値(最大値)である。しかし、そのクラスタにそのような予備選択幅が不要なコンテキストも含まれているならば、クラスタを分割し、予備選択幅がより小さくてもよいコンテキストのクラスタを作成する事ができる。ただし、ここでは安定した推定のために、クラスタリング基準に順位の最悪値を用いるのではなく、クラスタ内のサンプルの、ターゲットコストによる予備選択順位上での上位から97%の位置の順位を予備選択幅予測値とし、これをクラスタリング基準とする。

40

【0040】

クラスタの分割は、あるクラスタを分割した後の二つのクラスタのサンプル数、及びそれらクラスタから決まる予備選択幅予測値をそれぞれ c_1 、 c_2 、 k_1 、及び k_2 とするとき、次の式(1)

50

【 0 0 4 1 】

【 数 1 】

$$\sigma^2 = \frac{c_1(k_1 - \mu)^2 + c_2(k_2 - \mu)^2}{c_1 + c_2} \quad (1)$$

$$\text{ただし } \mu = \frac{c_1 k_1 + c_2 k_2}{c_1 + c_2}$$

の値が最大となる質問で分割を繰返す事で行なわれる。これは、上位 9 7 % 点の値を用いて定義された分布間距離を基準とするクラスタリングと考えられる。

10

【 0 0 4 2 】

なお、本実施の形態では、素片選択にテキストの情報を利用できない場合も想定し、コンテキストとしては音素環境のみを考慮している。

【 0 0 4 3 】

回帰木 6 0 のサイズを抑えるために、ノードの分割において以下の 3 つの条件を用いた。

【 0 0 4 4 】

- (1) 分割後のノードに属するサンプル数が制限値 C m i n 未満にならない事
- (2) 分割によって、少なくとも一方のノードの予備選択幅予測値が、分割前の予測値に対して 1 0 % 以上変化する事
- (3) 回帰木 6 0 の深さが 3 0 段を超えない事

20

図 4 に示す予測回帰木作成部 1 5 0 は、このクラスタリングを行なうためのものである。図 5 に、予測回帰木作成部 1 5 0 の機能をコンピュータ及びコンピュータプログラムで実現する場合のコンピュータプログラムの制御構造をフローチャート形式で示す。図 5 を参照して、この処理では、最初にステップ 1 7 0 で素片選択データ 1 4 4 (図 4 参照) を準備する。具体的には、素片選択データ 1 4 4 を格納したファイルをオープンする。以後、このファイルから読出された素片選択データ 1 4 4 の個々のデータを「サンプル」と呼ぶ。ステップ 1 7 4 では、質問データを準備する。具体的には、質問データを格納したファイルをオープンする。以後、クラスタリング処理が開始される。

【 0 0 4 5 】

30

ステップ 1 7 6 において、全サンプルを素片候補予備選択部 5 0 の最初の一つのノード (ルートノード) に属するサンプルとして分類する。すなわち、最初のクラスタが作成される。また、ルートノードの予備選択幅予測値 k を、ルートノードに属するサンプルの予備選択順位上での上位 9 7 % 点として算出し、ルートノードに情報として付加する。これ以後の処理は停止条件が充足されるまでの繰返処理である。

【 0 0 4 6 】

ステップ 1 7 8 において、停止条件を満たしていないノードがあるか否かが判定される。停止条件は、前述した 3 つの条件の裏返しである。すなわち、(1) 分割後のノードに属するサンプル数が制限値 C m i n 未満になるか、(2) ノードの分割によって得られる二つのノードのいずれの予備選択幅予測値も、分割前の予測値に対して 1 0 % 以上変化しないか、(3) 回帰木 6 0 の深さが 3 0 段を超えたか、という条件が成立するとそのノードに対するそれ以上の分割は行なわない。停止条件を満たしていないノードがあれば、それらノードの中のいずれかを選択してステップ 1 8 0 以下の処理を行なう。停止条件を満たしていないノードがなければ、得られた回帰木を出力して処理を終了する。

40

【 0 0 4 7 】

ステップ 1 8 0 では、処理対象のノードに分類されたサンプル (当該ノードにより示されるクラスタに分類されたサンプル) について、最初に準備した所定数の質問のうち、既にノードに割当て済みの質問以外の質問の全てに答える事で、それぞれ二つずつのクラスタ (クラスタ対) に分ける。ステップ 1 8 2 では、得られたクラスタ対のうち、前述した式 (1) で示される分布間距離が最大となる質問を、処理中のノードに割当てる。続いて

50

ステップ184では、処理中のノードに割当てられた質問により得られた二つのクラスタに対応する二つのノードを、現在処理中のノードの子ノードとして、回帰木に追加する。各ノードには、処理中のノードに割当てられた質問に対する答えがイエスかノーかによってサンプルを分類して割当てる。また、各ノードに割当てられたサンプルに基づき、各ノードの予備選択幅予測値kを算出し、各ノードに情報として付加する。この後、ステップ178に戻る。

【0048】

こうして、回帰木60中の全てのノードが停止条件を充足すると処理が終了し、回帰木60が完成する。

【0049】

<動作>

上記した音声合成装置30は以下の様に動作する。音声合成装置30の動作に先立ち、音声合成装置30で使用する回帰木60(図2参照)を作成する必要がある。従って、最初に図4及び図5を参照して回帰木60の動作を説明する。

【0050】

図4を参照して、学習用データ140を記憶媒体に記憶させる。学習用データ140は、前述した通り、多数の音声合成用テキストからなる。素片選択データ作成部142は、学習用データ140中のテキストを読み込み、コスト計算に基づく素片選択によって音声合成を行なう。ここでは、コストとしてターゲットコストと接続コストとの双方を用いる。ただし、ここでは、実際に選択された素片について、ターゲットコストによる順位を付ける処理も行なう。各ターゲット音素のコンテキストと、最終的に選択された素片について算出されたターゲットコストによる順位とを記憶させる事で、素片選択データ144を作成する。こうして、学習用データ140の全てについて音声合成(素片選択)が終了すると、素片選択データ144が完成する(図5のステップ170)。

【0051】

素片選択データ144が完成すると、予測回帰木作成部150が以下の様にして回帰木60を作成する。この処理に先立ち、質問データ格納部152に予め所定個数の質問がコンピュータ読取可能な形式で準備される(図5のステップ174)。

【0052】

図5を参照して、ステップ176以下の処理は、予測回帰木作成部150が行なう処理である。ステップ176において、まず素片選択データ144の全てが、最初の一つのノード(ルートノード)に分類される。またここでは、ルートノードのサンプルに基づき、ルートノードの予備選択幅予測値が算出され、ルートノードに付与される。

【0053】

次に予測回帰木作成部150は、停止条件を充足していないノードがあるか否かを判定する(ステップ178)。繰返しの最初ではノードはルートノードのみであり、この停止条件は充足されていない事が通常である。従って判定の結果は「YES」となり、制御はステップ180に進む。

【0054】

ステップ180では、当該ノードに属する素片選択データの全てについて、まだノードに割当てられていない質問の各々に対する答えによって分類する。この分類の結果、質問の数だけのクラスタ対候補が作成される。

【0055】

ステップ182では、それらクラスタ対候補の各々について、クラスタ対を構成するクラスタ間の分布間距離が式(1)により算出される。この分布間距離が最大となる質問を、ルートノードに割当てる。図2に示す回帰木60のルートノード70に割当てられた質問「半音素前半か」という質問はこうして選択されたものである。

【0056】

続いてステップ184では、ステップ182でルートノード70に割当てられた質問に従って分類されたクラスタ対に対応する二つの子ノードをルートノード70から分岐する

10

20

30

40

50

形で作成する。この処理により、図 2 に示すノード 7 2 及びノード 7 4 が作成される。これらノードに属するサンプルとして、ルートノード 7 0 に割当てられた質問によってクラスタに分類されたものがそれぞれ割当てられる。また、こうしてノード 7 2 及びノード 7 4 に割当てられたサンプルに基づき、これらの予備選択幅予測値が算出され、これらノードに付与される。ただしこれらノードには、まだ質問は割当てられていない。

【 0 0 5 7 】

次に再度ステップ 1 7 8 に進み、停止条件を満たしていないノードが存在するか否かが判定される。本例では、ノード 7 2 及びノード 7 4 のいずれもまだ停止条件を満たしていないものとする。従って処理はステップ 1 8 0 に進む。ステップ 1 8 0 以下の処理は、停止条件を満たしていないノードが複数個ある場合、それらのいずれかを所定の選択方式で選択して行なわれる。ここでは、例えばノード 7 2 が選択されたものとする。

10

【 0 0 5 8 】

ステップ 1 8 0 において、ノード 7 2 に属する素片選択データの全てを、残りの質問（まだノードに割当てられていない質問）の各々でクラスタ対候補に分類する。ここでも、分類の結果、残りの質問の数だけのクラスタ対候補が作成される。

【 0 0 5 9 】

ステップ 1 8 2 では、ノード 7 2 に対しステップ 1 8 0 で作成されたクラスタ対候補のうち、クラスタ対を構成するクラスタ間の式 (1) による分布間距離が最大となる質問を、ノード 7 2 に割当てる。図 2 に示す例では、「 / e / 又は / o / か」という質問がノード 7 2 に割当てられる。

20

【 0 0 6 0 】

ステップ 1 8 4 では、ステップ 1 8 2 でノード 7 2 に割当てられた質問により分類されたクラスタ対に従い、新たな二つの子ノード（図 2 におけるノード 7 6 及び 7 8 ）が作成される。これらノードには、ノード 7 2 に割当てられた質問によって分類されたクラスタ対に属する素片選択データがそれぞれ属する事になる。各ノードにおいて、当該ノードに属するサンプルに基づいて予備選択幅予測値が算出され、これらノードに付与される。この後、ステップ 1 7 8 に戻る。

【 0 0 6 1 】

こうして、停止条件を満たしていないノードが回帰木 6 0 内に存在する限り、ステップ 1 7 8 ~ ステップ 1 8 4 の処理が繰返し行なわれ、図 2 に示す回帰木 6 0 がルートノード 7 0 から順番に下側に枝分かれしていく態様で作成される。回帰木 6 0 内の全ノードが停止条件を充足すると、回帰木作成システム 1 2 0 による処理が終了する。作成された回帰木 6 0 は、所定の記憶装置に記憶される。

30

【 0 0 6 2 】

こうして作成された回帰木 6 0 は、音声合成装置 3 0 で使用できる様に、音声合成装置 3 0 を構成するコンピュータ（その具体的構成は後述する。）内の記憶装置、又は外部記憶装置に格納され、音声合成装置 3 0 の動作時に素片候補数予測部 4 8 により利用できる様にコンピュータ内で準備される。

【 0 0 6 3 】

次に、図 1 及び図 3 を特に参照して、音声合成装置 3 0 の動作について述べる。入力テキスト 3 2 が与えられると、音声合成装置 3 0 のテキスト処理部 4 0 はこの入力テキスト 3 2 に対して形態素解析、構文解析、単語辞書の参照などを行なう事により、入力テキスト 3 2 を音声合成の単位である音素単位に分割し出力する。ここでは、入力テキスト 3 2 に対する解析結果を用いて、各音素についての韻律情報が生成され各音素に付される。

40

【 0 0 6 4 】

合成パラメータ生成部 4 2 は、テキスト処理部 4 0 の出力する音素列に対し、合成すべき音声の韻律に対応するターゲットパラメータを生成し、素片選択部 4 4 に与える。

【 0 0 6 5 】

素片選択部 4 4 は、音素単位で素片候補を選択しながら音素列に対応する素片の系列を作成していく。この処理において素片選択部 4 4 は、ある時刻での音声合成に用いる音声

50

素片の選択のために、合成パラメータ生成部 4 2 から与えられる音素列のうち、合成対象となる音素を中心とする所定のコンテキスト（中心音素 ± 2 音素の音素列）を素片候補数予測部 4 8 に与える。

【 0 0 6 6 】

素片候補数予測部 4 8 は、与えられたコンテキストに対し、図 2 に示す回帰木 6 0 のルートノード 7 0 の質問に対する答えを判定する。そして、判定結果に従ってノード 7 2 及びノード 7 4 のいずれかを選択する。選択されたノードにおいて、同じくそのノードに割当てられた質問に対する答えを判定する。以下同様に、与えられたコンテキストに対する、各ノードに割当てられた質問の答えを判定しながら、回帰木 6 0 をたどる。最終的に到達した葉ノードには、予備選択幅予測値としてある値が付与されている。素片候補数予測部 4 8 は、最終的に到達した葉ノードの予備選択幅予測値を素片候補予備選択部 5 0 に与える。

10

【 0 0 6 7 】

図 3 を参照して、素片候補予備選択部 5 0 の素片抽出部 1 0 0 は、素片選択部 4 4 からターゲット音素のコンテキストが与えられると、その中心音素を音素ラベルに持つ音声素片全てを素片 DB 5 2 から抽出し、ターゲットコスト算出部 1 0 2 に与える。

【 0 0 6 8 】

ターゲットコスト算出部 1 0 2 は、与えられた音声素片の全てに対し、素片選択部 4 4 から与えられたコンテキスト中の中心音素に関するターゲットパラメータに基づいてターゲットコストを算出する。ターゲットコストの算出にはそれほどのリソースは必要ではない。ターゲットコスト算出部 1 0 2 は、各音声素片に対してターゲットコストを付して順位比較部 1 0 4 に与える。

20

【 0 0 6 9 】

順位比較部 1 0 4 は、与えられた音声素片をターゲットコストの低いものから昇順にソートする。順位比較部 1 0 4 はさらに、こうしてソートされた音声素片のうち、ターゲットコストの低いものから素片候補数予測部 4 8 により予測された数だけの音声素片を素片選択部 4 4 に返す。

【 0 0 7 0 】

図 1 を参照して、素片選択部 4 4 は、素片候補予備選択部 5 0 から与えられた素片候補に対し、ターゲットコストと接続コストとの双方を用いた素片選択を行なう。波形接続部 4 6 は、選択された音声素片の波形を接続する事により出力音声波形 3 4 を生成し出力する。

30

【 0 0 7 1 】

こうして入力テキスト 3 2 を構成する全ての形態素の音素について、音声素片が選択され波形接続部 4 6 により接続されると、音声合成装置 3 0 の処理が終了する。

【 0 0 7 2 】

< コンピュータによる実現 >

【 0 0 7 3 】

図 7 は、この音声合成装置 3 0 を実現するコンピュータシステム 5 3 0 の外観を示し、図 8 はコンピュータシステム 5 3 0 の内部構成を示す。

40

【 0 0 7 4 】

図 7 を参照して、このコンピュータシステム 5 3 0 は、メモリドライブ 5 5 2 及び DVD ドライブ 5 5 0 を有するコンピュータ 5 4 0 と、キーボード 5 4 6 と、マウス 5 4 8 と、モニタ 5 4 2 と、マイクロフォン 5 7 0 と、音声合成の結果を出力するための一対のスピーカ 5 7 2 とを含む。

【 0 0 7 5 】

図 8 を参照して、コンピュータ 5 4 0 は、メモリドライブ 5 5 2 及び DVD ドライブ 5 5 0 に加えて、CPU（中央処理装置）5 5 6 と、CPU 5 5 6、メモリドライブ 5 5 2 及び DVD ドライブ 5 5 0 に接続されたバス 5 6 6 と、ブートアッププログラム等を記憶する読出専用メモリ（ROM）5 5 8 と、バス 5 6 6 に接続され、プログラム命令、シス

50

テムプログラム、及び作業データ等を記憶するランダムアクセスメモリ（RAM）560と、バス566に接続された不揮発性の外部記憶装置であるハードディスクドライブ（HDD）554とを含む。

【0076】

コンピュータ540はさらにローカルエリアネットワーク（LAN）574への接続を提供するネットワークアダプタ576を含む。

【0077】

コンピュータシステム530に音声合成装置30としての動作を行なわせるためのコンピュータプログラムは、DVDドライブ550又はメモリドライブ552に挿入されるDVD562又は不揮発性メモリ564に記憶され、さらにハードディスク554に転送される。又は、プログラムはネットワーク574を通じてコンピュータ540に送信されハードディスク554に記憶されてもよい。プログラムは実行の際にRAM560にロードされる。DVD562から、不揮発性メモリ564から、又はネットワーク574を介して、直接にRAM560にプログラムをロードしてもよい。

10

【0078】

図1に示す素片DB52、及び図2に示す回帰木60は、ハードディスク554上に格納され、プログラムの実行の際にRAM560にロードされる。CPU556は、図示しないプログラムカウンタレジスタにより示される、RAM560上のアドレスから命令を読み出し、命令をデコードし、RAM560又はハードディスク554の、デコード結果により特定されるアドレスからデータを読み出して命令に従い処理し、デコード結果によって特定されるアドレスに格納する。CPU556はこうした処理を繰返す事により、入力テキスト32から出力音声波形34（図1を参照）を合成する処理を行なう。

20

【0079】

このプログラムは、コンピュータ540にこの実施の形態に係る音声合成装置30としての動作を行なわせる複数の命令を含む。この動作を行なわせるのに必要な基本的機能のいくつかはコンピュータ540上で動作するオペレーティングシステム（OS）若しくはサードパーティのプログラム、又はコンピュータ540にインストールされる各種ツールキットのモジュールにより提供される。従って、このプログラムはこの実施の形態のシステムを実現するのに必要な機能全てを必ずしも含まなくてよい。このプログラムは、命令のうち、所望の結果が得られる様に制御されたやり方で適切な機能又は「ツール」を呼出す事により、上記した音声合成装置30としての動作を実行する命令のみを含んでいればよい。コンピュータシステム530の動作は周知であるので、ここでは繰返さない。

30

【0080】

< 実験結果 >

次のテーブル1に、回帰木60の作成において、様々な制限値Cminを設定した場合のクラスタリングの結果、及びテストセットの予備選択幅を推定した場合の結果を示す。テーブル1において、Nは回帰木60のノード数、「mean」及び「RMSE」はそれぞれ、予測結果の平均値及び二乗平均平方根誤差、（A）は予測誤り率（必要な予備選択幅より小さく予測した割合）、（B）は予測誤り箇所のRMSEである。RMSEの値が全体に大きな値となっているのは、予測値と予備選択順位との差を評価したためである。

40

【0081】

【表1】
テーブル1

C_{min}	N	mean	RMSE	(A)	(B)
10	3373	382.4	553.3	0.068	76.6
20	1705	478.3	636.4	0.045	55.4
50	731	481.6	607.1	0.037	54.9
100	387	478.4	580.8	0.034	47.3
200	193	484.5	575.1	0.030	44.3
500	77	508.6	574.7	0.029	42.7
1000	35	495.8	554.4	0.027	47.8
2000	19	539.3	584.6	0.029	48.2
5000	7	535.1	554.8	0.028	54.6

10

図6は、図2と同様、回帰木60のルートノード付近の図であるが、これは上記した終了条件の一つに使用されている制限値 $C_{min} = 500$ のときのものである。ただし、 k は予備選択幅予測値、 Y 及び N の中のカッコの中の値はデータサンプル数を示す。

【0082】

必要な素片選択幅を予測する回帰木を用いて素片選択を行なった場合の結果を調べるため、素片選択実験を行なった。用いた素片DBは女声47.6時間のコーパスから作成されたもので、合成目標は、所定の53文からなるコーパスである。接続コストの計算が素片選択に必要な計算時間の多くを占めていることから、本実験では接続コストの計算回数を計算量の基準とした。

20

【0083】

まず最初に、接続コストの計算回数が所定の値となるような予備選択幅の上限値を各サンプルについて推定した。この際、ビーム幅は100に固定した。これは、予備選択幅推定結果を用いる場合も同様である。従って、計算回数削減の影響は、素片候補数が多い箇所に現れる事になる。また推定値の下限は10とし、必ず(素片候補が存在するならば)10個以上の素片が考慮される様にした。

30

【0084】

結果を図9に示す。図9における「constant K」(図9中、「+」で示す。)が、予備選択幅を一定とする従来法である。横軸は1ターゲット音素あたりの接続コスト計算回数であり、縦軸は正規化コスト(ほぼ1音素あたりのコストに相当する。)である。なお、予備選択幅推定を行なった場合に10000と20000の結果がほぼ一致しているが、これは、予備選択幅推定によって逆に計算回数を設定値まで増やす事ができなかった場合も区別せずに図示しているためである。実際に行なわれた計算回数はこれらの値よりも小さい。

【0085】

図9に示す結果より、計算回数が5000程度のとき、 C_{min} が200、500、1000(図9中、それぞれ「 \square 」「 \triangle 」及び「 \diamond 」で示す。)の回帰木において、予備選択幅推定の効果が得られている事が分かる。 C_{min} が200の場合、計算回数が50000のときのコストの値は、従来法における計算回数100000のときと同程度である。従って、この場合、従来法の半分の計算回数で同等の素片選択が得られた事になる。

40

【0086】

その他の領域で従来法よりも素片選択結果が悪い原因は、主として予測誤りによるものと考えられる。計算回数が多くても構わない場合は、もともと予備選択幅削減の効果は期待できない。一方、計算回数を少なく設定した場合、今回用いた計算回数を制御する方法では、予備選択幅上限値が極端に下がる。従って従来法との差異は小さくなる。

【0087】

50

以上の通り本実施の形態によれば、予め素片選択を行なった結果に基づき、どの程度の数の素片候補を予備選択すればその中に最適素片が入ると期待できるかをコンテキスト別に予測するための回帰木を作成した。この回帰木を用い、ターゲット音素のコンテキストが与えられると、そのコンテキストに対する予備選択幅を予測する。この予備選択幅により定まる数だけの素片をターゲットコストに基づいて予備選択する。予備選択された素片候補中から、ターゲットコスト及び接続コストに基づいて最終的な素片を選択する。実際の素片選択結果に基づいて、コンテキストごとに予備選択幅を動的に切替えて素片候補を予備選択するので、予備選択により選ばれた候補中に最適な素片が存在する可能性が高い。しかも、回帰木を使用するために、ごく負荷の低い処理によって効率的に予備選択を行なう事ができる。選択のための処理のうち、最も負荷が高いのは、接続コストによるコスト計算の部分であるので、本実施の形態によれば、精度を下げずに、処理量を下げながら素片選択を行なう事ができる。

10

【0088】

なお、上記した実施の形態では、予備選択幅を予測するために、回帰木を使用した。しかし本発明は回帰木を用いるものには限定されない。コンテキストデータが与えられると、当該コンテキストデータに対して最適と思われる予備選択幅を返す事ができるものであれば、どのようなものでも利用できる。例えばニューラルネットワークなど、実際の素片選択結果に基づいて学習を行なう事ができるものであれば、結果の信頼性も高く、本発明を実現するのに特に適している。

【0089】

20

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及び範囲内のすべての変更を含む。

【図面の簡単な説明】

【0090】

【図1】本発明の一実施の形態に係る音声合成装置30のブロック図である。

【図2】音声合成装置30で使用される回帰木60の一例のルートノード付近の概略構成を示す図である。

【図3】素片候補予備選択部50のより詳細な構成を示すブロック図である。

30

【図4】回帰木60を作成するための回帰木作成システム120のブロック図である。

【図5】図4に示す予測回帰木作成部150をコンピュータで実現するためのコンピュータプログラムの制御構造を示すフローチャートである。

【図6】回帰木60の一例のルートノード付近の概略構成を、各ノードの質問によるサンプルの分類数と、各ノードにおける予備選択幅予測値とともに示す図である。

【図7】本発明の一実施の形態に係る音声合成装置30を実現するコンピュータシステム530の外観を示す図である。

【図8】図7に示すコンピュータ540のブロック図である。

【図9】本発明の一実施の形態に係る音声合成装置30と同様の原理を用いた素片の予備選択の効果を説明するためのグラフである。

40

【符号の説明】

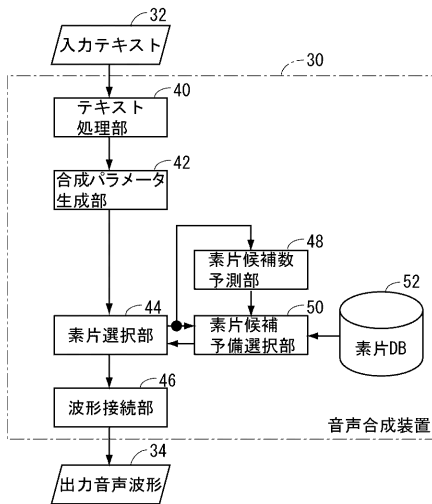
【0091】

- 30 音声合成装置
- 32 入力テキスト
- 34 出力音声波形
- 40 テキスト処理部
- 42 合成パラメータ生成部
- 44 素片選択部
- 46 波形接続部
- 48 素片候補数予測部

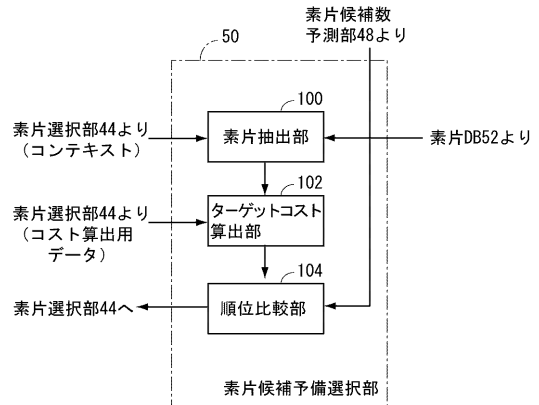
50

- 5 0 素片候補予備選択部
- 5 2 素片DB
- 1 0 0 素片抽出部
- 1 0 2 ターゲットコスト算出部
- 1 0 4 順位比較部
- 1 4 2 素片選択データ作成部
- 1 5 0 予測回帰木作成部

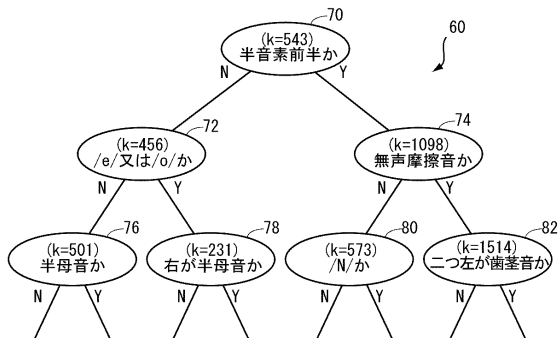
【図1】



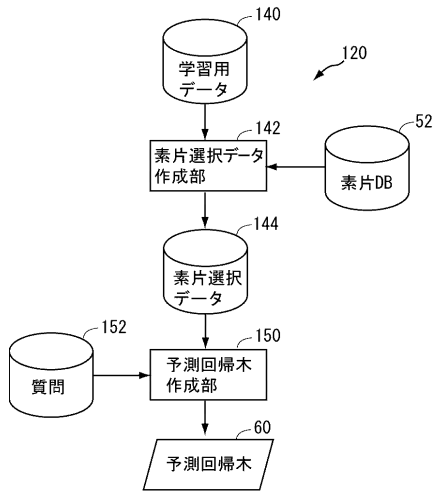
【図3】



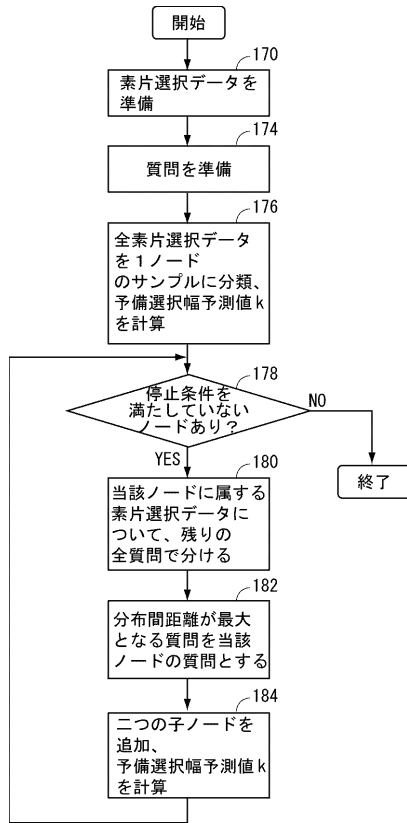
【図2】



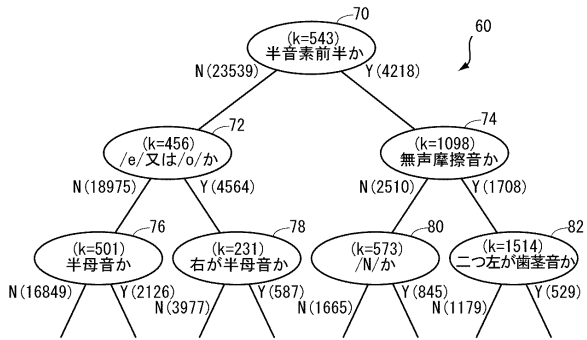
【図4】



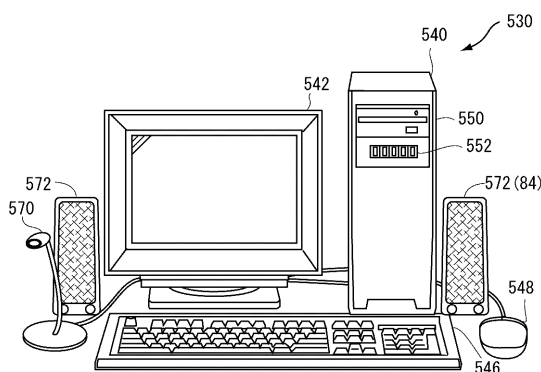
【図5】



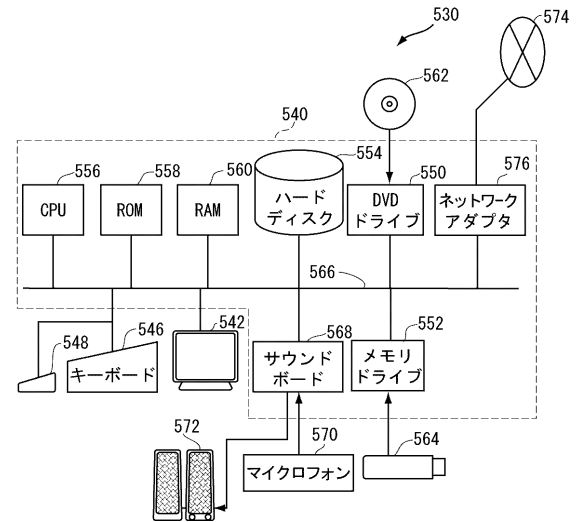
【図6】



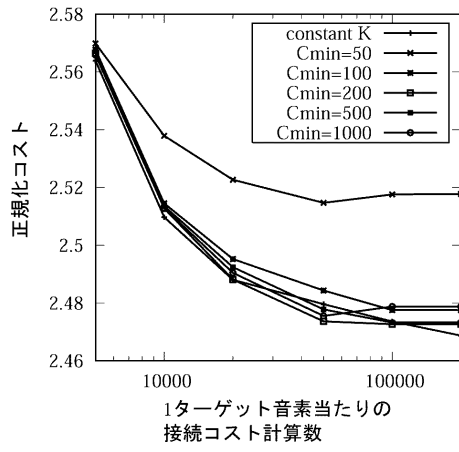
【図7】



【図8】



【図9】



フロントページの続き

(56)参考文献 特開2005 - 265895 (JP, A)
特開2005 - 241789 (JP, A)
特開2004 - 109535 (JP, A)

(58)調査した分野(Int.Cl., DB名)
G10L 13/06