

(19)日本国特許庁 (J P)

(12) 特 許 公 報 (B 2)

(11)特許番号

特許第3459600号
(P3459600)

(45)発行日 平成15年10月20日(2003.10.20)

(24)登録日 平成15年8月8日(2003.8.8)

(51)Int.Cl.⁷

識別記号

F I

G 1 0 L 13/06
13/08

G 1 0 L 5/04
3/00

E
H

請求項の数6(全24頁)

(21)出願番号 特願平11-274233
(22)出願日 平成11年9月28日(1999.9.28)
(65)公開番号 特開2001-100775(P2001-100775A)
(43)公開日 平成13年4月13日(2001.4.13)
審査請求日 平成13年1月17日(2001.1.17)

(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
(72)発明者 ニック・キャンベル
京都府相楽郡精華町大字乾谷小字三平谷
5番地 株式会社エイ・ティ・アール音
声翻訳通信研究所内
(72)発明者 北川 敏
京都府相楽郡精華町大字乾谷小字三平谷
5番地 株式会社エイ・ティ・アール音
声翻訳通信研究所内
(74)代理人 100062144
弁理士 青山 葆 (外2名)

審査官 渡邊 聡

最終頁に続く

(54)【発明の名称】 音声合成装置のための音声データ量削減装置及び音声合成装置

1

(57)【特許請求の範囲】

【請求項1】 音素ラベルに対応した音声波形信号の音声セグメントのデータからなる音声波形データベースを記憶する記憶装置を備え、上記自然発話の音声波形信号の音声セグメントを連結することにより任意の音素列を音声合成する音声合成装置のための音声データ量削減装置であって、
上記音声波形データベースに含まれる1対の音素のリストを生成する生成手段と、
上記生成された1対の音素のリストに基づいて各1対の音素に対する韻律的特徴パラメータと音響的特徴パラメータとに関する所定の類似度を計算し、上記計算された類似度が所定の第1のしきい値以上であるとき、当該各1対の音素のうち一方の1対の音素に係る音声波形信号の音声セグメントのデータを上記音声波形データベ

2

スから削除することにより音声データ量を削減するとともに、上記計算された類似度が所定の第1のしきい値以上であり、かつ上記1対の音素のリスト中の同一の1対の音素の数が所定の第2のしきい値以上であるときに、当該各1対の音素のうち一方の1対の音素に係る音声波形信号の音声セグメントのデータを上記音声波形データベースから削除する削減手段とを備えたことを特徴とする音声合成装置のための音声データ量削減装置。

【請求項2】 上記類似度は、それぞれ所定の重み係数で重み付けされた、上記韻律的特徴パラメータに関する類似度のスコアと、上記音響特徴パラメータに関する類似度のスコアとの線形結合の式を用いて計算されることを特徴とする請求項1記載の音声合成装置のための音声データ量削減装置。

【請求項3】 上記韻律的特徴パラメータは、音素時間

10

長と、音声基本周波数 F_0 と、経過時間に対する音声基本周波数 F_0 の傾きと、パワーとを含むことを特徴とする請求項 1 又は 2 記載の音声合成装置のための音声データ量削減装置。

【請求項 4】 上記音響的特徴パラメータは、スペクトラム情報を含むことを特徴とする請求項 1 乃至 3 のうちの 1 つに記載の音声合成装置のための音声データ量削減装置。

【請求項 5】 請求項 1 乃至 4 のうちの 1 つに記載の音声合成装置のための音声データ量削減装置によって音声データ量が削減された音声波形データベースに基づいて、入力された自然発話文の音素列に対して、音素候補を上記音声波形データベースから検索して連結することにより音声合成を行なう音声合成手段を備えたことを特徴とする音声合成装置。

【請求項 6】 請求項 1 乃至 4 のうちの 1 つに記載の音声合成装置のための音声データ量削減装置によって音声データ量が削減された音声波形データベースに基づいて、入力された自然発話文の音素列に対して、目標音素と音素候補との間の近似コストと、時間的に隣接して連結されるべき音素候補間の近似コストとを含むコストが最小となるように、音素候補を上記音声波形データベースから検索して連結することにより音声合成を行なう音声合成手段を備えたことを特徴とする音声合成装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、音素ラベルに対応した音声波形信号の音声セグメントのデータからなる音声波形データベースを用いて、自然発話の音声波形信号の音声セグメントを連結することにより任意の音素列を音声合成する音声合成装置のための音声データ量削減装置及び音声合成装置に関する。

【0002】

【従来の技術】例えば、特開平 10 - 049193 号公報において、音素ラベルに対応した音声波形信号の音声セグメントのデータからなる音声波形データベースを用いて、自然発話の音声波形信号の音声セグメントを連結することにより任意の音素列を音声合成する音声合成装置が開示されている。

【0003】

【発明が解決しようとする課題】この従来例の音声合成装置においては、最小限の信号処理を用いて単純に音声波形の連結を行うために、適切な音響的パラメータと韻律的パラメータとを有する音声セグメントを選択するため、大規模な音声波形データベースを必要とする。この大規模な音声波形データベースでは、当該データベースを記憶するメモリの容量が大きくなり、また、それに伴って探索空間が大きくなるために、適切な音声セグメントを探索するときの探索速度を高めることができないという問題点があった。

【0004】本発明の目的は以上の問題点を解決し、音声波形データベースを格納するメモリ容量を削減することができ、音声合成時の探索速度を高めることができる音声合成装置のための音声データ量削減装置及び音声合成装置を提供することにある。

【0005】

【課題を解決するための手段】本発明に係る音声合成装置のための音声データ量削減装置は、音素ラベルに対応した音声波形信号の音声セグメントのデータからなる音声波形データベースを記憶する記憶装置を備え、上記自然発話の音声波形信号の音声セグメントを連結することにより任意の音素列を音声合成する音声合成装置のための音声データ量削減装置であって、上記音声波形データベースに含まれる 1 対の音素のリストを生成する生成手段と、上記生成された 1 対の音素のリストに基づいて各 1 対の音素に対する韻律的特徴パラメータと音響的特徴パラメータとに関する所定の類似度を計算し、上記計算された類似度が所定の第 1 のしきい値以上であるとき、当該各 1 対の音素のうち一方の 1 対の音素に係る音声波形信号の音声セグメントのデータを上記音声波形データベースから削除することにより音声データ量を削減するとともに、上記計算された類似度が所定の第 1 のしきい値以上であり、かつ上記 1 対の音素のリスト中の同一の 1 対の音素の数が所定の第 2 のしきい値以上であるときに、当該各 1 対の音素のうち一方の 1 対の音素に係る音声波形信号の音声セグメントのデータを上記音声波形データベースから削除する削減手段とを備えたことを特徴とする。

【0006】

【0007】また、上記音声データ量削減装置において、上記類似度は、好ましくは、それぞれ所定の重み係数で重み付けされた、上記韻律的特徴パラメータに関する類似度のスコアと、上記音響特徴パラメータに関する類似度のスコアとの線形結合の式を用いて計算されることを特徴とする。

【0008】さらに、上記音声データ量削減装置において、上記韻律的特徴パラメータは、好ましくは、音素時間長と、音声基本周波数 F_0 と、経過時間に対する音声基本周波数 F_0 の傾きと、パワーとを含み、上記音響的

特徴パラメータは、スペクトラム情報を含む。

【0009】本発明に係る音声合成装置は、上記音声データ量削減装置によって音声データ量が削減された音声波形データベースに基づいて、入力された自然発話文の音素列に対して、音素候補を上記音声波形データベースから検索して連結することにより音声合成を行なう音声合成手段を備えたことを特徴とする。

【0010】また、本発明に係る音声合成装置は、上記音声データ量削減装置によって音声データ量が削減された音声波形データベースに基づいて、入力された自然発話文の音素列に対して、目標音素と音素候補との間の近

似コストと、時間的に隣接して連結されるべき音素候補間の近似コストとを含むコストが最小となるように、音素候補を上記音声波形データベースから検索して連結することにより音声合成を行なう音声合成手段を備えたことを特徴とする。

【0011】

【発明の実施の形態】以下、図面を参照して本発明に係る実施形態について説明する。

【0012】図1は、本発明に係る一実施形態である音声データ量削減処理装置のブロック図である。この実施形態の音声データ量削減装置は、図2の音声合成装置に提供される音素ラベルデータにてなるテキストデータベースと音声波形データベース内の音声データ量を削減するために、図3の音声データ量削減処理を用いて、各バイフォン（biphone；時間的に隣接する1対の音素をいう。）に対する評価韻律データとパイスpekトラム（詳細後述）とのデータを含む評価データ行列に基づいて、テキストデータベースと音声波形データベース内で所定の類似度以上の音声データを削除することにより音声データ量の削減を行う音声データ量削減処理部45を備えたことを特徴とする。

【0013】すなわち、本実施形態においては、韻律的特徴パラメータの領域及びスペクトル領域の両方における近似性の物理的尺度を使用して、出力される音声合成された音声の品質を維持しながら音声波形データベースにおける音声セグメントの数を減らして音声波形データベースにおける冗長さを小さくする方法について開示する。

【0014】本実施形態では、上記物理的尺度の音響的特徴パラメータとして、パイスpekトラムを用いる。パワースpekトラムは2次のデータによってのみ決定される*

$$B(\omega_1, \omega_2) = \sum_{m=-\infty}^{+\infty} \sum_{n=-\infty}^{+\infty} R(m, n) \exp\{-j(\omega_1 m + \omega_2 n)\}$$

【0018】ここで、R(m, n)は上述の3次の自己相関関数であり、3次のモーメントと累積係数（累積率ともいう。）が同一であるので、パイスpekトラムは3次の累積係数のスペクトラムとなる。パワースpekトラムとパイスpekトラムの物理的な重要性は、次式のX(k)のフーリエ・スティルチェスの表現式（クラメル

【0019】すべてのkに対して、

【数5】

$$X(k) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} e^{j\omega k} dZ(\omega)$$

であり、ここで、

【数6】E{dZ()} = 0

【数7】E{dZ()dZ()}

= 0, 1 2 のとき

= 2 P() d , 1 = 2 = のとき

*ため、より高次の情報は無視される。もし音声が高ス分布を示すならば、2次のデータのみで完全に音声を復元合成することができる。しかしながら、実際は音声が高ス分布ではない。以上の2つの理由により、本発明者は音声セグメントの類似度を測定するための音響特徴パラメータの尺度として、パイスpekトラムを使用して音声合成を評価する。

【0015】ここで、本実施形態において用いるパイスpekトラム（bispectrum）について、従来技術文献1

10 「C. L. Nikias et al., "Bispectrum Estimation: A Digital Signal Processing Framework", Proceedings of the IEEE 76, pp.869-891,1987」を参照して説明する。この実施形態において、実離散処理を仮定しており、{X(k)}を実離散ゼロ平均定常処理とすると、パワースpekトラムP()は次式で定義される。

【0016】

【数1】

$$P(\omega) = \sum_{\tau=-\infty}^{+\infty} r(\tau) \exp\{-j(\omega\tau)\}, |\omega| < \pi$$

20 ここで、

【数2】r() = E{X(k)X(k+)}

はその自己相関シーケンスである。もしR(m, n)が{X(k)}の3次のモーメントのシーケンスを示すならば、すなわち、

【数3】R(m, n) = E{X(k)X(k+m)X(k+n)}

であるとき、そのパイスpekトラムは次式で定義される。

【0017】

【数4】

並びに、

【数8】E{dZ()dZ()dZ()}

$$= B(,) d_1 d_2 , 1 + 2 = 3$$

$$= 0 , 1 + 2 = 3$$

である。

【0020】従って、パワースpekトラムP()は周波数が同一である2つのフーリエ成分の平均値の積に寄与することを表わす一方、パイスpekトラムB(,)は、1つの周波数が他の2つの周波数の和に等しいときの3つのフーリエ成分の平均値の積に寄与することを表わす。すなわち、パイスpekトラムは、2つのスペクトラム列に対する寄与度又は類似度を表わす。

【0021】本実施形態では、2次元のフーリエ変換を使用してパイスpekトラムを計算する。パイスpekトラムの次元は一般に高いため、2次元のDCTを使用してパイスpekトラムをより低次の係数に圧縮して計算する。な

50 お、本実施形態においては、音響特徴パラメータとして

バイスペクトラムを用いているが、本発明はこれに限らず、スペクトラム情報を用いてもよい。

【0022】本実施形態では、上記物理的尺度の韻律的特徴パラメータとして、各音声セグメントの持続時間（又は音素時間長）、最大信号振幅（又はパワー）、平均基本周波数（又はピッチ周波数） F_0 及び基本周波数 F_0 の傾き（又は傾斜）を用いる。1対の音声セグメント間の音声距離を決定するにはスペクトル領域における物理的尺度が必要であるが、音声波形データベースの多様性を維持しようとする場合には、1対の音声セグメント間の韻律的特徴パラメータの距離を決定することも必要である。スペクトル特性及び韻律的特性の両方で所定のしきい値内で類似する（又は重複する）音声セグメントは音声波形データベースから除去することが可能である。言い換えれば、音声波形データベースのある部分を削除しても、残った部分から削除前と実質的に同一の音声波形を合成できるならば、その部分を削除してもかまわないという技術的思想に基づいている。これにより音声波形データベースの量は減少するが、音声波形データ*

$$M = w_1 \times (w_{11} \times S_{F0} + w_{12} \times S_P + w_{13} \times S_D + w_{14} \times S_{F1}) + w_2 \times S_{BS}$$

【0026】ここで、

(a) w_1 : 韻律的特徴パラメータのスコアに対する重み係数であり、例えば0.5である。

(b) w_{11} : 基本周波数 F_0 のスコアに対する重み係数であり、例えば0.3である。

(c) S_{F0} : 1対のバイフォンに対する基本周波数 F_0 の差の絶対値である。

(d) w_{12} : パワーのスコアに対する重み係数であり、

例えば0.2である。

(e) S_P : 1対のバイフォンに対するパワーの差の絶対値である。

(f) w_{13} : 音素時間長のスコアに対する重み係数であり、例えば0.1である。

(g) S_D : 1対のバイフォンに対する音素時間長の比の値（又は比率）である。

(h) w_{14} : 基本周波数 F_0 の傾きのスコアに対する重み係数であり、例えば0.3である。

(i) S_{F1} : 1対のバイフォンに対する基本周波数 F_0 の傾きに係るスコアであり、2つのバイフォンの音声波形の変化パターンの組み合わせから求める。

< 1 > まず、基本周波数 F_0 の変化を上昇、水平（実質的に変化せずを意味する。）、下降の3つのパターンに分類する。

< 2 - 1 > 音声波形Aと音声波形Bのパターンが同じならば、スコア = 0とする。

< 2 - 2 > 音声波形Aと音声波形Bのパターンが次の表に該当するときは、スコア = 1とする。

【表1】

* としてのカバー範囲は実質的に変化しないといえる。本実施形態においては、韻律的特徴パラメータの類似度を測定するため、音素サイズの各音声セグメントの上述の韻律的特徴パラメータを用いる。

【0023】そして、本実施形態においては、各1対のバイフォンに対して、以下に示す類似度のスコアMを計算して所定のしきい値M_{th}以下のときに、1対のバイフォンは互いに類似していると判断し、かつこのバイフォンの数Nを予め決められたしきい値N_{th}未満とならないように、すなわち当該しきい値N_{th}以上であれば削除しても複製可能であり、それは冗長であると判断して、当該バイフォンの音素に関するデータをテキストデータベースメモリ22及び音声波形データベースメモリ21から削除することにより、音声データ量の削減を行う。

【0024】1対のバイフォンに対する類似度のスコアMは次式で表わされる。

【0025】

【数9】

波形 A	波形 B
上昇	水平
水平	上昇
水平	下降
下降	水平

< 2 - 2 > 音声波形Aと音声波形Bのパターンが次の表に該当するときは、スコア = 1とする。

【表2】

波形 A	波形 B
上昇	下降
下降	上昇

(j) w_2 : バイスペクトラムに対する重み係数であり、例えば、0.8である。

(k) S_{BS} : 1対のバイフォンのスペクトラム列のベクトルの距離、すなわち、各スペクトラム要素の2乗和の平方根である。

【0027】上記数9から明らかなように、類似度のスコアMが0に近いほど、類似度の度合いが高いといえる。

【0028】次いで、図1を参照して音声データ量削減処理装置の構成及び動作について説明する。図1において、テキストデータベースメモリ22は、自然発話の書

き下し文の音素ラベルデータにてなるテキストデータベースを記憶し、音声波形データベースメモリ 2 1 は、上記テキストデータベースにおける自然発話の書き下し文の音素ラベルデータに対応する音声波形信号の音声セグメントからなる音声波形データベースを記憶する。基本韻律データ生成部 4 0 は、これら 2 つのメモリ 2 1, 2 2 内のデータに基づいて、各音素ラベルに対応して、例えば L P C 分析法を用いて基本周波数 F_0 を検出するとともに、所定の音声波形分析を行うことによりパワーを検出することにより、基本韻律データとして生成して基本韻律データメモリ 5 0 に格納する。

【0029】次いで、評価韻律データ生成部 4 1 は、テキストデータベースメモリ 2 2 内の音素ラベルデータと、基本韻律データメモリ 5 0 内の基本韻律データとに基づいて、基本周波数 F_0 とパワーに加えて、基本周波数 F_0 の時間方向の傾きと、音素時間長を計算してこれら 4 つのデータを評価韻律データとして評価韻律データメモリ 5 1 に格納する。また、バイフォンリスト生成部 4 2 は、テキストデータベースメモリ 2 2 内の音素ラベルデータに基づいて、時間方向に隣接する 2 つの音素並びであるバイフォンを計数して、そのバイフォンと計数値をバイフォンリストメモリ 5 2 に格納する。さらに、バースペクトラムデータ生成部 4 3 は、テキストデータベースメモリ 2 2 内の音素ラベルデータと音声波形データベースメモリ 2 1 内の音声波形データベースとに基づいて、上述のバースペクトラムデータ、すなわち、1 対のバイフォンのスペクトラム列のベクトルの距離、すなわち、各スペクトラム要素の 2 乗和の平方根を計算してバースペクトラムデータメモリ 5 3 に格納する。

【0030】さらに、評価データ行列生成部 4 4 は、評価韻律データメモリ 5 1 内の評価韻律データと、バイフォンリストメモリ 5 2 内のバイフォンと、バースペクトラムデータメモリ 5 3 内のバースペクトラムデータとに基づいて、テキストデータベース内の各バイフォンの列に対して、上述の 4 つの評価韻律データと 1 つのバースペクトラムデータを行方向並置することにより、評価データ行列を生成して評価データ行列メモリ 5 4 に格納する。そして、音声データ量削減処理部 4 5 は、評価データ行列メモリ 5 4 内の評価データ行列と、バイフォンリストメモリ 5 2 内のバイフォンに基づいて、図 3 の音声データ量削減処理を実行することにより、テキストデータベースメモリ 2 2 内のテキストデータと音声波形データベースメモリ 2 1 内の音声波形データベースの音声データ量を削減して、削減後のデータをテキストデータベースメモリ 2 2 a と音声波形データベースメモリ 2 1 a にコピーする。そして、図 2 の音声合成装置は、音声データ量が削減された後のテキストデータベースメモリ 2 2 a 及び音声波形データベースメモリ 2 1 a 内のデータを用いて音声合成処理を行う。

【0031】図 3 は、図 1 のデータ量削減処理部によ

て実行されるデータ量削減処理を示すフローチャートである。図 3 において、まず、ステップ S 1 においてバイフォンリストメモリ 5 2 内のバイフォンリストから同一の 1 対のバイフォンを検索し、ステップ S 2 において存在するかが判断され、YES のときはステップ S 3 に進む一方、NO のときはステップ S 9 に進む。ステップ S 3 では、検索した 1 対のバイフォンに対する評価韻律データとバースペクトラムデータに基づいて、上記数 9 を用いて類似度のスコア M を計算する。

10 【0032】そして、ステップ S 4 において 1 対のバイフォンが類似しているかが判断され、具体的には M M_{th} であるかが判断され、YES のときは類似していると判断してステップ S 5 に進む一方、NO のときは類似していないと判断してステップ S 8 に進む。ここで、 M_{th} は予め決められたしきい値であり、例えば 0.1 である。次いで、ステップ S 5 においてバイフォンリストメモリ 5 2 内のバイフォンリストから同一のバイフォンの件数 N を計算し、ステップ S 6 において N N_{th} が判断され、YES のときはステップ S 7 20 に進む一方、NO のときはステップ S 8 に進む。ここで、 N_{th} は予め決められたしきい値であり、例えば 50 であり、これは元のテキストデータベースのバイフォンの数にも依存する。ステップ S 6 の判断では、音声波形データベースにおいて所定数の同一のバイフォンに対する音声波形データを確保して、音声合成後の音声の品質を所定以上に確保するために設けられる。さらに、ステップ S 7 において当該 1 対のバイフォンのうちの一方のバイフォンの音素をバイフォンリストから削除し、当該バイフォンの音素のラベルデータ及び音声波形データをそれぞれテキストデータベース及び音声波形データベースから削除して、ステップ S 8 に進む。ステップ S 8 30 では、バイフォンリストから別の組み合わせの 1 対のバイフォンを検索して、ステップ S 2 に戻る。

【0033】ステップ S 2 において NO であるときは、ステップ S 9 においてバイフォンリストから別の種類のバイフォンを検索し、ステップ S 10 において存在するかが判断され、YES のときはステップ S 1 に戻る一方、NO のときはステップ S 10 a に進む。ステップ S 10 a では、音声データ量が削減された後のテキストデータベースメモリ 2 2 及び音声波形データベースメモリ 2 1 内のデータをそれぞれ、テキストデータベースメモリ 2 2 a 及び音声波形データベースメモリ 2 1 a にコピーして当該音声データ量削減処理を終了する。

【0034】以上の図 3 の実施形態において、ステップ S 4 における判断とステップ S 6 における判断がともに YES であるときに、当該バイフォンの音素のラベルデータと音声波形データとを削除しているが、本発明はこれに限らず、ステップ S 4 における判断のみが YES であるときに、当該バイフォンの音素のラベルデータと音声波形データとを削除してもよい。

【0035】次いで、図1の音声データ量削減処理装置で音声データ量が削減された、テキストデータベースメモリ22a及び音声波形データベースメモリ21a内のデータを用いて音声合成を行う音声合成装置について以下に説明する。

【0036】図2は、本発明に係る一実施形態である自然発話音声波形信号接続型音声合成装置のブロック図である。本実施形態では、大きく分類すれば、次の4つの処理部に分類される。(1)音声波形信号データベースメモリ21a内の音声波形信号データベースの音声波形信号データの音声分析、具体的には、音素記号系列の生成、音素のアラインメント、特徴パラメータの抽出を含む処理を実行する音声分析部10。(2)最適重み係数を学習しながら決定する重み係数学習部11。(3)入力される音素列に基づいて音声単位の選択を実行して入力音素列に対応する音声波形信号データの索引情報を出力する音声単位選択部12。(4)音声単位選択部12から出力される索引情報に基づいて音声波形信号データベースメモリ21a内の音声波形信号データベースをランダムにアクセスして最適とされた各音素候補の音声波形信号を再生してスピーカ14に出力する音声合成部13。

【0037】具体的には、音声分析部10は、入力される自然発話の音声波形信号の音声セグメントと、上記音声波形信号に対応する音素列とに基づいて、音素HMMメモリ23を参照して、上記音声波形信号における音素毎の索引情報と、上記索引情報によって示された音素毎の第1の音響的特徴パラメータと、上記索引情報によって示された音素毎の第1の韻律的特徴パラメータとを抽出して出力する。特徴パラメータメモリ30は、上記音声分析部10から出力される索引情報と、上記第1の音響的特徴パラメータと、上記第1の韻律的特徴パラメータとを記憶する。次いで、重み係数学習部11は、特徴パラメータメモリ30に記憶された第1の音響的特徴パラメータと韻律的特徴パラメータとに基づいて、同一の音素種類の1つの目標音素とそれ以外の音素候補との間の第2の音響的特徴パラメータにおける音響的距離を計算し、上記計算した音響的距離に基づいて各音素候補に対して上記第2の音響的特徴パラメータ毎に所定の統計的解析を実行することにより、各音素候補に対する上記第2の音響的特徴パラメータにおける寄与度を表わす各目標音素毎の重み係数ベクトルを決定する。重み係数ベクトルメモリ31は、重み係数学習部11によって決定された上記第2の音響的特徴パラメータにおける各目標音素毎の重み係数ベクトルと、予め与えられた、各音素候補に関する第2の韻律的特徴パラメータにおける寄与度を表わす各目標音素毎の重み係数ベクトルとを記憶する。さらに、音声単位選択部12は、重み係数ベクトルメモリ31に記憶された各目標音素毎の重み係数ベクトルと、特徴パラメータメモリ30に記憶された第1の韻

律的特徴パラメータとに基づいて、入力される自然発話文の音素列に対して、目標音素と音素候補との間の近似コストを表わす目標コストと、隣接して連結されるべき2つの音素候補間の近似コストを表わす連結コストとを含むコストが最小となる、音素候補の組み合わせを検索して、検索した音素候補の組み合わせの索引情報を出力する。そして、音声合成部13は、音声単位選択部12から出力される索引情報に基づいて、当該索引情報に対応する音声波形信号の音声セグメントを音声波形信号データベースメモリ21aから逐次読み出して連結してスピーカ14に出力することにより、音声合成装置は、上記入力された音素列に対応する音声合成して出力する。

【0038】ここで、音声分析部10の処理は新しい音声波形信号データベースに対しては必ず一度行なう必要があるが、重み係数学習部11の処理は、一般に一度の処理でよく、重み係数学習部11によって求めた最適重み係数は異なる音声合成条件に対しても再利用が可能である。さらに、音声単位選択部12と音声合成部13の処理は、音声合成すべき入力音素列が変われば、その都度実行される。

【0039】本実施形態の音声合成装置は与えられたレベルの入力に基づいて必要とする、すべての特徴パラメータを予測し、所望の音声の特徴に最も近いサンプル(すなわち、音素候補の音声波形信号)を音声波形データベースメモリ21a内の音声波形信号データベースの中から選び出す。最低限、音素ラベルの系列が与えられれば処理は可能であるが、音声基本周波数F₀や音素時間長が予め与えられていれば、さらに高品質の合成音声を得られる。なお、入力として単語の情報だけが与えられた場合には、例えば音素隠れマルコフモデルメモリ23に格納された音素隠れマルコフモデル(以下、隠れマルコフモデルをHMMという。)などの辞書や規則に基づいて音素系列を予測する必要がある。また、韻律特徴が与えられなかった場合には音声波形信号データベース中のいろいろな環境における音素の既知の特徴を基に標準的な韻律を生成する。

【0040】本実施形態では、音声波形信号データベースメモリ21a内の録音内容を少なくとも正書法で記述されたテキストデータが例えば、テキストデータベースメモリ22a内のテキストデータベースのように存在するならば、あらゆる音声波形信号データベースが合成用の音声波形信号データとして利用可能であるが、出力音声の品質は録音状態、音声波形信号データベース中の音素のバランス等に大きく影響を受け、音声波形データベースメモリ21a内の音声波形信号データベースが豊富な内容であれば、より多様な音声合成でき、反対に音声波形信号データベースが貧弱であれば、合成音声は不連続感が強く、ブツブツしたものになる。

【0041】次いで、自然な発話音声に対する音素ラベ

10

20

30

40

50

ル付けについて説明する。音声単位の選択の善し悪しは音声波形信号データベース中の音素のラベル付けと検索の方法に依存する。ここで、好ましい実施例においては、音声単位は、音素である。まず、録音された音声に付与された正書法の発話内容を音素系列に変換し、さらに音声波形信号に割り当てる。韻律的特徴パラメータの抽出はこれに基づいて行なわれる。音声分析部 10 の入力テキストデータベースメモリ 22 a 内の音素表記を伴った音声波形データベースメモリ 21 a 内の音声波形信号データであり、出力は特徴ベクトル又は特徴パラメータである。この特徴ベクトルは音声波形信号データベース中で音声サンプルを表す基本単位となり、最適な音声単位の選択に用いられる。

【0042】音声分析部 10 の処理における第 1 段階においては、正書法で書かれた発話内容が実際の音声波形信号データでどのように発音されているかを記述するための正書法テキストから音素記号への変換である。次い*

* で、第 2 段階においては、韻律的及び音響的特徴を計測するために各音素の開始及び終了時点を決めるために、各音素記号を音声波形信号に対応付ける処理である（以下、当該処理を、音素のアラインメント処理という。）。さらに、第 3 段階においては、各音素の特徴ベクトル又は特徴パラメータを生成することである。この特徴ベクトルには、必須項目として音素ラベル、メモリ 30 内の音声波形信号データベース中の各ファイルにおける当該音素の開始時刻（開始位置）、音声基本周波数 F_0 、音素時間長、パワーの情報が記憶され、さらに、特徴パラメータのオプションとしてストレス、アクセント型、韻律境界に対する位置、スペクトル傾斜等の情報が記憶される。以上の特徴パラメータを整理すると、例えば、次の表のようになる。

【0043】

【表 3】

索引情報：

索引番号（1つのファイルに対して付与）

メモリ 30 内の音声波形信号データベース中の各ファイルにおける当該音素の開始時刻（開始位置）

第 1 の音響的特徴パラメータ：

1 2 次メルケプストラム係数

1 2 次メルケプストラム係数

音素ラベル

弁別素性：

母音性 (vocalic) (+) / 非母音性 (non-vocalic) (-)

子音性 (consonantal) (+) / 非子音性 (non-consonantal) (-)

中断性 (interrupted) (+) / 連続性 (continuant) (-)

抑止性 (checked) (+) / 非抑止性 (unchecked) (-)

粗擦性 (strident) (+) / 円熟性 (mellow) (-)

有声 (voiced) (+) / 無声 (unvoiced) (-)

集約性 (compact) (+) / 拡散性 (diffuse) (-)

低音調性 (grave) (+) / 高音調性 (acute) (-)

変音調性 (flat) (+) / 常音調性 (plain) (-)

嬰音調性 (sharp) (+) / 常音調性 (plain) (-)

緊張性 (tense) (+) / 弛緩性 (lax) (-)

鼻音性 (nasal) (+) / 口音性 (oral) (-)

第 1 の韻律的特徴パラメータ：

音素時間長

音声基本周波数 F_0

パワー

【0044】として代わって、第 1 の音響的特徴パラメータは、好ましくは、フォルマントパラメータと、声道音源パラメータであってもよい。上記索引情報内の開始時刻（開始位置）、第 1 の音響的特徴パラメータ及び第

1 の韻律的特徴パラメータは、各音素毎に特徴パラメータメモリ 30 に記憶される。ここで、音素ラベルに付与される、例えば 1 2 個の弁別素性の特徴パラメータは各項目別に (+) 又は (-) のパラメータ値が与えられ

る。さらに、音声分析部 10 の出力結果である特徴パラメータの一例を次の表に示す。ここで、索引番号は、音声波形信号データベースメモリ 21 a において、例えば複数の文からなる 1 つのパラグラフ又は 1 つの文のファイル毎に、索引番号が付与され、そして、1 つの索引番号が付与されたファイル中の任意の音素の位置を示すために当該ファイル内の開始時刻から計時された当該音素の開始時刻及びその当該音素の音素時間長とを付与することにより、当該音素の音声波形信号の音声セグメントを特定することができる。

【0045】

【表 4】音声分析部 10 の出力結果である特徴パラメータの一例
索引番号 X 0 0 0 5

音素	時間長	基本周波数	パワー
#	1 2 0	9 0	4 . 0
s	1 7 5	9 8	4 . 7
e i	9 5	1 0 2	6 . 5
d h	3 0	1 1 4	4 . 9
i h	7 5	1 4 3	6 . 9
s	1 5 0	1 4 0	5 . 7
p	8 7	1 3 7	5 . 1
l	3 4	1 0 7	4 . 9
i i	1 5 0	9 8	6 . 3
z	1 4 0	8 7	5 . 8
#	2 5 3	8 7	4 . 0

【0046】表 4 において、# はポーズを示す。音声単位を選択する場合に、音響的及び韻律的な各特徴パラメータがそれぞれの音素でどれだけの寄与をするかを予め調べておくことが必要であり、第 4 段階では、このために音声波形信号データベース中のすべての音声サンプルを用いて各特徴パラメータの重み係数を決定する。

【0047】音声分析部 10 における音素記号系列の生成処理においては、上述した通り、本実施形態では、少なくとも録音内容が正書法で記述されたものがあれば、あらゆる音声波形信号データベースが合成用の音声波形信号データとして利用可能である。入力として単語の情報だけが与えられた場合には辞書や規則に基づいて音素系列を予測する必要がある。また、音声分析部 10 における音素のアライメント処理においては、読み上げ音声の場合、各単語がそれぞれの標準の発音に近く発音されることが多く、躊躇したり、言い淀んだりすることもまれである。このような音声波形信号データの場合には簡単な辞書検索によって音素ラベリングが正しく行なわれ、音素アライメント用の音素 HMM の音素モデルの学習が可能となる。

【0048】音素アライメント用の音素モデルの学習

では完全な音声認識の場合と異なり、学習用の音声波形信号データとテスト用の音声波形信号データとを完全に分離する必要はなく、すべての音声波形信号データを用いて学習を行なうことができる。まず、別の話者用のモデルを初期モデルとし、すべての単語について標準発音が限られた発音変化のみを許し、適切なセグメンテーションが行なわれるように、全音声波形信号データを用いてピタビの学習アルゴリズムを用いて音素のアライメントを行ない、特徴パラメータの再推定を行なう。単語間のポーズは単語間ポーズ生成規則によって処理するが、単語内にポーズがあってアライメントが失敗した場合には人手により修正する必要がある。

10

20

【0049】どういった音素ラベルを音素表記として用いるかは選択が必要である。もし良く学習された HMM モデルが利用できるような音素セットが存在するならば、それを用いることが有利である。反対に、音声合成装置が完全な辞書を持っているならば、音声波形信号データベースのラベルを完全に辞書と照合する方法も有効である。我々は、重み係数の学習に対して選択の余地があるから、後で音声合成装置が予測したものと等価なものを音声波形信号データベースの中から照合できるかどうかを最も重要な基準とすればよい。発音の微妙な違いはその発音の韻律的環境によって自動的に把握されるため、特に手作業で音素のラベル付けを行なう必要はない。

【0050】前処理の次の段階として、個々の音素の調音的な特徴を記述するための韻律特徴パラメータの抽出を行なう。従来の音声学では、調音位置や調音様式といった素性で言語音を分類した。これに対して、ファース (Firth) 学派のような韻律を考慮した音声学では、韻律的文脈の違いから生ずる細かな音質の違いをとらえるために、明瞭に調音されている箇所や強調が置かれている箇所を区別する。これらの違いを記述する方法はいろいろなものがあるが、ここでは以下の 2 つの方法を用いる。まず低次のレベルでは、1 次元の特徴を求めるために、パワー、音素時間長の伸び及び音声基本周波数 F_0 を、ある音素について平均した値を用いる。一方、高次のレベルでは、韻律特徴における上記の違いを考慮した韻律境界や強調箇所をマークする方法を用いる。これらの 2 種類の特徴は相互に密接に関係しているため一方から他方を予測することができるが、両者は共に各音素の特徴に強い影響を与えている。

30

40

【0051】音声波形信号データベースを記述するための音素セットの規定法に自由度があるのと同様に、韻律的特徴パラメータの記述方法についても自由度があるが、これらの選び方は音声合成装置の予測能力に依存する。もし音声波形信号データベースが予めラベリングされているならば、音声合成装置の仕事は内部表現から音声波形信号データベース中の実音声をいかに行なうかを学習することである。これに対して、もし音声波形信号データベースが音素のラベル付けがなされていないならば、

50

どのような特徴パラメータを使えば音声合成装置が最も適切な音声単位を予測できるか否か、から検討することが必要となる。この検討及び最適な特徴パラメータの重みの決定学習は、各特徴パラメータに対する重み係数を学習しながら決定する重み係数学習部 11 において実行される。

【0052】次いで、重み係数学習部 11 によって実行される重み係数学習処理について述べる。与えられた目標音声の音響的及び韻律的な環境に最適なサンプルを音声波形信号データベースから選択するために、まずどの特徴がどれだけ寄与しているかを音素的及び韻律的な環境の違いによって決める必要がある。これは音素の性質によって重要な特徴パラメータの種類が変化するため、例えば、音声基本周波数 F_0 は有声音の選択には極めて有効であるが、無声音の選択にはほとんど影響がない。また、摩擦音の音響的特徴は前後の音素の種類によって影響が変わる。最適な音素を選択するためにそれぞれの特徴にどれだけの重みを置くかを最適重み決定処理、すなわち重み係数学習処理で自動的に決定する。

【0053】重み係数学習部 11 によって実行される最適重み係数の決定処理で、最初に行なわれることは音声波形信号データベース中で該当するすべての発話サンプルの中から最適なサンプルを選ぶときに使われる特徴をリストアップすることである。ここでは、調音位置や調音様式等の音素的特徴と先行音素、当該音素、及び後続音素の音声基本周波数 F_0 、音素時間長、パワー等の韻律的特徴パラメータ等を用いる。具体的には、詳細後述する第 2 の韻律的特徴パラメータを用いる。次いで、第 2 段階では各音素毎に、最適な候補を選ぶ際にどの特徴パラメータがどれだけ重要かを決定するために、1 つの音声サンプル（又は音素の音声波形信号）に着目し、他のすべての音素サンプルとの音素時間長の差をも含む音響的距離を求め、上位 N 2 個の最良の類似音声サンプル、すなわち N 2 ベストの音素候補の音声波形信号の音声セグメントを選び出す。

【0054】さらに、第 3 段階では線形回帰分析を行ない、それらの類似音声サンプルを用いて種々の音響的及び韻律的環境におけるそれぞれの特徴パラメータの重要度を示す重み係数を求める。当該線形回帰分析処理における韻律的特徴パラメータとして、例えば、次の特徴パラメータ（以下、第 2 の韻律的特徴パラメータという。）を用いる。

(1) 処理すべき当該音素から 1 つだけ先行する先行音素（以下、先行音素という。）の第 1 の韻律的特徴パラメータ；

(2) 処理すべき当該音素から 1 つだけ後続する後続音素（以下、後続音素という。）の音素ラベルの第 1 の韻律的特徴パラメータ；

(3) 当該音素の音素時間長；

(4) 当該音素の音声基本周波数 F_0 ；

(5) 先行音素の音声基本周波数 F_0 ；及び、

(6) 後続音素の音声基本周波数 F_0 。

ここで、先行音素は、当該音素から 1 つだけ先行する音素としているが、これに限らず、複数の音素だけ先行する音素を含んでもよい。また、後続音素は、当該音素から 1 つだけ後続する音素としているが、これに限らず、複数の音素だけ後続する音素を含んでもよい。さらに、後続音素の音声基本周波数 F_0 を除外してもよい。以上の実施形態においては、線形回帰分析を行って、重み係数を求めているが、本発明はこれに限らず、例えば、所定のニューラルネットワークを用いた統計的解析などの種々の統計的解析を用いて、重み係数を求めてもよい。

【0055】次いで、自然な音声サンプルの選択を行う音声単位選択部 12 の処理について説明する。従来例の音声合成装置では目的の発話に対して音素系列を決定し、さらに韻律制御のための F_0 と音素時間長の目標値が計算された。これに対して、本実施形態では最適の音声サンプルを選択するために韻律が計算されるだけで、直接韻律を制御することは行なわれない。

【0056】図 4 にはこの音声合成装置における図 2 の音声単位選択部 12 の処理を示す。この処理の入力は、目的発話の音素系列と、それぞれの音素毎に求めた各特徴に対する重みベクトル及び音声波形信号データベース中の全サンプルを表す特徴ベクトルである。一方、出力は音声波形信号データベース中の音素サンプルの位置を表す索引情報であって、音声波形信号の音声セグメントを接続するためのそれぞれの音声単位（具体的には音素、場合により複数の音素の系列が連続して選択され、一つの音声単位となることがある）の開始位置と音声単位時間長を示したものである。

【0057】最適な音声単位は目的発話との差の近似コストを表す目標コストと、隣接音声単位間での不連続性の近似コストを表す連結コストの和を最小化するパスとして求められる。経路探索には公知のビタビの学習アルゴリズムが利用される。目的とする目標音声 $t_1^n = (t_1, \dots, t_n)$ に対しては、目標コストと連結コストの和を最小化することで、各特徴が目的音声に近く、しかも音声単位間の不連続性が少ない音声波形信号データベース中の音声単位の組合せ $u_1^n = (u_1, \dots, u_n)$ を選ぶことができ、これらの音声単位の音声波形信号データベース内での位置を示すことにより、任意の発話内容の音声合成が可能になる。

【0058】音声単位の選択コストは、図 4 に示すように、目標コスト $C^1(u_i, t_i)$ と連結コスト $C^0(u_{i-1}, u_i)$ からなり、目標コスト $C^1(u_i, t_i)$ は、音声波形信号データベース中の音声単位（音素候補） u_i と、合成音声として実現したい音声単位（目標音素） t_i の間の差の予測値であり、連結コスト $C^0(u_{i-1}, u_i)$ は接続単位（接続する 2 つの音素） u_{i-1} と u_i との間の接続で起こる不連続の予測値である。例えば、本

出願人によって研究実用化された従来の A T R - T a l k 音声合成システムも目標コストと連結コストを最小化するという点では類似の考え方を取っていたが、韻律的な特徴パラメータを直接に単位選択に用いるということは本実施形態の音声合成装置の新しい特徴となっている。

【0059】次いで、コストの計算について述べる。目標コストは実現したい音声単位の特徴ベクトルと音波形信号データベース中から選ばれた候補の音声単位の特徴ベクトルの各要素の差の重み付き合計であり、各目標サブコスト $C^j(t_i, u_i)$ の重み係数 w^j が与えられた場合、目標コスト $C^i(t_i, u_i)$ は次式で計算することができる。

【0060】

【数10】

$$C^i(t_i, u_i) = \sum_{j=1}^p w^j C^j(t_i, u_i)$$

【0061】ここで、特徴ベクトルの各要素の差は p 個の目標サブコスト $C^j(t_i, u_i)$ (ただし、 j は 1 から p までの自然数である。) で表され、特徴ベクトルの次元数 p は、好ましい実施例においては、20 から 30 の範囲で可変としている。より好ましい実施例においては、次元数 $p = 30$ であり、目標サブコスト $C^j(t_i, u_i)$ 及び重み係数 w^j における変数 j の特徴ベクトル又は特徴パラメータは、上述の第 2 の韻律的特徴パラメータである。

【0062】一方、連結コスト $C^c(u_{i-1}, u_i)$ も同様に q 個の連結サブコスト $C^c_j(u_{i-1}, u_i)$ (ただし、 j は 1 から q までの自然数である。) の重み付き合計で表される。連結サブコストは接続する音声単位 u_{i-1} と u_i の音響的特徴から決定することができる。好ましい実施形態においては、連結サブコストとしては、(1) 音素接続点におけるケプストラム距離、(2) 対数パワーの差の絶対値、(3) 音声基本周波数 F_0 の差の絶対値の 3 種類を用いており、すなわち、 $q = 3$ である。これら 3 種類の音響的特徴パラメータと、先行音素の音素ラベルと、後続音素の音素ラベルとを、第 3 の音響的特徴パラメータという。各連結サブコスト $C^c_j(u_{i-1}, u_i)$ の重み w^j は予め経験的に (又は実験的に) *

$$C(t_1^n, u_1^n) = \sum_{i=1}^n \sum_{j=1}^p C^j(t_i, u_i) + \sum_{i=2}^n \sum_{j=1}^q C^c_j(u_{i-1}, u_i) + C^c(S, u_1) + C^c(u_n, S)$$

【0068】音声単位選択処理は上式で決まる全体のコストを最小にするような音声単位の組合せ u_1^n を決定するためのものである。ここで、日本出願の明細書で

* 与えられ、この場合、連結コスト $C^c(u_{i-1}, u_i)$ は次式で計算することができる。

【0063】

【数11】

$$C^c(u_{i-1}, u_i) = \sum_{j=1}^q w^j C^c_j(u_{i-1}, u_i)$$

【0064】もし、音素候補 u_{i-1} と u_i が音波形信号データベース中の連続する音声単位であった場合には、接続は自然であり、連結コストは 0 になる。ここで、好ましい実施例においては、連結コストは、特徴パラメータメモリ 30 内の第 1 の音響的特徴パラメータと第 1 の韻律的特徴パラメータに基づいて決定され、連続量である上記 3 つの第 3 の音響的特徴パラメータを取り扱うから例えば 0 から 1 までの任意のアナログ量をとる一方、目標コストは、それぞれの先行あるいは後続音素の弁別素性が一致するか否かなどを示す上記 30 個の第 2 の音響的特徴パラメータを取り扱うから、例えば 0 (特徴が一致しているとき) 又は 1 (特徴が一致していないとき) のデジタル量で表される要素を含む。そして、 N 個の音声単位の連結コストはそれぞれの音声単位の目標コストと連結コストの和となり、次式で表される。

【0065】

【数12】

$$C(t_1^n, u_1^n) = \sum_{i=1}^n C^i(t_i, u_i) + \sum_{i=2}^n C^c(u_{i-1}, u_i) + C^c(S, u_1) + C^c(u_n, S)$$

【0066】このとき、 S はポーズを表しており、 $C^c(S, u_1)$ 及び $C^c(u_n, S)$ はポーズから最初の音声単位へ及び最後の音声単位からポーズへの接続における連結コストを表している。この表現からも明らかのように、本実施形態ではポーズも音波形信号データベース中の他の音素とまったく同じ扱い方をしている。さらに上の式をサブコストで直接表現すると次式のようになる。

【0067】

【数13】

$$C(t_1^n, u_1^n) = \sum_{i=1}^n \sum_{j=1}^p C^j(t_i, u_i) + \sum_{i=2}^n \sum_{j=1}^q C^c_j(u_{i-1}, u_i) + C^c(S, u_1) + C^c(u_n, S)$$

は、オーバーラインを記述することができないために、オーバーラインの代わりに / を用いる。

【0069】

【数 1 4】
$$/ u_1^n = \min_{u_1, u_2, \dots, u_n} C(t_1^n, u_1^n)$$

【0 0 7 0】上記数 1 4 において、関数 \min は、当該関数の引数である $C(t_1^n, u_1^n)$ を最小にする音素候補の組み合わせ（すなわち、音素列候補） $u_1, u_2, \dots, u_n = / u_1^n$ を表わす関数である。

【0 0 7 1】図 2 の重み係数学習部 1 1 における重み係数の学習処理について以下説明する。目標サブコストの重みは音響的距離に基づく線形回帰分析を用いて決定する。重み係数の学習処理ではすべての音素毎に異なる重み係数を決めることもできるし、音素カテゴリ（例えば、すべての鼻音）毎に重み係数を決めることもできる。また、すべての音素について共通の重み係数を決めることもできるが、ここでは各音素で別々の重み係数を用いることとする。特徴パラメータメモリ 3 0 内のデータベースにおける各トークン（又は各音声サンプル）は、各トークンの音響的特徴に關係する第 1 の音響的特徴パラメータと第 1 の韻律的特徴パラメータの組で記述されている。重み係数は、第 1 の音響的特徴パラメータと第 1 の韻律的特徴パラメータの各パラメータと、トークン又はコンテキストにおける音素の第 2 の音響的特徴パラメータにおける差又は音響的距離との間の關係の強さ（寄与度）を決定するために学習される。以下に線形回帰分析における処理の流れを示す。

【0 0 7 2】< 1 > 現在学習を行なっている音素種類（又は音素カテゴリ）に属する音声波形信号データベース中のすべてのサンプルについて繰り返し以下の 4 つの処理（a）乃至（d）を実行する。

（a）取り上げた音声サンプルを目的の発話内容と見なす。

（b）音声波形信号データベース中の同一の音素種類（カテゴリ）に属する他のすべてのサンプルと当該音声サンプルとの音響的距離を計算する。

（c）目標音素に近いもの上位 $N 1$ 個（例えば、 $N 1 = 2 0$ 個である。）の最良の音素候補を選び出す。

（d）目標音素自身 t_i と上記（c）で選んだ上位 $N 1$ 個のサンプルについて目標サブコスト $C^1_j(t_i, u_i)$ を求める。

< 2 > すべての目標音素 t_i と上位 $N 1$ 個の最適サンプルについて音響的距離と目標サブコスト $C^1_j(t_i, u_i)$ を求める。

< 3 > p 個の目標サブコストに対して線形回帰分析を実行することにより、上記目標音素を表わす第 1 の音響的特徴パラメータと第 1 の韻律的特徴パラメータの各特徴パラメータにおける寄与度を予測して、当該音素種類（カテゴリ）に対する、 p 個の目標サブコストの線形重み係数を求める。この重み係数を用いて上記コストを計算する。そして、< 1 > から < 3 > までの処理をすべての音素種類（カテゴリ）について繰り返す。

【0 0 7 3】もし仮に目的音声単位の音響的距離が直接

求められた場合に最も近い音声サンプルを選び出すためにはそれぞれの目標サブコストにどのような重み係数をかければ良いのかを決定するのが、この重み係数学習部 1 1 の目的である。本実施形態の利点は音声波形信号データベース中の音声波形信号の音声セグメントを直接的に利用できることである。

【0 0 7 4】以上のように構成された図 2 の音声合成装置において、音声分析部 1 0 と、重み係数学習部 1 1 と、音声単位選択部 1 2 と、音声合成部 1 3 とは、例えば、マイクロプロセッシングユニット（MPU）などのデジタル計算機又は演算制御装置によって構成される一方、テキストデータベースメモリ 2 2 a と、音素 HMM メモリ 2 3 と、特徴パラメータメモリ 3 0 と、重み係数ベクトルメモリ 3 1 とは例えばハードディスクなどの記憶装置で構成される。ここで、好ましい実施例においては、音声波形信号データベースメモリ 2 1 a は、CD-ROM の形式の記憶装置である。以下、以上のように構成された図 2 の音声合成装置の各処理部 1 0 乃至 1 3 における処理について説明する。

【0 0 7 5】図 5 は、図 1 の音声分析部 1 0 によって実行される音声分析処理のフローチャートである。図 5 において、まず、ステップ S 1 1 で、音声波形信号データベースメモリ 2 1 a から自然発話の音声波形信号の信号を入力して A / D 変換してデジタル音声波形信号データに変換するとともに、当該音声波形信号の音声文を書き下したテキストデータをテキストデータベースメモリ 2 2 a 内のテキストデータベースから入力する。ここで、テキストデータはなくてもよく、ない場合は、音声波形信号から公知の音声認識装置を用いて音声認識してテキストデータを得てもよい。なお、A / D 変換した後のデジタル音声波形信号データは、例えば 1 0 ミリ秒毎の音声セグメントに分割されている。そして、ステップ S 1 2 で、音素列が予測されているか否かが判断され、音素列が予測されていないときは、ステップ S 1 3 で、例えば音素 HMM メモリ 2 3 内の音素 HMM を用いて音素列を予測して記憶した後、ステップ S 1 4 に進む。ステップ S 1 2 で音素列が予測されている又は予め与えられている、もしくは手作業で音素ラベルが付与されているときは、直接にステップ S 1 4 に進む。

【0 0 7 6】ステップ S 1 4 では、各音素セグメントに対する、音声波形信号の複数の文又は 1 つの文からなるファイルにおける開始位置と終了位置を記録し、当該ファイルに索引番号を付与する。次いで、ステップ S 1 5 では、各音素セグメントに対する上記第 1 の音響的特徴パラメータを例えば公知のピッチ抽出法を用いて抽出する。そして、ステップ S 1 6 では、各音素セグメントに対して音素ラベル付けを実行して、音素ラベルとそれに対する第 1 の音響的特徴パラメータを記録する。さらに、ステップ S 1 7 では、各音素セグメントに対する第 1 の音響的特徴パラメータと、音素ラベルと、音素ラベ

ルに対する上記第 1 の韻律的特徴パラメータを、ファイルの索引番号と、ファイル内の開始位置と時間長とともに、特徴パラメータメモリ 3 0 に記憶する。最後に、ステップ S 1 8 で、各音素セグメントに対して、ファイルの索引番号とファイル内の開始位置と時間長とを含む索引情報を付与して、当該索引情報を特徴パラメータメモリ 3 0 に記憶して、当該音声分析処理を終了する。

【0077】図 6 及び図 7 は、図 2 の重み係数学習部 1 1 によって実行される重み係数学習処理のフローチャートである。図 6 において、まず、ステップ S 2 1 で、特徴パラメータメモリ 3 0 から 1 個の音素種類を選択する。次いで、ステップ S 2 2 で、選択された音素種類と同一の音素種類を有する音素の第 1 の音響的特徴パラメータから第 2 の音響的特徴パラメータを取り出して目標音素の第 2 の音響的特徴パラメータとする。そして、ステップ S 2 3 で、同一の音素種類を有する目標音素以外の残りの音素と、第 2 の音響的特徴パラメータにおける目標音素との間の、音響的距離であるユークリッドケプストラム距離と、底を 2 とする対数音素時間長とを計算する。ステップ S 2 4 では、すべての残りの音素についてステップ S 2 2 及び S 2 3 の処理をしたか否かが判断され、処理が完了していないときは、ステップ S 2 5 で別の残りの音素を選択してステップ S 2 3 からの処理を繰り返す。

【0078】一方、ステップ S 2 4 で処理が完了しているときは、ステップ S 2 6 で、ステップ S 2 3 で得られた距離及び時間長に基づいて、上位 N 1 個の最良の音素候補を選択する。次いで、ステップ S 2 7 で選択された上位 N 1 個の最良の音素候補について 1 番目から N 1 番目までランク付けする。そして、ステップ S 2 8 で、ランク付けされた N 1 個の最良の音素候補に対して各距離から中間値を引いてスケール変換値を計算する。そして、ステップ S 2 9 において、すべての音素種類及び音素についてステップ S 2 2 から S 2 8 までの処理を完了したか否かが判断され、完了していないときは、ステップ S 3 0 で別の音素種類又は音素を選択した後、ステップ S 2 2 からの処理を繰り返す。一方、ステップ S 2 9 で処理が完了しているときは、図 7 のステップ S 3 1 に進む。

【0079】図 7 において、ステップ S 3 1 では、1 個の音素種類を選択する。次いで、ステップ S 3 2 では、選択された音素種類に対して各音素の第 2 の音響的特徴パラメータを抽出する。そして、ステップ S 3 3 で、選択された音素種類に対するスケール変換値に基づいて線形回帰分析を行うことにより、各第 2 の音響的特徴パラメータにおけるスケール変換値に対する寄与度を計算し、計算された寄与度を目標音素毎の重み係数として重み係数ベクトルメモリ 3 1 に記憶する。ステップ S 3 4 では、すべての音素種類について上記ステップ S 3 2 及び S 3 3 の処理を完了したか否かが判断され、完了して

いないときは、ステップ S 3 5 で別の音素種類を選択した後、ステップ S 3 2 からの処理を繰り返す。一方、ステップ S 3 4 で処理が完了しているときは、当該重み係数学習処理を終了する。なお、各第 2 の韻律的特徴パラメータにおける寄与度は経験的に（又は実験的に）予め与えられて、当該寄与度を目標音素毎の重み係数ベクトルとして重み係数ベクトルメモリ 3 1 に記憶する。

【0080】図 8 は、図 2 の音声単位選択部 1 2 によって実行される音声単位選択処理のフローチャートである。図 8 において、まず、ステップ S 4 1 で、入力された音素列のうち最初から 1 個目の音素を選択する。次いで、ステップ S 4 2 で、選択された音素と同一の音素種類を有する音素の重み係数ベクトルを重み係数ベクトルメモリ 3 1 から読み出し、目標サブコスト及び必要な特徴パラメータを特徴パラメータメモリ 3 0 から読み出してリストアップする。そして、ステップ S 4 3 ですべての音素について処理したか否かが判断され、完了していないときはステップ S 4 4 で次の音素を選択した後、ステップ S 4 2 の処理を繰り返す。一方、ステップ S 4 3 で完了していないときは、ステップ S 4 5 に進む。

【0081】ステップ S 4 5 では、入力された音素列に対して数 4 を用いて各音素候補における全体のコストを計算する。次いで、ステップ S 4 6 では、計算されたコストに基づいて、上位 N 2 個の最良の音素候補をそれぞれの目標音素に対して選択する。そして、ステップ S 4 7 では、数 5 を用いてビタビサーチにより、全体のコストを最小にする音素候補の組み合わせの索引情報と、その各音素の開始時刻と時間長とともに検索した後、音声合成部 1 3 に出力して、当該音声単位選択処理を終了する。

【0082】さらに、音声合成部 1 3 は、音声単位選択部 1 2 から出力される索引情報と、その各音素の開始時刻と時間長とに基づいて、音声波形信号データベースメモリ 2 1 a に対してアクセスして単位選択された音素候補のデジタル音声波形信号データを読み出して、逐次 D/A 変換して変換後のアナログ音声信号をスピーカ 1 4 を介して出力する。これにより、入力された音素列に対応する音声合成された音声はスピーカ 1 4 から出力される。

【0083】以上説明したように、本実施形態の音声合成装置においては、出力音声の自然性を最大にするために、大規模な自然音声のデータベースを用いて処理を最小に抑える方法について述べた。本実施形態は 4 つの処理部 1 0 乃至 1 3 から構成される。

<音声分析部 1 0> 正書法の書き起こしテキストを伴った任意の音声波形信号データを入力とし、この音声波形信号データベース中のすべての音素について、それらの性質を記述する特徴ベクトルを与える処理部。

<重み係数学習部 1 1> 音声波形信号データベースの特徴ベクトルと音声波形信号データベースの原波形を用い

て、目的の音声を合成する場合に最も適するように音声単位を選ぶための、各特徴パラメータの最適重み係数を重みベクトルとして決定する処理部。

<音声単位選択部 1 2> 音声波形信号データベースの全音素の特徴ベクトルと重みベクトルと目的音声の発話内容の記述から音声波形信号データベースメモリ 2 1 a の索引情報を作成する処理部。

<音声合成部 1 3> 作成された索引情報に従って、音声波形データベースメモリ 2 1 a 内の音声波形信号データベース中の音声波形信号データの音声セグメントに飛び飛びにアクセスし、目的の音声波形信号の音声セグメントを連結しかつ D / A 変換してスピーカ 1 4 に出力して音声を合成する処理部。

【0084】本実施形態においては、音声波形信号の圧縮や音声基本周波数 F₀ や音素時間長の修正は不要になったが、代わって音声サンプルを注意深くラベル付けし、大規模な音声波形信号データベースの中から最適なものを選択することが必要となる。本実施形態の音声合成方法の基本単位は音素であり、これは辞書やテキスト-音素変換プログラムで生成されるが、同一の音素であっても音声波形信号データベース中に音素の十分なバリエーションを含んでいることが要求される。音声波形信号データベースからの音声単位選択処理では目的の韻律的環境に適合し、しかも接続したときに隣接音声単位間での不連続性が最も低い音素サンプルの組合せが選ばれる。このために、音素毎に各特徴パラメータの最適重み係数が決定される。

【0085】本実施形態の音声合成装置の特徴は、次の通りである。

<単位選択規準としての韻律的情報の利用> スペクトルの特徴は韻律的特徴と不可分であるとの立場から、音声単位の選択規準に韻律的な特徴を導入した。

<音響的及び韻律的特徴パラメータの重み係数の自動学習> 音素環境や音響的特徴、韻律的特徴等の各種の特徴量が音声単位の選択にどれだけの寄与があるかを音声波形信号データベース中の全音声サンプルを利用することで自動的に決定し、テキストデータベース及び音声波形*

バイフォンリストの作成例

発話内容：あらゆる現実を全て自分の方へねじ曲げたのだ。

音素列：

```
# a r a y u r u g e N j i t s u o
# s u b e t e
# j i b u N n o h o o e
# n e j i m a g e t a n o d a #
```

バイフォンリスト：

```
# a a_r r_a a_y y_u u_r r_u u_g g_e e_N N_j j_i i_t s t_s_u u_o o_#
#_s_s_u u_b b_e e_t t_e e_#
```

* データベースを基本とする音声合成装置を構築した。

<音声波形信号の直接接続> 上記の自動学習により、大規模音声波形信号データベースから最適な音声サンプルを選び出すことにより、何らの信号処理も利用しない任意音声合成装置を構築した。

<音声波形信号データベースの外部情報化> 音声波形信号データベースを完全に外部情報として取り扱うことにより、単に CD-ROM 等に記憶した音声波形信号データを取り替えることで任意の言語、任意の話者に利用できる音声合成装置を構築した。

【0086】

【実施例】本発明者は、本実施形態の図 1 の音声データ量削減装置を用いて、音声波形データベースの音声データを削減し、かつ音声データ量が削減された音声波形データベースに基づいて、図 2 の音声合成装置を用いて音声合成した実験結果について以下に説明する。

【0087】まず、本発明者らによる音素バランス調査について説明する。本特許出願人が所有するテキストデータベース及び音声波形データベース（以下、CHAT R という。）で使用される日本語話者のラベルデータに基づいて音素バランス調査を行った。この調査では、音素ラベルデータから抽出した音素列から、2 音素ずつ組み合わせたバイフォンデータを作成した。例えば、ラベルデータの発話内容が「あらゆる」であった場合、抽出した音素列は、

【数 15】# a r a y u r u #

となる。ここで、# は発話開始記号又は発話終了記号を示す。この音素列から 2 音素ずつ組み合わせたバイフォンリストは、

【数 16】#_a a_r r_a a_y y_u u_r r_u u_#

の 8 種類となる。

【0088】次いで、バイフォンリストの作成例について説明する。この作成例の結果を次の表に示す。

【0089】

【表 5】

#_j j_i i_b b_u u_N N_n n_o o_h h_o o_o o_e e_#
 #_n n_e e_j j_i i_m m_a a_g g_e e_t t_a a_n n_o o_d d_a a_#

【0090】上記の場合、バイフォンリストは、50件 * 示す。
 44種類となり、6件が重複したバイフォンであること 【0091】
 が見える。ここで、件数別バイフォンリストを次の表に* 【表6】
 件数別バイフォンリスト

1件:#_a a_r r_a a_y y_u u_r r_u u_g e_N N_j i_ts ts_u u_o o_#
 #_s s_u u_b b_e e_t_e
 #_j i_b b_u u_N N_n n_o o_h h_o o_o o_e
 #_n n_e e_j j_i m_m_a a_g t_a a_n o_d d_a a_#

2件:g_e, e_t, e_#, n_o

3件:j_i

【0092】次いで、CHATR内のある日本語話者M	503文章	503ファイル
YAに対してのバイフォンリストを作成し、以下のよう	22文章	22ファイル
に出現頻度を集計することにより音素バランスを調べ	20 旅行会話	675ファイル
た。ここで、話者MYAの音声波形データベースの情報	合計	1200ファイル
を次の表に示す。		

【0093】

【表7】話者MYAの音声波形データベースの情報

【0094】日本語話者MYAに対する音素バランス調査の結果を以下に示す。

発話内容

ファイル数

【0095】

【表8】

バイフォンリスト集計結果

バイフォンの件数：60, 405件(ラベリングミス除く)

バイフォンの種類：415通り(ラベリングミス除く)

【0096】

【表9】出現頻度の多いバイフォンリスト10種類

バイフォン

件数

a__i	970
n__o	955
k__a	886
t__o	786
g__a	772
i__m	744

o__o 1436

a__# 1326

d__e 992

40 【0097】

m__a 988

【表10】

出現頻度1件のバイフォン：42通り

l_g l_n U_b U_r d_y f_o t_y l_kk a_ff a_gg ch_e dd_a e_dd e_gg e_hh
 e_ss ff_u
 gg_a gg_u hh_a hh_i hh_o i_hh i_jj jj_i o_hh o_zz pp_U ss_o ss_u tt_U
 u_dd
 zz_u U_tts cch_l cch_e cch_u e_cch o_tts ssh_l ssh_U tts_U

【0098】日本語話者MYAに対する音素バランス調査におけるバイフォンリストの詳細結果を以下に示す。

【0099】

* * 【表11】

出現頻度500件以上のバイフォンリスト
(バイフォン/件数)

o_o	1436	a_#	1326	d_e	992	m_a	988	a_i	970	n_o	955
k_a	886	t_o	786	g_a	772	i_m	744	a_r	738	t_a	716
w_a	696	n_i	678	a_s	651	o_#	604	o_k	603	n_a	594
s_U	586	r_u	582	k_u	581	o_n	573	y_o	563	a_k	551
r_i	544	i_#	536	o_r	517	k_o	507	sh_i	506	r_a	506
u_u	504	e_#	503	a_N	503						

【0100】

【表12】

出現頻度100件以上500件未満のバイフォンリスト
(バイフォン/件数)

#_k	490	t_e	482	m_o	465	u_n	452	r_e	444	i_n	440
sh_l	424	e_N	415	s_a	412	o_m	405	a_n	397	s_u	394
l_t	394	i_r	365	a_sh	364	#_h	363	ts_u	358	h_a	356
#_s	356	g_o	355	k_i	349	o_d	348	e_s	348	o_t	340
e_r	340	i_N	332	y_u	326	d_a	324	U_#	317	u_#	313
d_o	312	u_r	311	y_a	310	o_sh	309	ch_i	299	s_o	299
s_e	288	i_t	286	a_t	281	a_m	280	U_k	279	j_i	276
h_o	274	o_g	271	N_d	267	k_e	262	n_e	261	m_i	259
o_N	254	e_k	254	e_e	248	i_k	243	sh_o	242	r_o	239
o_y	238	N_n	238	#_n	237	i_g	234	#_o	233	u_g	232
u_k	230	i_d	227	u_d	222	e_t	218	a_d	218	i_i	217
o_h	216	o_s	214	a_g	210	o_i	205	e_n	205	m_e	200
#_#	199	l_k	198	#_m	192	tt_e	188	z_a	188	u_m	188
k_y	188	#_g	187	a_tt	181	e_d	180	e_g	178	#_i	173
b_a	170	u_N	169	e_sh	167	e_i	164	a_a	162	#_d	160
#_j	159	#_a	159	#_t	157	z_u	156	e_m	155	N_k	153
o_ch	152	j_u	152	#_y	152	i_sh	149	o_z	145	b_u	140
i_o	139	#_sh	137	i_h	137	u_i	132	N_g	131	o_w	129
i_s	129	tt_a	127	k_U	126	e_w	125	u_o	124	i_w	124
e_o	123	a_ts	118	j_o	115	a_w	114	g_e	112	o_b	111
o_j	110	a_o	110	u_t	109	U_t	109	k_l	108	a_e	106
N_t	106	u_s	102								

【0101】

【表13】

出現頻度50件以上100件未満のバイフォンリスト
(バイフォン/件数)

u_y	99	m_u	99	i_y	99	u_sh	98	h_i	97	h_l	97
a_ch	95	i_ch	93	ts_U	92	ch_l	92	o_a	92	N_j	92
f_u	91	u_b	90	g_u	90	g_i	90	U_s	90	#_ch	88
i_b	83	sh_a	82	f_U	81	b_e	81	a_h	81	o_tt	80
N_w	80	i_ts	78	u_j	77	N_sh	76	a_y	76	b_i	75
#_w	74	a_z	73	a_b	73	o_e	71	i_a	71	b_o	71
sh_u	70	i_tt	69	e_b	68	N_b	68	N_#	68	tt_o	67
i_z	67	#_b	67	a_j	66	i_j	65	e_a	65	#_r	65

31
 #_f 65 u_w 63 u_h 63 o_ts 62 u_z 62 h_y 61
 N_m 61 r_y 60 e_ts 57 ch_o 57 h_e 56 N_s 55
 e_y 51

32

【0102】

* * 【表14】

出現頻度10件以上50件未満のバイフォンリスト
 (バイフォン/件数)

z_e 47 u_e 47 o_f 47 N_r 47 N_o 47 l_ts 46
 ch_u 45 i_e 45 i_pp 42 z_o 42 N_h 42 #_z 41
 n_y 40 #_u 40 a_u 39 #_ts 38 o_u 34 n_u 34
 j_a 34 g_y 34 pp_a 33 u_ts 32 p_a 32 l_s 32
 #_e 32 kk_a 30 e_h 29 N_y 29 u_ch 28 e_tt 26
 U_sh 26 i_u 26 N_p 26 i_kk 25 b_y 25 e_z 23
 e_j 23 a_f 23 a_kk 22 U_n 21 u_tt 20 i_f 20
 #_p 20 kk_u 19 o_kk 18 e_kk 18 u_f 18 i_ssh 17
 pp_o 17 N_ch 17 l_ch 17 o_p 17 m_y 17 u_a 16
 p_u 16 N_z 16 N_i 16 kk_o 15 ch_a 15 p_e 15
 j_e 15 N_e 15 l_# 15 o_pp 14 t_i 14 U_ts 13
 p_o 13 ssh_i 12 pp_u 12 kk_i 12 w_o 12 p_i 12
 kk_y 11 a_pp 11 U_d 11

【0103】

【表15】

出現頻度1件以上10件未満のバイフォンリスト
 (バイフォン/件数)

u_kk 9 ss_e 9 i_p 9 e_u 9 N_f 9 ssh_a 8
 a_ssh 8 u_pp 8 pp_i 8 u_p 8 r_U 8 f_i 8
 pp_y 7 e_ch 7 U_tt 7 p_y 7 e_f 7 e_ssh 6
 cch_i 6 kk_e 6 i_ss 6 N_ts 6 e_p 6 d_i 6
 a_p 6 U_z 6 U_p 6 U_g 6 tts_u 5 ssh_o 5
 o_cch 5 l_tt 5 l_sh 5 f_e 5 u_ssh 4 i_tts 4
 i_cch 4 cch_a 4 u_ss 4 tt_i 4 p_U 4 U_h 4
 N_u 4 N_a 4 a_cch 3 ss_a 3 pp_e 3 kk_U 3
 dd_o 3 a_ss 3 U_kk 3 l_pp 3 j_l 3 f_a 3
 l_d 3 u_cch 2 o_ssh 2 cch_o 2 sh_e 2 i_dd 2
 hh_e 2 ff_e 2 e_pp 2 e_ff 2 a_hh 2 w_e 2
 p_l 2 d_u 2 U_j 2 l_p 2 l_m 2 l_h 2
 l_b 2 tts_U 1 ssh_U 1 ssh_l 1 o_tts 1 e_cch 1
 cch_u 1 cch_e 1 cch_l 1 U_tts 1 zz_u 1 u_dd 1
 tt_U 1 ss_u 1 ss_o 1 pp_U 1 o_zz 1 o_hh 1
 jj_i 1 i_jj 1 i_hh 1 hh_o 1 hh_i 1 hh_a 1
 gg_u 1 gg_a 1 ff_u 1 e_ss 1 e_hh 1 e_gg 1
 e_dd 1 dd_a 1 ch_e 1 a_gg 1 a_ff 1 l_kk 1
 t_y 1 f_o 1 d_y 1 U_r 1 U_b 1 l_n 1
 l_g 1

【0104】以上示したように、ある日本語話者に対して、バイフォンの分布にかなりバラツキがあることがわかる。これに基づいて、本発明者らは、音声合成後の品

質を低下させずに、図3の処理を用いて音声データ量の削減を行った。

50 【0105】音声波形データベースの縮小の実験におい

ては、音素的にバランスのとれた文章セット（文数 n = 5 2 5）及び旅行用対話テキストのセット（文数 n = 6 7 5）を読み上げる男性話者 M Y A（職業アナウンサーではない）による音声波形データを使用した。そして、本発明者らは、上述の冗長性尺度、すなわち類似度の判断基準に従って音声波形データベースを元のサイズの3分の2及び半分にまで、すなわち1 2 0 0発話から各々8 0 6及び6 2 5発話にまで音声データ量を削減した。ここで、それぞれのデータベースを第1減量 M Y Aと第*

音声波形データベースの縮小

* 2減量 M Y Aと呼ぶ。音声データ量を削減した後に、得られた音声波形データベースのスコアを、韻律的特徴パラメータ及び音響特徴パラメータにかかる距離に関して測定した。表 1 6 は、3つの音声波形データベースの各々における音声ファイルの計数値を示している。表 1 7 は韻律的特徴パラメータのスコアを示し、表 1 8 は音響特徴パラメータのスコアを示している。

【 0 1 0 6 】
【 表 1 6 】

データベース	音声セグメント数	音声波形 D B 内の文の数	テストされた文の数
M Y A	1 2 0 0	5 2 5 / 6 7 5	1 1 9 8
第 1 減量 M Y A	8 0 6	4 7 4 / 3 3 2	8 0 4
第 2 減量 M Y A	6 2 5	3 3 8 / 2 3 7	6 2 3

（注）D B はデータベースを示す。以下、同様である。

【 0 1 0 7 】

【 表 1 7 】

韻律的特徴パラメータに基づくスコア

データベース	全体のスコア	音声波形 D B 内の文の数	スコアの平均値
M Y A	9 5 6 0	1 1 9 8	* 7 . 9 8 0 5 1 8
第 1 減量 M Y A	6 9 8 3	8 0 4	8 . 6 8 5 6 2 7
第 2 減量 M Y A	5 6 0 5	6 2 3	8 . 9 9 7 5 5 7

（注）* は最低のスコアを示す。

【 0 1 0 8 】

【 表 1 8 】

パースペクトルに基づくスコア

データベース	全体のスコア	音声波形 D B 内の文の数	スコアの平均値
M Y A	1 5 0 2 0	1 1 9 8	1 2 . 5 3 7 7 7 0
第 1 減量 M Y A	9 7 5 7	8 0 4	* 1 2 . 1 3 6 1 1 9
第 2 減量 M Y A	7 5 7 1	6 2 3	1 2 . 1 5 2 6 8 7

（注）* は最低のスコアを示す。

【 0 1 0 9 】 表 1 7 及び表 1 8 から明らかなように、音声波形データベースの音声データ量を削減した結果、韻律的特徴パラメータの変動性ではスコアで約 0 . 7 ポイント低下しているものの、音響特徴パラメータにおける平滑性の面ではスコアで 0 . 4 ポイント向上していることが分かる。

【 0 1 1 0 】 以上説明したように、本実施形態によれば、スマートフォンリストに基づいて各 1 対のスマートフォンに対する韻律的特徴パラメータと音響的特徴パラメータとに関する所定の類似度を計算し、上記計算された類似度が所定の第 1 のしきい値以上であり、かつ上記スマートフォンのリスト中の同一のスマートフォンの数が所定の第 2 のし

きい値以上であるときに、当該 1 対のスマートフォンのうちの一方のスマートフォンの音素に係る音声波形信号の音声セグメントのデータを上記音声波形データベースから削除することにより音声データ量を削減する。従って、音声合成時の音声品質を実質的に低下させることなく、音声波形データベースを格納するメモリ容量を削減することができ、音声合成時の探索速度を高めることができる。

【 0 1 1 1 】

【 発明の効果 】 以上詳述したように本発明に係る音声合成装置のための音声データ量削減装置によれば、1 対の音素のリストに基づいて各 1 対の音素に対する韻律的特徴パラメータと音響的特徴パラメータとに関する所定の

類似度を計算し、上記計算された類似度が所定の第1のしきい値以上であるときに、当該各1対の音素のうちの一方の1対の音素に係る音声波形信号の音声セグメントのデータを上記音声波形データベースから削除することにより音声データ量を削減する。従って、音声合成時の音声品質を実質的に低下させることなく、音声波形データベースを格納するメモリ容量を削減することができ、音声合成時の探索速度を高めることができる。

【0112】また、上記計算された類似度が所定の第1のしきい値以上でありかつ上記パイフォンのリスト中の同一のパイフォンの数が所定の第2のしきい値以上であるときに、当該各1対の音素のうちの一方の1対の音素に係る音声波形信号の音声セグメントのデータを上記音声波形データベースから削除することにより音声データ量を削減する。従って、音声波形データベースを格納するメモリ容量を削減することができ、音声合成時の探索速度を高めることができるとともに、音声波形データベースにおいて所定数の同一のパイフォンに対する音声波形データを確保して、音声合成後の音声の品質を所定以上に確保することができる。

【図面の簡単な説明】

【図1】 本発明に係る一実施形態である音声データ量削減処理装置のブロック図である。

【図2】 本発明に係る一実施形態である自然発話音声波形信号接続型音声合成装置のブロック図である。

【図3】 図1の音声データ量削減処理部によって実行される音声データ量削減処理を示すフローチャートである。

【図4】 図2の音声単位選択部によって計算される音声単位選択コストの定義を示すモデル図である。

*【図5】 図2の音声分析部によって実行される音声分析処理のフローチャートである。

【図6】 図2の重み係数学習部によって実行される重み係数学習処理の第1の部分のフローチャートである。

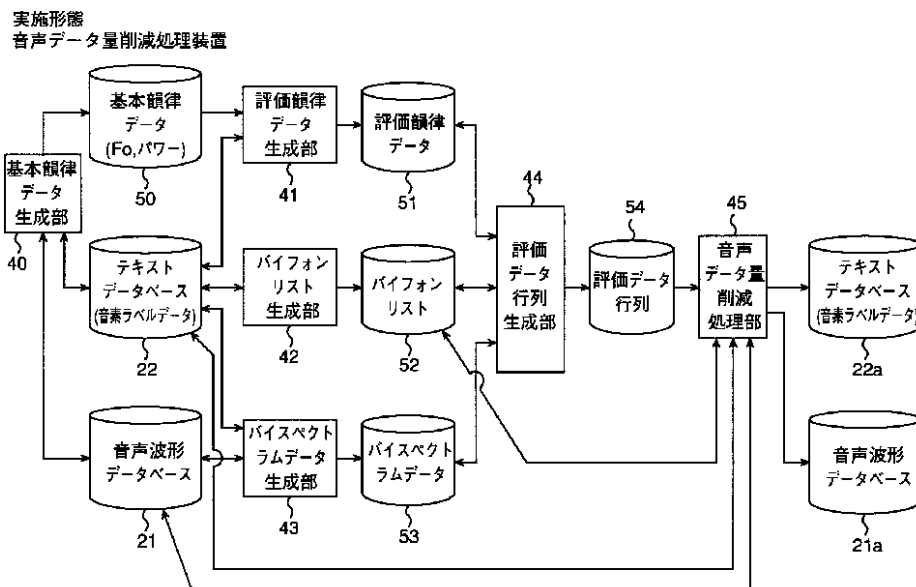
【図7】 図2の重み係数学習部によって実行される重み係数学習処理の第2の部分のフローチャートである。

【図8】 図2の音声単位選択部によって実行される音声単位選択処理のフローチャートである。

【符号の説明】

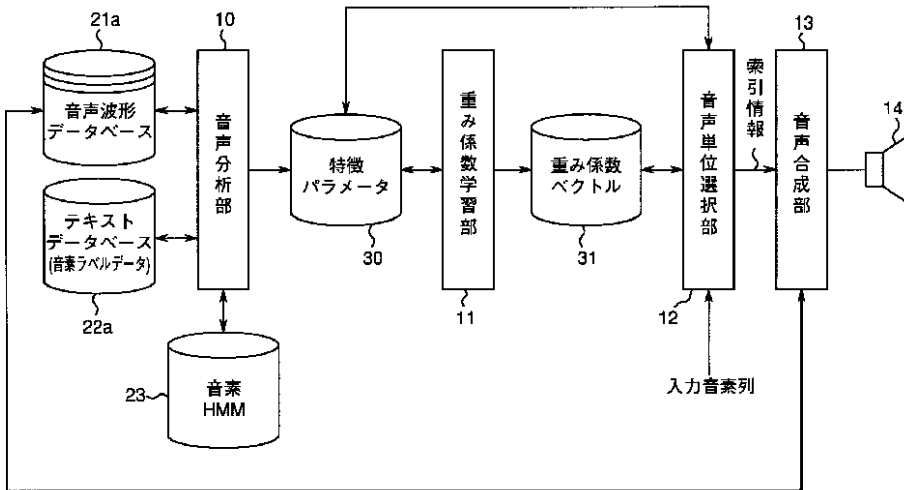
- 10 ... 音声分析部、
- 11 ... 重み係数学習部、
- 12 ... 音声単位選択部、
- 13 ... 音声合成部、
- 14 ... スピーカ、
- 21, 21a ... 音声波形信号データベースメモリ、
- 22, 22a ... テキストデータベースメモリ、
- 23 ... 音素HMMメモリ、
- 30 ... 特徴パラメータメモリ、
- 31 ... 重み係数ベクトル、
- 40 ... 基本韻律データ生成部、
- 41 ... 評価韻律データ生成部、
- 42 ... パイフォンリスト生成部、
- 43 ... バイスペクトラムデータ生成部、
- 44 ... 評価データ行列生成部、
- 45 ... 音声データ削減処理部、
- 50 ... 基本韻律データメモリ、
- 51 ... 評価韻律データメモリ、
- 52 ... パイフォンリストメモリ、
- 53 ... バイスペクトラムデータメモリ、
- 54 ... 評価データ行列メモリ、
- 55 ... 音声データ削減処理部メモリ。

【図1】

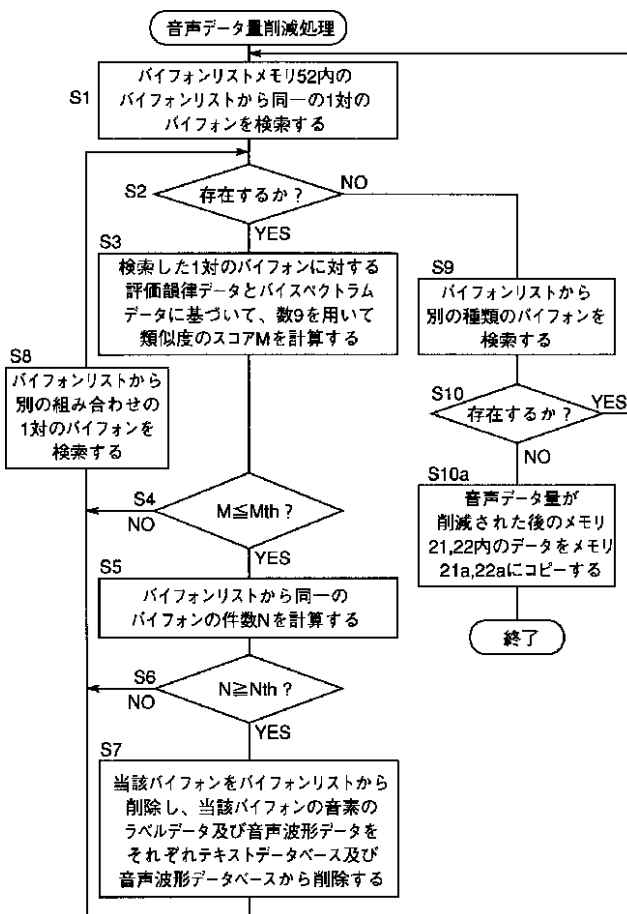


【図 2】

実施形態
自然発話音声波形信号接続型音声合成装置

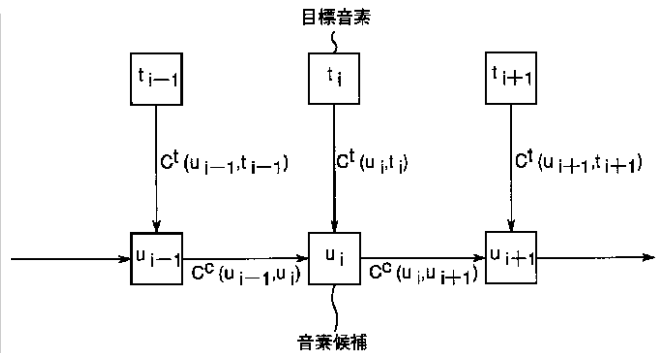


【図 3】

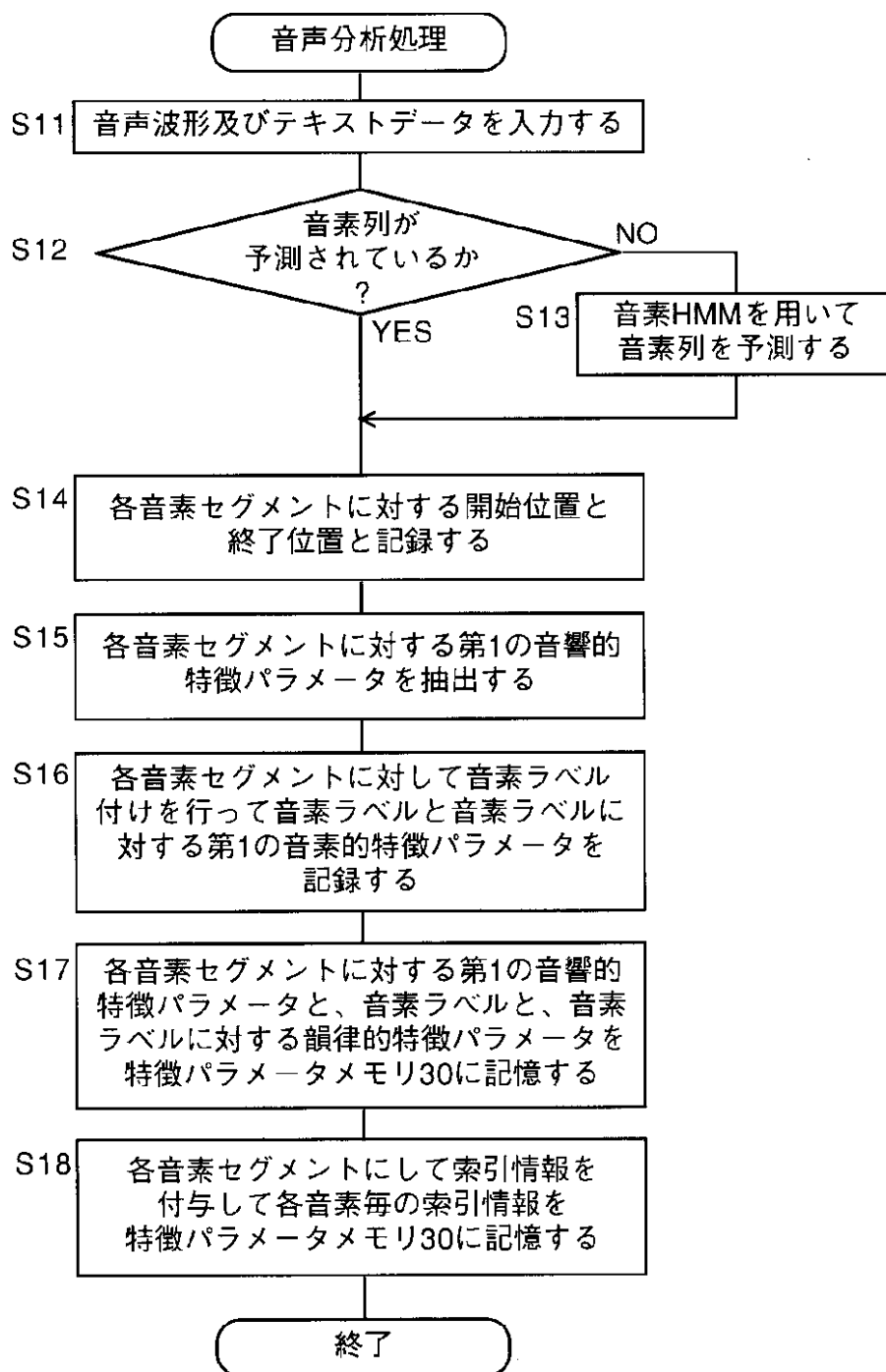


【図 4】

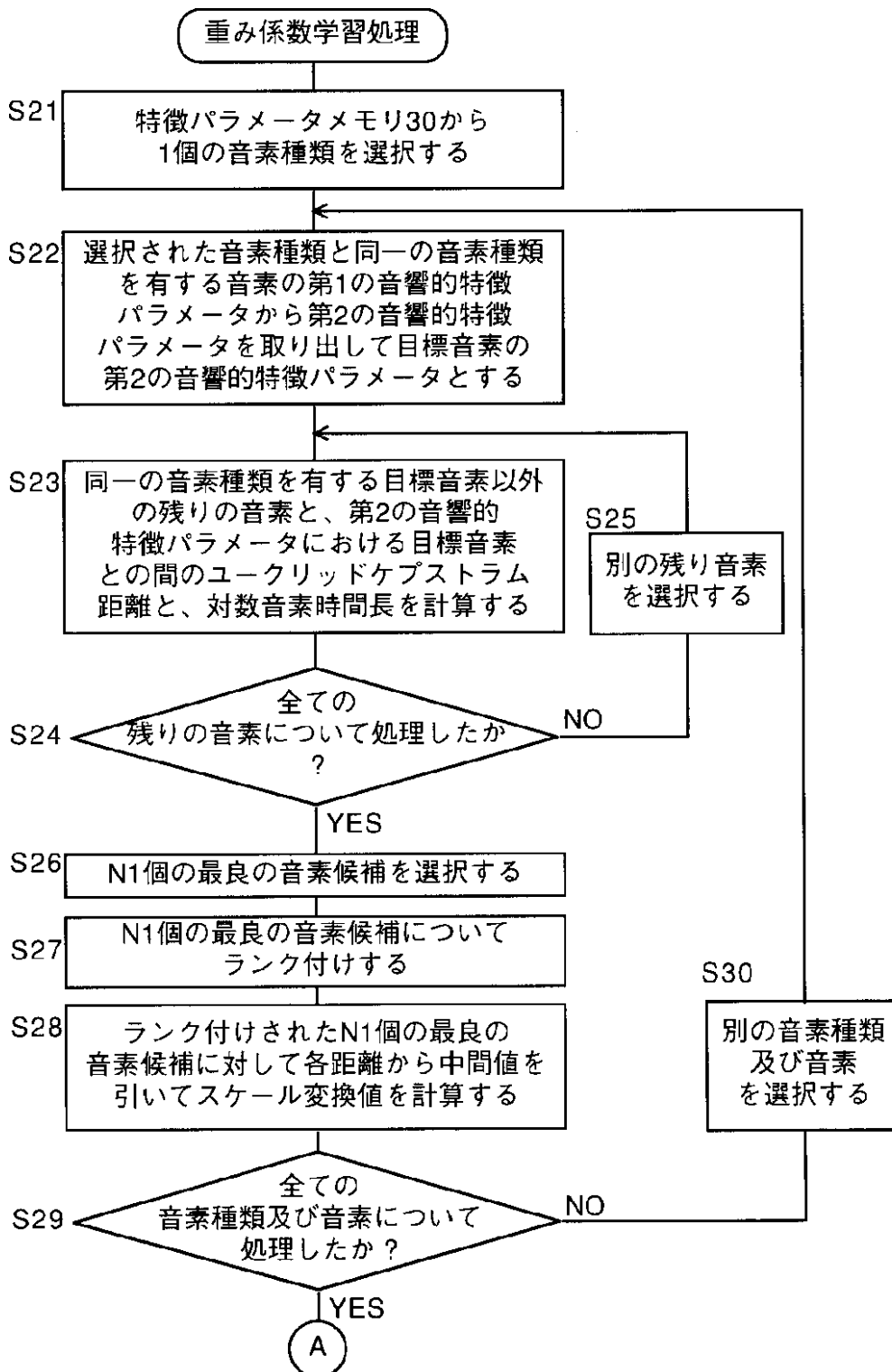
音声単位選択コストの定義



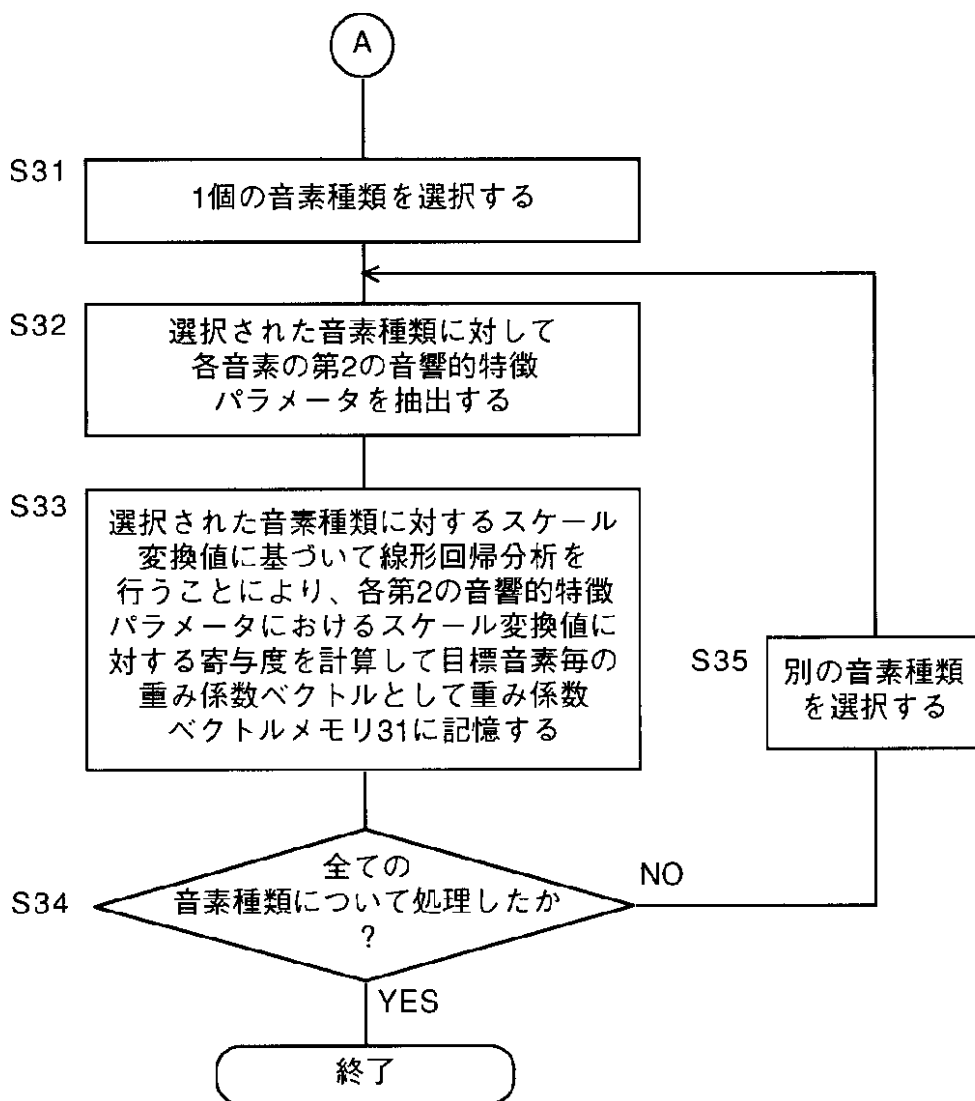
【図5】



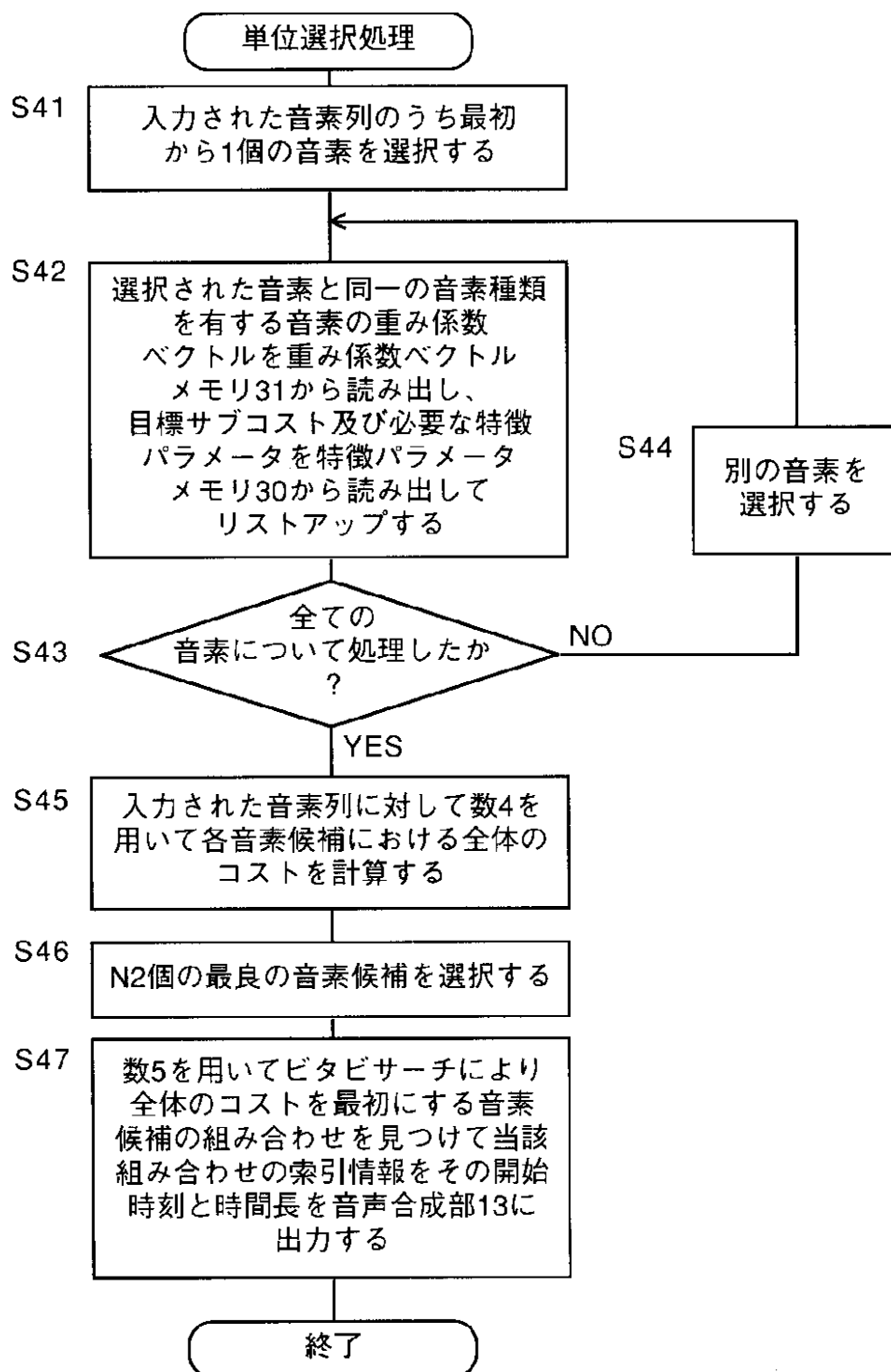
【図6】



【図7】



【図8】



フロントページの続き

(56)参考文献 特開 平11 - 95796 (J P , A)
特開 平10 - 39889 (J P , A)
特開 平 8 - 248975 (J P , A)
特開 昭56 - 51800 (J P , A)
特開 平 7 - 84590 (J P , A)
特開 平10 - 49193 (J P , A)

(58)調査した分野(Int.Cl.⁷, D B 名)
G10L 13/06