

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4720974号
(P4720974)

(45) 発行日 平成23年7月13日(2011.7.13)

(24) 登録日 平成23年4月15日(2011.4.15)

(51) Int.Cl. F I
G 1 0 L 21/04 (2006.01) G 1 0 L 21/04 1 2 0 D
G 1 0 L 13/02 (2006.01) G 1 0 L 13/02 1 2 2 Z

請求項の数 4 (全 23 頁)

<p>(21) 出願番号 特願2004-369108 (P2004-369108)</p> <p>(22) 出願日 平成16年12月21日(2004.12.21)</p> <p>(65) 公開番号 特開2006-178052 (P2006-178052A)</p> <p>(43) 公開日 平成18年7月6日(2006.7.6)</p> <p>審査請求日 平成19年11月16日(2007.11.16)</p> <p>(出願人による申告)平成16年度独立行政法人情報通信研究機構、研究テーマ「超高速知能ネットワーク社会に向けた新しいインタラクション・メディアの研究開発」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</p> <p>(74) 代理人 100099933 弁理士 清水 敏</p> <p>(72) 発明者 米澤 朋子 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>(72) 発明者 鈴木 紀子 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>(72) 発明者 小暮 潔 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 style="text-align: right;">最終頁に続く</p>
--	--

(54) 【発明の名称】 音声発生装置およびそのためのコンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

同じ内容の複数種類の音声に基づいて、声質を変化させながら音声を発生させるための音声発生装置であって、

複数の中間段階の数、及び、シグモイド関数決定のためのパラメータを受付ける受付手段と、

前記受付手段が受付けた前記パラメータにしたがってシグモイド関数を決定するシグモイド関数決定手段と、

前記シグモイド関数決定手段が決定した前記シグモイド関数を用いて、前記受付手段が受付けた数の複数の中間段階におけるモーフィング率をそれぞれ決定する第1のモーフィング率決定手段と、

前記複数種類の音声に対して、前記第1のモーフィング率決定手段が決定した前記複数種類のモーフィング率で音声モーフィングを行なうことによって複数の中間音声を作成するモーフィング実行手段と、

前記複数種類の音声と、前記モーフィング実行手段が作成した前記複数の中間音声とを記憶する音声記憶装置と、

前記複数種類の音声のモーフィング率を指定するためのモーフィング率指定手段と、前記モーフィング率指定手段によりモーフィング率が指定されたことに応答して、前記音声記憶装置から与えられる、前記複数種類の音声と前記複数の中間音声との中で、指定されたモーフィング率に最も近いモーフィング率の音声を選択する音声選択手段と、

前記音声選択手段から与えられる前記選択された音声を音声信号に変換する音声信号変換手段とを含み、

前記複数種類の音声には、予めそれぞれの声質を特定するラベルが付されており、前記モーフィング率指定手段は、

前記複数種類の音声の各々について、対応するラベルを表示するラベル表示手段と、前記複数種類の音声をそれぞれ表す複数の基準点をそれぞれ対応するラベルと関連付けて所定の位置に表示する基準点表示手段と、

前記基準点と所定の関係にある、予め定める領域内にユーザにより配置されたポイントの位置と、前記複数の基準点との間の距離にしたがって、前記複数種類の音声のモーフィング率を決定する第2のモーフィング率決定手段と、を含む、音声発生装置。

10

【請求項2】

前記モーフィング率指定手段はさらに、

所定空間内における、前記複数種類の音声の各々にそれぞれ対応する前記複数の基準点の位置を決定する基準点決定手段を含む、請求項1に記載の音声発生装置。

【請求項3】

前記音声選択手段による音声発生の基準時刻を定めるタイマと、

前記モーフィング率指定手段によりモーフィング率の指定がされたことに応答して、前記タイマを参照して音声再生の時刻を得て、前記音声選択手段により選択された音声を示す情報とともに音声再生シーケンスとして記録する選択音声記録手段と、

音声再生シーケンスの再生を指示する信号に応答して、当該信号により示される音声再生シーケンスを読み出して、当該音声再生シーケンスにより指定される時刻に、当該時刻に指定された音声を選択して前記音声信号変換手段に与える音声再生制御処理手段とをさらに含む、請求項1又は請求項2に記載の音声発生装置。

20

【請求項4】

コンピュータにより実行されると、請求項1～請求項3のいずれかに記載の音声発生装置として当該コンピュータを動作させる、コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は音声処理技術に関し、特に、音声に任意の表情付けを容易に行なうことが可能な音声処理技術に関する。

30

【背景技術】

【0002】

様々な音声表現においてユーザの所望する表情付けを可能にし、豊かな表情のついた音声を実現するための技術は、様々な用途に利用可能であると思われる。例えば、豊かな表情の歌声の合成などである。そのためには様々な表情付きの音声だけでなく、表情付けの程度が中間的な音声も必要だと考えられる。さらに、それらは自然な音声に近いことが望ましい。

【0003】

よって、多種多様な音声データを集めることが望まれる。しかし、そのためには所望の表情付けがされた音声を収録する作業が必要であるが、そのような作業は極めて困難である。その結果、音声に対しユーザの所望するような種々の表情付けを容易に行なうことができる従来技術は存在していない。

40

【非特許文献1】カワハラ、H. およびマツイ、H.、「無干渉時間 - 周波数表現における弾性的知覚的距離尺度に基づく聴覚的モーフィング」、ICASSP'2003 予稿集、第1巻、pp. 256 - 259、2003年 (Kawahara, H., and Matsui, H., "AUDITORY MORPHING BASED ON AN ELASTIC PERCEPTUAL DISTANCE METRIC IN AN INTERFERENCE-FREE TIME-FREQUENCY REPRESENTATION," Proc. ICASSP'2003, vol.1, pp.256-259, 2003.)

【非特許文献2】<http://www.wakayama-u.ac.jp/~kawahara/STRAIGHTadv/> (高品質音声分

50

析変換合成システム (S T R A I G H T)

【発明の開示】

【発明が解決しようとする課題】

【0004】

以上のように、音声に対しユーザの所望するような種々の表情付けを容易に行なうことができる技術は存在しておらず、そうした技術に対する需要が高まっている。

【0005】

したがって本発明の主たる目的は、音声の表情付けを自由に变化させながら音声の発生を行なうことができる音声発生装置を提供することである。

【0006】

本発明の他の目的は、音声の表情付けを自由に、かつ簡単な操作で变化させながら音声の発生を行なうことができる音声発生装置を提供することである。

【0007】

本発明のさらに他の目的は、多数の音声を収録する必要なく、音声の表情付けを自由に变化させながら音声の発生を行なうことができる音声発生装置を提供することである。

【課題を解決するための手段】

【0008】

本発明の第1の局面に係る音声発生装置は、同じ内容の複数種類の音声に基づいて、声質を变化させながら音声を発生させるための音声発生装置である。この音声発生装置は、複数種類の音声のモーフィング率を指定するためのモーフィング率指定手段と、モーフィング率指定手段によりモーフィング率が指定されたことに応答して、複数種類の音声と、複数種類の音声に対して複数種類のモーフィング率で音声モーフィングを行なって得た中間音声との中で、指定されたモーフィング率に最も近いモーフィング率の音声を選択するための音声選択手段と、音声選択手段により選択された音声を音声信号に変換するための音声信号変換手段とを含む。

【0009】

複数種類の音声と、それらの中間音声とを予め準備しておき、モーフィング率指定手段によりモーフィング率が指定されると、それに応答してこれら音声の中で指定されたモーフィング率に最も近いモーフィング率の音声を選択して音声信号に変換する。多くの種類の音声を準備しなくても、中間音声を用いることにより種々の表情付けがされた音声の発生をリアルタイムで行なうことができる。

【0010】

好ましくは、複数種類の音声には、予めそれぞれの声質を特定するラベルが付されている。そしてモーフィング率指定手段は、複数種類の音声の各々について、対応するラベルを表示するためのラベル表示手段と、複数種類の音声をそれぞれ表す複数の基準点をそれぞれ所定の位置に対応するラベルと関連付けて表示するための基準点表示手段と、基準点と所定の関係にある予め定める領域内にユーザにより配置されたポイントの位置と、複数の基準点との間の距離にしたがって、複数種類の音声のモーフィング率を決定するためのモーフィング率決定手段とを含む。

【0011】

複数種類の音声に付されたラベルが表示され、さらにそれら複数種類の音声に対応する基準点がラベルと関連付けて表示される。ユーザがそれら基準点と所定の関係にある領域、例えば基準点が二つの場合には基準点を結ぶ直線上、基準点が3つの場合にはそれら3点で囲まれる三角形内にポイントを配置すると、その位置と基準点との間の距離にしたがって、モーフィング率決定手段が音声のモーフィング率を決定する。視覚的で直感的に分りやすい、簡単な操作で音声のモーフィング率を指定することができる。

【0012】

モーフィング率指定手段は、所定空間内において、複数種類の音声の各々にそれぞれ対応する複数の基準点を決定するための基準点決定手段と、複数の基準点と所定の関係にある、予め定める領域内にユーザにより配置されたポイントの位置と、複数の基準点との間

10

20

30

40

50

の距離にしたがって、前記複数種類の音声のモーフィング率を決定するためのモーフィング率決定手段とを含んでもよい。

【0013】

複数種類の音声に対応する基準点が空間内において決定される。ユーザがそれら基準点と所定の関係にある領域、例えば基準点が4つの場合にはそれら4点で囲まれる三角錐内にポイントを配置すると、その位置と基準点との間の距離にしたがって、モーフィング率決定手段が音声のモーフィング率を決定する。ポイントの配置によってモーフィング率を指定できるため、直感的で分りやすく、簡単な操作で音声のモーフィング率を指定することができる。

【0014】

さらに好ましくは、音声発生装置は、音声選択手段による音声発生 of 基準時刻を定めるタイマと、モーフィング率指定手段によりモーフィング率の指定がされたことに応答して、タイマを参照して音声再生の時刻を得て、音声選択手段により選択された音声を示す情報とともに音声再生シーケンスとして記録するための選択音声記録手段と、音声再生シーケンスの再生を指示する信号に応答して、当該信号により示される音声再生シーケンスを読み出して、当該音声再生シーケンスにより指定される時刻に、当該時刻に指定された音声を選択して音声信号変換手段に与えるための音声再生制御処理手段とをさらに含む。

【0015】

音声を選択する操作をしながら音声を発声させると、その操作により選択された音声と、その音声を選択されたときの時刻とが音声再生シーケンスとして記録される。それを音声再生制御手段によって再生することにより、過去に行なった音声再生と同じ順序で、同じ音声を聞いた音声再生を再現することができる。

【0016】

音声発生装置は、複数種類の音声と中間音声とを記憶し、音声選択手段に与えるための音声記憶装置をさらに含んでもよい。

【0017】

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と、上記したいずれかの音声発生装置として当該コンピュータを動作させるものである。したがって、対応する音声発生装置と同様の作用効果を奏することができる。

【発明を実施するための最良の形態】

【0018】

異なる二つの音声を所望のモーフィング率で混合することにより中間的な声質を有する音声（以下「中間音声」と呼ぶ。）を作成する音声モーフィングと呼ばれる技術が存在する。この音声モーフィングに利用可能な音声分析変換合成ツールとしてSTRAIGHTと呼ばれるシステムが知られている（非特許文献1に記載）。

【0019】

この音声モーフィング技術を用いれば、予め録音した2種類の音声間の中間の声質を持つ音声を作成できる。本実施の形態では、実際の音声データだけではなく、音声モーフィングにより作成した中間音声も使用し、再生中の音声の声質を所望にしたがって変化させる。中間音声を作成する上では、元になる2種類の音声、例えば平坦（normal）な音声、暗い（dark）音声、ねっとりした（wet）音声の中の2種類により、同じ内容のテキストを朗読する、または同じ歌を歌うときの声を録音することが有効である。このとき、テキストの朗読速度、または歌の速さを同じにするようにするとよい。例えば歌の場合であれば予め録音された伴奏にあわせて歌を歌うようにすればよい。また、テキストの朗読の場合には、最初に朗読した音声をヘッドホンで話者に聞かせながら2回目の朗読を行なうようにしてもよい。

【0020】

なお、本明細書において「表情付け」とは、音声を聞く人が、その音声によりその音声に対して感じる主観的な印象のことをいう。また、本明細書では、そのような印象により表される声の性質を「声質」と呼ぶ。すなわち、本明細書においては「表情付け」と「声

10

20

30

40

50

質」とは同じ意味を表す。

【 0 0 2 1 】

なお、以下に記載する実施の形態の説明では、同じ部品には同じ参照符号を付す。それらの名称および機能も同一である。したがってそれらについての詳細な説明は繰返さない。

【 0 0 2 2 】

図 1 に、本発明の一実施の形態にかかる音声モーフィングシステム 2 0 の概略構成を示す。図 1 を参照して、この音声モーフィングシステム 2 0 は、予め準備された、基準となる第 1 の音声 3 8 および第 2 の音声 4 0 を記憶する基準音声記憶装置 2 6 と、基準音声記憶装置 2 6 に記憶された第 1 の音声 3 8 および第 2 の音声 4 0 を利用者により指定されたパラメータを用いて得られる複数種類のモーフィング率をそれぞれ用いてモーフィングすることにより、第 1 の音声 3 8 と第 2 の音声 4 0 との間での等間隔の中間音として知覚される 9 個のモーフィング後音声 2 8 を作成するための音声モーフィング装置 2 2 と、音声モーフィング装置 2 2 により作成されたモーフィング後音声 2 8 を格納するための記憶装置と、音声モーフィング装置 2 2 に対してモーフィング後音声 2 8 として作成される中間音声の数などのパラメータをユーザが入力する際に使用する入出力インタフェース 2 4 とを含む。

10

【 0 0 2 3 】

ここで「等間隔の中間音として知覚される」とは、聴者の主観的な印象として、発生された音声の声質が、一定の割合で一方の声質から他方の声質に変化していくように感じられることを指す。

20

【 0 0 2 4 】

音声モーフィング装置 2 2 は、実質的には前述した非特許文献 1 に記載の S T R A I G H T を使用する。

【 0 0 2 5 】

音声モーフィングシステム 2 0 はさらに、モーフィング後音声 2 8 と第 1 の音声 3 8 および第 2 の音声 4 0 とを記憶するための音声データ記憶装置 3 0 と、所与のモーフィング率に基づいて、音声データ記憶装置 3 0 に記憶された音声データのうち、与えられたモーフィング率に最も近い音声をリアルタイムで選択し音声信号として再生するためのモーフィング音声選択・再生装置 3 2 と、モーフィング音声選択・再生装置 3 2 に対してモーフィング率などの指示を与えるためにユーザが利用するユーザインタフェース 3 4 と、モーフィング音声選択・再生装置 3 2 により再生された音声信号を音声に変換するためのスピーカシステム 3 6 とを含む。

30

【 0 0 2 6 】

図 1 に示す音声モーフィングシステム 2 0 は、一般的にはコンピュータシステムのハードウェアと、当該ハードウェアにより実行されるプログラムとにより実現される。図 2 にこの音声モーフィングシステム 2 0 を実現するコンピュータシステム 5 0 の外観図を、図 3 にそのブロック図を、それぞれ示す。

【 0 0 2 7 】

図 2 を参照して、音声モーフィングシステム 2 0 を実現するコンピュータシステム 5 0 は、マイクロフォン 6 4 と、一組のスピーカ 3 6 と、C D - R O M (Compact Disc Read Only Memory)ドライブ 7 0 および F D (Flexible Disk)ドライブ 7 2 とを有するコンピュータ 7 8 と、いずれもコンピュータ 7 8 に接続された L C D (液晶表示装置) 7 4 とキーボード 6 6 とマウス 6 8 とを含む。

40

【 0 0 2 8 】

図 3 はコンピュータ 7 8 のハードウェアのブロック図である。図 3 を参照して、コンピュータ 7 8 は、C P U (Central Processing Unit: 中央処理装置) 8 0 と、C P U 8 0 に接続されたバス 8 2 と、バス 8 2 に接続された読出専用メモリ (R O M) 8 4 と、バス 8 2 に接続されたランダムアクセスメモリ (R A M) 8 6 と、バス 8 2 に接続されたハードディスク 8 8 と、C D - R O M 9 6 からデータを読出す C D - R O M ドライブ 7 0 と

50

、FD98からデータを読み出したりFD98にデータを書込んだりするためのFDドライブ72と、バス82に接続され、マイクロフォン64とスピーカ36とが接続されるサウンドボード90と、バス82に接続され、ローカルエリアネットワーク(LAN)等のネットワーク上でのデータコミュニケーション能力を提供するネットワークボード92と、ビデオキャプチャボード94とを含む。

【0029】

図2、図3に示すコンピュータシステム50上で所定の制御構造を有するソフトウェアを動作させることにより、図1に示す音声モーフィングシステム20を実現できる。

【0030】

図4は、図1に示す音声モーフィング装置22のブロック図である。図4を参照して、音声モーフィング装置22は、入出力インタフェース24を介して操作者から音声モーフィングの特徴を定めるパラメータを受取るためのパラメータ入力部100と、パラメータ入力部100により受取られたパラメータにしたがい、モーフィング後音声を知覚的に等間隔に並ぶように二つの音声のモーフィング率を定める際に用いるシグモイド関数を決定するためのシグモイド関数決定部102とを含む。

10

【0031】

ところで、音声モーフィングにおいては、音声のモーフィング率をどのように決めるかが問題となる。最も簡単な方法として、二つの音声の混合割合を一定割合ずつ増減させていく方法が考えられる。例えば第1の段階として第1の音声を90%と第2の音声を10%、第2の段階として第1の音声を80%と第2の音声を20%、のように一定の差分でモーフィング率を変えていく方法である。

20

【0032】

しかし、本発明の発明者らは、このように一定の差分でモーフィング率を変化させた場合、実際に聴者にとっては一定の割合で音声に変化していくようには感じられないことを実験により確認した。さらに実験により、以下に説明するようにシグモイド関数を用いてモーフィング率を決定すると、聴者にとって一定の割合で音声に変化していくように感じられることが分った。以下、本実施の形態で使用する、シグモイド関数を用いたモーフィング率の決定の仕方を説明する。

【0033】

図5に、シグモイド関数のグラフ130を示す。シグモイド関数は、一般的に以下の式により定義される。

30

【0034】

【数1】

$$f(x) = 1 / (1 + \exp(-ax + b)) + c$$

例えばこの式のパラメータa、bおよびcをパラメータ入力部100を介して操作者から受取ることにより、シグモイド関数決定部102は任意のシグモイド関数を定義できる。

【0035】

こうして決定されたシグモイド関数130のグラフのうち、中央の変曲点を含んだ左右対称な部分をモーフィング率の決定に用いる。まず、その横軸を中間音声の段階数に合わせて等分し、それぞれに段階番号を割当てる。図5には、その例として0~10までを示す。割当てられた0~10までの数字のうち、左端の「0」は第1の音声100%、第2の音声0%(すなわち第1の音声のみ)の音声を示す。右端の「10」は、第1の音声0%、第2の音声100%(すなわち第2の音声のみ)の音声を示す。中間の1~9までは、それぞれ中間の音声の段階を示す。段階0における第2の音声のモーフィング率を0%、段階10における第2の音声のモーフィング率を100%として、縦軸にモーフィング率をとることができる。

40

【0036】

図5に示すシグモイド関数130について、上記した各中間段階(1~9)での値を求

50

め、その値をその段階での第2の音声のモーフィング率とする。シグモイド関数の曲線は中央に変曲点があるため、上記した図5に示す例では段階5における第2の音声のモーフィング率は50%となる。しかしそれ以外の点では、モーフィング率は段階番号に対し非線形に変化する。

【0037】

図4に示す音声モーフィング装置22は、このようにして求めた各段階でのモーフィング率の値132を用いて、各段階での音声モーフィングを行なう。このため、音声モーフィング装置22はさらにモーフィング率決定部104を含む。

【0038】

音声モーフィング装置22はさらに、モーフィング率決定部104により決定された複数のモーフィング率の全てについて、予め準備された第1の音声38と第2の音声40との間の音声モーフィングを行なう繰返し制御を実行するための繰返し制御部106と、繰返し制御部106による制御にしたがい、第1の音声38と第2の音声40とを繰返し制御部106から指定されるモーフィング率で混合してモーフィング後の音声を生成するためのモーフィング実行部108と、モーフィング実行部108により作成されたモーフィング後音声28を、繰返し制御部106による繰返し制御にしたがって異なる名称を付したファイルとして、予め準備された音声記憶装置120内に保存するための保存処理部110とを含む。音声記憶装置120としては、図3に示すハードディスクを用いることができる。

【0039】

音声モーフィング装置22は、後述するようにコンピュータハードウェアと、コンピュータにより実行されるプログラムとにより実現される。図6に、繰返し制御部106、モーフィング実行部108、および保存処理部110を実現するためのプログラムの制御構造をフローチャート形式で示す。なお、図4に示すモーフィング率決定部104により決定された9つの中間段階のモーフィング率を $r(k)$ ($k=1\sim 9$)とする。

【0040】

図6を参照して、このプログラムは、起動されるとまず初期設定を行なう(ステップ140)。続いてステップ142において、繰返し制御変数 k に0を代入する。ステップ144で繰返し制御変数 k に1を加算する。ステップ146で繰返し制御変数 k の値が予め設定された中間段階数(本実施の形態では9)を越えたか否かを判定する。越えていれば処理を終了する。越えていなければステップ148に進む。

【0041】

ステップ148では、モーフィング率 $r(k)$ で第1の音声と第2の音声とをSTRAIGHTを用いて音声モーフィングする。ステップ150では、得られたモーフィング後の音声を「morph_k.wav」(k は1~9までの数字)というファイル名で音声記憶装置120に保存する。制御はステップ144に戻る。

【0042】

こうして、 $k=1\sim 9$ まで音声モーフィングと保存とを繰返すことにより、図4の音声記憶装置120には、9段階のモーフィング後音声28が記憶される。なお、第1の音声38および第2の音声40も予め音声記憶装置120に記憶しておくことにより、音声記憶装置120には後述するモーフィング音声選択・再生装置32において使用する音声資源が全て記憶されることになる。

【0043】

図7に、そのモーフィング音声選択・再生装置32の機能的ブロック図を示す。図7を参照して、モーフィング音声選択・再生装置32は、モーフィング音声選択・再生装置32をコンピュータにより実現するように予め準備されたプログラム160と、ユーザインタフェース34を用いたモーフィング音声の発生に関する入出力を行なうために、プログラム160に基づいて入出力画面164を作成し、モニタ74に表示させるための表示作成部162と、入出力画面164に対しユーザがユーザインタフェース34を用いて何らかの操作を行なった際、その操作によりシステム内に発生するイベントを検知し、プログ

10

20

30

40

50

ラム 160 内のオブジェクトのうち、適切なものに当該イベントを振分けるためのイベント検知部 166 とを含む。表示作成部 162 およびイベント検知部 166 としての基本的な機能は、コンピュータのオペレーティングシステム (OS) により提供される。

【0044】

モーフィング音声選択・再生装置 32 には、第 1 の音声 38、第 2 の音声 40、およびモーフィング後音声 28 を記憶し音声発生のためにモーフィング音声選択・再生装置 32 に与えるための、図 4 に示すものと同じ音声記憶装置 120 と、表示作成部 162 により発生される音声信号を音声に変換するためのスピーカ 36 と、表示作成部 162 により作成される音声再生シーケンスファイルを記憶するための記憶装置 178 とが接続される。記憶装置 178 は、音声記憶装置 120 と同様、図 3 に示すハードディスク 88 により実現できる。

10

【0045】

表示作成部 162 は、音声記憶装置 120 に記憶された第 1 の音声 38 および第 2 の音声 40 に付されたラベルを読み出し、入出力画面 164 に表示することでそれぞれの音声の声質をユーザに対し提示することができる。

【0046】

図 8 に入出力画面 164 の例を示す。図 8 上段を参照して、入出力画面 164 は、それぞれモーフィング対象となる第 1 の音声 38 および第 2 の音声 40 に付されたラベルを表示するためのラベルテキスト領域 210 および 212 と、音声をつまみ調節のいづれに設定するかをユーザが指定するように準備されたスライダ 214 と、音声発生を開始および停止をそれぞれ指示する際にユーザが使用するための開始ボタン 216 および停止ボタン 218 とを含む。

20

【0047】

図 8 に示す例において、ラベルテキスト領域 210 には「normal」(特に特徴をもたない、中立的な平坦な音)、ラベルテキスト領域 212 には「wet」(鼻にかかったような「ねっとり」した声色)というラベルがそれぞれ表示されている。これら以外にも例えば「dark」(母音が全体的に後舌母音に近くなるような「暗い」音)、子守唄のような「whisper」(ささやき声)などのラベルが考えられるが、一般的にはラベルは使用者(または録音者)の主観に基づいて自由に音声ファイルに付しておけばよい。

30

【0048】

スライダ 214 は、スライダ目盛 242 と、スライダつまみ 240 とを含む。スライダつまみ 240 を例えば図 8 の下段に示すようにマウスポインタ 244 でドラッグすることにより、スライダつまみ 240 がスライダ目盛 242 上を移動する。スライダ 214 は、スライダ目盛 242 上のスライダつまみ 240 の位置に対応する値をリアルタイムで検知し、属性値として保持する。同時にスライダ 214 が操作されたというイベントをイベント検知部 166 に与える機能を持つ。

【0049】

開始ボタン 216 が押され、音声発生が開始された後にユーザがスライダつまみ 240 をスライドさせることにより、音声の発生中に、その音声を第 1 の音声から第 2 の音声まで、中間段階を含めて 11 種類の音にリアルタイムで変更させることができる。モーフィング音声選択・再生装置 32 は、音声発生時のユーザ操作による声質変更のシーケンスを記憶し、記憶装置 178 (図 7 参照) にファイルとして出力する機能を持つ。

40

【0050】

再び図 8 を参照して、入出力画面 164 はさらに、記憶装置 178 に記憶された声質変更のシーケンスファイルからシーケンスを読み出して当該シーケンスにしたがって声質を変更して所定の音声を再生する際にユーザが使用する再生ボタン 220 を含む。

【0051】

プログラム 160 は、図 8 に示すスライダ 214 の実体を構成するスライダオブジェクト 190 と、それぞれ開始ボタン 216、停止ボタン 218、および再生ボタン 220 の

50

実体を構成する開始ボタンオブジェクト192、停止ボタンオブジェクト194、および再生ボタンオブジェクト196とを含む。

【0052】

これら各オブジェクトについて、プログラム160の実行開始とともにそれぞれインスタンスが生成されて入出力画面164の作成、イベント検知部166による適切なメソッドの実行、および各インスタンスに付随する属性値の取得と記憶などが実行される。

【0053】

以下、プログラム160に含まれるプログラムコードをコンピュータで実行することにより実現される機能を、それぞれ機能ブロックとして説明する。

【0054】

すなわち、モーフィング音声選択・再生装置32はさらに、プログラム160により実現される機能ブロックとして、イベント検知部166により開始ボタン216の押下というイベントが検知されたことに応答して計時を開始し、停止ボタン218の押下というイベントが検知されたことに応答して計時を終了するタイマ184と、イベント検知部166によりスライダオブジェクト190の操作に関するイベントが検知されたことに応答して、スライダオブジェクト190からその属性値としてスライダ値(区間[0, 10]の間の整数)を読み取り、音声記憶装置120に記憶された音声(第1の音声38、第2の音声40、および9種類の間音声)のうちいずれを使用するかを決定し、音声ファイルの名称を出力するための音声選択処理部168とを含む。

【0055】

プログラム160により実現される機能ブロックはさらに、第1および第2の二つの入力を持ち、第1の入力が音声選択処理部168の出力を受けると接続され、図8に示す開始ボタン216が押されたときには第1の入力を、再生ボタン220が押されたときには第2の入力を、それぞれ選択し、選択された入力に与えられた信号を出力するためのセレクタ170と、セレクタ170の出力を受け、音声記憶装置120に記憶された音声ファイル(第1の音声38、第2の音声40、および9種類の間音声)のうち、セレクタ170の出力により指定される音声ファイルを読み出してタイマ184により指定される再生位置から音声信号への変換を開始しスピーカ36に与えるための音声発生処理部172とを含む。

【0056】

プログラム160により実現される機能ブロックはさらに、開始ボタン216の押下というイベントが検知されたことに応答して、音声発生処理部172による音声の発生を開始させるための開始指示部173と、開始指示部173からの音声発生の開始指示と、音声選択処理部168による音声選択処理とに併せて、そのときの音声選択処理部168の出力をタイマ184の計時値とともにシーケンスとして記録するための選択音声記録部174と、停止ボタン218の押下というイベントに併せて、選択音声記録部174により記録されている再生シーケンスをファイルとして記憶装置178に保存させるためのシーケンス保存部176とを含む。

【0057】

プログラム160により実現される機能ブロックはさらに、再生ボタン220の押下というイベントに併せて、記憶装置178に保存されている1または複数の再生シーケンスファイルのいずれかをユーザに選択させるためのファイル選択処理部180と、ファイル選択処理部180により再生シーケンスファイルが選択されると、タイマ184を起動し、ファイル選択処理部180により選択された再生シーケンスファイルを記憶装置178から読み出して、タイマ184の計時に基づいて、再生開始後、選択された再生シーケンスにより指定された時刻になると再生シーケンスにより指定された音声ファイル名をセレクタ170の第2の入力に与えることにより、再生シーケンスに基づく音声再生を制御するための音声再生制御処理部182とを含む。

【0058】

[動作]

10

20

30

40

50

図1～図8を参照して、上記した音声モーフィングシステム20は以下のように動作する。動作は大きく3つのフェーズに分けられる。第1のフェーズでは、予め準備された、互いの声質の異なる第1の音声38と第2の音声40とから9個の中間段階のモーフィング後音声28を作成する。なお、これに先立ち、同じ話者により、声質（音声の表情）を変えて同一の文章を読んだり同一の歌を歌ったりすることによって二つの音声を収録しておき、これらをそれぞれ第1の音声38および第2の音声40として保存しておく。また、第1の音声38および第2の音声40の音声ファイルには、付属情報として声質を示すラベルを付しておく。

【0059】

第2のフェーズでは、このようにして作成されたモーフィング後音声28と、最初に準備された第1の音声38および第2の音声40とを用い、声質を自由に変更しながらこれら音声の発生を行なう。このとき、再生シーケンスが記憶装置178に保存される。第3のフェーズでは、記憶装置178にファイルとして保存された再生シーケンスを読み出し、その再生シーケンスにしたがって音声を選択し発生させることにより、再生シーケンスを再現する。以下、各フェーズでの音声モーフィングシステム20の動作を説明する。

【0060】

- 音声モーフィング -

図4を参照して、パラメータ入力部100は、入出力インタフェース24を用いてユーザから、中間段階の数と、シグモイド関数決定のためのパラメータとを受取る。パラメータ入力部100は、このパラメータをシグモイド関数決定部102に与える。

【0061】

シグモイド関数決定部102は、与えられたパラメータにしたがってシグモイド関数を決定する。決定されたシグモイド関数に関する情報はモーフィング率決定部104に与えられる。

【0062】

モーフィング率決定部104は、このシグモイド関数を用い、図5を参照して説明した方法にしたがって、各中間段階におけるモーフィング率 $r(k)$ ($k=1\sim 9$)を決定する。モーフィング率決定部104は、決定したモーフィング率を、中間段階数とともに繰返し制御部106に与える。

【0063】

繰返し制御部106は、指定されたモーフィング率で中間段階の音声を音声モーフィングにより作成するようモーフィング実行部108を制御し、作成させる。繰返し制御部106はまた、このようにして作成された9つの中間段階のモーフィング後音声28を音声記憶装置120に格納するよう保存処理部110を動作させる。このとき、保存処理部110は、各中間音声のファイルに前述したとおり「morph_k.wav」というファイル名を付す。音声記憶装置120には、予めモーフィングに用いられる第1の音声38および第2の音声40も準備されているものとする。

【0064】

以上で第1のフェーズの動作は完了である。音声記憶装置120を音声データ記憶装置30としてモーフィング音声選択・再生装置32に接続することにより、モーフィング音声選択・再生装置32による、中間音声を用いて表情を変化させた音声の発生が可能になる。

【0065】

- 中間音声を用いた任意の声質の音声の発生 -

図7を参照して、第2のフェーズにおいてはモーフィング音声選択・再生装置32は以下のように動作する。モーフィング音声選択・再生装置32の各機能ブロックは、予めプログラム160の形で準備されている。

【0066】

ユーザがモーフィング音声を用いた音声の発生をするようにモーフィング音声選択・再生装置32に指示すると、図7および図8に示す入出力画面164が表示作成部162に

10

20

30

40

50

より作成され、モニタ 74 上に表示される。このとき、表示作成部 162 は、音声記憶装置 120 に記憶された第 1 の音声 38 および第 2 の音声 40 に付されたラベルを読み込み、それぞれラベルテキスト領域 210 および 212 に表示する。この表示により、ユーザは二つの音声の声質がどのようなものであるかを知ることができる。

【0067】

ユーザは、まずスライダ 214 のスライダつまみ 240 をマウスポインタ 244 によりスライドさせることで音声発生開始時の中間音声の段階を指定する。この操作によりイベントが発生し、イベント検知部 166 はこのイベントをスライダオブジェクト 190 に与える。スライダオブジェクト 190 は、属性値として保持しているスライダつまみ 240 の位置を示す値を音声選択処理部 168 に与える。音声選択処理部 168 は、この値に基づいて、音声記憶装置 120 に記憶された音声ファイルのうちどの音声ファイルを選択するかを決定し、セレクトア 170 に与える。

10

【0068】

ユーザが開始ボタン 216 を押下すると開始ボタン 216 の押下イベントが発生する。イベント検知部 166 はこのイベントを開始ボタンオブジェクト 192 に与える。開始ボタンオブジェクト 192 は、このイベントに回答し開始指示部 173 を制御してタイマ 184 を起動する。セレクトア 170 は、第 1 の入力、すなわち音声選択処理部 168 から与えられたファイル名を選択して音声発生処理部 172 に与える。開始指示部 173 は、音声発生処理部 172 に対して音声の発生の開始を指示する。

【0069】

音声発生処理部 172 は、音声記憶装置 120 から、セレクトア 170 を介して音声選択処理部 168 から与えられたファイル名に対応するファイルを読み出し、タイマ 184 の計時にしたがった位置からタイマ 184 の計時に同期して再生を開始する。発生された音声信号はスピーカ 36 に与えられ、音声に変換される。

20

【0070】

一方、選択音声記録部 174 は、開始ボタン 216 が押下されたことに回答して、タイマ 184 の計時が 0 のときの音声選択処理部 168 の出力をタイマ 184 の値とともに記録する。

【0071】

このようにして音声ファイルが再生されている途中でユーザが図 8 に示すスライダつまみ 240 を操作して、スライダつまみ 240 を別の位置に移動させたものとする。このイベントはイベント検知部 166 により検知され、スライダオブジェクト 190 に与えられる。スライダオブジェクト 190 はこのイベントに回答して、属性値として保持しているスライダつまみ 240 の位置を示す値を音声選択処理部 168 に与える。音声選択処理部 168 は、この値に対応する音声ファイル名を選択し、セレクトア 170 に与える。セレクトア 170 はこの値を音声発生処理部 172 に与えるので、音声発生処理部 172 は指定された音声ファイルを新たに読み出し、タイマ 184 により示される位置からタイマ 184 に同期して音声の再生を開始する。

30

【0072】

また、音声選択処理部 168 が新たなファイル名を選択したことに回答し、選択音声記録部 174 はその新たなファイル名と、そのときのタイマ 184 の値とを組にして追加して記録する。

40

【0073】

図 9 の上段に、ユーザによるスライダつまみ 240 の操作の例を時系列で示す。図 9 上段では、縦軸にスライダつまみ 240 の位置をスライダ目盛の値で、横軸に時間を、それぞれ示す。図 9 上段に示すように、スライダの位置は時間的にある軌跡 250 を描く。これがモーフィング後音声を用いた再生シーケンスを示す。この再生シーケンスを選択音声記録部 174 によって記録しておけば、同じ再生シーケンスを再現することができる。

【0074】

ただし、スライダつまみ 240 の位置はスライド目盛の中間となることもあり得る。そ

50

うした場合には、スライダつまみ 240 の位置に最も近い目盛を選択し、その目盛に対応する音声ファイルを選択する。したがって、なだらかな線を描くスライダつまみ 240 の軌跡 250 は、図 9 下段に示すようにある中間音声から次の中間音声に、音声の種類としては不連続な形で、ただし時間的には連続して、再現されることになる。ただし、このままでは音声のつながり目で「ブツン」という雑音が入る。そこで、こうしたつながり目では、先の音声を徐々にフェードアウトし、後の音声を徐々にフェードインする形で音声を混合することで雑音の発生を防止する。

【0075】

中間段階の数としてある程度大きい値を選択しておけば、図 9 下段に示すような形で音声を再生しても、聴者には違和感を与えない。また、各中間音声のモーフィング率は、シグモイド関数を用いてできるだけ等しい間隔で相違した音声となるように設定されている。したがって、聴者には、こうして生じる声質の変化は、滑らかでかつ自然に感じられることになる。

10

【0076】

停止ボタン 218 が押されると、そのイベントはイベント検知部 166 により検知され、停止ボタンオブジェクト 194 に与えられる。停止ボタンオブジェクト 194 は、このイベントにตอบสนองして選択音声記録部 174 による音声シーケンスの記録を中止させる。またタイマ 184 も停止させる。さらに停止ボタンオブジェクト 194 は、シーケンス保存部 176 に指示を与え、選択音声記録部 174 により記録された再生シーケンスを記憶装置 178 に保存させる。このときの保存名は、ダイアログボックスを開いてユーザに指定させる。

20

【0077】

以上が声質を変化させて音声を発生させる際のモーフィング音声選択・再生装置 32 の動作である。

【0078】

記憶装置 178 に保存された再生シーケンスに基づいて音声発生を行なう際には、モーフィング音声選択・再生装置 32 は以下のように動作する。

【0079】

再生シーケンスに基づく音声発生を行なう際には、ユーザは再生ボタン 220 (図 8 参照) を押下する。このイベントはイベント検知部 166 により検知される。イベント検知部 166 はこのイベントを再生ボタンオブジェクト 196 に与える。

30

【0080】

再生ボタンオブジェクト 196 は、イベントが与えられたことにตอบสนองしてファイル選択処理部 180 を起動する。ファイル選択処理部 180 は、記憶装置 178 に保存された各ファイルのファイル名を読み出し、ファイル選択処理部 180 にダイアログボックスとしてファイル選択ダイアログを表示する。ユーザが所望のファイルを選択すると、ファイル選択処理部 180 は選択されたファイル名を音声再生制御処理部 182 に与える。

【0081】

音声再生制御処理部 182 はタイマ 184 を起動させる。音声再生制御処理部 182 はさらに、選択された再生ファイルを読み込み、最初に選択されていた音声ファイルを指定する信号をセクタ 170 に与え、同時に音声発生処理部 172 を起動する。セクタ 170 は、第 2 の入力を選択して音声発生処理部 172 に与える。

40

【0082】

音声発生処理部 172 は、セクタ 170 を介して音声再生制御処理部 182 から与えられた音声ファイルを音声記憶装置 120 から読み出し、タイマ 184 の計時にしたがって再生を開始する。

【0083】

音声再生制御処理部 182 は、再生ファイル中の音声ファイル名とタイマ計時との組のうち、タイマ計時の値をタイマ 184 による計時と常に照合し、タイマ 184 の計時と一致するタイマ計時を持つものがあればそのタイマ計時と組になっている音声ファイル名を

50

セクタ 170 に与える。したがって音声発生処理部 172 は、この新たな音声ファイル名により指定される音声ファイルを音声記憶装置 120 から読出し、タイマ 184 の計時にしたがって再生を行なう。

【0084】

こうして、ファイル選択処理部 180 によって指定された再生ファイルによる再生シーケンスにしたがって、音声発生処理部 172 が音声記憶装置 120 中の音声を随時切替えながら音声発生を行なう。

【0085】

以上のように本実施の形態に係る音声モーフィングシステム 20 によれば、所望の音声を全て収録しなくても、第 1 の音声 38 および第 2 の音声 40 の中間音声を音声モーフィングで準備し、さらにユーザの選択にしたがってリアルタイムで音声をそれらの中で切替えながら音声発生を行なうことができる。中間音声の間の相違が一定に知覚されるように予め音声のモーフィング率を決めて中間音声を作成しているため、音声発生の途中で音声の切替えを行なっても不自然には感じられない。また、リアルタイムで作成した音声シーケンスを記録しておくことで、いつでも同じ再生シーケンスを再現できる。

【0086】

また、このように音声モーフィングを使用して中間段階の音声を作成し、それらを切替えて音声発生を行なうと、中間段階の音声については第 1 の音声 38 および第 2 の音声 40 の声質に応じ、それらの中間の声質を表すものとして知覚される。

【0087】

したがって、本実施の形態に係る音声モーフィングシステム 20 によれば、多数の音声を収録しなくても、任意の時刻にユーザが選択した声質を用いて音声を発生させることで、豊かな表情を持つ音声の発生が可能になる。

【0088】

なお、上記実施の形態ではパーソナルコンピュータのユーザインタフェースを使用してスライダを実現したが、本発明はそのような実施の形態には限定されない。例えばシンセサイザ等に組み込む形で、シンセサイザのスライダを用いたコントロールを行なっても良い。

【0089】

また、予め音声ファイルに付されていたラベルを変更したいとユーザが考えることもあるので、ラベルテキスト領域 210 をテキスト入力可能な領域とし、ラベルをユーザが変更可能にしてもよい。例えば、図 8 において「normal」（平坦）とラベルが表示されている音声が、ユーザには「dark」（暗い）と感じられることもある。そうした場合には、図 10 に示すようにラベルテキスト領域 210 に「dark」と入力して音声ファイルとともに保存しておくことにより、次にこの音声ファイルを使用する場合には「dark」というラベルがラベルテキスト領域 210 に表示される。

【0090】

なお、図 8 に示すように「normal」というラベルを持つ音声と、「wet」というラベルを持つ音声との間で表情付けを変化させるということは、特に特徴をもたない平坦な（中立的な）音声に対し、「wet」な（ねっとりした）表情の強度を変化させながら付加させることであると考えられる。これに対し図 10 に示すように「dark」というラベルを持つ音声と、「wet」というラベルを持つ音声との間で表情付けを変化させるということは、「dark」という表情付けと、「wet」という表情付けとの間での、表情付けの種類を変化させる、ということであると考えられる。

【0091】

[変形例]

上記した実施の形態では、2種類の音声の間で音声モーフィングを行なって得た中間音声を用いた。しかし本発明はそのように2種類の音声の間でのモーフィングには限定されない。例えば3種類以上の音声の間でのモーフィングを行なうことも可能である。モーフィング自体はSTRAIGHTを使用して行なうことができる。問題は、3種類以上の音

10

20

30

40

50

声の間でのモーフィング率を定める方法である。

【0092】

図11を参照して、3種類の音声の間での音声モーフィングを行なう際のモーフィング率の決定の方法について説明する。今、3種類の音声A、BおよびCの間でのモーフィングを行なうものとする。図11に示すように、これら3つの音声に対応する頂点260、262および264を有する三角形を考える。

【0093】

この三角形の各辺を所定数に分割し、各辺と並行な線で分割点同士を結ぶことにより、図11においてメッシュ270を作成できる。このメッシュ270を構成する各点に対応したモーフィング音声は以下のようにして作成できる。

【0094】

例えば音声Aおよび音声Bの間での各分割点に対応する中間音声は、上記した実施の形態での方法と同様の音声モーフィングで行なうことができる。音声Aおよび音声Cの間、音声Bおよび音声Cの間での音声モーフィングもそれぞれ行なうことができる。さらに、メッシュ270の各交点(例えば交点272)での中間音声は、その交点を通る任意の線の両端(例えば点274、276)の中間音声を、その両端からその交点までの距離の比に応じたモーフィング率でモーフィングすることにより作成できる。したがって、メッシュ270の各点に対応する中間段階の音声を全て作成できる。

【0095】

音声発生時には、上記したメッシュ270を有する三角形をコンピュータモニタ上に表示し、メッシュ270中の交点をポイントにより指定する。具体的には、ポイントの座標を調べ、メッシュ270の交点のうちポイントにより表される点に最も近い座標を持つ交点に対応する中間音声を選択すればよい。例えば図12を参照して、ポイント280がメッシュの3つの交点290、292および294で形成される三角形の内部にあるものとする。このときには、ポイント280の位置と各交点290、292および294との間の距離d1、d2およびd3を調べ、距離がもっとも小さくなる点を選択する。

【0096】

このような中間音声の発生方法は、元となる音声は3種類の場合だけでなく、4種類以上の場合にも同様に適用できる。

【0097】

さらに、音声発生時には、上記のように作成したメッシュの二つの交点に対応する中間音声の間でさらに音声モーフィングを行なうようにしてもよい。この場合の例を図13に示す。図13を参照して、3種類の音声A、BおよびCに対応する3つの頂点260、262および264を有する三角形を考える。その中に、上記した方法と同様にしてメッシュ270を作成する。このメッシュ中の任意の交点、例えば二つの交点310と312とを選択し、この二つの交点を結ぶ線分314を任意の数に分割することにより、交点310と312とに対応する中間段階の音声の間での中間音声をモーフィングにより作成できる。こうして作成した音声を発生させるときには、実施の形態で説明したのと同様の方法を利用できる。

【0098】

また、例えば元となる音声は3種類の場合には、磁気センサ、光学センサ、画像処理技術を用いた物体検出など、対象物の位置を3次元的な座標系中で検出できるシステムを用い、音声のモーフィング率を指定することができる。例えば予め3種類の音声に対し種々のモーフィング率で中間音声を作成しておく。音声発生時には、図14に示すように空間にxyz座標系を設定する。利用者は、所定の3次元ポイントで空間内の1点322を指定する。この1点322に対し、座標値(X, Y, Z)が定まる。その座標値(X, Y, Z)に応じたモーフィング率($X / (X + Y + Z)$, $Y / (X + Y + Z)$, $Z / (X + Y + Z)$)に最も近いモーフィング率の中間音声を選択して音声を発生させる。

【0099】

このように3次元的な座標指定によって発生させる音声を切替えることにより、ダンス

10

20

30

40

50

などのパフォーマンスとそれに伴う歌などとを連動させることができる。もちろん、次元数は3次元に限定されず、4次元以上の任意の次元数を用いることも可能である。

【0100】

なお、本実施の形態による音声の表情付けの変化は、リアルタイムで実行できる。また、音声としては、予め一連の発話を別々の種類の表情付けがされた音声で実際に朗読したり歌ったりした場合だけではなく、予め別々の種類の表情付けがされた合成音声を準備しておいてもよい。したがって、本実施の形態による音声の表情付けを次のような場合にも利用できる。

【0101】

例えば、声を出ることができない人が自動プレゼンテーションにより合成音声を発する場合を考える。この場合には、予め複数種類の声質により表情付けがされた合成音声を作成しておき、合成音声によるプレゼンテーションでは、リアルタイムで音声の表情付けを変化させることができる。例えば強調したいところでは張りのある声質の音声でプレゼンテーションし、重要でないところはぼそぼそとした声質の音声でプレゼンテーションさせることができる。すなわち、観客の反応を見ながら最適と思われる声質で自動プレゼンテーションの音声を発生させることができ、プレゼンテーションをより効果的にすることができる。

10

【0102】

また、ディスクジョッキーのパフォーマンスにおいて、歌声を用いたモーフィングにより、歌唱にリアルタイムで変化する表情付けを行なうこともできる。例えばバックミュージックの方が歌よりも重要な場合には歌声は平坦な（表情のない）ものとし、歌詞に注目してほしいときには「ねっとりした」表情付け音声にモーフィングし、ゆっくり歌い終わるときには次第に「ささやき」に表情付けをモーフィングし、など、同一歌唱曲中で連続的に音声モーフィングを行なうことができる。こうすることで、ディスクジョッキーが思う通りの表情付けで歌唱を再生させることができる。

20

【0103】

3つ以上の音声のモーフィング率を指定する方法は、上記したようにモニタ上のポイントまたは3次元的なポイントだけでなく、最初に説明したのと同様、スライダによって指定することもできる。その場合のモニタ表示例を図15に示す。

【0104】

図15を参照して、モニタ表示340には、3つの音声に対応するスライダ350、352および354を表示する。これらスライダ350、352および354の左端には、各音声に付されたラベルを表示するラベルテキスト領域370、372および374が表示される。スライダ350、352および354のつまみ360、362、364を左右にスライドさせることにより、各音声のモーフィング率を調整できる。なおこの場合、スライダ目盛の数値そのものはモーフィング率に対応しない。モーフィング率は、3つのスライダのつまみが指す目盛の値の合計を基準（100%）とし、合計に対する各スライダの目盛の値の率で定めるようにすればよい。

30

【0105】

さらに、平坦な音声を基準として、任意の表情付けを行ないながら音声を発生させる場合にも本発明を適用することができる。平坦な音声も含めて音声の種類が4種類の場合について、図16～図18を参照して音声の選択方法について説明する。

40

【0106】

図16を参照して、原点380を持つ立体座標系を考える。原点をnormalというラベルを持つ音声に割当て、3軸をそれぞれ3種類の音声A、B、およびCに割当てる。これら3軸をそれぞれ音声A軸、音声B軸、および音声C軸と呼ぶことにする。

【0107】

図17を参照して、例えば原点380と音声A軸との間の中間表情付け390、原点380と音声B軸との間の中間表情付け392、および原点380と音声C軸との間の中間表情付け394を、それぞれ第1の実施の形態と同様に行なうことができる。これらの中

50

間表情付けによって発生される音声は、原点 3 8 0 に対応する平坦な音声から、表情付けの種類強度を音声 A、B、C 軸に沿って変化させて発生したものとなる。

【 0 1 0 8 】

一方、音声 A 軸と音声 B 軸との間の中間表情付け 4 0 0、音声 B 軸と音声 C 軸との間の中間表情付け 4 0 4、および音声 C 軸と音声 A 軸との間の中間表情付け 4 0 2 も考えられる。これらは、ある種類の表情付けがされた軸上の音声を、他軸上の別種類の表情付けがされた音声に変えることであるから、音声に対する表情付けの「種類」の変化に相当すると考えることができる。

【 0 1 0 9 】

図 1 6 および図 1 7 を参照して説明したような方向の中間音声の発生だけでなく、それらを所定の割合で混合した音声も発生可能である。その方法は図 1 1 を参照して説明した方法と同様である。

【 0 1 1 0 】

例えば、図 1 8 に示すように、音声 A 軸、音声 B 軸、および音声 C 軸の各々において、各音声の割合が 1 0 0 % となる点 4 2 0、4 2 2、および 4 2 4 を決める。次にこれら 3 点 4 2 0、4 2 2、および 4 2 4 を互いに結ぶことによって得られる空間上の 3 つの半直線 4 1 2、4 1 4 および 4 1 6 を考える。これら 3 つの半直線 4 1 2、4 1 4 および 4 1 6 は空間上で一つの三角形を規定する。この三角形を図 1 1 に示す方法と同様に分割することでメッシュ 4 1 0 が得られる。

【 0 1 1 1 】

さらに、音声 A 軸、音声 B 軸および音声 C 軸のそれぞれにおいて、原点 3 8 0 と前述した点 4 2 0、4 2 2 および 4 2 4 との間を 1 0 分割する。この分割により、各軸上には 0 % から 1 0 0 % まで、1 1 個の点が規定される。それら点のうち、各軸上でそれぞれ 1 0 % に相当する 3 点を互いに結んで 3 角形を考えることで、メッシュ 4 1 0 と同様のメッシュが形成できる。同様にして、2 0 % から 9 0 % までの各々の率について、メッシュ 4 1 0 と同様のメッシュが形成できる。

【 0 1 1 2 】

中間音声として、平坦な音声と、音声 A、B および C とを準備し、これら音声を用いて予め上記のように得られたメッシュの各交点に相当する混合割合の中間音声を作成しておく。音声の発生時には、ユーザが 3 次元空間上のある点（図 1 8 において原点 3 8 0 と 3 点 4 2 0、4 2 2 および 4 2 4 により形成される三角錐内のある点）を指定すると、上記したメッシュの交点のうち、指定された点に最も近い交点を定めることができる。その交点に対応する中間音声で音声を発生させる。

【 0 1 1 3 】

このようにすることにより、normal 音声を中心として、3 種類の表情付けのうち、任意のものを選択し、任意の強度で normal 音声に対しそれらの表情付けを行なうことができる。さらにそれだけでなく、3 種類の音声の種類を互いに入れ替えたり、それら音声の持つ表情を任意の割合で混合した音声を作成したりすることができる。

【 0 1 1 4 】

もちろん、中心におく音声は normal 音声に限定されるわけではなく、利用者の意図に応じて任意の表情付けを持つ音声を中心としてもよい。もっとも、中心に置く音声を normal 音声とすると、中間音声の表情がどのようなものになるか直感的に分りやすいと思われる。したがって中心に normal 音声を置くことが実用的である。

【 0 1 1 5 】

なお、複数（例えば二つ）の表情付音声の間でモーフィング率 = 0 . 5 としてモーフィングを行なった場合でも、得られる音声は normal なものとは異なるものとして知覚されることが実験によって確認されている。3 つ以上の音声についても同様で、全ての表情付音声について互いにモーフィング率が等しくなるような条件でモーフィングを行なったとしても、得られる音声は normal とは異なって知覚されると思われる。したがって、図 1 6 ~ 図 1 8 に示すように normal 音声を中心におき、この音声を基準として

10

20

30

40

50

様々な表情を付ける形でモーフィングを行なうようにすることが好ましい。

【0116】

このようにnormal音声を基準として、他に例えば3種類の表情付音声をモーフィングする場合、すなわち4種類の音声の間でのモーフィングを行なう場合でも、図15に示す3個のスライダ350、352および354を用いてモーフィング率を指定できる。各スライダには、normal音声を基準とした3個の表情付音声のモーフィング率を指定する。したがって、3個のスライダにより指定された値が全て0の場合にはnormal音声が発生されることになる。

【0117】

さらに、上記したように表情付けの強度を任意に変化させたり、種類の変化を任意に行なわせたりする場合、音声の種類は3種類または4種類に限定されるわけではない。理論的には、5種類以上の音声の間でも中間音声を同様にして定め、利用することができる。

10

【0118】

なお、上記した実施の形態では同一話者による複数声質の音声をを用いて音声モーフィングを行なっている。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。別の話者による音声の間での音声モーフィングを行なってもよい。発声される内容はテキスト朗読でもよいし、ある歌の歌声でもよい。

【0119】

さらに、上記した実施の形態では、2種類の音声の間で音声モーフィングを行なって得る中間音声の数を9個としたが、中間音声の数が9個に限定されるわけではないことはもちろんである。一般的には、基準となる音声の数が多くなると、音声の変化を滑らかにするためには中間音声の数を多くする必要がある。

20

【0120】

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のすべての変更を含む。

【図面の簡単な説明】

【0121】

【図1】本発明の一実施の形態に係る音声モーフィングシステム20の全体の概略ブロック図である。

30

【図2】音声モーフィングシステム20を実現するコンピュータシステム50の外観図である。

【図3】コンピュータシステム50のブロック図である。

【図4】図1に示す音声モーフィング装置22の概略ブロック図である。

【図5】モーフィング率を決定するためのシグモイド関数を示すグラフである。

【図6】モーフィング後音声作成処理の制御の流れを示すフローチャートである。

【図7】図1に示すモーフィング音声選択・再生装置32のブロック図である。

【図8】モーフィング音声選択・再生装置32における入出力画面164を示す図である。

40

【図9】スライド目盛の軌跡および音声シーケンスの記録例を示す図である。

【図10】入出力画面164のラベルテキスト領域210において音声ラベルを変更した状態を示す図である。

【図11】3種類の音声から中間音声を作成する方法を説明するための図である。

【図12】3種類の音声から作成された中間音声のうち、いずれを選択するかに関する方法を説明するための図である。

【図13】3種類の音声から作成した二つの中間音声の間でさらに中間音声を作成する方法を説明するための図である。

【図14】3次元センサによる中間音声の選択方法を説明するための図である。

【図15】3種類の音声から得られた中間音声を選択する際のスライド表示を説明するた

50

めの図である。

【図 1 6】声質の種類を選択および各音声種類における声質の強さの指定のために使用する座標系を説明するための図である。

【図 1 7】平坦な音声と 3 種類の音声との間の声質の種類と強度との指定方法を説明するための図である。

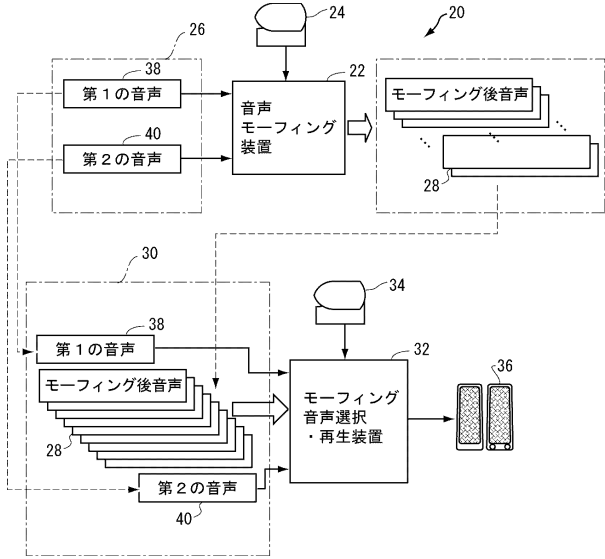
【図 1 8】平坦な音声と 3 種類の音声との間での中間音声の指定方法を説明するための図である。

【符号の説明】

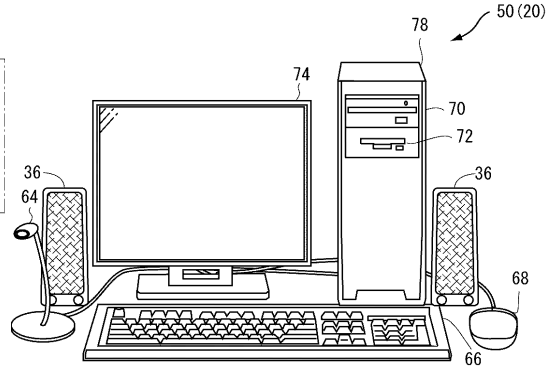
【 0 1 2 2 】

2 0	音声モーフィングシステム	10
2 2	音声モーフィング装置	
2 4	入出力インタフェース	
2 6	基準音声記憶装置	
2 8	モーフィング後音声	
3 0	音声データ記憶装置	
3 2	モーフィング音声選択・再生装置	
3 4	ユーザインタフェース	
3 8	第 1 の音声	
4 0	第 2 の音声	
5 0	コンピュータシステム	20
1 0 0	パラメータ入力部	
1 0 2	シグモイド関数決定部	
1 0 4	モーフィング率決定部	
1 0 6	繰返し制御部	
1 0 8	モーフィング実行部	
1 1 0	保存処理部	
1 2 0	音声記憶装置	

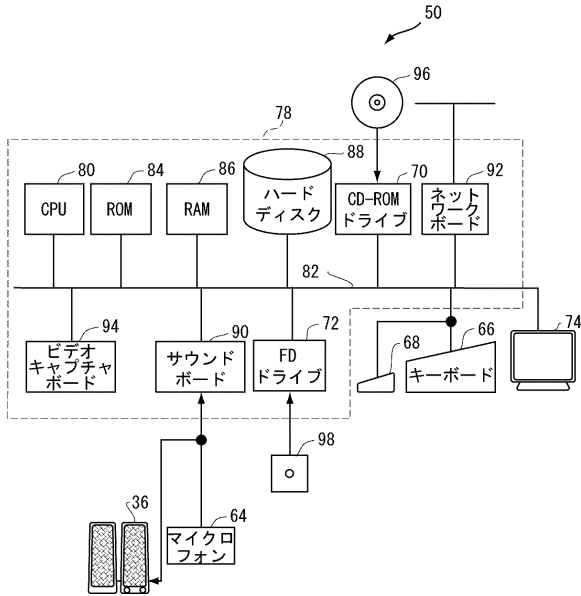
【図1】



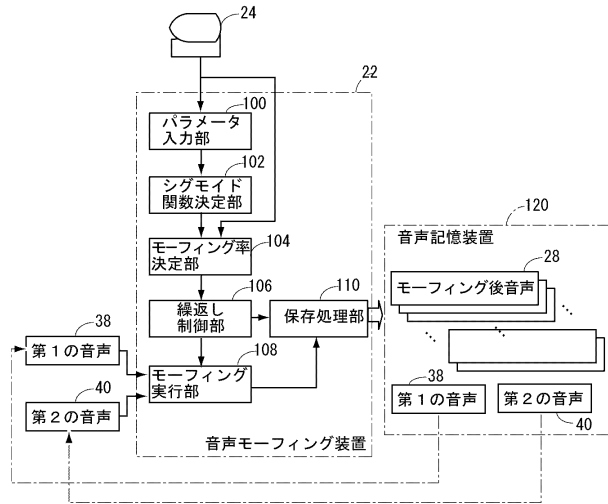
【図2】



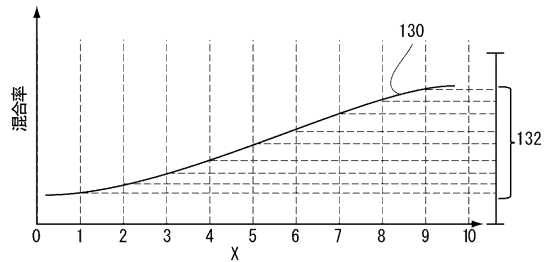
【図3】



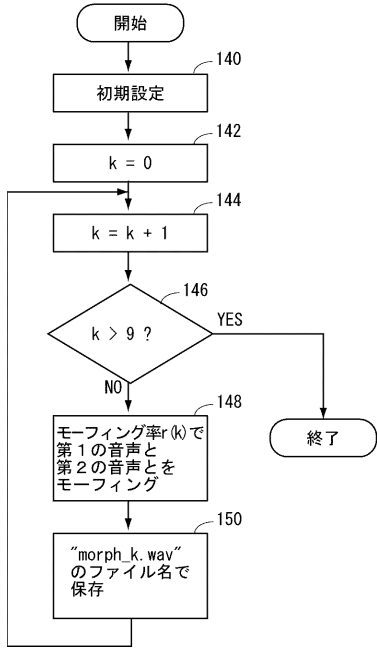
【図4】



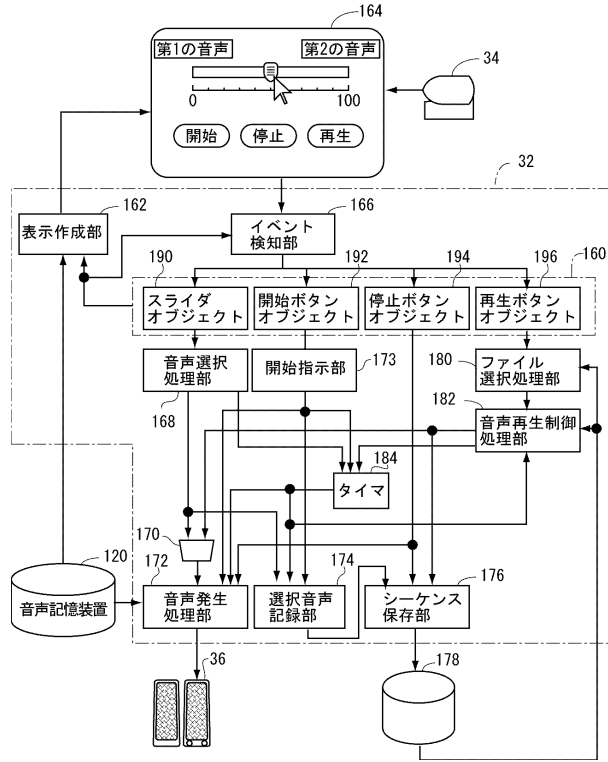
【図5】



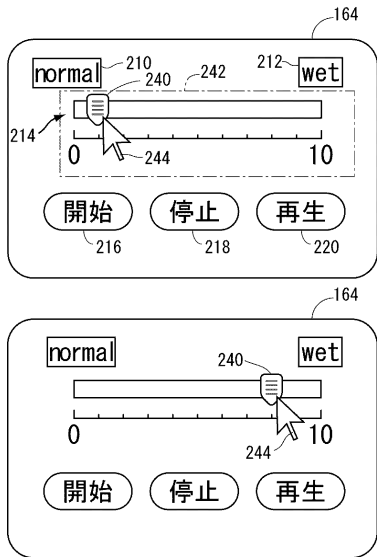
【図6】



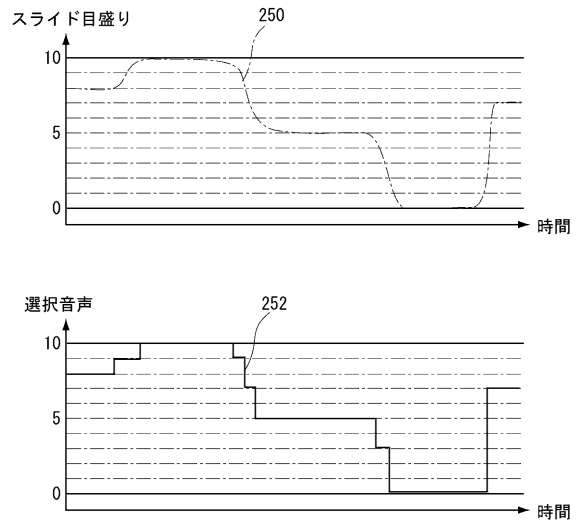
【図7】



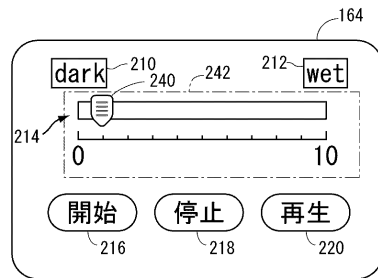
【図8】



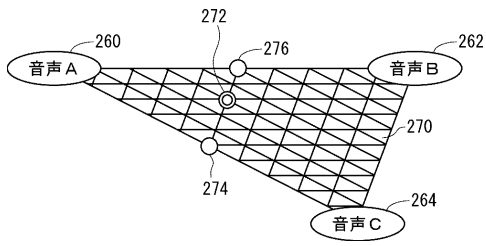
【図9】



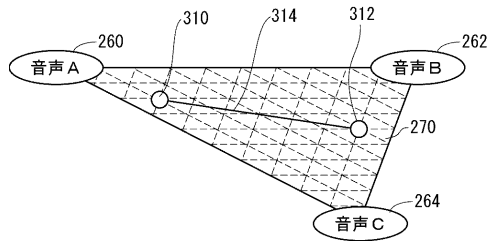
【図10】



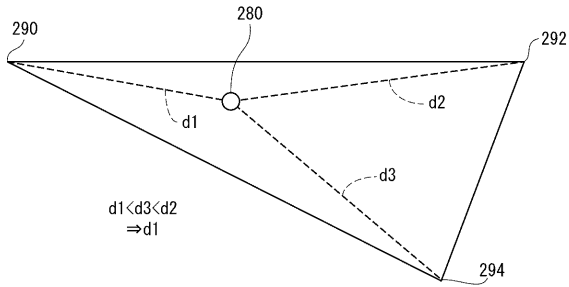
【図 1 1】



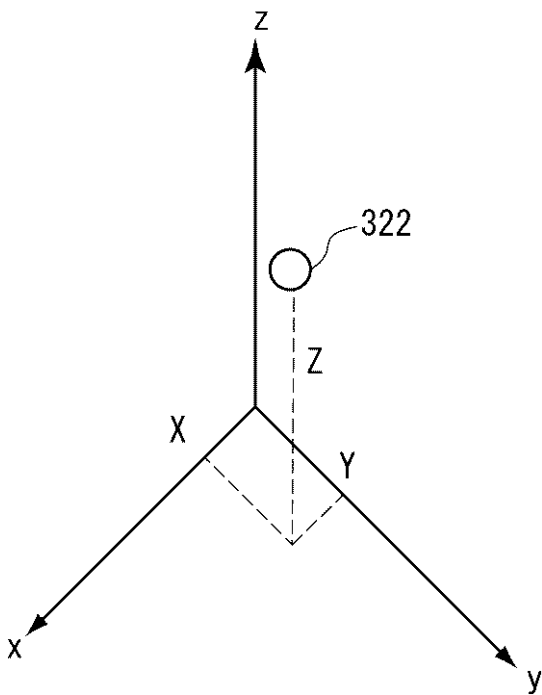
【図 1 3】



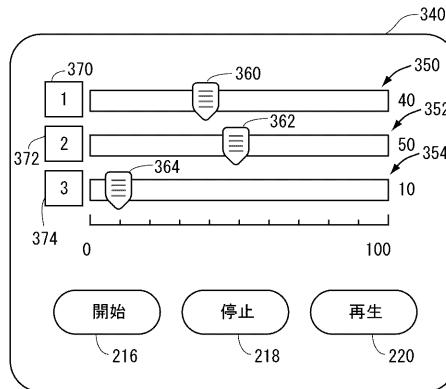
【図 1 2】



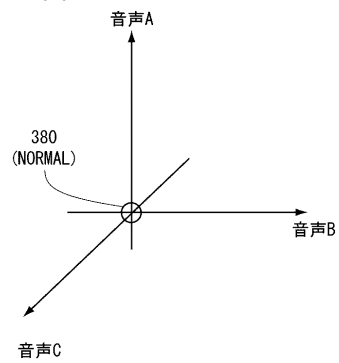
【図 1 4】



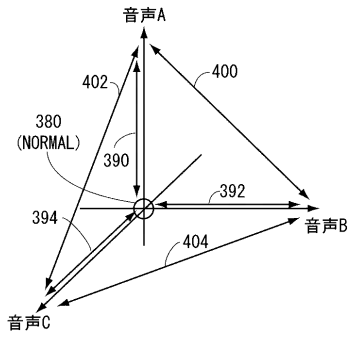
【図 1 5】



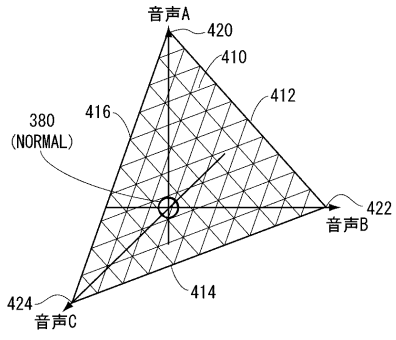
【図 1 6】



【 図 1 7 】



【 図 1 8 】



フロントページの続き

審査官 清水 正一

- (56)参考文献 特開2001-109901(JP,A)
特開2002-094881(JP,A)
特開2004-005265(JP,A)
特開2003-295882(JP,A)
特開2003-219262(JP,A)
特開平09-050295(JP,A)
特開平09-152892(JP,A)
特開2003-283613(JP,A)
特開2001-117564(JP,A)
特開平10-254500(JP,A)
特開平09-244693(JP,A)
特開平09-146597(JP,A)
特開2004-102118(JP,A)
特開2002-333897(JP,A)
坂野 秀樹, 武田 一哉, 鹿野 清宏, 板倉 文忠, 包絡と音源の独立操作による音声モーフィング, 電子情報通信学会論文誌, 日本, 社団法人電子情報通信学会, 1998年 2月, 第J81-A巻
鈴木 紀子, 米澤 朋子, 片桐 恭弘, 男性 - 女性間モーフィング音声の印象評定, 日本音響学会2004年秋季研究発表会講演論文集 - I -, 日本, 社団法人日本音響学会, 2004年 9月21日, p.409-410

(58)調査した分野(Int.Cl., DB名)

G10L 21/00 - 21/06
G10L 13/00 - 13/08