

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4543263号  
(P4543263)

(45) 発行日 平成22年9月15日(2010.9.15)

(24) 登録日 平成22年7月9日(2010.7.9)

(51) Int.Cl.		F I			
<b>G06T 13/00</b>	<b>(2006.01)</b>	G06T 13/00		B	
<b>G10L 13/00</b>	<b>(2006.01)</b>	G10L 13/00	I00V		
<b>G06T 15/70</b>	<b>(2006.01)</b>	G06T 15/70		B	

請求項の数 7 (全 34 頁)

(21) 出願番号	特願2006-230543 (P2006-230543)	(73) 特許権者	393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2
(22) 出願日	平成18年8月28日(2006.8.28)	(74) 代理人	100099933 弁理士 清水 敏
(65) 公開番号	特開2008-52628 (P2008-52628A)	(72) 発明者	四倉 達夫 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
(43) 公開日	平成20年3月6日(2008.3.6)	(72) 発明者	川本 真一 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
審査請求日	平成19年12月18日(2007.12.18)	(72) 発明者	中村 哲 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

最終頁に続く

(54) 【発明の名称】 アニメーションデータ作成装置及びアニメーションデータ作成プログラム

(57) 【特許請求の範囲】

【請求項1】

視覚素コーパスを記憶した第1の記憶手段を備えたコンピュータにおいて、入力される音声データに基づき、前記音声データに対応して動く口を含む顔のアニメーションデータを作成するためのアニメーションデータ作成プログラムであって、

前記視覚素コーパスは、音声付の発話時の顔の映像から作成した複数の視覚素ユニットを含み、

各視覚素ユニットは、視覚素ラベルと、当該視覚素ユニットに対応する顔の動きを示す動きデータと、当該視覚素ユニットに対応する音声から得られた、当該視覚素ユニットに対応する音素の継続長を含む韻律情報とを含み、

前記プログラムは、前記音声データを、音声データにより表される音素を特定する音素データ列に変換するための第1の変換手段として前記コンピュータを機能させ、

前記音素データ列は、音素ラベルと、前記音声データ中の当該音素部分の継続長を含む韻律情報とからなる音素データを含み、

前記プログラムはさらに、前記第1の変換手段の出力する前記音素データ列中の音素データに含まれる音素ラベルの各々を、対応の視覚素ラベルに変換することにより、視覚素データ列を出力するための第2の変換手段として前記コンピュータを機能させ、

前記第2の変換手段の出力する視覚素データ列は、視覚素ラベルと、前記音声データにおける、当該視覚素データに対応する部分から得られる、少なくとも当該視覚素データに対応する音素の継続長を含む韻律情報とからなる視覚素データを含み、

前記プログラムはさらに、

前記視覚素データ列に含まれる視覚素データの各々について、前記視覚素コーパス内の視覚素ユニットの内、当該視覚素データに含まれる視覚素ラベルと同じ視覚素ラベルを持ち、かつ当該視覚素データに含まれる韻律情報と、前記視覚素コーパスに含まれる各視覚素が有する韻律情報とにより音声の類似度を評価する評価関数により、当該視覚素データに含まれる音声と最も類似した音声を持つと評価された視覚素ユニットを前記視覚素コーパスから選択するための第1の選択手段と、

前記第1の選択手段により選択された視覚素ユニットに含まれる動きデータを視覚素データ列の順序にしたがい時間軸上で連結することにより、前記入力される音声データに対応する口のアニメーションデータを作成するための連結手段として前記コンピュータを機能させる、アニメーションデータ作成プログラム。

10

【請求項2】

前記視覚素コーパスに含まれる視覚素ユニットの各々に含まれる音声の韻律情報は、当該視覚素ユニットに対応する音声の継続長に加えて当該継続期間中の音声の平均パワーを含み、

前記第1の変換手段は、前記音声データを、音素データ列に変換するための手段を含み、前記音素データ列は、音素ラベルと、前記音声データ中の当該音素部分の継続長及び平均パワーとからなる音素データを含み、

前記第1の選択手段は、

前記視覚素データ列に含まれる視覚素データの各々について、前記視覚素コーパス内の視覚素ユニットの内、当該視覚素データに含まれる視覚素ラベルと同じ視覚素ラベルを持つ視覚素ユニットの各々について、当該視覚素データに含まれる継続長及び平均パワーと、当該視覚素ユニットが有する継続長及び平均パワーとにより音声の類似度を評価する評価関数の値を評価するための評価手段と、

20

前記評価手段により、当該視覚素データに含まれる音声と最も類似した音声を持つと評価された視覚素ユニットを前記視覚素コーパスから選択するための第2の選択手段とを含む、請求項1に記載のアニメーションデータ作成プログラム。

【請求項3】

前記コンピュータは、音素ラベルと、視覚素ラベルとの対応関係を記憶した音素 - 視覚素変換テーブルを記憶するための第2の記憶手段をさらに含み、

30

前記第2の変換手段は、前記第1の変換手段の出力する前記音素データ列の音素データに含まれる音素ラベルの各々を、前記音素 - 視覚素変換テーブルを参照することによって対応の視覚素ラベルに変換して、視覚素データ列を出力するための手段を含む、請求項1又は請求項2に記載のアニメーションデータ作成プログラム。

【請求項4】

前記視覚素コーパスの各視覚素ユニットは、前記音声付の発話時の顔の映像から前記複数の視覚素ユニットを作成した際の、前記各視覚素に先行する第1の数の視覚素ユニットの視覚素ラベル、及び前記各視覚素に後続する第2の数の視覚素ユニットの視覚素ラベルをさらに含み、前記先行する前記第1の数の視覚素ユニットの視覚素ラベルと、前記各視覚素ユニットの視覚素ラベルと、前記後続する前記第2の数の視覚素ユニットの視覚素ラベルとは、視覚素ラベルの組を構成し、

40

前記第2の変換手段は、

前記第1の変換手段の出力する前記音素データ列中の音素データの各々に対し、当該音素データに含まれる音素ラベルと、その前の前記第1の数の音素データに含まれる音素ラベルと、その後の前記第2の数の音素データに含まれる音素ラベルとの各々を、対応の視覚素ラベルに変換し、音素データの順番に組合せて視覚素ラベルの組を作成するための手段と、

前記第1の変換手段の出力する前記音素データ列中の音素データの各々に対し、前記第1の変換手段の出力する前記音素データ列中の音素データに含まれる音素ラベルを、前記視覚素ラベルの組を作成するための手段により得られた視覚素ラベルの組で置換すること

50

により、前記視覚素ラベルデータを作成し、出力するための手段とを含み、

前記第1の選択手段は、前記視覚素データ列に含まれる視覚素データの各々について、前記視覚素コーパス内にある、処理対象の視覚素データに含まれる視覚素ラベルの組と同じ視覚素ラベルの組を持ち、かつ当該処理対象の視覚素データに含まれる韻律情報と、前記視覚素コーパスに含まれる各視覚素ユニットが有する韻律情報とにより音声の類似度を評価する評価関数により、当該視覚素データに含まれる音声と最も類似した音声を持つと評価された視覚素ユニットを前記視覚素コーパスから選択するための第2の選択手段を含む、請求項1に記載のアニメーションデータ作成プログラム。

【請求項5】

前記第2の選択手段は、

前記視覚素データ列に含まれる視覚素データの各々について、前記視覚素コーパス内に、当該処理対象の視覚素データに含まれる視覚素ラベルの組と同じ視覚素ラベルの組を持つ視覚素ユニットが存在するか否かを判定するための判定手段と、

前記判定手段により、前記視覚素コーパス内に、当該処理対象の視覚素データに含まれる視覚素ラベルの組と同じ視覚素ラベルの組を持つ視覚素ユニットが存在すると判定されたことに応答して、それら視覚素ユニットの各々に関し、当該視覚素データに含まれる韻律情報と、前記視覚素コーパスに含まれる各視覚素ユニットが有する韻律情報とにより音声の類似度を評価する評価関数の値を算出するための第1の算出手段と、

前記判定手段により、前記視覚素コーパス内に、当該視覚素データに含まれる視覚素ラベルの組と同じ視覚素ラベルの組を持つ視覚素ユニットが存在しないと判定されたことに  
20 応答して、処理対象の視覚素データの視覚素ラベルの組のうち、処理対象の視覚素データの視覚素ラベルを含む一部からなる部分的視覚素ラベルのみを基準として、前記視覚素コーパス内から、当該一部と位置及び内容が一致する視覚素ラベルの組を持つ視覚素ユニットを選択するための手段と、

前記選択するための手段により選択された視覚素ユニットの各々について、処理対象の視覚素データに含まれる韻律情報との間で前記評価関数の値を算出するための第2の算出手段と、

前記第1の算出手段又は前記第2の算出手段により算出された評価関数の値が最も小さな視覚素ユニットを選択するための手段とを含む、請求項4に記載のアニメーションデータ作成プログラム。  
30

【請求項6】

前記連結手段は、前記選択手段により選択された視覚素ユニットに含まれる動きデータのうち、時間軸上で連続する二つの視覚素ユニットの動きデータについて、先行する視覚素ユニットの動きデータの最後の一部分の動きデータと、後続する視覚素ユニットの先頭の一部分の動きデータとの各々を、時間に応じた重み付けをして加算することにより、視覚素ユニットの動きデータを時間軸上で連結するための加重加算手段を含む、請求項1又は請求項2に記載のアニメーションデータ作成プログラム。

【請求項7】

複数の三つ組視覚素ユニットを含む視覚素コーパスを用い、入力される音声データに対応する顔の動きを示すアニメーションデータを作成するためのアニメーションデータ作成装置であって、  
40

前記三つ組視覚素ユニットの各々は、三つ組視覚素ラベルと、当該三つ組視覚素ユニットに対応する視覚素の継続時間と、当該視覚素を収録したときに発話されていた音声の平均パワーと、当該視覚素を収録したときの発話者の顔の特徴点の動きデータとを含み、

入力される音声データに対して音声分析を行なうことにより、音素ラベル、音素の継続長、及び当該音素の発話時の平均パワーからなる音素データ列を作成するための音素変換手段と、

音素ラベルと視覚素ラベルとの対応関係を示すテーブルを記憶するための手段と、

前記音素データ列に含まれる音素ラベルを、前記テーブルを参照して対応する視覚素ラベルに変換することにより、視覚素データ列を作成するための第1の変換手段と、  
50

前記第1の変換手段の出力する前記視覚素データ列中の視覚素データの各々について、視覚素ラベルを前後の視覚素データの視覚素ラベルと組合せた三つ組視覚素ラベルに変換し、三つ組視覚素データ列を出力するための第2の変換手段と、

前記第2の変換手段の出力する前記三つ組視覚素データ列に含まれる三つ組視覚素データの各々について、前記視覚素コーパスから、前記三つ組視覚素データの有する三つ組視覚素ラベルと一致する三つ組視覚素ラベルを持つ三つ組視覚素ユニットであって、当該三つ組視覚素ユニットの持つ継続長及びパワーと、前記三つ組視覚素データの持つ継続長及び平均パワーとの間の類似度を評価する評価関数によって前記三つ組視覚素データの継続長及び平均パワーと類似する継続長及び平均パワーを持つと評価される三つ組視覚素ユニットを選択するための選択手段と、

10

前記三つ組視覚素ユニット選択手段により選択された三つ組視覚素ユニットに含まれる顔の動きデータを、前記三つ組視覚素データの時系列にしたがって時間軸上で連結することにより、顔のアニメーションデータを作成するための手段とを含む、アニメーションデータ作成装置。

【発明の詳細な説明】

【技術分野】

【0001】

この発明はアニメーションデータ作成技術に関し、特に、予め準備した顔モデルを用い、音声から、音声と同期した顔画像のアニメーションを作成するためのアニメーションデータ作成装置及びプログラムに関する。

20

【背景技術】

【0002】

アニメーション作品の制作にコンピュータ・グラフィックス(CG)が用いられることが多くなり、従来のセルアニメーション等では制作者の高度な技能を要していたようなアニメーションが、単純な作業によって実現できるようになった。CGを用いる技術の中には例えば、3次元モデルを用いてアニメーションを制作する技術がある。この技術では、アニメーションの各フレームにおいて、オブジェクトの形状・位置・方向等を仮想空間上のポリゴンによって定義する。そしてその定義に基づきオブジェクトの画像を合成し、それら画像からアニメーションを構成する。オブジェクトの形状が一度定義されると、その形状について、あらゆる視点からの画像を何度でも合成できる。

30

【0003】

フレームごとにオブジェクトを変形させて画像化することにより、キャラクタの表情の変化等も表現できる。キャラクタの声として別途音声を用意し、キャラクタの口の形及び表情などをその音声に合わせて変化させると、あたかもキャラクタが発話しているようなアニメーションを制作できる。本明細書では、音声に合わせてキャラクタの口の形や表情を変化させることを、「リップシンク」と呼ぶ。また、本明細書では、リップシンクが実現しているアニメーションを「リップシンクアニメーション」と呼ぶ。

【0004】

リップシンクを実現するには、キャラクタの声と各フレームの画像で表現されるキャラクタの表情とを同期させなければならない。リップシンクを実現するための手法として従来から広く用いられている手法は、次の二つに分類される。一つの手法は、予め制作された映像に合わせて後から音声を録音する手法(アフターレコーディング:いわゆる「アフレコ」)である。もう一つの手法は、音声を先に録音しておき、その音声に合わせて映像を後から制作する方法(プレレコーディング:これを以下「プレレコ」と呼ぶ。)である。アフレコでは、アニメーションの制作者が、発話中のキャラクタの表情変化を予測しながら各フレームの画像を制作し、アニメーションを構成する。キャラクタの声を担当する発話者(又は声優)は、アニメーション上でのキャラクタの表情を見ながらタイミングを調整してセリフを発話する。これに対しプレレコでは、発話者は自由にセリフを発話する。制作者は、その音声に合わせて表情を調整しながら、各フレームの画像を制作する。

40

【0005】

50

CGを用いてリップシンクアニメーションを生成するための技術として、後掲の非特許文献1では、発話時の音声を録音することにより得られる収録音声データと、当該収録音声データの収録時に同時に収録される発話者の顔の複数個の特徴点に関するモーションキャプチャデータとからなるデータセットから、リップシンクアニメーション作成用の統計確率モデルを作成するための統計確率モデル作成装置が開示されている。この統計確率モデルは、入力される音素ラベル列又は視覚素列ラベルに対する、各特徴点の位置の確率を与えるモデルである。

【0006】

なお、本明細書では、「視覚素」とは、音素と同様、顔（主として口）の基本的な形状のことをいう。視覚素は複数個存在するが、それらは視覚素を識別する名称により区別される。本明細書では視覚素の名称を視覚素ラベルと呼ぶ。

10

【0007】

この統計確率モデルを用い、入力音声から得られた音素ラベル列又は視覚素ラベル列に対して最も尤度が高くなるような特徴点の位置データの系列を推定することができる。推定された特徴点の位置データの系列により、入力音声と同期した顔モデルの特徴点の軌跡、すなわち顔画像のアニメーションのフレームごとのワイヤフレームモデルが得られる。各フレームにおけるワイヤフレームモデルに対するレンダリングによってアニメーション画像を得ることができる。

【0008】

非特許文献1の開示によると、統計確率モデルの学習の際に、顔の特徴点の位置データだけでなく、その速度及び加速度までモデル学習用のパラメータに加えることにより、位置データのみを用いた場合と比較してより自然な動きをする顔アニメーションを得ることができる。

20

【非特許文献1】T. ヨツクラ他、「動的特徴を用いたHMMからのリップシンクアニメーション」、ACM SIGGRAPH 2006 予稿集CD、2006年7月30日(T. Yotsukura et al., "Lip-sync Animation from HMM Using Dynamic Features", ACM SIGGRAPH 2006, 30 July 2006, Boston, Massachusetts)

【発明の開示】

【発明が解決しようとする課題】

【0009】

上記した非特許文献1による手法は、位置データという静的データのみを用いた場合と比較してよりスムーズで自然な動きを持つ顔アニメーションを作成するために有効である。しかし、モデルの学習に特徴点の動的データを用いるために、モデル学習時のパラメータ数が静的データのみを用いる場合と比較して3倍になる。そのため、モデル学習に時間を要するという問題がある。特に、より精密なアニメーションを作成するために特徴点の数を増加させたりすると、モデル学習の時間がそれだけ増加してしまう。

30

【0010】

また、HMM（隠れマルコフモデル）による統計的処理により顔の特徴点の位置データを推定するため、実際の顔の特徴点の動きと比較すると、推定された位置データには、わずかではあるがずれが生ずるといった問題点がある。

40

【0011】

そのため、アニメーション作成のための準備がより短時間で可能で、しかも動きが自然で実際の顔の動きをよく反映したリップシンクアニメーションを作成できる技術が望まれている。

【0012】

それ故に本発明の目的は、アニメーション作成のための準備が短時間で可能で、実際の顔の動きをよく反映した自然な動きを実現できるリップシンクアニメーションを作成可能なアニメーションデータ作成装置及びそのためのプログラムを提供することである。

【課題を解決するための手段】

【0013】

50

本発明の第1の実施の形態に係るアニメーションデータ作成プログラムは、視覚素コーパスを記憶した第1の記憶手段を備えたコンピュータにおいて、入力される音声データに基づき、音声データに対応して動く口を含む顔のアニメーションデータを作成するためのアニメーションデータ作成プログラムである。視覚素コーパスは、音声付の発話時の顔の映像から作成した複数の視覚素ユニットを含む。各視覚素ユニットは、視覚素ラベルと、当該視覚素ユニットに対応する顔の動きを示す動きデータと、当該視覚素ユニットに対応する音声から得られた、当該視覚素ユニットに対応する音素の継続長を含む韻律情報とを含む。このプログラムは、音声データを、音声データにより表される音素を特定する音素データ列に変換するための第1の変換手段としてコンピュータを機能させる。音素データ列は、音素ラベルと、音声データ中の当該音素部分の継続長を含む韻律情報とからなる音素データを含む。このプログラムはさらに、第1の変換手段の出力する音素データ列中の音素データに含まれる音素ラベルの各々を、対応の視覚素ラベルに変換することにより、視覚素データ列を出力するための第2の変換手段としてコンピュータを機能させる。第2の変換手段の出力する視覚素データ列は、視覚素ラベルと、音声データ中における、当該視覚素データに対応する部分から得られる、少なくとも当該視覚素データに対応する音素の継続長を含む韻律情報とからなる視覚素データを含む。このプログラムはさらに、視覚素データ列に含まれる視覚素データの各々について、視覚素コーパス内の視覚素ユニットの内、当該視覚素データに含まれる視覚素ラベルと同じ視覚素ラベルを持ち、かつ当該視覚素データに含まれる韻律情報と、視覚素コーパスに含まれる各視覚素が有する韻律情報とにより音声の類似度を評価する評価関数により、当該視覚素データに含まれる音声と最も類似した音声を持つと評価された視覚素ユニットを視覚素コーパスから選択するための第1の選択手段と、第1の選択手段により選択された視覚素ユニットに含まれる動きデータを視覚素データ列の順序にしたがい時間軸上で連結することにより、入力される音声データに対応する口のアニメーションデータを作成するための連結手段としてコンピュータを機能させる。

#### 【0014】

予め、視覚素コーパスを第1の記憶手段に記憶させておく。視覚素コーパスは、音声付の発話時の顔の映像から作成した複数の視覚素ユニットを含む。視覚素ユニットに含まれる動きデータは、発話時の実際の顔の動きを反映している。第1の変換手段は、入力される音声データを、音素データ列に変換する。第2の変換手段は、音素データ列に含まれる音素ラベルを対応の視覚素ラベルに変換し、視覚素データ列として出力する。第1の選択手段は、評価関数により、視覚素データに含まれる音声と最も類似した音声を持つと評価された視覚素ユニットを視覚素コーパスから選択する。連結手段は、こうして選択された視覚素ユニットの動きデータを時間軸上で連結し、アニメーションデータを作成する。

#### 【0015】

アニメーションデータの作成時に使用される動きデータは、実際の顔の動きから得られたものである。したがって、それらを連結したとき、少なくとも各視覚素データに対応する部分で得られる顔アニメーションの動きは、実際の顔の動きをよく反映した自然なものとなる。視覚素コーパスの作成には、非特許文献1で挙げられたような多数のデータを用いた学習処理は必要ない。したがって、アニメーション作成のための準備が短時間で可能で、実際の顔の動きをよく反映した自然な動きを実現できるリップシンクアニメーションを作成可能なアニメーションデータ作成プログラムを提供できる。

#### 【0016】

好ましくは、視覚素コーパスに含まれる視覚素ユニットの各々に含まれる音声の韻律情報は、当該視覚素ユニットに対応する音声の継続長に加えて当該継続期間中の音声の平均パワーを含み、第1の変換手段は、音声データを、音素データ列に変換するための手段を含み、音素データ列は、音素ラベルと、音声データ中の当該音素部分の継続長及び平均パワーとからなる音素データを含み、第1の選択手段は、視覚素データ列に含まれる視覚素データの各々について、視覚素コーパス内の視覚素ユニットの内、当該視覚素データに含まれる視覚素ラベルと同じ視覚素ラベルを持つ視覚素ユニットの各々について、当該視覚

10

20

30

40

50

素データに含まれる継続長及び平均パワーと、当該視覚素ユニットが有する継続長及び平均パワーとにより音声の類似度を評価する評価関数の値を評価するための評価手段と、評価手段により、当該視覚素データに含まれる音声と最も類似した音声を持つと評価された視覚素ユニットを視覚素コーパスから選択するための第2の選択手段とを含む。

【0017】

視覚素ユニットの選択における評価に、継続長だけでなく音声の平均パワーも使用される。顔の各部の動きは、発話時の声の大きさにより影響される。したがって、このように音声の平均パワーも用いて、選択すべき視覚素ユニットを評価することにより、顔の各部の動きに大きな不連続がない視覚素ユニットを選択できる。

【0018】

その結果、アニメーション作成のための準備が短時間で可能で、実際の顔の動きをよく反映した自然で滑らかな動きを実現できるリップシンクアニメーションを作成可能なアニメーションデータ作成プログラムを提供できる。

【0019】

より好ましくは、コンピュータは、音素ラベルと、視覚素ラベルとの対応関係を記憶した音素-視覚素変換テーブルを記憶するための第2の記憶手段をさらに含む。第2の変換手段は、第1の変換手段の出力する音素データ列の音素データに含まれる音素ラベルの各々を、音素-視覚素変換テーブルを参照することによって対応の視覚素ラベルに変換して、視覚素データ列を出力するための手段を含む。

【0020】

音声データから音素データ列への変換という確立した技術を用いて音素データ列を得て、その後に音素ラベルを対応する視覚素ラベルに変換する。したがって、既存の技術を用いて効率的にシステムを構築できる。

【0021】

さらに好ましくは、第2の変換手段による変換により得られる視覚素ラベルの数は、第1の変換手段により出力される音素ラベルの数よりも少ない。

【0022】

音声と比較して、視覚素の数は少なくてもよい。そこで、このように視覚素ラベルの数を音素ラベルの数より少なくすることで、処理を安定させることができる。

【0023】

視覚素コーパスの各視覚素ユニットは、音声付の発話時の顔の映像から複数の視覚素ユニットを作成した際の、各視覚素に先行する第1の数の視覚素ユニットの視覚素ラベル、及び各視覚素に後続する第2の数の視覚素ユニットの視覚素ラベルをさらに含んでもよい。先行する第1の数の視覚素ユニットの視覚素ラベルと、各視覚素ユニットの視覚素ラベルと、後続する第2の数の視覚素ユニットの視覚素ラベルとは、視覚素ラベルの組を構成する。第2の変換手段は、第1の変換手段の出力する音素データ列中の音素データの各々に対し、当該音素データに含まれる音素ラベルと、その前の第1の数の音素データに含まれる音素ラベルと、その後の第2の数の音素データに含まれる音素ラベルとの各々を、対応の視覚素ラベルに変換し、音素データの順番に組合せて視覚素ラベルの組を作成するための手段と、第1の変換手段の出力する音素データ列中の音素データの各々に対し、第1の変換手段の出力する音素データ列中の音素データに含まれる音素ラベルを、視覚素ラベルの組を作成するための手段により得られた視覚素ラベルの組で置換することにより、視覚素ラベルデータを作成し、出力するための手段とを含む。第1の選択手段は、視覚素データ列に含まれる視覚素データの各々について、視覚素コーパス内にある、処理対象の視覚素データに含まれる視覚素ラベルの組と同じ視覚素ラベルの組を持ち、かつ当該処理対象の視覚素データに含まれる韻律情報と、視覚素コーパスに含まれる各視覚素ユニットが有する韻律情報とにより音声の類似度を評価する評価関数により、当該視覚素データに含まれる音声と最も類似した音声を持つと評価された視覚素ユニットを視覚素コーパスから選択するための第2の選択手段を含む。

【0024】

10

20

30

40

50

視覚素ラベルをこのように視覚素ラベルの組で置換することにより、視覚素ユニットに対応する発話時の前後の顔の形まで考慮した形で視覚素ユニットを選択できる。したがって、実際の顔の動きを、その前後の顔の形まで考慮した形で反映した、自然で滑らかな動きを実現できるリップシンクアニメーションを作成可能なアニメーションデータ作成プログラムを提供できる。

【0025】

第1の数は1でもよく、第2の数も1でよい。

【0026】

好ましくは、第2の選択手段は、視覚素データ列に含まれる視覚素データの各々について、視覚素コーパス内に、当該処理対象の視覚素データに含まれる視覚素ラベルの組と同じ視覚素ラベルの組を持つ視覚素ユニットが存在するか否かを判定するための判定手段と、判定手段により、視覚素コーパス内に、当該処理対象の視覚素データに含まれる視覚素ラベルの組と同じ視覚素ラベルの組を持つ視覚素ユニットが存在すると判定されたことに応答して、それら視覚素ユニットの各々に関し、当該視覚素データに含まれる韻律情報と、視覚素コーパスに含まれる各視覚素ユニットが有する韻律情報とにより音声の類似度を評価する評価関数の値を算出するための第1の算出手段と、判定手段により、視覚素コーパス内に、当該視覚素データに含まれる視覚素ラベルの組と同じ視覚素ラベルの組を持つ視覚素ユニットが存在しないと判定されたことに応答して、処理対象の視覚素データの視覚素ラベルの組のうち、処理対象の視覚素データの視覚素ラベルを含む一部分からなる部分的視覚素ラベルのみを基準として、視覚素コーパス内から、当該一部と位置及び内容が一致する視覚素ラベルの組を持つ視覚素ユニットを選択するための手段と、選択するための手段により選択された視覚素ユニットの各々について、処理対象の視覚素データに含まれる韻律情報との間で評価関数の値を算出するための第2の算出手段と、第1の算出手段又は第2の算出手段により算出された評価関数の値が最も小さな視覚素ユニットを選択するための手段とを含む。

【0027】

前後の視覚素ラベルまで含んだ視覚素ラベルの組と一致するような視覚素ラベルを持つ視覚素ユニットを視覚素コーパスから選択しようとする場合、特に視覚素コーパスに含まれる視覚素のバリエーションが十分大きくないときには、条件を満たす視覚素ユニットが存在しないこともあり得る。そこで、そうした場合には、前半のみ、又は後半のみの視覚素ラベルの組が一致するような視覚素ユニットを視覚素コーパスから選択することにより、確実に適切な視覚素ユニットを選択することができる。

【0028】

さらに好ましくは、連結手段は、選択手段により選択された視覚素ユニットに含まれる動きデータのうち、時間軸上で連続する二つの視覚素ユニットの動きデータについて、先行する視覚素ユニットの動きデータの最後の一部分の動きデータと、後続する視覚素ユニットの先頭の一部分の動きデータとの各々を、時間に応じた重み付けをして加算することにより、視覚素ユニットの動きデータを時間軸上で連結するための加重加算手段を含む。

【0029】

視覚素コーパスから選択した視覚素ユニットは、通常は互いに連続して収録されたものではない。したがって、顔の動きに多少の不連続が生じ得る。そこで、このように連続する二つの視覚素ユニットの動きデータを、その境界部分で加重加算することによって、滑らかに両者を連結することができる。その結果、アニメーション作成のための準備が短時間で可能で、実際の顔の動きをよく反映した自然で滑らかな動きを実現できるリップシンクアニメーションを作成可能なアニメーションデータ作成プログラムを提供できる。

【0030】

本発明の第2の局面にかかる記録媒体は、上記したいずれかのアニメーションデータ作成プログラムを記録した、コンピュータ読取可能な記録媒体である。

【0031】

本発明の第3の局面に係るアニメーションデータ作成装置は、複数の三つ組視覚素ユニ

10

20

30

40

50

ットを含む視覚素コーパスを用い、入力される音声データに対応する顔の動きを示すアニメーションデータを作成するためのアニメーションデータ作成装置である。三つ組視覚素ユニットの各々は、三つ組視覚素ラベルと、当該三つ組視覚素ユニットに対応する視覚素の継続時間と、当該視覚素を収録したときに発話されていた音声の平均パワーと、当該視覚素を収録したときの発話者の顔の特徴点の動きデータとを含む。アニメーションデータ作成装置は、入力される音声データに対して音声分析を行なうことにより、音素ラベル、音素の継続長、及び当該音素の発話時の平均パワーからなる音素データ列を作成するための音素変換手段と、音素ラベルと視覚素ラベルとの対応関係を示すテーブルを記憶するための手段と、音素データ列に含まれる音素ラベルを、テーブルを参照して対応する視覚素ラベルに変換することにより、視覚素データ列を作成するための第1の変換手段と、第1の変換手段の出力する視覚素データ列中の視覚素データの各々について、視覚素ラベルを前後の視覚素データの視覚素ラベルと組合せた三つ組視覚素ラベルに変換し、三つ組視覚素データ列を出力するための第2の変換手段と、第2の変換手段の出力する三つ組視覚素データ列に含まれる三つ組視覚素データの各々について、視覚素コーパスから、三つ組視覚素データの有する三つ組視覚素ラベルと一致する三つ組視覚素ラベルを持つ三つ組視覚素ユニットであって、当該三つ組視覚素ユニットの持つ継続長及びパワーと、三つ組視覚素データの持つ継続長及び平均パワーとの間の類似度を評価する評価関数によって三つ組視覚素データの継続長及び平均パワーと類似する継続長及び平均パワーを持つと評価される三つ組視覚素ユニットを選択するための選択手段と、三つ組視覚素ユニット選択手段により選択された三つ組視覚素ユニットに含まれる顔の動きデータを、三つ組視覚素データの時系列にしたがって時間軸上で連結することにより、顔のアニメーションデータを作成するための連結手段とを含む。

#### 【0032】

予め、視覚素コーパスを作成しておく。視覚素コーパスの視覚素ユニットに含まれる動きデータは、実際の顔の動きを反映している。第1の変換手段は、入力される音声データを、音素データ列に変換する。第2の変換手段は、音素データ列に含まれる音素ラベルを対応の視覚素ラベルに変換し、視覚素データ列として出力する。選択手段は、評価関数により、視覚素データに含まれる音声と最も類似した音声を持つと評価された視覚素ユニットを視覚素コーパスから選択する。連結手段は、こうして選択された視覚素ユニットの動きデータを視覚素データの時系列にしたがって時間軸上で連結し、アニメーションデータを作成する。

#### 【0033】

アニメーションデータの作成時に使用される動きデータは、実際の顔の動きから得られたものである。したがって、それらを連結したとき、各視覚素データに対応する部分で得られる顔アニメーションの動きは、実際の顔の動きをよく反映した自然なものとなる。視覚素コーパスの作成には、非特許文献1で挙げられたような多数のデータを用いた学習処理は必要ない。したがって、アニメーション作成のための準備が短時間で可能で、実際の顔の動きをよく反映した自然な動きを実現できるリップシンクアニメーションを作成可能なアニメーションデータ作成プログラムを提供できる。

#### 【発明を実施するための最良の形態】

#### 【0034】

以下、本発明の一実施の形態に係るリップシンクアニメーション作成装置について説明する。後述するように、このリップシンクアニメーション作成装置は、コンピュータハードウェアと、コンピュータハードウェアにより実行されるプログラムと、コンピュータの記憶装置に格納される音響モデルなどのデータとにより実現される。

#### 【0035】

最初に、以下の説明で使用される用語について説明する。

#### 【0036】

「視覚素」とは、英語の「viseme」の訳語である。「口形素」とも呼ばれる。視覚素は、音声における音素と同じく、顔の動きの中に存在する基本的な顔（特に口）の形

10

20

30

40

50

状を表す情報の組のことをいう。

【 0 0 3 7 】

「視覚素ラベル」とは、視覚素を識別するために各視覚素に付与される名称のことをいう。音素における「音素ラベル」と同様に使用される。

【 0 0 3 8 】

「視覚素コーパス」とは、発話しているときの発話者の顔の動きをモーションキャプチャ装置によって収録し、視覚素別に分割して保持したデータベースのことをいう。本実施の形態では、視覚素コーパスは複数の視覚素ユニットを含む。各視覚素ユニットは、顔の特徴点の位置ベクトルの時系列データと、視覚素名と、位置ベクトルの時系列データのうち、各視覚素に対応する部分の時間情報と、各視覚素に対応する部分の音声のパワーとを含んでいる。なお、本実施の形態では、視覚素コーパスに、最初に収録された音声データも付してある。これを「音声 - 視覚素コーパス」と呼ぶ。

10

【 0 0 3 9 】

「視覚素データ」とは、入力される音声から得られる、視覚素コーパス中から視覚素を選択するための基準となるデータのことをいう。本実施の形態では、視覚素データは、選択されるべき視覚素の視覚素ラベルと、その継続長と、視覚素に対応する入力音声の平均パワーとを含む。視覚素の継続長も、その視覚素に対応する入力音声の音素の継続長から得られる。

【 0 0 4 0 】

「三つ組視覚素ラベル」とは、ある視覚素の視覚素ラベルと、その視覚素の直前の視覚素の視覚素ラベルと、その視覚素の直後の視覚素の視覚素ラベルとを、時間軸上での順序にしたがって組合せたもののことをいう。本実施の形態では、視覚素コーパス中の各視覚素ユニットには、この三つ組視覚素のラベルが付されている。これらを本明細書では三つ組視覚素ユニットと呼ぶ。

20

【 0 0 4 1 】

[ 構成 ]

以下、本発明の一実施の形態に係るプログラムにより実現されるリップシンクアニメーション作成装置の機能的構成について説明する。図 1 に、このリップシンクアニメーション作成装置 4 0 のブロック図を示す。図 1 を参照して、リップシンクアニメーション作成装置 4 0 は、所定のテキストを発話しているときの発話者 5 0 の顔の特徴点の動きをその音声とともに収録し、音声 - 視覚素コーパスを作成するための収録システム 6 0 と、収録システム 6 0 により作成された音声 - 視覚素コーパスを記憶するための音声 - 視覚素コーパス記憶部 6 2 と、入力される音声データ 4 2 から、音声データ 4 2 と同期して動く、顔の特徴点の動きベクトル列をアニメーションデータとして合成するためのアニメーションデータ合成装置 4 4 と、アニメーションデータ合成装置 4 4 により合成されたアニメーションデータを記憶するためのアニメーションデータ記憶部 4 6 とを含む。

30

【 0 0 4 2 】

音声 - 視覚素コーパス記憶部 6 2 に記憶される音声 - 視覚素コーパスは、発話時の発話者 5 0 の映像から得られた三つ組視覚素ユニット列を含む。

【 0 0 4 3 】

リップシンクアニメーション作成装置 4 0 はさらに、実際のアニメーションの作成時に、アニメーションデータ記憶部 4 6 に記憶されたアニメーションデータを読み出し、予め準備されたワイヤフレームからなる、アニメーションのキャラクタの顔モデルに対してこのアニメーションデータを適用することにより、入力される音声データ 4 2 と同期して動く顔モデルの時系列データを作成し、さらに顔モデルに対し顔のテクスチャを適用してレンダリングをすることによって、所定フレーム / 秒のレートで表示されるキャラクタの顔のアニメーションを作成するためのアニメーション作成装置 4 8 と、アニメーション作成装置 4 8 により作成されたアニメーションを音声データ 4 2 とともに記憶するためのアニメーション記憶部 9 8 とを含む。

40

【 0 0 4 4 】

50

リップシンクアニメーション作成装置40はさらに、アニメーションの表示時に、アニメーション記憶部98に記憶されているアニメーションを読み出して所定フレームレートで図示しないフレームメモリに書込むためのアニメーション読出部100と、アニメーション読出部100によりフレームメモリに書込まれたアニメーションをその音声とともに再生し表示するための表示部52とを含む。

【0045】

図2に、収録システム60の構成を示す。図2を参照して、収録システム60は、発話者50による発話音声と発話時における発話者50の動画像とを収録するための録画・録音システム112と、発話時における発話者50の顔の各部位の位置及びその軌跡を計測するためのモーションキャプチャ(Motion Capture。以下「MoCap」と呼ぶ。)システム114と、録画・録音システム112により収録された音声・動画データ116及びMoCapシステム114により計測されたデータ(以下、このデータを「MoCapデータ」と呼ぶ。)118から、音声のデータ、発話時の発話者の顔の各部位の三次元の動きベクトル、視覚素ラベル、視覚素の継続長、及びその視覚素の発話時の音声の平均パワー等の系列からなるデータセット120を作成し、音声-視覚素コーパス記憶部62に音声-視覚素コーパスとして格納するためのデータセット作成装置122とを含む。なお、発話者の顔の特徴点の三次元データは、後述するように頭部の動きを除去した動きベクトルとなるように加工される。本明細書ではこの処理を正規化処理と呼び、正規化された後の顔の特徴点の三次元動きベクトル系列を顔パラメータと呼ぶ。

【0046】

録画・録音システム112は、発話者50により発せられた音声を受けて音声信号に変換するためのマイクロホン130A及び130Bと、発話者50の動画像を撮影しその映像信号とマイクロホン130A及び130Bからの音声信号とを同時に記録して音声・動画データ116を生成するためのカムコーダ132とを含む。

【0047】

カムコーダ132は、MoCapシステム114に対してタイムコード134を供給する機能を持つ。カムコーダ132は、音声信号及び映像信号を所定の形式でデータ化し、さらにタイムコード134と同じタイムコードを付与して図示しない記録媒体に記録する機能を持つ。

【0048】

本実施の形態に係るMoCapシステム114は、高再帰性光学反射マーカ(以下、単に「マーカ」と呼ぶ。)の反射光を利用して計測対象の位置を計測する光学式のシステムを含む。MoCapシステム114は、発話者50の頭部の予め定める多数の部位にそれぞれ装着されるマーカからの赤外線反射光の映像を、所定の時間間隔のフレームごとに撮影するための複数の赤外線カメラ136A, ..., 136Fと、赤外線カメラ136A, ..., 136Fからの映像信号をもとにフレームごとに各マーカの位置を計測し、カムコーダ132からのタイムコード134を付与して出力するためのデータ処理装置138とを含む。

【0049】

図3に、発話者50の頭部110に装着されるマーカの装着位置の例を模式的に示す。図3を参照して、発話者50の頭部110に近い顔、首、及び耳の多数の箇所160にそれぞれマーカが装着される。マーカの形状は半球状又は球状であり、その表面は光を再帰反射するように加工されている。マーカの大きさは直径数ミリメートル程度である。音声-視覚素コーパス62を充実したものにするには、複数日にわたり又は複数の発話者50について計測を行なうことが必要となる。そのため、マーカの装着順序を予め定めておき、装着位置として、顔器官の特徴的な位置又は装着済みのマーカとの相対的な関係によって定められる位置を予め定めておく。こうして定められる装着位置を、本明細書では「特徴点」と呼ぶ。

【0050】

顔の物理的な構造上、発話者50の顔の表面上には、頭自体の動きに追従して移動する

10

20

30

40

50

が発話者50の表情変化の影響をほとんど受けない箇所がある。例えばこめかみ160A及び160B,鼻の先端160Cがこのような特徴を持つ。本実施の形態では、このような箇所を特徴点として予め定めておく。以下、このような特徴点を不動点と呼ぶ。モーションキャプチャでは、顔の特徴点の三次元的位置が計測されるが、その位置の変動は発話者50の頭部110自体の移動による変動も含む。顔の動きを得るためには、各特徴点の位置データから、頭部の動きを差引く必要がある。この処理を正規化と呼ぶ。その詳細については後述する。不動点は正規化処理で用いられる。正規化処理のためには4点以上の不動点を定めることが望ましい。

#### 【0051】

再び図2を参照して、データ処理装置138は、各マーカの位置の計測データ(以下、「マーカデータ」と呼ぶ。)をフレームごとにまとめてMoCapデータ118を生成し、データセット作成装置122に出力する。MoCapシステム114には、市販の光学式MoCapシステムを利用できる。市販の光学式MoCapシステムにおける赤外線カメラ及びデータ処理装置の機能及び動作については周知であるので、これらについての詳細な説明はここでは繰返さない。

#### 【0052】

データセット作成装置122は、音声・動画データ116を取込んで記憶するための音声・動画記憶部140と、音声・動画記憶部140に記憶された音声・動画データ116を読み出し、三つ組視覚素データ列124を作成して出力するための三つ組視覚素データ列作成部144と、MoCapデータ118を取込んで記憶するためのMoCapデータ記憶部142と、MoCapデータ記憶部142に記憶されたMoCapデータを読み出し、MoCapデータ152を正規化して、顔の各特徴点の顔パラメータの系列126に変換するための正規化処理部146と、三つ組視覚素データ列作成部144からの三つ組視覚素データ列124及び正規化処理部146からの顔パラメータの系列126を、それらのタイムスタンプを利用して同期させて結合することによりデータセット120を生成し、音声・視覚素コーパス記憶部62に音声・視覚素コーパスとして格納させるための結合部148を含む。

#### 【0053】

正規化処理部146は、MoCapデータ152の各フレームにおいて、前述の不動点の位置変化が0になるよう、当該フレームの各マーカデータを変換することによって、当該フレームの顔パラメータを生成する機能を持つ。本実施の形態では、この変換にアフィン変換を用いる。

#### 【0054】

時刻 $t = 0$ のフレームのMoCapデータ152におけるマーカデータを同次座標系で $P = P_x, P_y, P_z, 1$ 、時刻 $t = 0$ におけるマーカデータを $P' = P'_x, P'_y, P'_z, 1$ と表すと、マーカデータ $P$ とマーカデータ $P'$ との関係は、アフィン行列 $M$ を用いて次の式(1)のように表現される。

#### 【0055】

##### 【数1】

$$P' = MP \quad M = \begin{bmatrix} M_{11} & M_{12} & M_{13} & M_{14} \\ M_{21} & M_{22} & M_{23} & M_{24} \\ M_{31} & M_{32} & M_{33} & M_{34} \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (1)$$

顔パラメータの系列126の各フレームにおいて不動点の位置データがすべて同じ値となれば、不動点の位置変化が0になり、それ以外の特徴点の位置を不動点の位置を基準として正規化できる。そこで、本実施の形態では、フレームごとに、 $t = 0$ のフレームにおける各不動点のマーカデータと、処理対象のフレームにおける当該不動点のマーカデータとから、当該フレームにおけるアフィン行列 $M$ を算出する。このアフィン行列 $M$ を用いて

、各マーカデータをアフィン変換する。変換後のマーカデータはそれぞれ、 $t = 0$ での頭の位置のまま発話を行なった状態での顔の特徴量の位置を表すものとなる。

【0056】

本実施の形態ではさらに、無表情の発話者の顔の画像から得られた各特徴点のマーカデータを、上記正規化により得られた各特徴点のマーカデータから差し引くことによって、各フレームでの特徴点の位置を動きベクトルで表す。こうすることで、顔モデルのアニメーションを作成する際には次のような処理をすればよいことになる。

【0057】

図4を参照して、アニメーションキャラクタの顔モデル170が予め準備されているものとする。この顔モデル170に対し、3つの連続するフレーム180、182及び184からなる顔画像のアニメーション172を作成するときには、顔モデル170の各特徴点のマーカデータに、上記した処理で得られた動きベクトル $V_{180}$ 、 $V_{182}$ 及び $V_{184}$ をそれぞれ加算する。この処理により、3つのフレーム180、182及び184の各々における顔モデルの各特徴点の位置が得られる。実際には、顔モデルはワイヤフレームで与えられ、特徴点の位置がワイヤフレームのノードの位置とは必ずしも一致しないので、顔モデル170のノードに、特徴点をマッピングしておく必要がある。顔モデルの変形の詳細については後述する。

【0058】

図2に示す三つ組視覚素データ列作成部144の詳細について図5を参照して説明する。図5を参照して、三つ組視覚素データ列作成部144は、音声・動画記憶部140から音声・動画データ116を読み出し、音声を音響処理のための所定のフレーム長及びフレーム間隔でフレーム化するためのフレーム化処理部200と、フレーム化処理部200により出力される各フレームの音声データから後述するビタピアライメントで使用する特徴量230を抽出するための特徴抽出部201と、発話者50(図1参照)の音声による学習によって得られた統計的音響モデルを記憶するための音響モデル記憶部202と、収録システム60による発話データの収録時の発話テキストを記憶するための発話テキスト記憶部204と、特徴抽出部201により出力される特徴量の系列から、音響モデル記憶部202に記憶された音響モデル及び発話テキスト記憶部204に記憶された発話テキストを用いたビタピアライメントにより、発話テキストに対応する各音素のラベルとその継続長とからなる音素データの系列であって尤度最大となるもの(音素データ列232)を出力するためのビタピアライメント部206と、ビタピアライメント部206により出力された音素データ列232を記憶するための音素データ列記憶部208とを含む。なお、本実施の形態では、音響モデルとしては音響HMMからなるものを用いる。

【0059】

三つ組視覚素データ列作成部144はさらに、音素ラベルと視覚素ラベルとの間の対応関係を示す音素-視覚素変換テーブルを記憶するための音素-視覚素変換テーブル記憶部210と、音素データ列記憶部208に記憶された音素データ列を読み出し、各音素データに含まれる音素ラベルを、音素-視覚素変換テーブル記憶部210に記憶された音素-視覚素変換テーブルを参照して、対応する視覚素ラベルに変換して、視覚素ラベルとその継続長とからなる視覚素データとし、視覚素データ列234を出力するための音素-視覚素変換部212と、音素-視覚素変換部212から出力される視覚素データ列234を記憶するための視覚素データ列記憶部214と、視覚素データ列記憶部214に記憶された視覚素データ列を読み出し、各視覚素データに含まれる視覚素ラベルを、その前の視覚素データの視覚素ラベル、処理対象の視覚素データの視覚素ラベル、及びその直後の視覚素データの視覚素ラベルをこの順番で組合せた三つ組視覚素ラベルに変換し、三つ組視覚素データ列236として出力するための視覚素-三つ組視覚素変換部216と、視覚素-三つ組視覚素変換部216により出力される三つ組視覚素データ列236を記憶するための三つ組視覚素データ列記憶部218とを含む。音素-視覚素変換部212は、音素ラベルを視覚素ラベルに変換した結果、同一の視覚素ラベルが連続するときには、それらをまとめて一つの視覚素データとし、その継続長も合計する。音素-視覚素変換部212はさらに、

各視覚素データに対応する音声の平均パワーも算出し、視覚素データに韻律的情報として付与する。

【 0 0 6 0 】

図 6 に、ピタビライメント部 2 0 6 が行なう処理の概略を示す。本実施の形態では、特徴抽出部 2 0 1 は、音声の各フレームから特徴量 2 3 0 として M F C C (メル周波数ケプストラム係数) を算出し、ピタビライメント部 2 0 6 に与える。ピタビライメント部 2 0 6 は、音響モデル記憶部 2 0 2 に記憶された多数の音素 H M M と、発話テキスト記憶部 2 0 4 に記憶された発話テキストとを用い、発話テキストに対応した音素列の分割として最も尤度の高くなるような分割方法にしたがって音声を音素列に分割し、各音素のラベルとその継続長とからなる音素データ列 2 3 2 を出力する。

10

【 0 0 6 1 】

図 7 に、視覚素 - 三つ組視覚素変換部 2 1 6 から出力され、三つ組視覚素データ列記憶部 2 1 8 に記憶される三つ組視覚素データ列 2 3 6 の例を示す。図 7 に示すように、三つ組視覚素データの各々は、三つ組視覚素ラベルと、ミリ秒単位の継続長と、その継続長全体での音声の平均パワーとを含む。図 7 において、三つ組視覚素の中央にある記号がその視覚素データ本来の視覚素ラベルである。その左側に記号「 - 」をはさんで付されているのがその直前の視覚素データの視覚素ラベルであり、右側に記号「 + 」をはさんで付されているのがその直後の視覚素データの視覚素ラベルである。なお、図中、「 s i l 」は無音状態に対応する視覚素ラベルを示し、「 s p 」は短いポーズに対応する視覚素ラベルを示し、A, R, Y 等はそれぞれ所定の音素に対応する視覚素ラベルを示す。

20

【 0 0 6 2 】

図 8 に、図 5 に示す音素 - 視覚素変換テーブル記憶部 2 1 0 に記憶された音素 - 視覚素変換テーブルの一例を示す。音素 - 視覚素変換テーブルの構成はこれ以外にも種々に考えられる。基本的には、音素ラベルを、その音素を発音しているときの口の形を示す視覚素ラベルに関連付けたものが音素 - 視覚素変換テーブルである。図 8 に示すように、本実施の形態では、一つの視覚素ラベルには 1 以上の音素ラベルが対応付けられている。これは、発音している音が異なっても、口の形がよく似ている場合があること、そのような場合には、異なる音に対し同じ口の形状でアニメーションを作成しても違和感を与えないこと、に基づく。

【 0 0 6 3 】

なお、図 8 において「 s i l B 」は発話の直前の無音状態を、「 s i l E 」は発話の直後の無音状態を、それぞれ表す。

30

【 0 0 6 4 】

以上から、音声 - 視覚素コーパス記憶部 6 2 に記憶される音声 - 視覚素コーパスの構成を示すと図 9 のようになる。図 9 を参照して、音声 - 視覚素コーパスは、音声波形データ 2 4 0 と、動きベクトル列 2 4 2 と、三つ組視覚素ユニット列 2 4 4 とを含む。音声波形データ 2 4 0 及び動きベクトル列 2 4 2 にはいずれもタイムコードが付されている。本実施の形態では音声波形データ 2 4 0 は使用しない。

【 0 0 6 5 】

三つ組視覚素ユニット列 2 4 4 中の各ユニットは、ユニットを識別するためのユニット I D (識別番号) と、そのユニットの三つ組視覚素ラベルと、そのユニットの視覚素の継続長と、その視覚素に対応する音声の平均パワーと、その視覚素に対応する動きベクトルの、動きベクトル列 2 4 2 における開始位置を示す時間とを含む。本実施の形態では、三つ組視覚素ユニットには動きベクトル列は含まれていないが、開始位置と、継続長とで動きベクトル列 2 4 2 を参照することにより、その視覚素ユニットに属する動きベクトル系列が動きベクトル列 2 4 2 中のどこにあるかを知ることができる。

40

【 0 0 6 6 】

再び図 1 を参照して、アニメーションデータ合成装置 4 4 は、入力される音声データ 4 2 から三つ組視覚素データ列を作成するための、図 2 に示す三つ組視覚素データ列作成部 1 4 4 と同様の機能を実現する三つ組視覚素データ列作成部 8 0 と、三つ組視覚素データ

50

列作成部 80 により作成された三つ組視覚素データ列に含まれる視覚素データの各々について、入力される音声データ 42 に同期したアニメーションを作成するために最適と評価される三つ組視覚素ユニットを音声 - 視覚素コーパス記憶部 62 の中から選択するための三つ組視覚素ユニット選択部 82 と、三つ組視覚素ユニット選択部 82 により選択された三つ組視覚素ユニットに含まれる顔の特徴点の三次元動きベクトルを時間軸に沿って互いに連結することにより、アニメーションデータを作成するための三つ組視覚素ユニット連結部 84 とを含む。

【 0067 】

図 10 に、三つ組視覚素データ列作成部 80 の構成の詳細を示す。三つ組視覚素データ列作成部 80 は、図 5 に示す三つ組視覚素データ列作成部 144 と基本的に同じ構成である。

10

【 0068 】

図 10 を参照して、三つ組視覚素データ列作成部 80 は、入力される音声データ 42 を所定フレーム長及び所定フレーム間隔のフレームによってフレーム化するためのフレーム化処理部 280 と、フレーム化処理部 280 により出力される音声データの各フレームから、MFCC を特徴量として抽出し、特徴量からなる系列を出力するための特徴量抽出部 282 と、音声データ 42 の発話者の音声により学習を行なった音響モデルを記憶するための音響モデル記憶部 284 と、入力される音声データ 42 の発話テキストを記憶するための発話テキスト記憶部 286 と、特徴量抽出部 282 により抽出された特徴量の系列に対し、音響モデル記憶部 284 に記憶された音響モデルと、発話テキスト記憶部 286 に記憶された発話テキストとを用いたビタピアライメントを行ない、発話テキストにしたがった音素の音素ラベル及びその継続長を含む音素データの系列（音素データ列）を出力するためのビタピアライメント部 288 とを含む。

20

【 0069 】

三つ組視覚素データ列作成部 80 はさらに、図 5 に示す音素 - 視覚素変換テーブル記憶部 210 に記憶されたものと同じの音素 - 視覚素変換テーブルを記憶するための音素 - 視覚素変換テーブル記憶部 290 と、ビタピアライメント部 288 により出力される音素データ列に含まれる音素データの各々の音素ラベルを、音素 - 視覚素変換テーブル記憶部 290 に記憶された音素 - 視覚素変換テーブルを参照して対応する視覚素ラベルに変換し、視覚素データ列を出力するための音素 - 視覚素変換部 292 と、音素 - 視覚素変換部 292 により出力される視覚素データ列を記憶するための視覚素データ列記憶部 293 と、視覚素データ列記憶部 293 に記憶された視覚素データ列を読み出し、視覚素データの各々に対し、その視覚素ラベルをその前後の視覚素データの視覚素ラベルと順番に結合して得られる三つ組視覚素ラベルに置換することによって、三つ組視覚素データ列を出力するための視覚素 - 三つ組視覚素変換部 294 と、視覚素 - 三つ組視覚素変換部 294 により出力される三つ組視覚素データ列を記憶するための三つ組視覚素データ列記憶部 295 とを含む。図 1 に示す三つ組視覚素ユニット選択部 82 は、三つ組視覚素データ列記憶部 295 から三つ組視覚素データ列 296 を読み出すことになる。

30

【 0070 】

図 11 に、図 1 に示す三つ組視覚素データ列作成部 80 から三つ組視覚素ユニット選択部 82 に渡される三つ組視覚素データ列 296 の構成を示す。図 11 を参照して、この三つ組視覚素データ列 296 に含まれる三つ組視覚素データの各々は、入力される音声データ 42 中の視覚素データの順序を示すシーケンス番号と、三つ組視覚素ラベルと、視覚素の継続長と、この視覚素データに対応する音声の平均パワーとを含む。

40

【 0071 】

図 12 に、三つ組視覚素ユニット選択部 82 より出力される三つ組視覚素データ列 300 の構成を示す。図 12 を参照して、この三つ組視覚素データ列 300 に含まれる三つ組視覚素データは、図 11 に示す三つ組視覚素データ列と同様の構成を持つが、アニメーションデータを生成するために最適であると三つ組視覚素ユニット選択部 82 により評価され、音声 - 視覚素コーパス記憶部 62 から選択された視覚素ユニットを識別するための選

50

扱ユニットIDをさらに含んでいる。この選択ユニットIDは、図9に示す三つ組視覚素ユニット列244の左端の「ユニットID」に相当する。このユニットIDがあれば、音声-視覚素コーパス記憶部62を参照して、三つ組視覚素ユニット列244の中の対応する三つ組視覚素ユニットの「開始時間」及び「継続長」のデータを用いて動きベクトル列242からこのユニットに属する動きベクトル系列を抽出できる。

【0072】

図13に、コンピュータを三つ組視覚素ユニット選択部82として機能させるため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13を参照して、この機能ブロックでは、ステップ310において、三つ組視覚素データ列作成部80により出力される三つ組視覚素データ列のうち、読出ポインタ位置にある三つ組視覚素データを読み 10  
読む。ステップ312では、読出ポインタ位置が、読込むべき三つ組視覚素データ列の終了位置に達したか否かを判定する。達していれば処理を終了する。終了位置に達していなければステップ314に進む。

【0073】

ステップ314では、ステップ310で読んだ三つ組視覚素データに含まれる三つ組視覚素ラベルと一致する三つ組視覚素ラベルを持つ視覚素ユニットが音声-視覚素コーパス記憶部62に記憶された音声-視覚素コーパス内に存在しているか否かを判定する。そのような視覚素ユニットが音声-視覚素コーパス内に存在していればステップ318に進み、なければステップ316に進む。

【0074】

ステップ318では、ステップ314で見つめられた三つ組視覚素ユニットを全て音声-視覚素コーパスから読み出す。この後ステップ320に進む。 20

【0075】

一方、ステップ316では、ステップ310で読んだ三つ組視覚素データに含まれる三つ組視覚素ラベルのうち、前半の二つ組視覚素ラベル、又は後半の二つ組視覚素ラベルと一致するような二つ組視覚素ラベルを三つ組視覚素ラベルの前半又は後半を持つ三つ組視覚素ユニットを音声-視覚素コーパスから全て読み出す。この後、ステップ320に進む。

【0076】

ステップ320では、ステップ316又はステップ318で読み出された三つ組視覚素ユニットの全てについて、以下の式によりコストCを計算する。 30

【0077】

【数2】

$$C = w_{TD} \cdot |TD - UD| + w_{TP} \cdot |TP - UP| \quad (2)$$

ただし、TD及びTPは、ステップ310で読んだ三つ組視覚素データに含まれる視覚素継続時間及び平均パワーであり、UD及びUPは、コスト計算の対象となっている三つ組視覚素ユニットに含まれる視覚素の継続時間及び平均パワーであり、 $w_{TD}$ 及び $w_{TP}$ はそれぞれ継続時間の差及び平均パワーの差に対して割当てられる重みである。重み $w_{TD}$ 及び $w_{TP}$ は、話者の相違などの条件によって異なるため、主観的テストによって決定する必要があるが、例えば $w_{TD} = w_{TP}$ としてもよい。また、この状態から $w_{TD}$ 及び $w_{TP}$ の値を少しずつ変えることにより、これらの値の好ましい組合せを徐々に求めるようにしてもよい。 40

【0078】

ここで算出するコストCは、処理中の三つ組視覚素データに対する顔の特徴点の動きベクトルを与える三つ組視覚素ユニットとして最適なものを、視覚素に対応する音声の韻律的特徴を用いて評価するための評価関数である。上の式から分かるように、本実施の形態では、コスト関数として、視覚素の継続長（これは視覚素に対応する音声の継続長に等しい。）の差の絶対値と、視覚素に対応する音声の平均パワーの差の絶対値の線形和を用いる。この他にもコスト関数としては種々のものが考えられる。三つ組視覚素データ及び三 50

つ組視覚素ユニットの構成を定める際には、コスト関数としてどのような情報を用いるかを検討し、必要なデータを保存するようにしなければならない。

【 0 0 7 9 】

ステップ 3 2 2 では、ステップ 3 2 0 で計算されたコストの最小値を求め、最小値を与えた三つ組視覚素ユニットを選択する。この三つ組視覚素ユニットが、処理中の三つ組視覚素データに対する最適な動きベクトル列を与えるものとして選択される。この後、制御はステップ 3 1 0 に戻り、次の三つ組視覚素データに対する処理を実行する。

【 0 0 8 0 】

図 1 3 に示すフローチャートに対応する制御構造を有するコンピュータプログラムにより、図 1 に示す三つ組視覚素ユニット選択部 8 2 を実現することができる。

10

【 0 0 8 1 】

次に、図 1 に示す三つ組視覚素ユニット連結部 8 4 の機能について説明する。以上述べたように、図 1 に示す三つ組視覚素ユニット選択部 8 2 により、入力される音声データ 4 2 に対応する三つ組視覚素ユニット列が選択される。これら三つ組視覚素ユニット列をそのまま時間軸上で連結すると、ユニットとユニットとの間で各特徴点の位置のずれが生じたり、ユニットとユニットとの間で時間軸上でのギャップ又は重複が生じたりするために、画像が不自然なものになってしまう。三つ組視覚素ユニット連結部 8 4 は、そのような特徴点の位置のずれを解消させながら三つ組視覚素ユニットを時間軸上で連結する機能を持つ。

【 0 0 8 2 】

20

図 1 7 に、三つ組視覚素ユニット連結部 8 4 による動きベクトルの連結方法を示す。図 1 7 ( A ) を参照して、ある三つ組視覚素ユニットにおけるある特徴点 M 1 の軌跡 4 3 0 と、後続する三つ組視覚素ユニットにおける対応する特徴点 M 1 ' の軌跡 4 3 2 とを、その両端で時間 T だけ重複させる。そして、この時間におけるこの特徴点の軌跡を、以下の式に従い平滑化して算出し、なめらかな軌跡 4 4 0 を生成する。

【 0 0 8 3 】

【数 3】

$$M = M_1 \cdot (1 - f(t)) + M_1' \cdot f(t) \quad (3)$$

$$(0 \leq t \leq T, \quad f(t) = t \cdot 1.0 / T)$$

30

ただし、M は平滑化後の特徴点の動きベクトル、M<sub>1</sub> 及び M<sub>1</sub>' はそれぞれ平滑化前の、先行及び後続する三つ組視覚素ユニットの特徴点の動きベクトル、t は重複区間 T の先頭からの経過時間を示す。三つ組視覚素ユニット連結部 8 4 は、これ以外の区間では、その三つ組視覚素ユニットの動きベクトル列をそのまま出力する。

【 0 0 8 4 】

なお、このような連結を行なうと、各三つ組視覚素ユニットの継続長は実質的に T だけ短縮されることになるので、それを防ぐため、各三つ組視覚素ユニットの一端（例えば後端）を T だけ延長する。図 9 に示すような音声 - 視覚素コーパス記憶部 6 2 の構造を採用することにより、そのような三つ組視覚素ユニットの延長は簡単に行なえる。

【 0 0 8 5 】

40

三つ組視覚素ユニット連結部 8 4 は、このような連結を、三つ組視覚素ユニット選択部 8 2 から出力される三つ組視覚素ユニット列内の連結部の全てについて、全ての特徴点に対して行なう。その結果、所定の周期ごとに、顔の特徴点の全てについての動きベクトルを持ったデータ系列が得られる。このデータを本明細書ではアニメーションデータと呼ぶ。アニメーションデータはアニメーションデータ記憶部 4 6 により格納される。

【 0 0 8 6 】

アニメーション作成装置 4 8 は、アニメーションデータ記憶部 4 6 に記憶されたアニメーションデータと、予め準備された、アニメーションキャラクタの顔モデルとからアニメーションを作成する機能を持つ。

【 0 0 8 7 】

50

図1を参照して、アニメーション作成装置48は、アニメーションキャラクタの顔モデルを記憶するための顔モデル記憶部90を含む。本実施の形態では、顔モデル記憶部90に記憶された顔モデルは、多数の多角形(ポリゴン)によって、静止状態における所定の顔の形状を表現した形状モデルを利用する。この顔モデルに基づき、アニメーションデータ記憶部46に格納されたアニメーションデータを利用してアニメーションを作成するためには、この顔モデルと、アニメーションデータに対応する特徴点の位置との対応付け(マッピング)を予め行なっておく必要がある。本実施の形態では、顔モデルに手作業でこの特徴点の位置をマッピングするものとし、顔モデル上の特徴点の位置を「仮想マーカ」と呼ぶマーカにより示すものとする。

#### 【0088】

図14に、顔モデル330及び仮想マーカの一例を示す。図14を参照して、顔モデル330を構成するポリゴンの辺(図14の三角形の辺を構成する黒い線)をエッジ、エッジ同士の交点を顔モデル330におけるノードと呼ぶ。図14には、仮想ノードのマッピング例を、記号と+マークとを組合せた記号332として示してある。

#### 【0089】

顔には、目・口・鼻の穴のように顔面を構成しない切れ目がある。一般に、これらの切れ目は、顔モデル330の一部としてはモデリングされない。すなわち、切れ目にはポリゴンを定義しないか、切れ目は、顔モデル330とは別のオブジェクトとして定義される。したがって、切れ目と顔面との間は境界エッジで仕切られる。境界エッジとは、二つのポリゴンによって共有されていないようなエッジのことを言う。

#### 【0090】

再び図1を参照して、アニメーション作成装置48はさらに、アニメーションデータ記憶部46に格納されたアニメーションデータのうち、アニメーションの各フレームに相当する時刻のデータを読み出し、顔モデル記憶部90に格納された顔モデルを、読み出されたアニメーションデータ内の動きベクトルにしたがって変形させて出力するための顔モデル変形部92と、顔モデルに対するレンダリングによりキャラクタのアニメーションを作成するための、顔のテクスチャデータ、照明位置、カメラ位置などの設定を記憶するためのレンダリングデータ記憶部94と、顔モデル変形部92により出力される各フレームの顔モデルに対し、レンダリングデータ記憶部94に記憶されたレンダリングのためのデータを用いてレンダリングを行ない、アニメーションのフレームごとに出力しアニメーション記憶部98に記憶させるためのレンダリング部96とを含む。

#### 【0091】

顔モデル330により表現される顔の形状は、アニメーションのキャラクタの顔の基本形状を示すものであり、ユーザにより創作される任意のものでよい。ただし、前述したとおり、動きベクトルを用いて顔モデル330に表情を付与するには、顔モデル330により表現される形状のどの部分が特徴点に対応しているかを定義する必要がある。仮想マーカ332によってそうした対応が示される。図1に示す顔モデル記憶部90には、顔モデルの各ノードの3次元位置データだけでなく、各仮想マーカの位置と、それら仮想マーカと特徴点との対応関係も記憶されている。

#### 【0092】

顔モデル変形部92は、以下のようにして顔モデル記憶部90に記憶された顔モデルの変形を行なう。基本的には顔モデル変形部92は、読み出したアニメーションデータごとに、顔モデルを構成する全てのノードに対して以下の処理(マーカラベリング処理と呼ぶ。)を行なう。すなわち、顔モデル変形部92は、顔モデルのノードの各々に対し、そのノードからの距離が最も近い仮想マーカを、仮想マーカの座標に基づき選択する。顔モデル変形部92は、選択された仮想マーカが、処理中のノードに対応付ける仮想マーカとして適切か否かを判定する。適切であれば選択マーカをこのノードに対応するマーカとして採用し、不適切であれば棄却する。このような処理を繰返し、顔モデルの一つのノードに対し所定数n(例えばn=3)の仮想マーカを採用する。本明細書では、あるノードに対し採用された仮想マーカを、当該ノードの「対応マーカ」と呼ぶ。

10

20

30

40

50

## 【 0 0 9 3 】

なお、本実施の形態では、選択マーカの対応マーカとしての適/不適を判断する際の基準に、顔モデルの境界エッジを利用する。

## 【 0 0 9 4 】

このマーカラベリング処理により、顔モデルの各ノードに対応マーカが関係付けられると、アニメーションデータから得られる、対応マーカに対応する特徴点の動きベクトルの値の内挿により、そのノードの三次元位置座標が計算される。この計算方法については後述する。

## 【 0 0 9 5 】

図 1 5 に、顔モデル変形部 9 2 により実行されるマーカラベリング処理のプログラムの制御構造をフローチャートで示す。図 1 5 を参照して、マーカラベリング処理では、ステップ 3 4 0 A とステップ 3 4 0 B とで囲まれた、ステップ 3 4 2 からステップ 3 5 4 までの処理を、顔モデル 3 3 0 の各ノードに対して実行する。

10

## 【 0 0 9 6 】

ステップ 3 4 2 では、処理対象のノードから仮想マーカまでの距離をそれぞれ算出する。さらに仮想マーカをこの距離の昇順でソートしたものをリストにする。

## 【 0 0 9 7 】

ステップ 3 4 4 では、以下の繰返しを制御するための変数  $i$  及び対応マーカとして採用したマーカの数を表す変数  $j$  に 0 を代入する。ステップ 3 4 6 では、変数  $i$  に 1 を加算する。

20

## 【 0 0 9 8 】

ステップ 3 4 7 では、変数  $i$  の値が仮想マーカの数  $M_{max}$  を超えているか否かを判定する。変数  $i$  の値が仮想マーカの数  $M_{max}$  を超えていればエラーとし、処理を終了する。これは、全ての仮想マーカを調べても、対応マーカとして採用されたものが 3 つに満たなかった場合に生ずる。普通このようなことはないが、念のためにこのようなエラー処理を設けておく。変数  $i$  の値が仮想マーカの数  $M_{max}$  以下であれば制御はステップ 3 4 8 に進む。

## 【 0 0 9 9 】

ステップ 3 4 8 では、リストの先頭から変数  $i$  で示される位置に存在する仮想マーカ（以下これを「マーカ ( $i$ ) 」と呼ぶ。）と処理対象のノードとを結ぶ線分が、顔モデル 3 3 0 におけるいずれの境界エッジも横切らない、という制約条件を充足しているか否かを判定する。当該線分が境界エッジのいずれかを横切るものであれば、ステップ 3 4 4 に戻る。さもなければステップ 3 5 0 に進む。

30

## 【 0 1 0 0 】

ステップ 3 5 0 では、この時点でのマーカ ( $i$ ) を処理対象のノードの対応マーカの一つに指定する。そしてマーカ ( $i$ ) を示す情報を、処理対象のノードのマーカ・ノード対応情報として保存する。この後制御はステップ 3 5 2 に進む。ステップ 3 5 2 では、変数  $j$  に 1 を加算する。ステップ 3 5 4 では、変数  $j$  の値が 3 となっているか否かを判定する。変数  $j$  の値が 3 であればステップ 3 4 0 B に進む。さもなければステップ 3 4 4 に進む。

40

## 【 0 1 0 1 】

上記したように、処理対象のノードと仮想マーカとを結ぶ線分が顔モデルの境界エッジを横切るものは、処理対象のノードに対応する仮想マーカから除外される。これは以下の理由による。例えば上唇と下唇とのように、間に境界エッジが存在する場合を考える。この場合、実際の顔では、上唇に位置するノードと、下唇に位置するノードとに相当する位置は互いに異なる動きをする。したがって、例えば上唇のノードの移動量を算出する際に、下唇に存在するマーカの移動量を用いることは適当ではない。線分がある境界エッジを横切っているか否かは、例えば、その境界エッジが顔モデルを構成するポリゴンのうち二つによって共有されているか、一つのみに属しているかによって判定する。

## 【 0 1 0 2 】

50

図16に、顔モデル330における唇周辺のポリゴンと仮想マーカとを示す。以下、図16を参照して、あるノードの対応マーカを特定する方法について具体例を用いて説明する。図16を参照して、顔モデル330(図14参照)の唇周辺には、多数の三角形ポリゴンが存在する。各ポリゴンは、三つのエッジに囲まれている。上唇と下唇の間には境界エッジ400が存在する。境界エッジ400は、顔モデル330と切れ目との境界、又は顔モデル330の外縁にあたる。そのため、境界エッジ以外のエッジは二つのポリゴンに共有されるが、境界エッジ400に該当するエッジは共有されない。

#### 【0103】

既に説明したように、顔モデル変形部92は、顔モデル330を構成するノードの中から処理対象のノードを一つ選択する。図16において、ノード410が処理対象のノードとして選択されたものとする。ノード410の近隣には、仮想マーカ412A, ..., 412Eが存在するものとする。顔モデル変形部92は、ノード410の座標と、仮想マーカ412A, ..., 412Eの座標とをもとに、ノード410と仮想マーカとの間の距離をそれぞれ算出する。そして、仮想マーカの中から、ノード410に最も近い位置にある仮想マーカ412Aを選択する。

#### 【0104】

続いて、顔モデル変形部92は、ノード410と仮想マーカ412A, ..., 412Eとを結ぶ線分414A, ..., 414Eが境界エッジ400を横切るか否かを検査する。すなわち、まずノード410と仮想マーカ412Aとを結ぶ線分414Aが境界エッジ400を横切るか否かを検査する。図16に示す例では、この線分414Aは、境界エッジ400を横切らない。そのため顔モデル変形部92は、仮想マーカ412Aをノード410の対応マーカの一つとする。そして、仮想マーカの中から、仮想マーカ412Aの次にノード410に近い位置にある仮想マーカ412Bを選択し検査を行なう。ノード410と仮想マーカ412Bとを結ぶ線分414Bは、境界エッジ400を横切っている。そのため、仮想マーカ412Bはノード410の対応マーカからは除外される。

#### 【0105】

顔モデル変形部92は、以上のような動作を所定数(3個)の対応マーカが選択されるまで繰返し、ノード410の対応マーカ(図16に示す例では仮想マーカ412A、412D、及び412E)を選択する。

#### 【0106】

再び図14を参照して、顔モデル変形部92は、顔モデル記憶部90に記憶された、特徴点と仮想マーカとの対応関係に基づき、あるフレームの三つ組視覚素ユニットにおける各特徴点の動きベクトルをそれぞれ対応の仮想マーカ332に付与する。さらに顔モデル変形部92は、顔モデル330の各ノードに、対応する仮想マーカ332の動きベクトルにより示される変化量から所定の内挿式により算出される変化量ベクトル $v$ を割当てることにより、顔モデル330の変形を行なう。顔モデル変形部92は、変形後の顔モデル330を、そのフレームにおける形状モデルとして出力する。

#### 【0107】

基本となる顔モデル330のうちの、あるノードの座標ベクトルを $N$ 、基本となる顔モデル330において、当該ノードと対応関係にある $i$ 番目の仮想マーカの座標を $M_i$ ( $1 \leq i \leq 3$ )、変形後の顔モデルにおける対応するマーカの座標を $M'_i$ とすると、顔モデル変形部92は、このノードの座標の変化量ベクトル $v$ を次の内挿式によって算出する。なお、 $M'_i - M_i$ が特徴点の動きベクトルに相当する。

#### 【0108】

##### 【数4】

$$v = \sum_i^n (M'_i - M_i) \cdot \left( \frac{1.0}{\|N - M_i\|} \right) \cdot \left( \sum_i^m \frac{1.0}{\|N - M_i\|} \right)^{-1} \quad (4)$$

レンダリング部96は、ポリゴンにより表された形状モデルに対するレンダリングを行

なうことができるものであればよく、市販のレンダリングエンジンを用いることもできる。アニメーションのフレームレートにしたがい、1フレームごとの顔モデルを上記式にしたがって生成し、レンダリングを行なうことにより、このレンダリングによりえられた画像のシーケンスとしてアニメーションが得られる。アニメーションはアニメーション記憶部98に記憶される。

【0109】

アニメーション読出部100は、アニメーション記憶部98から1フレームごとにアニメーションを読出して画像化し、フレーム間隔ごとに表示部52のフレームメモリに書き込む機能を持つ。

【0110】

[動作]

本実施の形態に係る顔アニメーションの作成システム40は以下のように動作する。リップシンクアニメーション作成装置40の動作は大きく四つのフェーズに分けることができる。第1のフェーズは音声・視覚素コーパス記憶部62を作成するフェーズである。第2のフェーズは音声・視覚素コーパス記憶部62を用いて入力される音声データ42からアニメーションデータ記憶部46を作成するフェーズである。第3のフェーズは、顔モデルを用い、アニメーションデータ記憶部46からアニメーションを作成しアニメーション記憶部98に格納するフェーズである。最後のフェーズは、アニメーション記憶部98に記憶されたアニメーションを表示部52に表示するフェーズである。以下、各フェーズにおけるリップシンクアニメーション作成装置40の動作について説明する。

【0111】

第1のフェーズ：音声・視覚素コーパス記憶部62の作成

以下に、収録システム60が収録を行ない、音声・視覚素コーパス記憶部62を生成する動作について説明する。図2及び図3を参照して、発話者50の頭部110の各特徴点160には、マーカを予め装着しておく。その状態で、発話者は発話を行なう。音声・視覚素コーパスを充実したものにするために、又は、各音素がバランスよく含まれるようにするために、発話の内容を事前に決めておき、発話者50にその内容で発話を行なってもらふ。この発話の内容は、図5に示す発話テキスト記憶部204に記憶される。

【0112】

収録が開始され、発話者50が発話すると、録画・録音システム112が、発話時の音声と顔の動画を収録し、音声・動画データ116を生成する。音声・動画データ116は音声・動画記憶部140に記憶される。この際、カムコード132は、MoCapシステム114に対してタイムコード134を供給するとともに、音声・動画データ116に、タイムコード134と同じタイムコードを付与する。

【0113】

同時に、発話時における特徴点160の位置が、MoCapシステム114により次のようにして三次元データとして計測される。マーカはそれぞれ、対応する特徴点の動きに追従して移動する。赤外線カメラ136A, ..., 136Fはそれぞれ、マーカによる赤外線反射光を、所定のフレームレート(例えば毎秒120フレーム)で撮影しその映像信号をデータ処理装置138に出力する。データ処理装置138は、それらの映像信号の各フレームにタイムコード134を付与し、当該映像信号をもとに、各マーカの三次元座標をフレームごとに算出する。データ処理装置138は、各マーカの三次元座標をフレームごとにまとめてMoCapデータ118として蓄積する。

【0114】

以上の収録プロセスにより収録された音声・動画データ116及びMoCapデータ118は、データセット作成装置122に与えられる。データセット作成装置122は、音声・動画データ116を音声・動画記憶部140に蓄積し、MoCapデータ118を、MoCapデータ記憶部142に蓄積する。

【0115】

正規化処理部146は、MoCapデータ記憶部142から、 $t = 0$ のフレームにおけ

10

20

30

40

50

る M o C a p データを読み出す。このときの不動点の M o C a p データが後の正規化処理の基準となる。正規化処理部 1 4 6 はさらに、各フレームでの各特徴点の座標を、不動点として指定された複数の特徴点の三次元座標を用いて以下の様に正規化する。

【 0 1 1 6 】

すなわち、正規化処理部 1 4 6 は、M o C a p データ 1 5 2 の各フレームにおいて、当該フレームの不動点の三次元座標と、 $t = 0$  のフレームにおける不動点の三次元座標とから、前述の式 ( 1 ) のアフィン行列を求め、当該アフィン行列を用いて、各特徴点の三次元座標をアフィン変換する。この変換により、変換後の特徴点の三次元座標はそれぞれ、 $t = 0$  での位置に頭を固定して発話を行なった状態での顔の特徴点の位置を表すものとなる。すなわち、各特徴点の三次元座標が正規化される。これら座標から、 $t = 0$  のときの各特徴点の座標から減算することで、その特徴点のその時点での動きベクトルが得られる。その結果、M o C a p データ 1 5 2 から顔パラメータの系列 1 2 6 が得られる。顔パラメータの系列 1 2 6 は、結合部 1 4 8 に与えられる。

10

【 0 1 1 7 】

図 5 を参照して、三つ組視覚素データ列作成部 1 4 4 のフレーム化処理部 2 0 0 は、音声・動画記憶部 1 4 0 に記憶された音声・動画データ 1 1 6 の音声データを所定フレーム長及び所定フレーム間隔でフレーム化し、特徴抽出部 2 0 1 に与える。

【 0 1 1 8 】

特徴抽出部 2 0 1 は、フレーム化処理部 2 0 0 から与えられた各フレームから、ピタピアライメント部 2 0 6 の処理で使用される音響特徴量 ( M F C C ) を算出し、特徴量 2 3 0 としてピタピアライメント部 2 0 6 に与える。このとき、各フレームの音声データもピタピアライメント部 2 0 6 に与えられる。

20

【 0 1 1 9 】

ピタピアライメント部 2 0 6 は、音響モデル記憶部 2 0 2 に記憶された音響モデルと、発話テキスト記憶部 2 0 4 に記憶された発話テキストとを用いて、特徴量 2 3 0 の系列に対するピタピアライメントを行ない、アライメントの結果得られた音素のラベル ( 音素ラベル ) 列を、各音素の継続長とともに音素データ列 2 3 2 として音素データ列記憶部 2 0 8 に格納させる。このとき、音素データ列 2 3 2 には各フレームの音声データも付される。

【 0 1 2 0 】

音素 - 視覚素変換部 2 1 2 は、音素データ列記憶部 2 0 8 から音素データを順次読み出し、各音素データに含まれる音素ラベルを、音素 - 視覚素変換テーブル記憶部 2 1 0 に記憶された音素 - 視覚素変換テーブルを参照して視覚素ラベルに変換する。音素ラベルにかえて視覚素ラベルを格納した音素データは視覚素データを構成する。音素 - 視覚素変換部 2 1 2 はこのとき、同一の視覚素ラベルが連続しているときにはそれらを一つの視覚素データにまとめ、その継続長も合計する。さらに音声の平均パワーを算出し、各視覚素データに付与する。音素 - 視覚素変換部 2 1 2 は、こうして得られた視覚素データ列 2 3 4 を、フレーム化された音声データとともに視覚素データ列記憶部 2 1 4 に格納させる。

30

【 0 1 2 1 】

視覚素 - 三つ組視覚素変換部 2 1 6 は、視覚素データ列記憶部 2 1 4 から視覚素データを順次読み出し、以下のような処理を行なう。すなわち、視覚素 - 三つ組視覚素変換部 2 1 6 は、各視覚素データの視覚素ラベルを、その直前の視覚素データに含まれる視覚素ラベルと、当該視覚素データの視覚素ラベルと、その直後の視覚素ラベルとをこの順で結合した三つ組視覚素ラベルに変換する。このようにして視覚素ラベルに代えて三つ組視覚素ラベルを格納した視覚素データは、三つ組視覚素データとなる。視覚素 - 三つ組視覚素変換部 2 1 6 は、こうして得られた三つ組視覚素データ列 2 3 6 を三つ組視覚素データ列記憶部 2 1 8 に音声データとともに格納させる。

40

【 0 1 2 2 】

結合部 1 4 8 は、三つ組視覚素データ列記憶部 2 1 8 に記憶された三つ組視覚素データ列と、正規化処理部 1 4 6 から与えられる顔パラメータの系列 1 2 6 とをそれらに付され

50

ている時間情報を用いて同期させて結合して、音声データ42とともに音声 - 視覚素コーパスを生成し、音声 - 視覚素コーパス記憶部62に格納する。

【0123】

第2のフェーズ：アニメーションデータ記憶部46の合成

第2のフェーズはアニメーションデータ合成装置44による。キャラクタの声を表す音声データ42が準備され、三つ組視覚素データ列作成部80に与えられる。この音声データ42は、事前に、キャラクタの声を担当する発話者（又は声優）によって発話されたものを録音することにより得られる。音声データ42の発話テキストは図10に示す発話テキスト記憶部286に格納される。

【0124】

図10を参照して、フレーム化处理部280は、入力される音声データ42を図5に示すフレーム化处理部200と同一のフレーム長及びフレーム間隔でフレーム化し、特徴量抽出部282に与える。

【0125】

特徴量抽出部282は、図5に示す特徴抽出部201と同様の処理により、音声の各フレームごとに、所定の音響特徴量（MFCC）を抽出し、ビタビライメント部288に与える。

【0126】

ビタビライメント部288は、音響モデル記憶部284及び発話テキスト記憶部286を用いて特徴量抽出部282に対するビタビライメントを行なって、音素ラベル及び各音素の継続長を含む音素データからなる音素データ列を音素 - 視覚素変換部292に与える。

【0127】

音素 - 視覚素変換部292は、この音素データ列に含まれる各音素データに対し、その中の音素ラベルを、音素 - 視覚素変換テーブル記憶部290に格納された音素 - 視覚素変換テーブルを参照して視覚素ラベルに変換する。音素ラベルに代えて視覚素ラベルを格納した音素データは視覚素データとなり、視覚素データ列として音素 - 視覚素変換部292から出力され視覚素データ列記憶部293に記憶される。各視覚素データは、視覚素ラベルと、元の音素の継続長とを含む。同一の視覚素ラベルが連続する場合、それらはまとめられ、継続長も合計される。また、視覚素データごとに、対応する音声の平均パワーが算出される。

【0128】

視覚素 - 三つ組視覚素変換部294は、視覚素データ列記憶部293に記憶された各視覚素データを順番に読み出し、各視覚素データに含まれる視覚素ラベルを、その直前及び直後の視覚素データの視覚素ラベルと結合することにより得られる三つ組視覚素ラベルで、視覚素データの視覚素ラベルを置換し、三つ組視覚素データ列として三つ組視覚素データ列記憶部295に記憶させる。

【0129】

図1を参照して、三つ組視覚素ユニット選択部82は、図10に示す三つ組視覚素データ列記憶部295から三つ組視覚素データ列296を読み出し、以下の処理を行なう。すなわち、三つ組視覚素ユニット選択部82は、三つ組視覚素データ列296に含まれる三つ組視覚素データごとに、音声 - 視覚素コーパス記憶部62に含まれる、同一の三つ組視覚素ラベルを持つ三つ組視覚素ユニットを探す（図13のステップ314）。そのようなユニットがあればそれら全てとの間で、その三つ組視覚素データに含まれる視覚素継続長及び平均パワーを用い、先に示した式（2）によってコスト計算を行なう（図13のステップ318）。そのようなユニットがなければ、三つ組視覚素ラベルのうちで前半の二つ組視覚素ラベル、又は後半の二つ組視覚素ラベルが一致する三つ組視覚素ユニットを音声 - 視覚素コーパス記憶部62から読み出し、それら全てとの間で、その三つ組視覚素データに含まれる視覚素の継続長及び平均パワーを用い、先に示した式（2）によってコスト計算を行なう（図13のステップ316）。

10

20

30

40

50

## 【 0 1 3 0 】

そして、このようにして計算されたコストの最小値を与える三つ組視覚素ユニットを処理対象の三つ組視覚素データに対する最適な三つ組視覚素ユニットとして選ぶ（図 1 3 のステップ 3 2 2）。

## 【 0 1 3 1 】

三つ組視覚素ユニット選択部 8 2 はこのようにして得られた三つ組視覚素ユニットからなる三つ組視覚素ユニット列を三つ組視覚素ユニット連結部 8 4 に与える。

## 【 0 1 3 2 】

三つ組視覚素ユニット連結部 8 4 は、三つ組視覚素ユニット選択部 8 2 から与えられた三つ組視覚素ユニット列中の各三つ組視覚素ユニットについて、その動きベクトル列を、先行する三つ組視覚素ユニットの動きベクトル列、及び後続する三つ組視覚素ユニットの動きベクトル列と時間軸上で連結する。なお、このとき、図 1 7 を参照して説明したように、各ユニットの動きベクトル列の最後部を時間 T だけ延長し、後続するユニットの動きベクトル列の先頭の時間 T の部分との間で、各特徴点ごとに上記した式（3）による平滑化処理を行なう。

## 【 0 1 3 3 】

以上の処理により、アニメーションデータが作成される。作成されたアニメーションデータはアニメーションデータ記憶部 4 6 に格納される。

## 【 0 1 3 4 】

第 3 のフェーズ：モデルを用いたアニメーションの作成

アニメーションの作成は、図 1 に示すアニメーション作成装置 4 8 により行なわれる。図 1 を参照して、顔モデル変形部 9 2 は、顔モデル記憶部 9 0 に記憶された顔モデルを読み出す。この顔モデルについては、音声・視覚素コーパス記憶部 6 2 を作成したときの特徴点と仮想マーカとの対応付けが既に行なわれており、さらに顔モデルを構成する各ノードに対応する仮想マーカも既に定められているものとする。

## 【 0 1 3 5 】

顔モデル変形部 9 2 は、アニメーションデータ記憶部 4 6 に記憶されているアニメーションデータのうちから、アニメーションのフレームレートにしたがった時間に最も近い時刻を持つフレームのアニメーションデータを順番に読み出し、各フレームについて以下の処理を行なう。

## 【 0 1 3 6 】

顔モデル変形部 9 2 は、顔モデルの各仮想マーカに、読み出されたアニメーションデータ内に含まれる対応する特徴点の三次元の動きベクトルを割り当てる。顔モデル変形部 9 2 はさらに、式（4）にしたがい、顔モデル記憶部 9 0 の各ノードの三次元位置座標を与える変化量ベクトル  $v$  を算出する。変化量ベクトルを全てのノードに対し算出することにより、そのフレームにおける顔モデルが完成する。この顔モデルはレンダリング部 9 6 に与えられる。

## 【 0 1 3 7 】

レンダリング部 9 6 は、顔モデル変形部 9 2 から与えられた顔モデルに対し、レンダリングデータ記憶部 9 4 に記憶されたレンダリングデータ及び設定にしたがったレンダリングを行なってアニメーションの一フレームに相当する画像を作成し、アニメーション記憶部 9 8 に格納させる。

## 【 0 1 3 8 】

顔モデル変形部 9 2 及びレンダリング部 9 6 の以上の動作を繰り返し、アニメーションデータ記憶部 4 6 に記憶されたアニメーションデータの末尾まで到達したところでアニメーション作成装置 4 8 は処理を終了させる。

## 【 0 1 3 9 】

以上の処理により、アニメーション記憶部 9 8 には、入力される音声データ 4 2 に対応した顔アニメーションを表す、所定のフレームレートでの一連の顔画像が記憶されていることになる。

10

20

30

40

50

## 【 0 1 4 0 】

第 4 のフェーズ：アニメーションの表示

アニメーションの表示はアニメーション読出部 1 0 0 及び表示部 5 2 により行なわれる。表示処理では、アニメーション読出部 1 0 0 がアニメーション記憶部 9 8 に格納された画像を先頭から順次読出し、規定の時間間隔で、表示部 5 2 内の図示されないフレームメモリに書込む。表示部 5 2 はこのフレームメモリに書き込まれた画像を所定時間間隔で読出し、画面に表示する。その結果、表示部 5 2 の画面上には、入力される音声データ 4 2 に対応して変化する、顔モデル記憶部 9 0 に記憶されたアニメーションキャラクタの顔モデルにより規定された顔のアニメーションが表示される。

## 【 0 1 4 1 】

以上のように、本実施の形態に係るリップシンクアニメーション作成装置 4 0 によれば、発話者の顔の多数の特徴点と、顔モデル 1 7 0 の各ノードとを予め対応付ける。さらに、発話時の音声から得た音素ラベルを、対応する視覚素ラベルに変換し、さらに三つ組視覚素ラベルに変換して、その継続長、その継続長中の音声の平均パワー、及びモーションキャプチャにより得た発話時の顔の特徴点の三次元動きベクトル列と組合せて、三つ組視覚素ユニットとして音声 - 視覚素コーパス記憶部 6 2 に記憶させることで、音声 - 視覚素コーパスを作成しておく。

## 【 0 1 4 2 】

音声データ 4 2 が与えられると、音声データ 4 2 から得た音素ラベル、音素の継続長、及び平均パワーからなる音素データ列を作成し、さらに各音素ラベルを音声 - 視覚素コーパスの作成時と同様の方法により三つ組視覚素ラベルに変換して三つ組視覚素データ列を得る。各三つ組視覚素ラベルの継続長及びその音声の平均パワーをもとに、所定のコストの評価関数を用い、コストが最小となる三つ組視覚素ユニットを音声 - 視覚素コーパス記憶部 6 2 から選択し、三つ組視覚素ユニット連結部 8 4 によってそれらの動きベクトル列を連結する。この連結の際に、隣接する三つ組視覚素ユニットの動きベクトル列のうち、先行する三つ組視覚素ユニットの動きベクトル列を延長し、この部分で後続する三つ組視覚素ユニットの動きベクトルとの間で平滑化処理を行なう。

## 【 0 1 4 3 】

こうした処理により、顔モデル 1 7 0 の各ノードの軌跡を表す動きベクトルが、音声 - 視覚素コーパス記憶部 6 2 に格納された実際の発話時の顔の特徴点の動きベクトルに基づいて算出される。したがって、ノードの集合としての顔モデルの時間的变化が、実際の動きに近い自然な動きを表すアニメーションデータとして得られる。アニメーションデータを構成する各特徴点の動きベクトルと、顔モデル 1 7 0 の各ノードとの対応関係とに基づいて顔モデル 1 7 0 の各ノードの動きベクトルを算出することで、フレームごとに各ノードの集合としての顔モデルを作成することができ、変形された顔モデルの系列が得られる。これら顔モデルに対するレンダリングを行なうことで、アニメーションを作成できる。

## 【 0 1 4 4 】

アニメーションデータ記憶部 4 6 に記憶された顔パラメータの系列 1 2 6 は、音声データ 4 2 により表される音声が発話されるとき顔の各特徴点の非線形的な軌跡を、実際にモーションキャプチャにより得られた測定データに基づいて表現する。したがって、発話中の表情の非線形的な変化を忠実に再現した、自然なアニメーションを作成することができる。

## 【 0 1 4 5 】

リップシンクアニメーション作成装置 4 0 は、モデルベースでアニメーションを作成する。ユーザは、音声 - 視覚素コーパス記憶部 6 2 が作成された後は、キャラクタの声に相当する音声データ 4 2 と、静止状態でのキャラクタの顔の形状を定義した顔モデルと、音声データ 4 2 に対するピタピアライメントを行なうための音響モデルと、音声データ 4 2 に対応する発話テキストとを用意し、顔モデル上に、特徴点に対応する仮想マーカを指定するだけで、キャラクタの声に合わせて表情の変化する、自然なリップシンクアニメーションを作成できる。キャラクタの顔のデザインが制限されることなく、顔モデル 4 4 が表す

10

20

30

40

50

キャラクターの顔の形状は任意のものでよい。そのため、ユーザによるアニメーション制作のバリエーションを狭めることなく、リップシンクアニメーションを作成できる。

【 0 1 4 6 】

なお、上記した実施の形態では、リップシンクアニメーション作成装置 4 0 は収録システム 6 0、アニメーションデータ合成装置 4 4、アニメーション作成装置 4 8、及びアニメーション読出部 1 0 0 の全てを含んでいる。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。これらが全て別々の装置、又はプログラムにより実現されてもよい。また、これらが物理的に同じコンピュータ上で実現される必要もないし、例えばアニメーションデータ合成装置 4 4 を構成する各部が単一のコンピュータ上で実現される必要もない。これらを別々のコンピュータ上で動作するプログラムにより実現し、それらの間のデータの移動を、ネットワーク経由又はリムーバブル記録媒体を介して実現するようにしてもよい。

10

【 0 1 4 7 】

[ コンピュータによる実現及び動作 ]

本実施の形態のリップシンクアニメーション作成装置 4 0 の各機能部は、収録システム 6 0 ( 図 1 及び図 2 参照 ) の録画・録音システム 1 1 2 及び M o C a p システム 1 1 4 に含まれる一部の特殊な機器を除き、いずれもコンピュータハードウェアと、そのコンピュータハードウェアにより実行されるプログラムと、コンピュータハードウェアに格納されるデータとにより実現される。図 1 8 はこのコンピュータシステム 4 5 0 の外観を示し、図 1 9 はコンピュータシステム 4 5 0 の内部構成を示す。

20

【 0 1 4 8 】

図 1 8 を参照して、このコンピュータシステム 4 5 0 は、DVD ( D i g i t a l V e r s a t i l e D i s k ) ドライブ 4 7 2 及びリムーバブルなメモリを装着可能なメモリポート 4 7 0 を有するコンピュータ 4 6 0 と、キーボード 4 6 6 と、マウス 4 6 8 と、モニタ 4 6 2 と、マイクロフォン 4 9 0 と、一対のスピーカ 4 5 8 とを含む。マイクロフォン 4 9 0 は、このコンピュータシステム 4 5 0 において音声データ 4 2 ( 図 1 参照 ) を収録する際に使用される。スピーカ 4 5 8 はアニメーションを表示する際の音声の再生に用いられる。

【 0 1 4 9 】

図 1 9 を参照して、コンピュータ 4 6 0 は、メモリポート 4 7 2 及び DVD ドライブ 4 7 0 に加えて、ハードディスク 4 7 4 と、CPU ( 中央処理装置 ) 4 7 6 と、CPU 4 7 6、ハードディスク 4 7 4、メモリポート 4 7 2、及び DVD ドライブ 4 7 0 に接続されたバス 4 8 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 ) 4 7 8 と、バス 4 8 6 に接続され、プログラム命令、システムプログラム、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 ) 4 8 0 と、バス 4 8 6 に接続され、マイクロフォン 4 9 0 からの音声信号をデジタル信号化したり、CPU 4 7 6 より出力されるデジタル音声信号をアナログ化してスピーカ 4 5 8 を駆動したりするためのサウンドボード 4 8 8 とを含む。コンピュータシステム 4 5 0 はさらに、プリンタを含んでもよい。

30

【 0 1 5 0 】

コンピュータ 4 6 0 はさらに、ローカルエリアネットワーク ( LAN ) 4 5 2 への接続を提供するネットワークインターフェイス ( I / F ) 4 9 6 を含む。

40

【 0 1 5 1 】

コンピュータシステム 4 5 0 にリップシンクアニメーション作成装置 4 0 の各機能部を実現させるためのコンピュータプログラムは、DVD ドライブ 4 7 0 又はメモリポート 4 7 2 に挿入される DVD 4 8 2 又はメモリ 4 8 4 に記憶され、さらにハードディスク 4 7 4 に転送される。又は、プログラムは図示しないネットワークを通じてコンピュータ 4 6 0 に送信されハードディスク 4 7 4 に記憶されてもよい。プログラムは実行の際に RAM 4 8 0 にロードされる。DVD 4 8 2 から、メモリ 4 8 4 から、又はネットワークを介して、直接に RAM 4 8 0 にプログラムをロードしてもよい。

【 0 1 5 2 】

50

このプログラムは、コンピュータ460にこの実施の形態のリップシンクアニメーション作成装置40の各機能部を実現させるための複数の命令を含む。この機能を実現させるのに必要な基本的機能のいくつかは、コンピュータ460にインストールされる各種ツールキットのモジュール、又はコンピュータ460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により提供される。したがって、このプログラムはこの実施の形態のシステム及び方法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ように制御されたやり方で適切な機能又は「ツール」を呼出すことにより、上記した顔アニメーションの作成システム40の各機能部が行なう処理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450の動作は周知であるので、ここでは繰返さない。

10

## 【0153】

[様々な変形例]

なお、上記した実施の形態では、視覚素コーパス及び視覚素データの視覚素ラベルとして、三つ組視覚素ラベルを用いている。こうすることにより、ある視覚素に対応する発話の前後の発話における顔の動きまで反映した形で、適切な視覚素ユニットを選択できる。しかし本発明はそのような実施の形態には限定されない。例えば視覚素ラベルとして一つだけを使用してもよい。この場合、得られる視覚素ユニットの不連続部分が大きくなる可能性があるが、上記した加重加算処理によって滑らかに連結することができる。

## 【0154】

また、上記した実施の形態では、三つ組視覚素ラベルとして、ある視覚素を中心に、その前後の視覚素ラベルを一つずつ採用し、中央の視覚素ラベルと組合せたものを用いた。しかし本発明はそのような実施の形態には限定されない。例えば、ある視覚素の前又は後の視覚素の視覚素ラベルと、対象の視覚素の視覚素ラベルとからなる、二つ組視覚素ラベルを用いてもよい。また、四つ以上の視覚素ラベルからなるものを採用してもよい。四つ以上の視覚素ラベルの組を採用した場合には、適切な視覚素ラベルの組を持つ視覚素ユニットが視覚素コーパスで見つからない可能性が高くなる。そうした場合、視覚素コーパスをより大きくしてもよいし、上記した実施の形態におけるように、それより少ない数の視覚素ラベルの組を用いて代替的な視覚素ユニットを探すようにしてもよい。

20

## 【0155】

また、上記した実施の形態では、視覚素ユニットの連結の際に、先行する視覚素ユニットの動きベクトルのみを時間Tだけ拡張している。しかし本発明はそのような実施の形態には限定されない。例えば、どの視覚素ユニットも、その前後にT/2だけ動きベクトルを拡張するようにしてもよい。また、拡張する時間についても主観的テストによって適切と思われる値に設定すればよい。

30

## 【0156】

上記した実施の形態では、どの視覚素ユニットを選択すべきかを決定するために、視覚素に対応する音素の発話継続長と、その間の音声の平均パワーとを用いている。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。これら以外の韻律的特徴、例えば音の高さ(基本周波数)を用いてもよいし、これらの任意の組合せを用いてもよい。顔画像が時間とともに変化するというアニメーションの特性上、発話継続長については評価の対象として採用することが望ましいが、発話継続長以外の韻律的特徴を用いて視覚素ユニットを選択したのち、発話継続長を視覚素データの継続長にあわせて調整するようにしてもよい。

40

## 【0157】

また、上記した実施の形態では、式(2)によって、韻律的に似た特徴を持つ三つ組視覚素ユニットを探し、そうしたユニットがアニメーション合成の上で適したものとして取り扱っている。しかし、使用できる式は式(2)に限らない。韻律的特徴の差の二乗和を式(2)に代えて用いてもよいし、それ以外の式で、視覚素ユニットと視覚素データとの韻律的な特徴の相違を的確に表すものがあればそうしたものを用いてもよい。

## 【0158】

50

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味および範囲内でのすべての変更を含む。

【図面の簡単な説明】

【0159】

【図1】本発明の実施の形態に係るリップシンクアニメーション作成装置40のブロック図である。

【図2】収録システム60の詳細な構成を示すブロック図である。

【図3】頭部110に装着されるマーカの配置例を示す図である。

10

【図4】アニメーションキャラクタの顔モデル170及びフレームごとの動きベクトルから顔画像のアニメーション172を作成する手順を示す模式図である。

【図5】三つ組視覚素データ列作成部144のブロック図である。

【図6】ビタピアライメントの概略を示す模式図である。

【図7】三つ組視覚素データ列234の構成を示す図である。

【図8】音素 - 視覚素変換テーブルの構成を示す図である。

【図9】音声 - 視覚素コーパス記憶部62の構成を示す図である。

【図10】三つ組視覚素データ列作成部80のより詳細な構成を示す図である。

【図11】三つ組視覚素データ列296の構成を示す図である。

20

【図12】三つ組視覚素データ列300の構成を示す図である。

【図13】三つ組視覚素ユニット選択部82を実現する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14】顔モデルを示す模式図である。

【図15】顔モデル変形部92により実行されるマーカラベリング処理を実現する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16】顔モデル変形部92により実行される、顔モデルにおける唇周辺のノードと仮想マーカとの対応付を説明するための図である。

【図17】三つ組視覚素ユニット連結部84による動きベクトルの連結方法を説明するための図である。

【図18】本発明の一実施の形態に係るリップシンクアニメーション作成装置40の主要な機能を実現するコンピュータシステムの外觀の一例を示す図である。

30

【図19】図18に示すコンピュータシステムのブロック図である。

【符号の説明】

【0160】

40 リップシンクアニメーション作成装置

44 アニメーションデータ合成装置

48 アニメーション作成装置

60 収録システム

62 音声 - 視覚素コーパス記憶部

80 三つ組視覚素データ列作成部

40

82 三つ組視覚素ユニット選択部

84 三つ組視覚素ユニット連結部

92 顔モデル変形部

96 レンダリング部

100 アニメーション読出部

112 録画・録音システム

114 M o C a pシステム

116 音声・動画データ

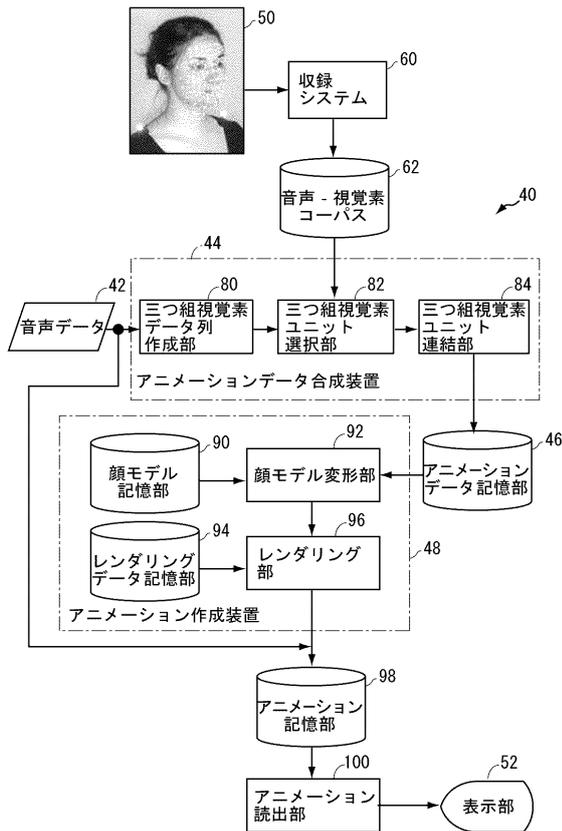
122 データセット作成装置

138 データ処理装置

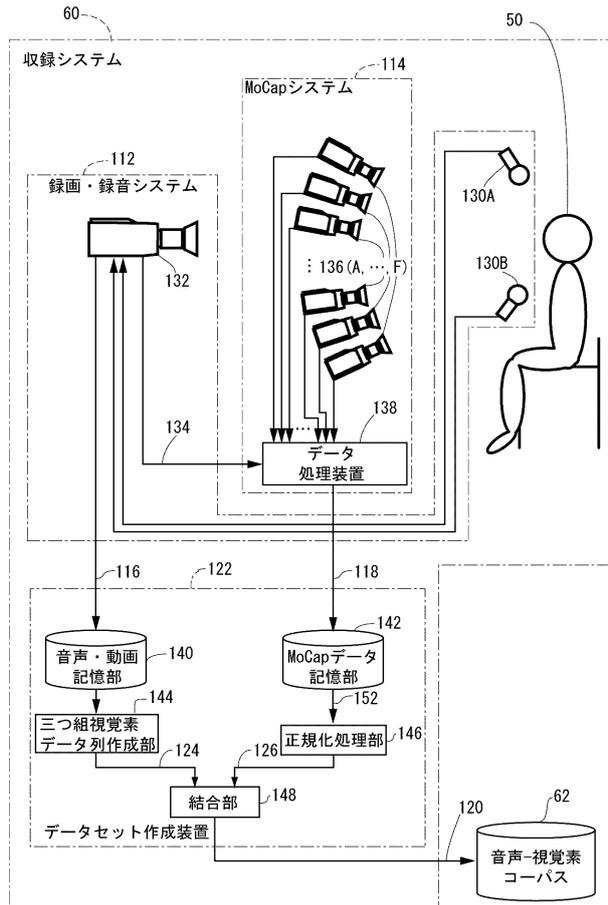
50

- 1 4 4 三つ組視覚素データ列作成部
- 1 4 6 正規化处理部
- 1 4 8 結合部
- 2 0 0 , 2 8 0 フレーム化处理部
- 2 0 1 , 2 8 2 特徴抽出部
- 2 0 2 , 2 8 4 音響モデル記憶部
- 2 0 4 , 2 8 6 発話テキスト記憶部
- 2 0 6 , 2 8 8 ビタピアライメント部
- 2 1 0 , 2 9 0 音素 - 視覚素変換テーブル記憶部
- 2 1 2 , 2 9 2 音素 - 視覚素変換部
- 2 1 4 , 2 9 3 視覚素データ列記憶部
- 2 1 6 , 2 9 4 視覚素 - 三つ組視覚素変換部
- 2 1 8 , 2 9 5 三つ組視覚素データ列記憶部
- 2 3 4 , 2 9 6 三つ組視覚素データ列
- 2 4 0 音声波形データ
- 2 4 2 動きベクトル列
- 2 4 4 三つ組視覚素ユニット列
- 3 0 0 三つ組視覚素データ列

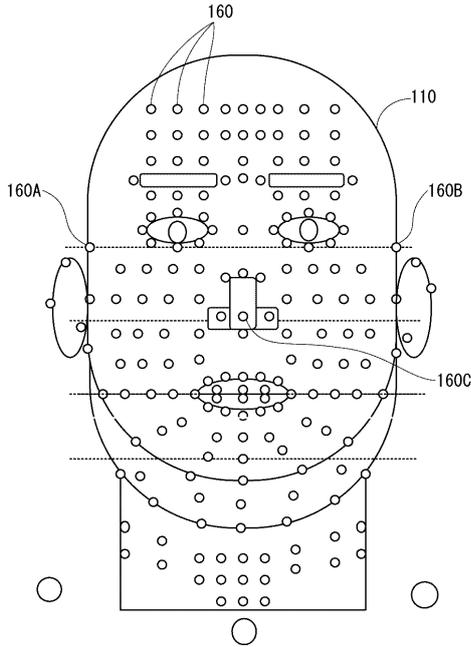
【図1】



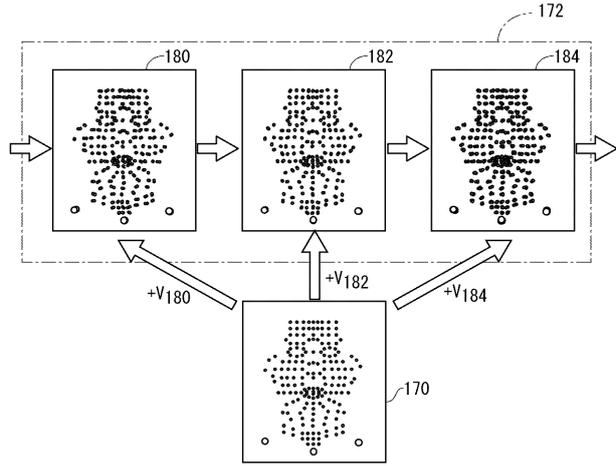
【図2】



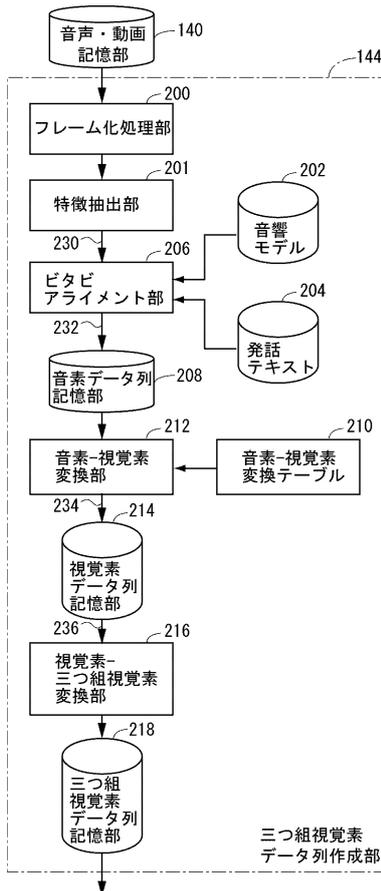
【図3】



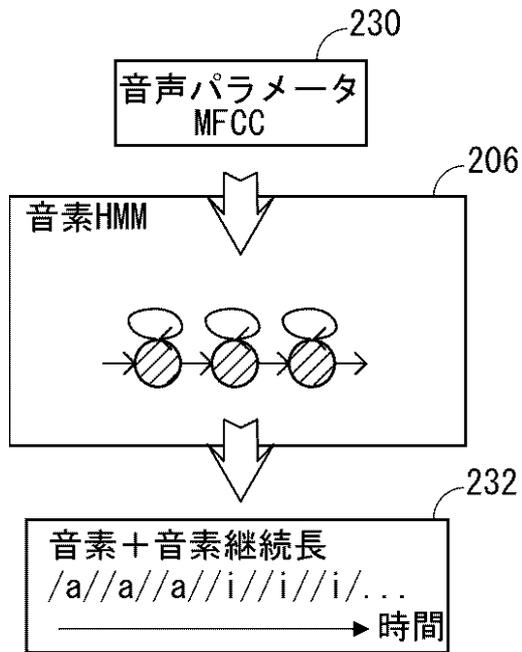
【図4】



【図5】



【図6】



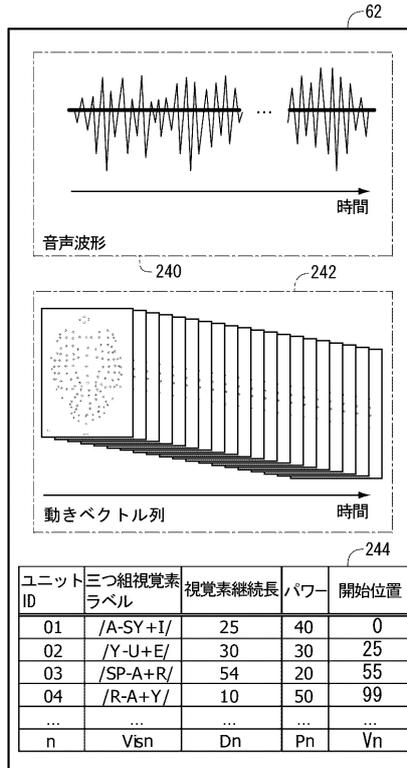
【図7】

番号	三つ組視覚素	継続長 (MSEC)	パワー
0:	sil-sp+A	119	P <sub>0</sub>
1:	sp-A+R	11	P <sub>1</sub>
2:	A-R+A	8	P <sub>2</sub>
3:	R-A+Y	11	P <sub>3</sub>
...	...	...	...

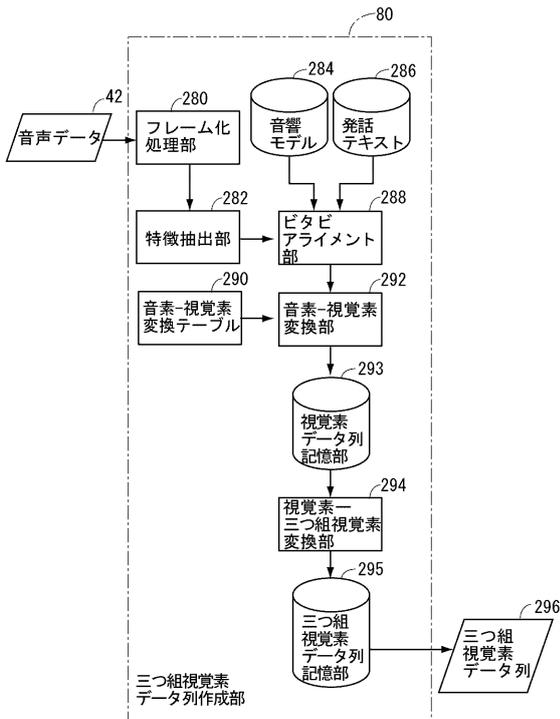
【図8】

番号	視覚素	音素					
1	a	a	A				
2	i	i	I				
3	u	u	U				
4	e	e	E				
5	o	o	O				
6	p	m	b	p			
7	sy	j	my	ky	by	gy	ny
		hy	ry	py	ch	dy	sh
8	w	w	f				
9	t	t	d	n			
10	s	ts	z	s			
11	y	y					
12	r	r					
13	vf	k	h	g			
14	sp	silB	silE	sp	q		
15	N	N					

【図9】



【図10】



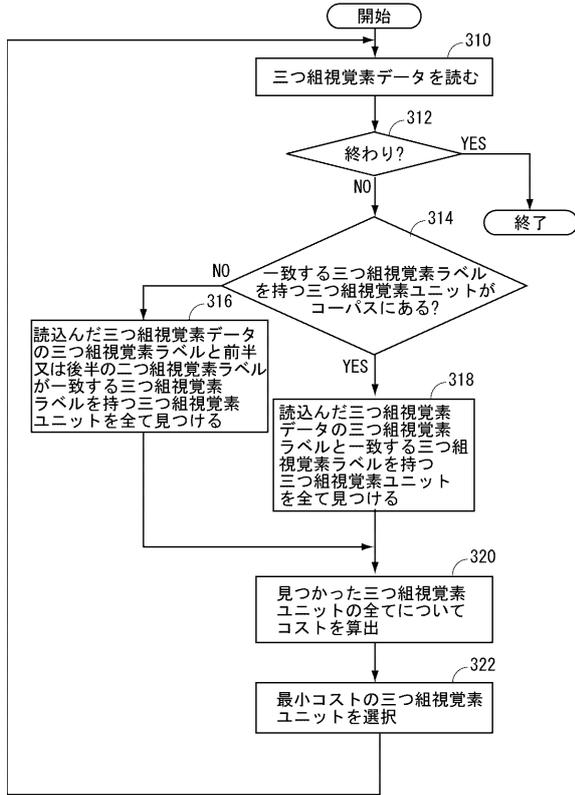
【図11】

シーケンス番号	三つ組視覚素ラベル	視覚素継続長	パワー
0	/SP-K+ O/	120	10
1	/K-O+ N/	40	40
2	/N-N+I /	10	20
...	...	...	...
m	/W-A+S P/	20	30

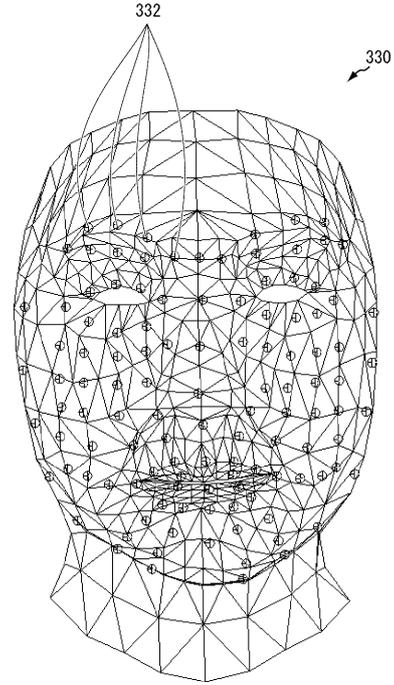
【図12】

シーケンス番号	三つ組視覚素ラベル	視覚素継続長	パワー	選択ユニットID
0	/SP-K+ 0 /	120	10	49201
1	/K-O+ N /	40	40	180129
2	/N-N+ I /	10	20	138002
...	...	...	...	...
m	/W -A+ SP/	20	30	29039

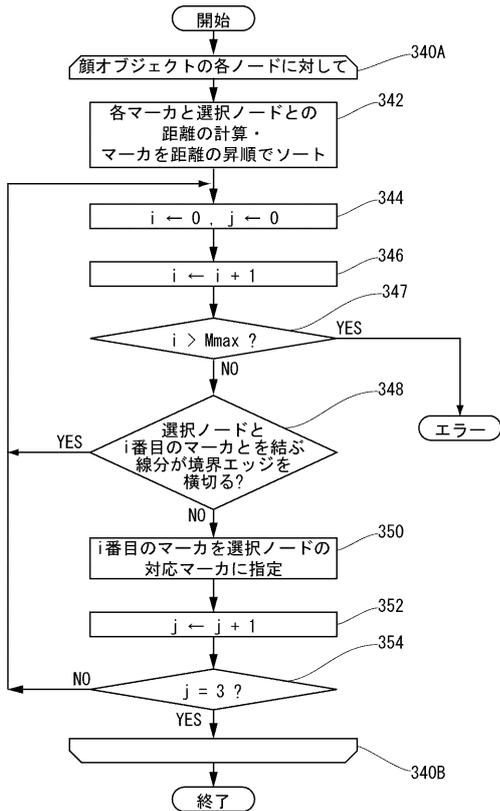
【図13】



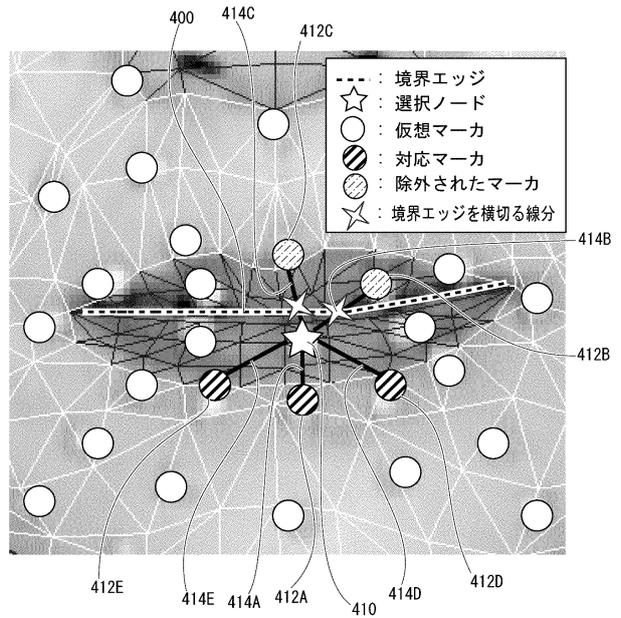
【図14】



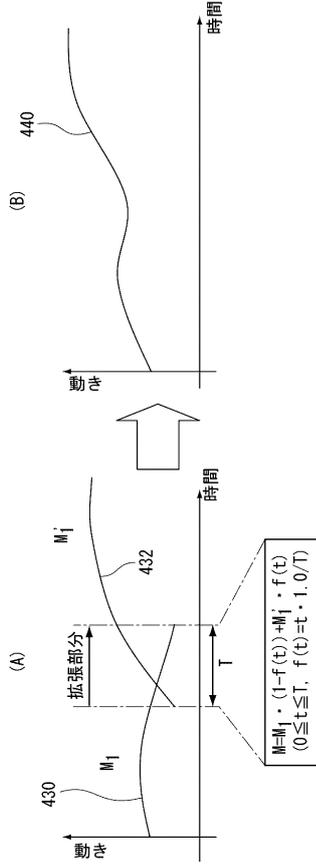
【図15】



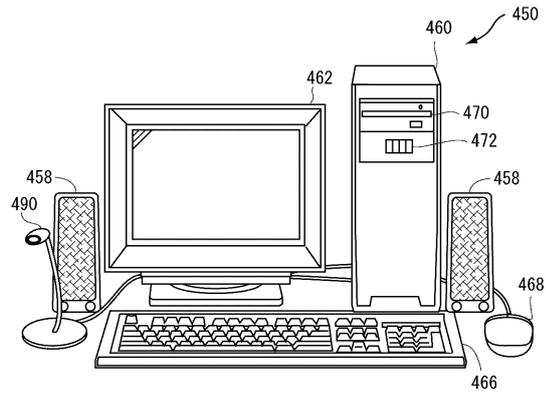
【図16】



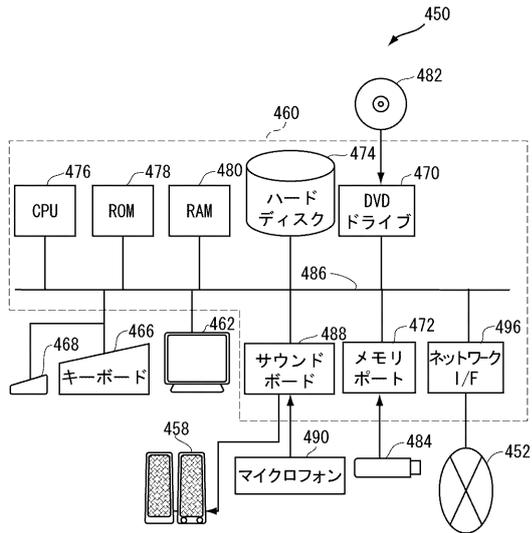
【図17】



【図18】



【図19】



## フロントページの続き

審査官 加内 慎也

- (56)参考文献 特開2000-011200(JP,A)  
特開2002-244689(JP,A)  
特開2003-208188(JP,A)  
カスタマイズ性を考慮した擬人化音声対話ソフトウェアツールキットの設計, 情報処理学会論文誌, 2002年 7月15日, 2249-2263  
HMMを用いた自然な発話動画像合成, 電子情報通信学会論文誌, 2000年11月25日, 2498-2506  
顔の分析・合成とその応用, 情報処理学会研究報告, 2003年 7月 4日, 107-114  
森島繁生, デジタルメディア作品の制作を支援する基盤技術, CREST研究年報 平成17年度 戦略的創造研究推進事業, 科学技術振興機構, 2005年, <http://www.jst.go.jp/kisoken/crest/report/heisei17/pdf/a07/f01/s004.pdf>

## (58)調査した分野(Int.Cl., DB名)

G06T 13/00

G06T 15/70