

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4631078号
(P4631078)

(45) 発行日 平成23年2月16日(2011.2.16)

(24) 登録日 平成22年11月26日(2010.11.26)

(51) Int. Cl.		F I	
G06T 13/40	(2011.01)	G06T 15/70	B
G06T 1/00	(2006.01)	G06T 1/00	340A
G10L 15/00	(2006.01)	G10L 15/00	200P
G10L 15/28	(2006.01)	G10L 15/28	500

請求項の数 11 (全 38 頁)

(21) 出願番号	特願2006-201026 (P2006-201026)	(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7月24日(2006.7.24)	(74) 代理人	100099933 弁理士 清水 敏
(65) 公開番号	特開2007-58846 (P2007-58846A)	(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19年3月8日(2007.3.8)	(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
(31) 優先権主張番号	特願2005-217860 (P2005-217860)		
(32) 優先日	平成17年7月27日(2005.7.27)		
(33) 優先権主張国	日本国(JP)		

最終頁に続く

(54) 【発明の名称】 リップシンクアニメーション作成用の統計確率モデル作成装置、パラメータ系列合成装置、リップシンクアニメーション作成システム、及びコンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

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予め定める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であって、

前記モーションキャプチャデータは複数のフレームを含み、前記複数のフレームの各々は当該フレームにおける前記複数個の特徴点の位置データを含み、前記複数のフレームと前記収録音声との間には時間的対応関係が付けられており、

前記統計確率モデル作成装置は、

前記音声の特徴量と音素とに関して予め準備された所定の音素統計確率モデルを用いて、前記データセットに含まれる収録音声データに含まれる音素列、及び当該音素列を構成する各音素に関する音素継続長を推定するための音素列推定手段と、

前記音素列推定手段により推定された音素列及び音素継続長に基づき、前記フレームの各々に対し、予め定義された所定のラベルセットに属するラベルによるラベリングを行なうためのラベリング手段と、

前記ラベリング手段によりラベリングされたモーションキャプチャデータからの統計的学習により、前記リップシンクアニメーション作成用の統計確率モデルとして、前記ラベル間の遷移確率と前記各特徴点の位置の出力確率とに関する統計確率モデルの学習を行なうための学習手段とを含む、統計確率モデル作成装置であって、

前記ラベルセットは、各々が発話時の口の形状を表す、複数個の所定の視覚素ラベルを含み、

前記ラベリング手段は、

音素と視覚素との間の所定の対応関係にしたがい、前記音素列推定手段により推定された前記音素列を前記視覚素ラベルの系列に変換し、前記音素継続長をもとに当該系列を構成する前記視覚素ラベルの各々の継続長を決定するための手段と、

前記決定するための手段により決定された視覚素ラベルの系列と継続長とをもとに、前記フレームの各々に対し、前記視覚素ラベルによるラベリングを行なうための視覚素ラベリング手段とを含む、統計確率モデル作成装置。

【請求項 2】

前記ラベルセットに含まれる視覚素ラベルの数は、前記音素列推定手段により推定される音素セットに含まれる音素の種類の数より少ない、請求項 1 に記載の統計確率モデル作成装置。

【請求項 3】

前記学習手段は、前記ラベリング手段によりラベリングされたモーションキャプチャデータから、連続する三つのラベルの組を学習単位として学習を行なうことにより、前記リップシンクアニメーション作成用の統計確率モデルとして、前記ラベル間の遷移確率と前記各特徴点の位置の出力確率とに関する統計確率モデルの学習を行なうための手段を含む、請求項 1 又は請求項 2 に記載の統計確率モデル作成装置。

【請求項 4】

前記統計確率モデル作成装置はさらに、前記モーションキャプチャデータ中の各フレームにおいて、当該フレームと、当該フレームに隣接するフレームとにおける前記複数個の特徴点の位置データから、前記複数個の特徴点の予め定められた動的特徴データを算出し対応する位置データに付加するための動的特徴データ算出手段を含み、

前記学習手段は、前記ラベリング手段によりラベリングされ、前記動的特徴データが付加された位置データを含むモーションキャプチャデータからの統計的学習により、前記リップシンクアニメーション作成用の統計確率モデルとして、前記ラベル間の遷移確率と前記各特徴点の位置の出力確率とに関する統計確率モデルの学習を行なうための手段を含む、請求項 1 ~ 請求項 3 のいずれかに記載の統計確率モデル作成装置。

【請求項 5】

前記動的特徴データ算出手段は、前記モーションキャプチャデータ中の各フレームにおいて、当該フレームの前記複数の特徴点の位置データと、当該フレームに隣接するフレームにおける前記複数個の特徴点の位置データとから、当該フレームにおける、前記複数個の特徴点の速度パラメータ及び加速度パラメータを前記動的特徴データとして算出し、対応する位置データに付加するための手段を含む、請求項 4 に記載の統計確率モデル作成装置。

【請求項 6】

コンピュータにより実行されると、当該コンピュータを請求項 1 ~ 請求項 5 のいずれかに記載の統計確率モデル作成装置として動作させる、コンピュータプログラム。

【請求項 7】

発話時における発話者の顔の複数個の特徴点の軌跡を時系列で表すパラメータ系列を合成するためのパラメータ系列合成装置であって、

発話により発生した音声の入力を受けて、音声の特徴量と音素とに関し予め学習を行なって得られた第 1 の統計確率モデルに基づき、当該音声を出力する音素列と当該音素列を構成する各音素の音素継続長とを推定するための音素列推定手段と、

前記音素列推定手段により推定された音素列と音素継続長とをもとに、予め定義された所定のラベルセットに属するラベルからなる系列を生成し、当該系列を構成する当該ラベルの各々の継続長を決定するためのラベル列生成手段と、

前記ラベル間の遷移確率と前記各特徴点の位置の出力確率とに関し予め学習することにより得られた第 2 の統計確率モデルに基づき、前記ラベル列生成手段により生成された系

10

20

30

40

50

列と継続長とを入力パラメータとして前記複数個の特徴点の軌跡を推定することにより、前記パラメータ系列を生成するための軌跡推定手段とを含む、パラメータ系列合成装置であって、

前記ラベルセットは、各々が発話時の口の形状を表す、複数個の所定の視覚素ラベルを含み、

前記第2の統計確率モデルは、前記視覚素ラベル間の遷移確率と前記各特徴点の位置の出力確率とに関し予め学習され、

前記ラベル列生成手段は、音素と前記視覚素ラベルとの間の所定の対応関係にしたがい、前記音素列推定手段により推定された音素列を前記視覚素ラベルの系列に変換し、前記音素継続長をもとに、当該系列を構成する各視覚素ラベルの継続長を決定するための変換手段を含む、パラメータ系列合成装置。

10

【請求項8】

前記ラベルセットに含まれる視覚素ラベルの数は、前記音素列推定手段により推定される音素セットに含まれる音素の種類の数より少ない、請求項7に記載のパラメータ系列合成装置。

【請求項9】

前記第2の統計確率モデルは、前記視覚素ラベル間の遷移確率と、前記各特徴点の位置パラメータ及び当該特徴点に関する動的特徴パラメータの出力確率とに関し予め学習された動的特徴による統計確率モデルを含み、

前記軌跡推定手段は、

前記ラベル間の遷移確率と前記各特徴点の位置パラメータ及び動的特徴パラメータの出力確率とに関し予め学習することにより得られた前記動的特徴による統計確率モデルに基づき、前記ラベル列生成手段により生成された系列と継続長とを入力パラメータとして、前記複数個の特徴点に対する前記位置パラメータ及び前記動的特徴パラメータの系列として最尤となる位置パラメータ及び動的特徴パラメータの系列を出力するための手段と、

20

前記位置パラメータ及び動的特徴パラメータの系列に対し、当該パラメータが得られた前記統計確率モデルに固有の変換によって、前記位置パラメータを前記動的特徴パラメータを用いて補正し、前記複数個の特徴点の各々の前記軌跡を推定するための手段とを含む、請求項7又は請求項8に記載のパラメータ系列合成装置。

【請求項10】

コンピュータにより実行されると、当該コンピュータを請求項7～請求項9のいずれかに記載のパラメータ系列合成装置として動作させる、コンピュータプログラム。

30

【請求項11】

第1の座標空間における複数のノードの座標値を用いて顔の形状を定義した所定の顔オブジェクトをもとに、音声に同期する前記顔のアニメーションを作成するためのリップシンクアニメーション作成システムであって、

請求項7～請求項10のいずれかに記載のパラメータ系列合成装置と、

前記音声の入力に対して前記パラメータ系列合成装置により合成される、発話者の顔の複数個の特徴点の軌跡を表すパラメータ系列に基づき、前記顔オブジェクトにおける前記ノードの座標値を変更することにより、前記顔の形状を定義するオブジェクトを、前記アニメーションのフレームごとに生成するための変形オブジェクト生成手段と、

40

前記アニメーションの前記各フレームについて、前記変形オブジェクト生成手段により生成されるオブジェクトから、当該フレームにおける前記顔の画像を合成するための画像化手段とを含む、リップシンクアニメーション作成システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、CG (Computer Graphics) を用いたアニメーションの作成技術に関し、キャラクタの発話時の表情を表現したリップシンクアニメーションを作成するための統計確率モデル作成装置、パラメータ系列合成装置、及びコンピュータプログラム、並びにそれ

50

らを用いたリップシンクアニメーション作成システムに関する。

【背景技術】

【0002】

アニメーション作品の制作に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。

10

【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。すなわち一つの手法は、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（ブレレコーディング：これを以下「ブレレコ」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。これに対しブレレコでは、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。

20

30

【0005】

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様々な技術が提案されている。後掲の非特許文献1には、キーフレーム法と呼ばれる手法によってリップシンクを実現するための技術が開示されている。この手法では、キャラクタの典型的な表情を表現したオブジェクトを予め複数用意しておく。そして、これら用意されたオブジェクトを用いて、発話中のキャラクタの表情を次のようにして指定する。まず、アニメーションを構成するフレームの中から、用意されたオブジェクトを使用してキャラクタの表情を表現するフレーム（キーフレーム）を定める。続いて、キーフレームで使用する表情のパラメータを指定する。この指定が完了すると、アニメーションの各フレームにおけるキャラクタの表情を表すオブジェクトを、フレームごとに生成する。この際、キーフレームについては、上記の指定により指定されたオブジェクトをそのまま当てはめる。二つのキーフレームの間にあるその他のフレーム（中間フレーム）については、その中間フレームの前後にある二つのキーフレームに使用されているオブジェクトからの、時間軸による線形補間によって、オブジェクトを生成する。

40

【0006】

後掲の非特許文献2には、物理モデルに基づく顔のシミュレーションによって、リップシンクを実現する技術が開示されている。この技術では、顔の筋肉と皮膚と骨格とを3層構造のばねモデルによって物理モデル化する。発話時の筋肉の動きに基づきモデル上で筋肉を操作し、筋肉が移動・変形した場合の皮膚の動きをシミュレートする。

【0007】

50

後掲の非特許文献3及び非特許文献4には、統計確率的な手法によって発話中における顔の動画像を合成する技術が開示されている。この手法では、予め発話時の顔の画像をデータベース(以下単に「DB」と書く。)化しておく。発話内容に適した特徴を備える画像をデータベース中の顔の画像から選び再構成する。

【0008】

このうち、非特許文献3に記載の技術では、写真画像がDB化される。合成されるアニメーションは、それら写真画像を再構成したものである。したがって、大規模かつ適切なDBを用意すれば、実写の動画に近い自然な映像でリップシンクを実現できる。

【0009】

非特許文献4に記載の技術では、3次元の顔のオブジェクトがDB化される。この技術では、発話中における顔の所定の複数の点についての位置計測と音声の収録とを同時に行なう。位置計測のデータについて主成分分析を行ない、顔のパラメータを生成する。顔のパラメータと音声の収録データとから、予め用意された音素隠れマルコフモデル(Hidden Markov Model: HMM)における状態に対応する顔のパラメータを選び、状態ごとに平均をとる。この平均されたパラメータを用い、音素HMMの各状態に対応するオブジェクトを生成しておく。このようにして生成されたオブジェクトと音素HMMとを用いて、プレレコでアニメーションを合成する。すなわち、まず、予め発話音声を用意しておき、当該発話音声から音素HMMを用いて音素列を合成する。この音素列に基づき、アニメーションの各フレームに対し、音素を指定する。指定された音素に対応するオブジェクトを当該フレームのオブジェクトに定め、オブジェクトの系列を作成し画像化する。

【非特許文献1】コーエン, M. M., マッサロ, D. W. 1993年. 「視覚的に合成された発話における同時調音のモデル」, コンピュータアニメーションのモデルと技法, 139 - 156頁(Cohen, M. M., Massaro, D. W. 1993. "Modeling coarticulation in synthetic visual speech", Models and Techniques in Computer Animation, pp.139-156)

【非特許文献2】ウォーターズ, K. 1987年. 「3次元の顔の表現をアニメーション化するための筋肉モデル」, ACM シーグラフ '87 17 - 24頁(Waters, K., 1987. A muscle model for animating three-dimensional facial expressions. ACM SIGGRAPH '87 pp.17-24)

【非特許文献3】エザット, T., ガイガー, G., ポッジョ, T. 2002年. 「学習可能なビデオリアリスティック発話アニメーション」, ACM シーグラフ 2002 (Ezzat, T., Geiger, G. and Poggio, T. "Trainable Videorealistic Speech Animation", Proceedings of ACM SIGGRAPH 2002)

【非特許文献4】K カキハラ, S ナカムラ, K シカノ 「HMMに基づく音声からの顔の動きの合成」, 米国電気電子学会(IEEE)マルチメディアの国際会議及び博覧会予稿集, 2000年7 - 8月 第1巻 427 - 430頁(K Kakahara, S Nakamura, K Shikano, "Speech-To-Face Movement Synthesis Based on HMMs", Proceedings of IEEE International Conference on Multimedia and Expo, July-August, 2000 Vol.1, pp.427-430)

【非特許文献5】徳田 恵一、「HMMによる音声合成の基礎」、電子情報通信学会技術研究報告、第100巻第392号、SP2000-74, pp. 43 - 50, 2000年10月

【発明の開示】

【発明が解決しようとする課題】

【0010】

アフレコであれプレレコであれ、手作業でのアニメーション制作によってリップシンクを実現するには、膨大な量の作業とそのための高度な技能とを要する。アフレコでリップシンクを実現するには、発話時の各フレームにおける表情を制作者が的確に予測しなければならない。しかし、この予測にも限度がある。また、アフレコでリップシンクを実現するには、発話者が発話のタイミングを調整しなければならない。しかし、発話のタイミン

グ等をフレーム単位で調整することは困難である。そのため、高度なリップシンクを実現するのに、制作者・発話者の双方に極めて高い技能が要求される。これに対しプレレコでは、予め収録された音声に合わせて各フレームの画像が制作される。画像は音声と異なり、フレーム単位での修正が可能であるため、高精度にタイミングの調整を行なうことができる。したがって高度なリップシンクが実現可能となる。しかしながらこの方法では、アニメーション画像の制作者がフレームごとに画像を調整しなければならない。又は制作者が、音声と画像とを照合して画像を修正しなければならない。そのため、制作者に過酷な作業を強いることになる。

【 0 0 1 1 】

リップシンクを実現するための作業に関する上記のような問題は、3次元モデルを用いたCGによるアニメーション制作においても同様に発生する。3次元のオブジェクトを用いて表情などを表現するには、仮想空間上でオブジェクトを変形させなければならない。すなわち、ポリゴンの頂点（ノード）の位置についていちいち再定義しなければならない。オブジェクトの変形によってアニメーションを制作するには、フレームごとにこの作業を行わなければならない。現在のアニメーションに用いられる形状モデルは、膨大な数のポリゴンにより構成されているため、再定義を要するノードの数もまた膨大である。そのため、制作に要する作業量及びコストは莫大なものとなる。

【 0 0 1 2 】

非特許文献1に記載の技術では、典型的な表情のオブジェクトが、そのままキーフレームにおける画像の合成に用いられる。したがって、あるキャラクタ用のオブジェクトは、他のキャラクタに転用できない。すなわち、キャラクタごとに典型的な表情のオブジェクトを用意しなければならない。また、この技術では、中間フレームにおける表情を表現するオブジェクトが予め用意されたオブジェクトの線形補間により生成される。しかし、人間の表情の変化はこのような線形的なものではない。したがって、この手法では、表情の変化を忠実に表現できず、リップシンクは不完全なものとなる。

【 0 0 1 3 】

非特許文献2に記載の技術は、顔の物理的構造を考慮した手法であり、シミュレーションを適切に行なえば、表情の変化を忠実に表現することができるかもしれない。しかし、この技術で意図した表情を表現するには、各筋肉組織の収縮量を解剖学的な知識に基づいていちいち設定しなければならない。そのため、この技術を用いてリップシンクアニメーションを作成するのは極めて困難である。

【 0 0 1 4 】

非特許文献3に記載の技術では、発話時の表情の特徴量を動画像から得ている。しかしこの技術では、次のような問題が発生する。すなわち、顔及びその表情は立体的（3次元）であるのに対し、動画像は2次元の情報である。3次元での形状変化に関する特徴量を2次元の動画像から得るのは困難である。したがってこの技術では、表情の変化についての情報を得るのが困難であるという問題が発生する。また、動画像の情報としての質はその画像を撮影するためのカメラの性能に依存する。したがって、動画像から求める特徴量に誤差が生じる恐れがあるという問題も発生する。

【 0 0 1 5 】

非特許文献4に記載の手法では、アニメーションとして作成可能な顔の表情は、DBに格納されたオブジェクトで表現される表情に限定されてしまう。多様な容顔のキャラクタの多彩な表情を表現するには、キャラクタごとに顔のオブジェクトを用意しDB化する必要がある。これは事実上不可能である。

【 0 0 1 6 】

それゆえに、本発明の目的は、任意のキャラクタについて、高度なリップシンクを実現するとともに、リップシンクアニメーションの制作作業を省力化する統計確率モデル作成装置、パラメータ系列合成装置、及びそれらを用いたリップシンクアニメーション作成システムを提供することである。

【課題を解決するための手段】

10

20

30

40

50

【0017】

本発明の第1の局面に係る統計確率モデル作成装置は、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予め定める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である。モーションキャプチャデータは複数のフレームを含み、複数のフレームの各々は当該フレームにおける複数個の特徴点の位置データを含み、複数のフレームと収録音声との間には時間的対応関係が付けられている。統計確率モデル作成装置は、音声の特徴量と音素とに関して予め準備された所定の音素統計確率モデルを用いて、データセットに含まれる収録音声データに含まれる音素列、及び当該音素列を構成する各音素に関する音素継続長を推定するための音素列推定手段と、音素列推定手段により推定された音素列及び音素継続長に基づき、フレームの各々に対し、所定のラベルセットに属するラベルによるラベリングを行なうためのラベリング手段と、ラベリング手段によりラベリングされたモーションキャプチャデータからの統計的学習により、リップシンクアニメーション作成用の統計確率モデルとして、ラベル間の遷移確率と各特徴点の位置の出力確率とに関する統計確率モデルの学習を行なうための学習手段とを含む。

10

【0018】

発話時の音声から音素列とその継続長が推定される。この音素列及び音素継続長に基づき、音声及びモーションキャプチャデータの各フレームについて、ラベリングが行なわれる。このラベリングがされたモーションキャプチャデータを学習データとして学習手段が統計的学習データを行なうことにより、統計確率モデルが得られる。この統計確率モデルを使用すると、音声を構成する音素についてラベリングがされた音声を与えられると、そのラベル系列に基づいて、音声に対応する顔の特徴点の位置の確率を出力できる。この確率に基づき、それらの特徴点の軌跡のうちで最尤となるものを定めることにより、音声から顔の動きを推定することができる。特徴点の軌跡が与えられるので、学習データを収録したときの発話者とは異なる顔モデルであっても、特徴点の対応付けがされていれば、音声に基づいてその顔モデルの動きを推定することができる。そのために、莫大な労力を要する作業は不要である。その結果、任意のキャラクタについて、高度なリップシンクを実現するとともに、リップシンクアニメーションの制作作業を省力化する統計確率モデル作成装置を提供することができる。

20

30

【0019】

ラベルセットは、各々発話時の口の形状を表す、複数個の所定の視覚素ラベルを含んでもよい。ラベリング手段は、音素と視覚素との間の所定の対応関係にしたがい、音素列推定手段により推定された音素列を視覚素ラベルの系列に変換し、音素継続長をもとに当該系列を構成する視覚素ラベルの各々の継続長を決定するための手段と、決定するための手段により決定された視覚素ラベルの系列と継続長とをもとに、フレームの各々に対し、視覚素ラベルによるラベリングを行なうための視覚素ラベリング手段とを含む。

【0020】

音声から視覚素への変換がされ、その視覚素とモーションキャプチャデータとを学習データとしてモデルの学習が行なわれる。モデルから推定される特徴点の軌跡は、視覚素という形で予め定められたラベルセットとの関連で推定されるので、推定を効率的に行なうことができる。

40

【0021】

ラベルセットに含まれる視覚素ラベルの数は、音素列推定手段により推定される音素セットに含まれる音素の種類の数より少なくてもよい。

【0022】

音素の種類より視覚素ラベルの数が少ないので、最終的な特徴点の位置の推定が効率的に行なえる。

【0023】

ラベルセットは、各々が一つの音素を表す複数個の音素ラベルを含んでもよい。ラベリ

50

ング手段は、音素列推定手段により推定された音素列をもとに、音素ラベルの系列を生成し、音素継続長をもとに音素ラベルの各々の継続長を決定するための手段と、決定するための手段により決定された音素ラベルの系列と継続長とをもとに、フレームの各々に対し、音素ラベルによるラベリングを行なうための音素ラベリング手段とを含む。

【0024】

音声から得られた音素ラベルをそのまま学習に使用する。モデルから顔の特徴点の位置データの系列を推定する場合、その組み合わせは膨大となるが、出力として音素ラベルの形をとれば、組み合わせが音素ラベルの数によって限定される。その結果、このモデルを用いた顔の形状の推定が効率的に行なえる。

【0025】

学習手段は、ラベリング手段によりラベリングされたモーションキャプチャデータから、連続する三つのラベルの組を学習単位として学習を行なうことにより、リップシンクアニメーション作成用の統計確率モデルとして、ラベル間の遷移確率と各特徴点の位置の出力確率とに関する統計確率モデルの学習を行なうための手段を含んでもよい。

【0026】

顔の表情は、発話中の音素だけではなく、その前後の音素にも影響される。そこで、このように連続する三つのラベルの組を学習単位とすることにより、モデルによる顔特徴点の位置データの系列を推定する際に、実際の音声の連続に合致した形での推定を行なうことができ、アニメーションの動きが自然なものになる。

【0027】

統計確率モデル作成装置はさらに、モーションキャプチャデータ中の各フレームにおいて、当該フレームと、当該フレームに隣接するフレームとにおける複数個の特徴点の位置データから、複数個の特徴点の予め定められた動的特徴データを算出し対応する位置データに付加するための動的特徴データ算出手段を含み、学習手段は、ラベリング手段によりラベリングされ、動的特徴データが付加された位置データを含むモーションキャプチャデータからの統計的学習により、リップシンクアニメーション作成用の統計確率モデルとして、ラベル間の遷移確率と各特徴点の位置の出力確率とに関する統計確率モデルの学習を行なうための手段を含む。

【0028】

このように動的特徴データを学習に用い、音声から顔の特徴点の位置を推定する際にも同様の動的特徴データを使用するようにすると、特徴点の軌跡が実際の軌跡に類似した、滑らかなものとなる。

【0029】

動的特徴データ算出手段は、モーションキャプチャデータ中の各フレームにおいて、当該フレームの複数の特徴点の位置データと、当該フレームに隣接するフレームにおける複数個の特徴点の位置データとから、当該フレームにおける、複数個の特徴点の速度パラメータ及び加速度パラメータを動的特徴データとして算出し、対応する位置データに付加するための手段を含んでもよい。

【0030】

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、当該コンピュータを本発明の第1の局面に係るいずれかの統計確率モデル作成装置として動作させる。

【0031】

本発明の第3の局面に係るパラメータ系列合成装置は、発話時における発話者の顔の複数個の特徴点の軌跡を時系列で表すパラメータ系列を合成するためのパラメータ系列合成装置である。パラメータ系列合成装置は、発話により発生した音声の入力を受けて、音声の特徴量と音素とに関し予め学習を行なって得られた第1の統計確率モデルに基づき、当該音声を出力する音素列と当該音素列を構成する各音素の音素継続長とを推定するための音素列推定手段と、音素列推定手段により推定された音素列と音素継続長とをもとに、予め定義されたラベルセットに属する所定のラベルからなる系列を生成し、当該系列を構成

10

20

30

40

50

する当該ラベルの各々の継続長を決定するためのラベル列生成手段と、ラベル間の遷移確率と各特徴点の位置の出力確率とに関し予め学習することにより得られた第2の統計確率モデルに基づき、ラベル列生成手段により生成された系列と継続長とを入力パラメータとして、複数個の特徴点の軌跡を推定することにより、パラメータ系列を生成するための軌跡推定手段とを含む。

【0032】

音声に含まれる音素列から所定のラベル列への変換がされ、そのラベル列とモーションキャプチャデータとを学習データとしてモデルの学習が行なわれる。モデルから推定される特徴点の軌跡は、予め定められたラベルセット内のラベルにより限定されるので、推定を効率的に行なうことができる。

10

【0033】

ラベルセットは、各々が発話時の口の形状を表す、複数個の所定の視覚素ラベルを含んでもよい。第2の統計確率モデルは、視覚素ラベル間の遷移確率と各特徴点の位置の出力確率とに関し予め学習される。ラベル列生成手段は、音素と視覚素ラベルとの所定の対応関係にしたがい、音素列推定手段により推定された音素列を視覚素ラベルの系列に変換し、音素継続長をもとに、当該系列を構成する各視覚素ラベルの継続長を決定するための変換手段を含む。

【0034】

ラベルセットに含まれる視覚素ラベルの数は、音素列推定手段により推定される音素セットに含まれる音素の種類の数より少ないとよい。

20

【0035】

音素の種類より視覚素ラベルの数が少ないので、最終的な特徴点の位置の推定が効率的に行なえる。

【0036】

ラベルセットは、各々が一つの音素を表す複数個の音素ラベルを含んでもよい。第2の統計確率モデルは、音素ラベル間の遷移確率と各特徴点の位置の出力確率とに関し予め学習することにより得られる。ラベル列生成手段は、音素列推定手段により推定された音素列をもとに、音素ラベルの系列を生成し、音素継続長をもとに当該系列を構成する音素ラベルの各々の継続長を決定するための手段を含む。

【0037】

モデルから顔の特徴点の位置データの系列を推定する場合、その組み合わせは膨大となるが、出力として音素ラベルの形をとれば、組合せが音素ラベルの数によって限定される。その結果、このモデルを用いた顔の形状の推定時にも、音素ラベル系列を得るようすることで、推定が効率的に行なえる。

30

【0038】

第2の統計確率モデルは、視覚素ラベル間の遷移確率と、各特徴点の位置パラメータ及び当該特徴点に関する動的特徴パラメータの出力確率とに関し予め学習された動的特徴による統計確率モデルを含んでもよい。軌跡推定手段は、ラベル間の遷移確率と各特徴点の位置パラメータ及び動的特徴パラメータの出力確率とに関し予め学習することにより得られた前記動的特徴による統計確率モデルに基づき、前記ラベル列生成手段により生成された系列と継続長とを入力パラメータとして、複数個の特徴点に対する位置パラメータ及び動的特徴パラメータの系列として最尤となる位置パラメータ及び動的特徴パラメータの系列を出力するための手段と、位置パラメータ及び動的特徴パラメータの系列に対し、当該パラメータが得られた統計確率モデルに固有の変換によって、位置パラメータを動的特徴パラメータを用いて補正し、複数個の特徴点の各々の軌跡を推定するための手段とを含む。

40

【0039】

このように動的特徴パラメータまで含んで学習したモデルを用い、位置パラメータ系列と動的特徴パラメータの系列とを得た後に、位置パラメータ系列を動的特徴パラメータ系列を用いて補正すると、推定された後の特徴点の動きは滑らかでかつ自然なものとなる。

50

【 0 0 4 0 】

本発明の第 4 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、当該コンピュータを本発明の第 3 の局面に係るいずれかのパラメータ系列合成装置として動作させる。

【 0 0 4 1 】

本発明の第 5 の局面に係るリップシンクアニメーション作成システムは、第 1 の座標空間における複数のノードの座標値を用いて顔の形状を定義した所定の顔オブジェクトをもとに、音声に同期する顔のアニメーションを作成するためのリップシンクアニメーション作成システムである。リップシンクアニメーション作成システムは、本発明の第 3 の局面に係るいずれかのパラメータ系列合成装置と、音声の入力に対してパラメータ系列合成装置により合成される、発話者の顔の複数個の特徴点の軌跡を表すパラメータ系列に基づき、顔オブジェクトにおけるノードの座標値を変更することにより、顔の形状を定義するオブジェクトを、アニメーションのフレームごとに生成するための変形オブジェクト生成手段と、アニメーションの各フレームについて、変形オブジェクト生成手段により生成されるオブジェクトから、当該フレームにおける顔の画像を合成するための画像化手段とを含む。

10

【 発明を実施するための最良の形態 】

【 0 0 4 2 】

以下、図面を参照しつつ、本発明の実施の形態に係る顔アニメーションの作成システムについて説明する。なお、以下の説明に用いる図面では、同一の部品及びデータ等には同一の符号を付してある。それらの名称及び機能も同一である。したがって、それらについての説明は繰返さない。

20

【 0 0 4 3 】

< 第 1 の実施の形態 >

【 0 0 4 4 】

[構成]

図 1 に、本実施の形態に係る顔アニメーションの作成システム全体の構成をブロック図形式で示す。図 1 を参照して、この顔アニメーションの作成システム 4 0 は、キャラクタの声となる音声のデータ（以下、単に「音声データ」と呼ぶ。）4 2 と、キャラクタの無表情な顔の形状を定義するためのデータである顔オブジェクト 4 4 とから、キャラクタの声に同期してキャラクタの表情が変化する（すなわちリップシンクする）アニメーション 4 6 を作成するシステムである。

30

【 0 0 4 5 】

顔アニメーションの作成システム 4 0 は、学習用の音声の収録とその音声の発話中に発話者の顔の各器官に生じる位置変化（以下、この位置を「顔パラメータ」と呼ぶ。）の計測とを同時に行なうための収録システム 6 0 と、収録システム 6 0 により収録された学習用のデータを蓄積するための音声 - 顔パラメータ DB 6 2 と、音声と音素との関係をモデル化した音素 HMM 6 4 と、発話時の口の形状を表す最小単位である視覚素（viseme）と音素との対応関係を表す視覚素対応表 6 6 とを含む。

【 0 0 4 6 】

顔アニメーションの作成システム 4 0 はさらに、音素 HMM 6 4 及び視覚素対応表 6 6 を用いて、音声 - 顔パラメータ DB 6 2 から、発話時の口の形状（視覚素）と顔パラメータとの関係をモデル化した統計確率モデルの学習を行なうための学習システム 6 8 と、学習システム 6 8 による学習の結果得られる統計確率モデルである顔パラメータ HMM 5 0 と、音声データ 4 2 及び顔オブジェクト 4 4 をもとに、音素 HMM 6 4、視覚素対応表 6 6、及び顔パラメータ HMM 5 0 を用いてアニメーション 4 6 を作成するためのアニメーション作成システム 8 0 とを含む。顔アニメーションの作成システム 4 0 はさらに、ユーザがアニメーション作成システム 8 0 を操作するための表示装置 9 6 及び入力装置 9 8 を含む。

40

【 0 0 4 7 】

50

アニメーション作成システム 80 は、音声データ 42 をもとに、音素 HMM 64、視覚素対応表 66、及び顔パラメータ HMM 50 を用いて音声データ 42 に対応する顔パラメータの系列 84 を合成するための顔パラメータ合成部 82 と、合成された顔パラメータの系列 84 及び顔オブジェクト 44 をもとに、発話時のキャラクタの顔の形状モデル 92 をフレームごとに生成するためのマッピング部 90 と、マッピング部 90 により生成されたフレームごとの形状モデル 92 を画像に変換して、アニメーション 46 を生成するための画像化部 94 とを含む。

【0048】

収録システム 60

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 110 による発話音声と発話時における発話者 110 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 110 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、音声のデータ及びその発話時の顔パラメータのデータからなるデータセット 120 を作成し、音声 - 顔パラメータ DB 62 に格納するためのデータセット作成装置 122 とを含む。

【0049】

録画・録音システム 112 は、発話者 110 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 110 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。

【0050】

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。

【0051】

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 110 の頭部の予め定める多数の箇所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。

【0052】

図 3 に、発話者 110 に装着されるマーカの装着位置を模式的に示す。図 3 を参照して、発話者 110 の顔、首、及び耳の多数の箇所 160A, ..., 160M 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加工されている。マーカの大きさは数 mm 程度である。音声 - 顔パラメータ DB 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 110 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。図 3 に示す例では、181 箇所の特徴点 160A, ..., 160M にそれぞれマーカが配置される。

【0053】

顔の物理的な構造上、発話者 110 の顔の表面上には、頭自体の動きに追従して移動するが発話者 110 の表情変化の影響をほとんど受けない、という特徴を持つ箇所がある。

例えばこめかみ，鼻の先端がこのような特徴を持つ。本実施の形態では、このような箇所も特徴点として定めておく。以下、このような特徴点を不動点と呼ぶ。後述する正規化処理のために4点以上の不動点を定めることが望ましい。

【0054】

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。

10

【0055】

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、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、音声・動画データ116及びMoCapデータ118をそれらに付されたタイムコードに基づいて切出し、互いに同期する音声のデータ（以下、「収録音声データ」と呼ぶ。）150及びMoCapデータ152を出力するための切出処理部144とを含む。

【0056】

データセット作成装置122はさらに、切出されたMoCapデータ152における頭の動きの成分をキャンセルするように当該MoCapデータ152を正規化して、顔の各器官の変化を表す顔パラメータの系列154に変換するための正規化処理部146と、収録音声データ150及び顔パラメータの系列154を同期させて結合してデータセット120を生成し、音声・顔パラメータDB62に格納するための結合部148とを含む。

20

【0057】

正規化処理部146は、切出された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。

【0058】

ここに、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 $P = (P_x, P_y, P_z, 1)$ と表現する。また時刻 $t = 0$ におけるマーカデータを $P' = (P'_x, P'_y, P'_z, 1)$ と表現する。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。

30

【0059】

【数1】

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

40

顔パラメータの系列154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。そして、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。

【0060】

50

音声 - 顔パラメータ D B 6 2

図 4 に、音声 - 顔パラメータ D B 6 2 (図 1 参照) に格納されるデータセット 1 2 0 の構成を模式的に示す。図 4 を参照して、データセット 1 2 0 は、切出された収録音声データ 1 5 0 と顔パラメータの系列 1 5 4 とを含む。顔パラメータの系列 1 5 4 は、複数フレーム分の顔パラメータ 1 7 0 A, ..., 1 7 0 N を含む。顔パラメータ 1 7 0 A, ..., 1 7 0 N はそれぞれ、収録音声データにより表される音声が発話されていた期間内のいずれかの時刻に対応している。すなわち、収録音声データ 1 5 0 と顔パラメータ 1 7 0 A, ..., 1 7 0 N とを相互参照することにより、ある特徴を持つ発話が行なわれたときの特徴点の位置変化について情報を得ることができる。

【 0 0 6 1 】

音素 H M M 6 4

図 1 に示す音素 H M M 6 4 は、音素ごとに設けられた音声の特徴に関する H M M である。図 5 に、音素 H M M 6 4 の概要を示す。図 5 を参照して、音素 H M M 6 4 は、音声の特徴を表す所定のパラメータ (以下、「音声パラメータ」と呼ぶ。) 1 8 0 が与えられると、音声中に対応する音素が存在する尤度を出力する機能を持つ。したがって、音素 H M M 6 4 を用いることにより、音声パラメータ 1 8 0 から音素列 1 8 2 及び音素列を構成する各音素の音素継続長が推定できる。なお本実施の形態では、音声パラメータ 1 8 0 として、M F C C (Mel-Frequency Cepstral Coefficient) を使用するものとする。

【 0 0 6 2 】

視覚素対応表 6 6

図 1 に示す視覚素対応表 6 6 は、音素と視覚素との対応関係を示す表である。図 6 に視覚素対応表 6 6 の構成を示す。図 6 を参照して、視覚素対応表 6 6 は、発話中の口の形を表す 1 0 種類の視覚素と、4 3 種類の音素との対応関係を表す。例えば視覚素「A」は、音素「a」又は「A」が発話されるとき口の形を表す。音素「h」を発話する際の口の形は、前後の音素を発話する際の口の形に依存する。そのため、この音素に対応する 1 0 種の視覚素とは別に記号「***」によって表している。

【 0 0 6 3 】

学習システム 6 8

図 7 に、学習システム 6 8 (図 1 参照) の構成をブロック図で示す。図 7 を参照して、学習システム 6 8 は、音声 - 顔パラメータ D B 6 2 内のデータセット 1 2 0 から顔パラメータ H M M 5 0 の学習に用いるデータセットを生成するための前処理部 2 0 2 と、学習用のデータセット 2 0 0 を蓄積するための学習用 D B 2 0 4 と、学習用 D B 2 0 4 内に蓄積された学習用のデータセット 2 0 0 から顔パラメータ H M M 5 0 の学習を行なうための H M M 学習部 2 0 6 とを含む。

【 0 0 6 4 】

前処理部 2 0 2 は、音声 - 顔パラメータ D B 6 2 から処理対象のデータセット 1 2 0 を選択するためのデータセット選択部 2 1 0 と、データセット選択部 2 1 0 により選択されたデータセット 1 2 0 内の収録音声データ 1 5 0 (図 4 参照) から、データセット 1 2 0 が収録された際の発話内容に対応する視覚素の系列と各視覚素の継続長とを、音素 H M M 6 4 及び視覚素対応表 6 6 を用いて推定するための視覚素列生成部 2 1 2 と、選択されたデータセット 1 2 0 内の顔パラメータの系列 1 5 4 に含まれる顔パラメータ 1 7 0 A, ..., 1 7 0 N (図 4 参照) に、視覚素を表すラベルによるラベルリングを行ない、学習用のデータセット 2 0 0 を生成するためのラベリング部 2 1 4 とを含む。

【 0 0 6 5 】

視覚素列生成部 2 1 2 は、収録音声データ 1 5 0 から、音声パラメータ 1 8 0 (図 5 参照) を抽出するための特徴量抽出部 2 2 0 と、抽出された音声パラメータ 1 8 0 をもとに、発話に対応する音素として最尤のものをフレームごとに、音素 H M M 6 4 を用いて推定するための音素列推定部 2 2 2 と、音素列推定部 2 2 2 により推定された音素列 1 8 2 を構成する各音素を視覚素対応表 6 6 に基づき視覚素に変換するための音素・視覚素変換部 2 2 4 とを含む。

10

20

30

40

50

【 0 0 6 6 】

特徴量抽出部 2 2 0 は、収録音声データから、音声パラメータ 1 8 0 として各フレームにおける MFCC ベクトルを抽出する機能を持つ。音素列推定部 2 2 2 は、音素 HMM 6 4 から、最尤の音素列 1 8 2 及び音素継続長を推定する機能を持つ。本実施の形態の音素列推定部 2 2 2 は、この推定にビタビアルゴリズムを用いる。すなわち、音素列推定部 2 2 2 は、与えられた MFCC ベクトルの系列を出力する音素のビタビ系列を推定する。音素・視覚素変換部 2 2 4 は、推定された音素のビタビ系列を構成する音素をそれぞれ、視覚素に変換する機能を持つ。ラベリング部 2 1 4 は、発話時の視覚素の時間変化と顔パラメータとの対応付けを、データセット 1 2 0 の各々について行なう。

【 0 0 6 7 】

HMM 学習部 2 0 6 は、視覚素によりラベリングされた顔パラメータ 1 7 0 A, ..., 1 7 0 N を用いて、所定の視覚素列が与えられた場合の顔パラメータ 1 7 0 A, ..., 1 7 0 N の系列とその尤度とを学習する機能を持つ。ただし、発話中の表情は、音素と音声パラメータとの関係における調音結合と同様に、前後の視覚素に依存して変化することがある。音声認識用の音素 HMM の学習においては、調音結合に対処すべくトライフォンを認識の処理単位として用いることがある。そこで、本実施の形態では、三つの視覚素を 1 組とする三つ組視覚素 (TriViseme) を処理単位として、顔パラメータ HMM 5 0 の学習を行なう。

【 0 0 6 8 】

アニメーション作成システム 8 0

(顔パラメータ合成部 8 2)

図 8 に、顔パラメータ合成部 8 2 (図 1 参照) の構成をブロック図で示す。図 8 を参照して、顔パラメータ合成部 8 2 は、音素 HMM 6 4 及び視覚素対応表 6 6 を用いて音声データ 4 2 を視覚素列に変換するための視覚素列生成部 2 4 0 と、音声データ 4 2 により表される音声が発話されているときの顔パラメータの系列 8 4 を、視覚素列生成部 2 4 0 により生成された視覚素列と顔パラメータ HMM 5 0 とを用いて合成するための HMM マッチング部 2 4 2 とを含む。

【 0 0 6 9 】

視覚素列生成部 2 4 0 は、特徴量抽出部 2 5 0 と、音素列推定部 2 5 2 と、音素・視覚素変換部 2 5 4 とを含む。これらの機能は、特徴量抽出部 2 5 0 が音声データ 4 2 の入力を受け、それぞれ、図 7 に示す学習システム 6 8 の特徴量抽出部 2 2 0、音素列推定部 2 2 2、及び音素・視覚素変換部 2 2 4 の機能と同一である。したがってこれらの機能の詳細については、説明は繰返さない。

【 0 0 7 0 】

HMM マッチング部 2 4 2 は、視覚素列生成部 2 4 0 より視覚素列とその継続長とを受け、当該視覚素列と継続長とにより表される発話全体で尤度最大となる顔パラメータの系列 8 4 を、顔パラメータ HMM 5 0 を用いて合成する機能を持つ。

【 0 0 7 1 】

(マッピング部 9 0)

図 9 に、マッピング部 9 0 (図 1 参照) の構成をブロック図で示す。図 9 を参照して、マッピング部 9 0 は、入力装置 9 8 及び表示装置 9 6 に接続され、顔オブジェクト 4 4 上に特徴点 1 6 0 A, ..., 1 6 0 M (図 3 参照) に対応する仮想のマーカ (以下、単に「仮想マーカ」と呼ぶ。) を、ユーザの操作にしたがい配置するための仮想マーカ配置部 2 7 0 と、顔オブジェクト 4 4 内の各ノードを、各ノードに近接する仮想マーカによってラベリングするためマーカラベリング部 2 7 2 と、マーカラベリング部 2 7 2 によるラベリングにより形成されたノードと仮想マーカとの対応関係を表す、マーカラベリングデータを記憶するためのマーカラベリングデータ記憶部 2 7 4 とを含む。

【 0 0 7 2 】

マッピング部 9 0 はさらに、マーカラベリングデータ記憶部 2 7 4 に記憶されたマーカラベリングデータと顔パラメータ合成部 8 2 により合成された顔パラメータの系列 8 4 と

10

20

30

40

50

を用いて、顔オブジェクト 4 4 が表現する顔の形状から、変形した顔オブジェクト 9 2 を順次作成するための顔オブジェクト変形部 2 7 6 を含む。

【 0 0 7 3 】

仮想マーカ配置部 2 7 0 は、入力装置 9 8 及び表示装置 9 6 を用いてユーザにより行なわれる仮想マーカの配置操作にしたがい、顔オブジェクト 4 4 を規定する座標系上での各特徴点の座標を設定する。このようにして特徴点の座標が設定されることにより、各特徴点の各マーカデータを顔オブジェクト 4 4 上の各仮想マーカの位置に割当てることができる。なおこの際、顔パラメータの座標系と顔オブジェクトの座標系との間の変換も行なわれる。

【 0 0 7 4 】

図 1 0 に、顔オブジェクト 4 4 及び仮想マーカの一例を示す。図 1 0 を参照して、顔オブジェクト 4 4 は、この図における黒い線分でそれぞれ示すエッジにより囲まれた多数の多角形（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である。ポリゴンの頂点（エッジ同士の交点）が、顔オブジェクト 4 4 におけるノードである。一般に顔には、目・口・鼻の穴のように顔面を構成しない切れ目がある。これらの切れ目は一般に、顔オブジェクト 4 4 の一部としてモデリングされることはない。すなわち切れ目にポリゴンを定義しない。又は顔オブジェクト 4 4 とは別のオブジェクトとして定義される。よって、切れ目と顔面との間は境界エッジで仕切られている。

【 0 0 7 5 】

顔オブジェクト 4 4 により表現される顔の形状は、ユーザにより創作される任意のものでよい。ただし、顔パラメータを用いて顔オブジェクト 4 4 に表情を付与するには、顔オブジェクト 4 4 により表現される形状のどの部分が顔の各器官であるかを定義する必要がある。そのために、仮想マーカ配置部 2 7 0（図 9 参照）は、顔オブジェクト 4 4 上に、仮想マーカ 3 0 0 A, ..., 3 0 0 M をそれぞれ、ユーザの操作にしたがって配置する。

【 0 0 7 6 】

この際、収録システム 6 0（図 2 参照）におけるモーションキャプチャデータの収録に用いられたマーカの装着順序に従い仮想マーカ 3 0 0 A, ..., 3 0 0 M が配置されるよう、ユーザに対し誘導を行なう。したがって、ユーザの主観を反映しつつ、適切な位置に仮想マーカを配置することができる。図 9 に示す仮想マーカ配置部 2 7 0 は、顔オブジェクト 4 4 を規定する座標系における各仮想マーカの座標をマーカラベリング部 2 7 2 に出力する。

【 0 0 7 7 】

マーカラベリング部 2 7 2 は、顔オブジェクト 4 4 のノードの中から、処理の対象となるノードを選択し、選択したノード（以下、「選択ノード」と呼ぶ。）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。そして、選択された仮想マーカ（以下、「選択マーカ」と呼ぶ）が、この選択ノードに対応付ける仮想マーカとして適切であるかを判定する。適切であれば選択マーカを選択ノードの対応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、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。

【 0 0 7 8 】

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、顔オブジェクトの境界エッジを利用する。

【 0 0 7 9 】

図 1 1 に、マーカラベリング部 2 7 2 により実行されるマーカラベリング処理の構造をフローチャートで示す。図 1 1 を参照して、処理が開始されると、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、顔オブジェクト 4 4 における各ノードに対して実行する。

【 0 0 8 0 】

ステップ 3 4 2 では、選択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

10

20

30

40

50

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。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 及び採用された対応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。

【 0 0 8 1 】

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 M_{max} を超えていればエラーとし、処理を終了する。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 M_{max} 以下であれば制御はステップ 3 4 8 に進む。

【 0 0 8 2 】

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。) と選択ノードとを結ぶ線分が、顔オブジェクト 4 4 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。

【 0 0 8 3 】

ステップ 3 5 0 では、この時点でのマーカ (i) を選択ノードの対応マーカの一つに指定する。すなわちマーカ (i) を示す情報を、選択ノードのマーカ・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。

【 0 0 8 4 】

上記したように、選択ノードと仮想マーカとを結ぶ線分が顔オブジェクトの境界エッジを横切るものは、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がある。この場合、上唇に位置するノードと、下唇に位置するノードと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オブジェクト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、一つに属しているかによって判定する。

【 0 0 8 5 】

図 1 2 に、顔オブジェクト 4 4 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 1 2 を参照しつつ、当該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 1 2 を参照して、顔オブジェクト 4 4 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、三つのエッジに囲われている。そして上唇と下唇の間には境界エッジ 4 0 0 が存在する。境界エッジは、顔オブジェクト 4 4 と切れ目との接線、又は顔オブジェクト 4 4 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 4 0 0 に該当するエッジは共有されない。

【 0 0 8 6 】

マーカラベリング部 2 7 2 はまず、顔オブジェクト 4 4 を構成するノードの中からノードを一つ選択する。このノードが選択ノードである。ここに、図 1 2 に示すノード 4 1 0 が選択ノードであるものとする。選択ノード 4 1 0 の近隣には、仮想マーカ 4 1 2 A, ..., 4 1 2 E が存在する。マーカラベリング部 2 7 2 は、ノード 4 1 0 の座標と、仮想マーカの座標とをもとに、選択ノード 4 1 0 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 4 1 0 に最も近い位置にある仮想マーカ 4 1 2 A を選択する。

【 0 0 8 7 】

続いて、マーカラベリング部 2 7 2 は、選択ノード 4 1 0 と仮想マーカ 4 1 2 A とを結ぶ線分 4 1 4 A が境界エッジ 4 0 0 を横切るか否かを検査する。この線分 4 1 4 A は、境界エッジ 4 0 0 を横切らない。そのためマーカラベリング部 2 7 2 は、仮想マーカ 4 1 2 A を選択ノード 4 1 0 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

10

20

30

40

50

ーカ 4 1 2 A の次にノード 4 1 0 に近い位置にある仮想マーカ 4 1 2 B を選択し検査を行なう。選択ノード 4 1 0 と仮想マーカ 4 1 2 B とを結ぶ線分 4 1 4 B は、境界エッジ 4 0 0 を横切っている。そのため、仮想マーカ 4 1 2 B は選択ノード 4 1 0 の対応マーカからは除外される。

【 0 0 8 8 】

マーカラベリング部 2 7 2 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 4 1 0 の対応マーカ (図 1 2 に示す例では仮想マーカ 4 1 2 A 、 4 1 2 D 、 及び 4 1 2 E) を選択する。

【 0 0 8 9 】

再び図 9 を参照して、顔オブジェクト変形部 2 7 6 は、あるフレームの顔パラメータにおける各マーカデータをそれぞれ仮想マーカに付与する。さらに顔オブジェクト変形部 2 7 6 は、マーカラベリングデータ記憶部 2 7 4 のマーカラベリングデータに基づき、顔オブジェクト 4 4 の各ノードに、対応する仮想マーカの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、顔オブジェクト 4 4 の変形を行なう。そして、変形後の顔オブジェクト 4 4 を、形状モデル 9 2 として出力する。顔オブジェクト 4 4 のノードの座標を N 、当該ノードと対応関係にある仮想マーカの座標を M_i 、変形後の顔オブジェクトである形状モデル 9 2 におけるマーカの座標を M'_i とすると、顔オブジェクト変形部 2 7 6 は、ノードの座標の変化量ベクトル v を次の内挿式 (2) によって算出する。

【 0 0 9 0 】

【 数 2 】

$$v = \sum_i^n (M'_i - M_i) \cdot \left(\frac{1.0}{\|N - M_i\|} \right) \cdot \left(\sum_i^m \frac{1.0}{\|N - M_i\|} \right)^{-1} \quad (2)$$

【 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 4 0 は以下のように動作する。

【 0 0 9 1 】

収録システムの動作

以下に、収録システム 6 0 が収録を行ない、データセット 1 2 0 を生成する動作について説明する。図 2 を参照して、発話者 1 1 0 の頭部の各特徴点 1 6 0 A , ... , 1 6 0 M (図 3 参照) には事前に、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 - 顔パラメータ D B 6 2 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 1 1 0 にその内容で発話を行なってもらうようにしてもよい。

【 0 0 9 2 】

収録が開始され、発話者 1 1 0 が発話すると、録画・録音システム 1 1 2 が、発話時の音声と顔の動画を次のようにして収録する。すなわち、マイクロホン 1 3 0 A 及び 1 3 0 B は、発話者 1 1 0 の音声を受音して音声信号を発生する。カムコーダ 1 3 2 は、発話中の発話者 1 1 0 の動画を撮影し、その映像信号をマイクロホン 1 3 0 A 及び 1 3 0 B からの音声信号を同時に記録して音声・動画データ 1 1 6 を生成する。この際、カムコーダ 1 3 2 は、M o C a p システム 1 1 4 に対してタイムコード 1 3 4 を供給するとともに、音声・動画データ 1 1 6 に、タイムコード 1 3 4 と同じタイムコードを付与する。

【 0 0 9 3 】

この際、同時に、発話時における特徴点 1 6 0 A , ... , 1 6 0 M の位置が、M o C a p システム 1 1 4 により次のようにして計測される。マーカはそれぞれ、対応する特徴点の動きに従って移動する。赤外線カメラ 1 3 6 A , ... , 1 3 6 F 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 1 2 0 フレーム) で撮影しその映像信号をデータ処理装置 1 3 8 に出力する。データ処理装置 1 3 8 は、それらの映像信号の各フレームにタイムコード 1 3 4 を付与し、当該映像信号をもとに、各マーカの位置を

10

20

30

40

50

フレームごとに算出する。データ処理装置 138 は、各マーカの位置のデータをフレームごとにまとめて MoCap データ 118 として蓄積する。

【0094】

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。

【0095】

切出処理部 144 はまず、MoCap データ記憶部 142 から、 $t = 0$ のフレームにおける MoCap データを読み出して正規化処理部 146 に与える。このフレームのデータは、正規化処理部 146 による正規化に用いられる。続いて切出処理部 144 は、音声・動画記憶部 140 に記憶される音声・動画データ 116 から、1 発話分など所定の単位で収録音声データ 150 を切出す。そして、切出した収録音声データ 150 に付与されているタイムコードを参照して、収録音声データ 150 の当該タイムコード上での位置を特定し、収録音声データ 150 を結合部 148 に与える。続いて切出処理部 144 は、MoCap データ 118 から、タイムコード上、収録音声データ 150 の位置と同じ位置で MoCap データ 152 を切出し、正規化処理部 146 に与える。

【0096】

正規化処理部 146 は、MoCap データ 152 の各フレームにおいて、当該フレームの不動点のマーカデータと、予め与えられている $t = 0$ のフレームにおける不動点のマーカデータとから、アフィン行列を求め、当該アフィン行列を用いて、各マーカデータをアフィン変換する。この変換により、変換後のマーカデータはそれぞれ、頭を $t = 0$ 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。その結果、MoCap データ 152 は、顔パラメータの系列 154 になる。顔パラメータの系列 154 は、結合部 148 に与えられる。

【0097】

結合部 148 は、収録音声データ 150 及び顔パラメータの系列 154 を同期させて結合してデータセット 120 (図 4 参照) を生成し、音声 - 顔パラメータ DB 62 に格納する。

【0098】

顔パラメータ HMM 50 の学習

以下に、学習システム 68 が顔パラメータ HMM を学習する動作について説明する。図 7 を参照して、音声 - 顔パラメータ DB 62 内のデータセット 120 (図 4 参照) の各々は、学習システム 68 の前処理部 202 により、次のようにして学習用データセット 200 に変換される。

【0099】

すなわちまず、データセット選択部 210 が、音声 - 顔パラメータ DB 62 から処理対象のデータセット 120 (図 4 参照) を選択する。そして、当該データセット 120 に含まれる収録音声データ 150 と顔パラメータの系列 154 とをそれぞれ、視覚素列生成部 212 とラベリング部 214 とに与える。

【0100】

視覚素列生成部 212 に収録音声データ 150 が与えられると、特徴量抽出部 220 が、収録音声データ 150 から、その音声の特徴量のベクトル系列 180 として、フレームごとに MFCC を抽出する。音素列推定部 222 は、抽出された MFCC ベクトルの系列に対応する音素列 182 (図 5 参照) を、音素 HMM 64 に基づきピタビアルゴリズムによって推定する。すなわち、与えられたベクトル系列から、発話全体で尤度最大となる音素列 182 及び当該音素列 182 を構成する各音素の音素継続長を推定する。音素・視覚素変換部 224 は、推定された音素列 182 を構成する音素をそれぞれ、視覚素に変換する。これにより 43 種類の音素は、10 種類の視覚素にグループ化される。したがって、視覚素変換部 224 により出力される視覚素列 208 として可能な組合せの数は、視覚素

10

20

30

40

50

変換部 224 に入力されうる音素列 182 の組合せの数より少なくなる。音素・視覚素変換部 224 により出力されるデータは、データセット 120 の各時刻において発話者 110 が発話する際の口の形に対応する視覚素を表す。

【0101】

ラベリング部 214 は、この視覚素列に基づき、顔パラメータの系列 154 内の各顔パラメータ 170A, ..., 170N に対するラベリングを行なう。発話時の視覚素の時間変化と顔パラメータとの対応付けを、データセット 120 の各々について行なうことになる。ラベリング部 214 は、視覚素によりラベリングされた顔パラメータ 170A, ..., 170N からなる学習用データセット 200 を生成し、これを学習用 DB 204 に格納する。

10

【0102】

HMM 学習部 206 は、作成された学習用 DB 204 に格納された、学習用データセット 200 を用いて、顔パラメータ HMM 50 の学習を行なう。この際 HMM 学習部 206 は、三つの視覚素を 1 組とする三つ組視覚素を処理単位として、顔パラメータ HMM 50 の学習を行ない、視覚素間の遷移確率と、顔パラメータ 170A, ..., 170N の出力確率に関する学習を行ない、顔パラメータ HMM 50 を形成する。

【0103】

以上のようにして顔パラメータ HMM 50 を学習することにより、顔パラメータ HMM 50 に基づき、視覚素列から顔パラメータの系列を合成することが可能になる。顔パラメータは、各フレームにおける顔の多数の特徴点 160A, ..., 160M (図 3 参照) の位置を表すものである。また、視覚素は発話時の口の形を表すものである。そのため、アニメーション上でのキャラクタの声に対応する各フレームのキャラクタの視覚素が特定されれば、当該視覚素からなる視覚素列と、顔パラメータ HMM 50 とを用いて、各フレームにおける顔の多数の特徴点 160A, ..., 160M の位置情報を合成することが可能になる。すなわち、視覚素列から、発話時の特徴点 160A, ..., 160M の軌跡を推定することができる。よって、発話時の口の形のみならず顔の表情の変化に関して、情報を得ることが可能になる。

20

【0104】

また、視覚素の種類は音素の種類より少ない。したがって、音素ごとに状態が設けられた HMM より、視覚素ごとに状態が設けられた HMM の方が、少ない状態数のモデルとなる。発話中の発話者の表情は、音素よりむしろ発話中の口の形に依存して変化すると考えられる。そのため、視覚素列から学習された顔パラメータ HMM 50 の品質が、音素列から、又は MFC の系列から学習された顔パラメータ HMM 50 の品質より劣ることはない。同一の量の学習データからの学習を行なう場合、状態数の少ないモデルを学習する方が、データのスパースネスな学習の問題（一部の領域で学習に用いるデータが不足しているため、正確な推定を行なうことが不可能となる問題）も生じず、効率的である。したがって、視覚素列から顔パラメータ HMM 50 を学習することにより、効率的で高い品質の顔パラメータ HMM を得ることができる。さらに、三つ組視覚素を処理単位として、HMM 学習を行なうため、前後の視覚素に依存した顔の表情の変化に対しても精度の高い学習を行なうことができる。

30

40

【0105】

(顔パラメータの合成)

以下、図 1 に示すアニメーション作成システム 80 の動作について説明する。キャラクタの声を表す音声データ 42 が準備され、図 8 に示す顔パラメータ合成部 82 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。又は、音声合成技術により合成された音声のデータであってもよい。顔パラメータ合成部 82 に音声データ 42 が入力されると、視覚素列生成部 240 が、音素 HMM 64 及び視覚素対応表 66 を用いて、音声データ 42 から視覚素列及び当該視覚素列を構成する各視覚素の継続長を推定する。この動作は、学習システム 68 の視覚素生成部 212 (図 7 参照) の動作と同様である。これによ

50

り、音声データ42により表される音声の発話時における口の形の変化が特定される。

【0106】

HMMマッチング部242は、視覚素列生成部240により生成された視覚素列と顔パラメータHMM50とのマッチングを行ない、発話全体で最尤の顔パラメータの系列84を合成する。

【0107】

以上のようにして顔パラメータ合成部82により合成された顔パラメータの系列84は、音声データ42により表現される音声の発話中における口の形の変化から得られたものである。よってこの系列84は、当該音声の発話時における顔の特徴点160A, ..., 160Mの軌跡を表すものとなる。したがって、発話時の口の形のみならず顔の各特徴点の位置の非線形的な変化を、合成された顔パラメータの系列84によって特定できる。

10

【0108】

また顔パラメータ合成部82は、音声データ42から、音素HMM64と顔パラメータHMM50とに用いた2段階の推定により顔パラメータの系列84を合成する。すなわち、音声データ42の音声パラメータ180の入力に対し出力されうる顔パラメータの系列84は、音素HMM64に基づく音素列182の推定により絞込まれることになる。さらに、音素を視覚素に変換することにより、出力され得る顔パラメータの系列84は、さらに絞込まれる。そのため、特徴点が多数存在する場合であっても、効率的に顔パラメータの系列84を合成することができる。

【0109】

20

ただし、上記の顔パラメータ合成部82により合成される顔パラメータは、図1に示す音声-顔パラメータDB62に格納された顔パラメータの系列154に基づき合成されるものである。すなわち、音声データ42により表される音声と等価な音声を、図2に示す収録システム60における発話者110が発話した場合の顔の表情変化を表すものである。そこで、本実施の形態に係るマッピング部90は、キャラクタの顔の形状を表す顔オブジェクト44と顔パラメータの系列84とから、発話時の各フレームに対応する形状モデル92を、以下のようにして生成する。

【0110】

(マッピングによる形状モデル92の生成)

図9を参照して、マッピング部90に顔オブジェクト44(図4参照)が与えられると、まず、顔オブジェクト44は、仮想マーカ配置部270、マーカラベリング部272、及び顔オブジェクト変形部276に与えられる。

30

【0111】

仮想マーカ配置部270は、顔オブジェクト44に仮想マーカ300A, ..., Mを、ユーザの操作にしたがい配置する。これにより、無表情な状態での顔オブジェクト44における、当該顔オブジェクト44の座標系上での特徴点160A, ..., 160M(図3参照)の位置が特定される。すなわち、仮想マーカ配置部270は、まず顔オブジェクト44を画像化して表示装置96に出力し、さらにユーザから当該初期顔モデル上における仮想マーカの位置の指定を入力装置98を介して受ける。顔オブジェクト44での仮想マーカの位置は、収録システム60における発話者へのマーカの配置と同様のルールにしたがって指定される。そのため、顔オブジェクト44と各仮想マーカとの位置関係は、発話者110(図2参照)の顔と当該発話者110に装着されたマーカとの位置関係に対応する。

40

【0112】

仮想マーカ配置部270は、各マーカのマーカデータに対しモーションキャプチャデータの座標系から顔モデルの座標系に対する座標変換を行ない、初期顔モデルの座標系における各仮想マーカの座標を特定する。仮想マーカ配置部270は、当該各仮想マーカの座標を、マーカラベリング部272に与える。

【0113】

マーカラベリング部272は、顔オブジェクト44と仮想マーカの座標とを受けて、顔オブジェクト44の各ノードに対して、当該ノードの3個の対応マーカを図11及び図1

50

2を参照して前述したようにして特定する。マーカラベリング部272は、全てのノードに対して対応マーカを決定し、ノードに対する対応マーカを表すマーカラベリングデータを作成し、各仮想マーカの座標とともに、マーカラベリングデータ記憶部274に記憶させる。

【0114】

顔オブジェクト変形部276は、顔パラメータの系列84と、顔オブジェクト44と、マーカラベリングデータとをもとに、次のようにして、各フレームにおける形状モデル92を作成する。

【0115】

顔オブジェクト変形部276は、顔パラメータの系列から84の1フレーム分が与えられると、マーカラベリングデータ記憶部274からマーカラベリングデータを読み出し、当該顔パラメータにおける各特徴点の位置に基づき、当該フレームの形状モデル92における各ノードの位置を次のようにして算出する。

【0116】

すなわち、顔オブジェクト変形部276はまず、顔オブジェクト44上における仮想マーカの座標を、マーカラベリングデータ記憶部274から取得する。仮想マーカはそれぞれ、顔パラメータにおける特徴点と対応関係にある。そこで、顔オブジェクト変形部276は、顔パラメータの系列84における1フレーム分のデータをもとに、仮想マーカの各々に、当該仮想マーカに対応する特徴点の位置を付与し、当該1フレーム分の変化後の各仮想マーカの座標を算出する。

【0117】

さらに顔オブジェクト変形部276は、一つのノードの変化量を、ノードに対し指定されたn個の対応マーカの座標をもとに、上記の内挿式(2)によって算出する。顔オブジェクト変形部276は、フレームごとに、顔オブジェクト44の各ノードに対しこの処理を実行する。これにより、各ノードの座標は変更され、変形した顔の形状モデル92がフレームごとに生成される。顔オブジェクト変形部276は、変形した顔の形状モデル92の各々を、画像化部94に与える。

【0118】

(画像化によるアニメーションの作成)

画像化部94は、フレームごとの変形した顔モデルを受けると、それらにテクスチャなどを付与するなど、所定のレンダリング処理を行なう。この処理により生成される画像が、アニメーション46における各フレームの画像となる。これら各フレームの画像により形成された動画が、アニメーション46となる。

【0119】

以上のように、本実施の形態に係るマッピング部90は、発話者の顔の多数の特徴点と、顔オブジェクト44の各ノードとを対応付ける。さらに、各特徴点についての計測データをもとに、顔オブジェクト44の軌跡を算出する。したがって、ノードの集合としての顔オブジェクトの時間的な変化が顔パラメータの系列84として得られ、これにより、アニメーション46を作成することができる。顔パラメータの系列84は、音声データ42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。

【0120】

本実施の形態のアニメーション作成システム80は、モデルベースでアニメーションを作成する。ユーザは、キャラクタの声に相当する音声データ42と、静止状態でのキャラクタの顔の形状を定義した顔オブジェクト44とを用意し、顔オブジェクト44上に特徴点をルールにしたがい指定するだけで、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。また、キャラクタの顔のデザインが制限されることなく、顔オブジェクト44が表すキャラクタ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

10

20

30

40

50

アニメーションを作成できる。

【 0 1 2 1 】

[コンピュータによる実現及び動作]

本実施の形態の顔アニメーションの作成システム 4 0 の各機能部は、収録システム 6 0 (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 3 はこのコンピュータシステム 5 0 0 の外観を示し、図 1 4 はコンピュータシステム 5 0 0 の内部構成を示す。

【 0 1 2 2 】

図 1 3 を参照して、このコンピュータシステム 5 0 0 は、F D (フレキシブルディスク) ドライブ 5 2 2 及び C D - R O M (コンパクトディスク読出専用メモリ) ドライブ 5 2 0 を有するコンピュータ 5 1 0 と、キーボード 5 1 6 と、マウス 5 1 8 と、モニタ 5 1 2 とを含む。

【 0 1 2 3 】

図 1 4 を参照して、コンピュータ 5 1 0 は、F D ドライブ 5 2 2 及び C D - R O M ドライブ 5 2 0 に加えて、ハードディスク 5 2 4 と、C P U (中央処理装置) 5 2 6 と、C P U 5 2 6、ハードディスク 5 2 4、F D ドライブ 5 2 2、及び C D - R O M ドライブ 5 2 0 に接続されたバス 5 3 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 (R O M) 5 2 8 と、バス 5 3 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 (R A M) 5 3 0 とを含む。コンピュータシステム 5 0 0 はさらに、プリンタ 5 1 4 を含んでいる。

【 0 1 2 4 】

ここでは示さないが、コンピュータ 5 1 0 はさらにローカルエリアネットワーク (L A N) への接続を提供するネットワークアダプタボードを含んでもよい。

【 0 1 2 5 】

コンピュータシステム 5 0 0 に顔アニメーションの作成システム 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、C D - R O M ドライブ 5 2 0 又は F D ドライブ 5 2 2 に挿入される C D - R O M 5 3 2 又は F D 5 3 4 に記憶され、さらにハードディスク 5 2 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 5 1 0 に送信されハードディスク 5 2 4 に記憶されてもよい。プログラムは実行の際に R A M 5 3 0 にロードされる。C D - R O M 5 3 2 から、F D 5 3 4 から、又はネットワークを介して、直接に R A M 5 3 0 にプログラムをロードしてもよい。

【 0 1 2 6 】

このプログラムは、コンピュータ 5 1 0 にこの実施の形態の顔アニメーションの作成システム 4 0 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 5 1 0 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 5 1 0 上で動作するオペレーティングシステム (O S) 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 4 0 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 5 0 0 の動作は周知であるので、ここでは繰返さない。

【 0 1 2 7 】

なお、上記の実施の形態では、学習システム 6 8 (図 7 参照) において、ラベリング部 2 1 4 は、音素・視覚素変換部 2 2 4 により変換された視覚素列及び各視覚素の継続長に基づくラベリングを行なった。しかし、本発明はこのような実施の形態には限定されない。例えば、ラベリング部 2 1 4 は、音素列推定部 2 2 2 により推定される音素列 1 8 2 及び音素継続長に基づくラベリングを行なうようにしてもよい。この場合、HMM学習部 2

10

20

30

40

50

06は、音素列182及び音素継続長から、顔パラメータHMMの学習を行なうことになる。また、この場合、図8に示す顔パラメータ合成部82のHMMマッチング部242は、顔パラメータ合成部82の音素列推定部252により推定される音素列182及び音素継続長をもとに、顔パラメータHMM50とのマッチングを行なうことになる。

【0128】

また、本実施の形態に係るシステムにおいて、顔の特徴点160A, ..., 160Mの位置及び数は、図3に示すようなものには限定されない。ただし、マッピングに用いる特徴点の数が多くなるほど、アニメーション46における顔の表情変化を正確かつ自然に表現するものとなる。また、特徴点の数が多くなるほど、リップシンクの同期性も向上する。アニメーション作成システム80は、アニメーション46を出力する代わりに、各フレームにおける形状モデル92を出力するようにしてもよい。このようにすると、形状モデル92と別のオブジェクト等とを組合わせてアニメーションを生成することも可能になる。

10

【0129】

<第2の実施の形態>

第1の実施の形態に係る顔アニメーションの作成システム40によれば、音声から自動的にアニメーションを作成することが可能になった。しかし、以下に述べるように、それだけでは例えば口の動きが不自然になるなど、さらに改良すべき点が存在している。

【0130】

図15(A)に、発話時の顔の画像から得た本来の口の動きを示し、図15(B)に、第1の実施の形態に係る顔アニメーションの作成システム40を用いて顔の画像のアニメーションを自動的に作成したときの口の動きを示す。図15(A)では、口の動きは滑らかである。一方、図15(B)に示すアニメーションでの口の動きは、図15(A)に示すものと概略で一致しているものの、詳細な点では多くのステップ状の段差が存在していることが分かる。これは、アニメーション画像上では口の開き方がステップ状に変化していることを示す。そのため、このアニメーションを見た場合、やや不自然な感じを受ける。こうした問題は、顔パラメータHMM50によって得られる顔の画像の各位置を決めるベクトルが、マッチングにより定まる顔パラメータHMM50の各状態における平均ベクトルからなるために生ずると考えられる。

20

【0131】

このようなステップ状のアニメーション画像の動きをより滑らかにするために、例えば顔の画像のパラメータ系列にローパスフィルタを適用したり、パラメータ系列にスプライン曲線による近似を行ったりすることも考えられる。しかしそのような方策をとる場合、得られる画像からはメリハリが失われてしまい、やはり自然な動きが得られないという問題がある。

30

【0132】

第2の実施の形態では、こうした問題を解決するために、顔の特徴点の位置だけでなく、それらの速度及び加速度という、動的特徴パラメータをも用いて顔パラメータHMMの学習を行なう。後に示すように、このように動的特徴パラメータまで含めて学習した顔パラメータHMMを用いることにより、元の顔の画像の動きによく似た、スムーズでメリハリがあり、かつ自然な顔の動きが得られる。なお、動的特徴パラメータは、音声認識の分野では広く用いられている特徴量である。

40

【0133】

以下、第2の実施の形態における顔パラメータHMMの学習の原理と、その顔パラメータHMMを用いた顔の画像の各特徴点の位置の決定方法とについて説明する。なお、以下に記載した、動的特徴を用いるHMMの学習及びHMMによるマッチング後の位置ベクトルの算出方法は、非特許文献5に教示されたものと同様である。

【0134】

学習に用いる顔の特徴点として、第1の実施の形態に用いたものと同数の特徴点を採用する場合、位置ベクトルに加えて速度及び加速度の情報を用いるので、一つの特徴点当たりのパラメータ数(ベクトル数)は第1の実施の形態におけるベクトル数の3倍となる。

50

ある時刻 t における、ある特徴点の静的位置ベクトル（無表情な顔における特徴点の位置を基準としたもの）を c_t 、サンプリング間隔を τ とする。この場合、時刻 t におけるこの特徴点の速度ベクトル Δc_t 及び加速度ベクトル $\Delta^2 c_t$ は一般に以下のように近似される。

【 0 1 3 5 】

【数 3】

$$\Delta c_t = \sum_{\tau=L_+^{(1)}}^{\tau=L_+^{(1)}} w_1(\tau) c_{t+\tau} \quad (1)$$

$$\Delta^2 c_t = \sum_{\tau=L_+^{(2)}}^{\tau=L_+^{(2)}} w_2(\tau) \Delta c_{t+\tau} \quad (2)$$

10

ただし $L^{(1)}$ 及び $L^{(2)}$ はそれぞれ、時刻 t における速度及び加速度の算出において、時刻 t の前後で考慮すべき位置ベクトル及び速度ベクトルを含む時間幅をサンプリング時間 τ を単位として表したものであり、 w_1 及び w_2 はそれぞれ、各時刻での速度ベクトル及び加速度ベクトルを算出するために使用する、位置ベクトル及び速度ベクトルに割当てる重みを示す。本実施の形態では、 $L^{(1)} = L^{(2)} = 1$ とし、また重み w_1 としては、連続する3つの重みとして $w_1 = [-0.5, 0, 0.5]$ という値を用い、重み w_2 としては同様に $w_2 = [0.25, -0.5, 0.25]$ を用いる。

20

【 0 1 3 6 】

また、このとき、HMMの出力ベクトル o_t を次のように表すものとし、出力ベクトル o_t の系列を O で表すものとする。

【 0 1 3 7 】

【数 4】

$$o_t = [c'_t, \Delta c'_t, \Delta^2 c'_t]^T$$

式 (1) (2) は、行列形式で表すと次のように書くことができる。

【 0 1 3 8 】

【数 5】

$$O = WC \quad (3)$$

ただし

$$C = [c_1, c_2, \dots, c_T]^T$$

ベクトル c_t が M 次元とすれば、 C 、 O はそれぞれ、 TM 次元及び $3TM$ 次元である。行列 W は、 $3TM$ 行 TM 列の重み行列である。行列 W の要素の一部は係数 1 、 $w_1(\tau)$ 、及び $w_2(\tau)$ であり、他の大部分の要素は 0 である。

【 0 1 3 9 】

ここでは式 (3) の条件の下で、出力ベクトル o_t からなる出力ベクトル系列 O が得られる確率を最大にするような行列 C を求めることが目的となる。一般に、そのような C は、次の線形方程式を解く事により定められることが分かっている。

【 0 1 4 0 】

【数 6】

$$W^T U^{-1} W C = W^T U^{-1} M^T \quad (4)$$

ただし、行列 M 、 U はそれぞれ、以下のように表される。

【 0 1 4 1 】

30

40

【数 7】

$$M = [\mu'_{q1}, \mu'_{q2}, \dots, \mu'_{qT}]$$

$$U = \text{diag}[U_{q1}, U_{q2}, \dots, U_{qT}]$$

μ_{qt} と U_{qt} とはそれぞれ、HMMの状態 q_t の平均ベクトルと共分散行列とである。

【0142】

式(4)はコレスキー分解又はQR分解を用いて $O(TM^3L^2)$ の演算量で解くことができる。ただし、ここでは $L=1$ である。

【0143】

このようにして、出力ベクトル o_t の系列 O から C を算出する演算式を求めることにより、出力ベクトルの系列が得られれば、それに対応する静的ベクトル C 、すなわち顔の特徴点の位置ベクトルを得ることができる。この位置ベクトルの値の算出には、元の顔の画像の位置ベクトルだけでなく、速度ベクトル及び加速度ベクトルという動的特徴が使用されているため、後述するように第1の実施の形態と比較してより滑らかな顔の画像を得ることができる。

【0144】

[構成]

図16を参照して、第2の実施の形態に係る顔アニメーションの作成システム640は、第1の実施の形態に係る顔アニメーションの作成システム40の収録システム60、音声・顔パラメータDB62、学習システム68、顔パラメータHMM50、アニメーション作成システム80に代えて、それぞれ上記したように顔の特徴点の位置ベクトルに加え、それらの速度ベクトル及び加速度ベクトルまでをデータとして処理する能力を持つ収録システム660、音声・顔パラメータDB662、学習システム668、顔パラメータHMM650、及びアニメーション作成システム680を含み、それによって図1に示すアニメーション46よりも自然で、かつ滑らかなアニメーション646を作成する。図16において、図1に示す顔アニメーションの作成システム40の各部品と同一の部品には同一参照符号を付してある。それらの名称及び機能も同一である。したがって、ここではそれらについての詳細な説明は繰返さない。

【0145】

図16から分かるように、アニメーション作成システム680が図1に示すアニメーション作成システム80と異なるのは、図1の顔パラメータ合成部82に代えて、音声データ42をもとに、音素HMM64、視覚素対応表66、及び顔パラメータHMM650を用いて、音声データ42に対応する顔パラメータであって、速度ベクトル及び加速度ベクトルまで考慮して定められたパラメータの系列684を合成しマッピング部90に与えるための顔パラメータ合成部682を含むことである。その他の点においては、アニメーション作成システム680は、図1に示すアニメーション作成システム80と同様の構成を持っている。

【0146】

図17に、第2の実施の形態における収録システム660の詳細な構成を示す。図17を参照して、収録システム660は、図2に示す収録システム60と同様の構成を有する。異なるのは、図2のデータセット作成装置122に代えて、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MoCapデータ118から、音声のデータ及びその発話時の動的特徴パラメータを含む顔パラメータのデータからなるデータセット720を作成し、音声・顔パラメータDB662に格納するためのデータセット作成装置722を含む点である。

【0147】

データセット作成装置722の構成も、図2に示すデータセット作成装置122の構成とほぼ同一である。ただしデータセット作成装置722は、図2に示す正規化処理部146の後に、正規化処理部146の出力する顔パラメータの系列154を受け、顔パラメー

10

20

30

40

50

タの系列154に含まれる各時刻における各特徴点の静止位置ベクトルから、各時刻における各特徴点の速度ベクトル及び加速度ベクトルを前述した重み w_1 及び w_2 を用いて算出し、動的特徴を含む顔パラメータの系列754を出力する動的特徴算出部746を含む点と、図2に示す結合部148に代えて、動的特徴算出部746から動的特徴を含む顔パラメータの系列754を受け、切出処理部144から受ける収録音声データ150と、動的特徴算出部746から受ける動的特徴を含む顔パラメータの系列754とを同期させて結合して動的特徴を含むデータセット720を生成し、音声・顔パラメータDB662に格納するための結合部748を含む点とにおいて、データセット作成装置122とは異なっている。

【0148】

図18に、結合部748が出力する動的特徴を含むデータセット720の構成を概略的に示す。図18を参照して、動的特徴を含むデータセット720は、図4に示す単なる顔パラメータの系列154に代えて、動的特徴と顔パラメータの系列とを組合せた動的特徴を含む顔パラメータの系列754を含む点でデータセット120と異なる。

【0149】

動的特徴を含む顔パラメータの系列754は、図4に示す複数フレーム分の顔パラメータ170A, ..., 170Nに加え、これらフレームの顔パラメータから算出される速度ベクトルパラメータ(以下「速度パラメータ」と呼ぶ。)772A, ..., 772Nと、同じくこれらフレームの顔パラメータから算出される加速度ベクトルパラメータ(以下「加速度パラメータ」と呼ぶ。)774A, ..., 774Nを含む。これら複数フレーム分の顔パラメータ170A, ..., 170N、速度パラメータ772A, ..., 772N、及び加速度パラメータ774A, ..., 774Nは、各フレームごとに対応付けられている。これらは図4のデータセット120においてと同様、収録音声データ150と同期されている。すなわち、収録音声データ150と顔パラメータ170A, ..., 170N、速度パラメータ772A, ..., 772N、及び加速度パラメータ774A, ..., 774Nとを相互参照することにより、ある特徴を持つ発話が行なわれたときの特徴点の位置、その速度、及び加速度についての情報を得ることができる。

【0150】

図19に、学習システム668のブロック図を示す。図19を参照して、学習システム668は、概略、図7に示す学習システム68と同様の構成を持つ。異なるのは、図7に示す前処理部202に代えて、前処理部202と同様ではあるが、動的特徴を含むデータセット720を処理して学習用のデータセット700を出力することができる前処理部802を含む点と、図7の学習用DB204に代えて、動的特徴を含む学習用のデータセット700を蓄積するための学習用DB804を含む点と、図7に示すHMM学習部206に代えて、学習用DB804に格納された動的特徴を含む学習用のデータセットを用い、顔パラメータHMM650の学習を行なうためのHMM学習部806を含む点とである。

【0151】

前処理部802は、図7に示す前処理部202と同様の構成を持つが、データセット選択部210に代えて、音声・顔パラメータDB662から処理対象のデータセット720を選択する機能を持つデータセット選択部810を含む点と、選択されたデータセット720内の動的特徴を含む顔パラメータの系列754に含まれる顔パラメータ170A, ..., 170N、速度パラメータ772A, ..., 772N及び加速度パラメータ774A, ..., 774N(図18参照)に対し、音素・視覚素変換部224が出力する視覚素のラベルによるラベルリングを行ない、学習用のデータセット800を生成するためのラベルリング部814を含む点とで前処理部202と異なっている。

【0152】

図20に、HMM学習部80による学習が行なわれた後の、一つの視覚素に対応する顔パラメータHMM780の構成を簡単に示す。図20に示すように、この顔パラメータHMM780は3状態 $S_1 \sim S_3$ のHMMであって、各状態 $S_1 \sim S_3$ はそれぞれ、出力 $o_i = (c_i, c_i, c_i)$ ($i = 1 \sim 3$)の出力確率を与える確率分布と、遷移確

10

20

30

40

50

率を与える確率分布とを含んでいる。与えられる出力 o_i の系列と、顔パラメータ HMM 780 とのマッチングによって、そうした出力系列を与える尤度が最大となるような顔パラメータ HMM 780 の系列を求めることにより、各時刻における顔パラメータが、その時刻に対応する HMM によって定まる。その顔パラメータから、前述した式 (4) を用いて行列 C を算出することで、動的特徴量を考慮した、滑らかな変化をする顔の特徴点の座標を得ることができる。

【0153】

図 21 に、図 16 に示す顔パラメータ合成部 682 のより詳細な構成を示す。図 21 を参照して、顔パラメータ合成部 682 は、図 8 に示す第 1 の実施の形態の顔パラメータ合成部 82 とよく似た構成を持つ。異なる点は、図 8 の HMM マッチング部 242 に代えて、視覚素列生成部 240 により生成された視覚素列と顔パラメータ HMM 650 とをマッチングすることにより、音声データ 42 により表される音声が発話されているときの、顔パラメータ HMM 650 からの出力パラメータの系列 844 を生成し出力するための HMM マッチング部 842 を含む点と、HMM マッチング部 842 から出力される動的特徴量を含む出力パラメータの系列 844 に対し、前述した式 (4) を用いた変換を行ない、動きベクトル及び加速度ベクトルまで考慮した特徴点の位置ベクトル系列、すなわち顔パラメータの系列 684 (式 (4) における行列 C) を出力するための変換部 846 をさらに含む点とである。

【0154】

HMM マッチング部 842 は、視覚素列生成部 240 より視覚素列とその継続長とを受け、当該視覚素列と継続長とにより表される発話全体で尤度最大となるような、動的特徴量を含む出力パラメータの系列 844 を、顔パラメータ HMM 650 を用いて合成する機能を持つ。

【0155】

[動作]

この第 2 の実施の形態に係る顔アニメーションの作成システム 640 の各部のうち、第 1 の実施の形態の顔アニメーションの作成システム 40 内の部品と同一か又は対応する部品の動作は、その部品と同様である。ただし、扱うデータに動的特徴量が含まれている点が異なる。以下、第 1 の実施の形態に係るシステム 40 の動作とは異なる点に重点をおき、顔アニメーションの作成システム 640 の動作について説明する。

【0156】

収録システムの動作

図 17 を参照して、発話者 110 の頭部の各特徴点 160A, ..., 160M (図 3 参照) には事前に、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。収録が開始され、録画・録音システム 112 が、発話時の音声と顔の動画像を収録する。

【0157】

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 722 に与えられる。データセット作成装置 722 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。

【0158】

切出処理部 144 はまず、MoCap データ記憶部 142 から、 $t = 0$ のフレームにおける MoCap データを読み出して正規化処理部 146 に与える。このフレームのデータは、正規化処理部 146 による正規化に用いられる。続いて切出処理部 144 は、音声・動画記憶部 140 に記憶される音声・動画データ 116 から、1 発話分など所定の単位で収録音声データ 150 を切出す。そして、切出した収録音声データ 150 に付与されているタイムコードを参照して、収録音声データ 150 の当該タイムコード上での位置を特定し、収録音声データ 150 を結合部 748 に与える。続いて切出処理部 144 は、MoCap データ 118 から、タイムコード上、収録音声データ 150 の位置と同じ位置で MoCap データ 152 を切出し、正規化処理部 146 に与える。

【 0 1 5 9 】

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マーカデータと、予め与えられている $t = 0$ のフレームにおける不動点のマーカデータとから、アフィン行列を求め、当該アフィン行列を用いて、各マーカデータをアフィン変換する。この変換により、変換後のマーカデータはそれぞれ、頭の位置を $t = 0$ での位置に保った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。その結果、M o C a p データ 1 5 2 は、顔パラメータの系列 1 5 4 になる。顔パラメータの系列 1 5 4 は、動的特徴算出部 7 4 6 に与えられる。

【 0 1 6 0 】

動的特徴算出部 7 4 6 は、前述した式 (1) (2) と、重み $w_1 = [- 0 . 5 , 0 , 0 . 5]$ 、及び重み $w_2 = [0 . 2 5 , - 0 . 5 , 0 . 2 5]$ とを使用して、各時刻における動的特徴量 (速度ベクトル及び加速度ベクトル) を算出して顔パラメータの系列 1 5 4 とあわせ、動的特徴を含む顔パラメータの系列 7 5 4 を結合部 1 4 8 に与える。

【 0 1 6 1 】

結合部 7 4 8 は、収録音声データ 1 5 0 及び動的特徴を含む顔パラメータの系列 7 5 4 を同期させて結合して動的特徴を含むデータセット 7 2 0 を生成し、音声 - 顔パラメータ DB 6 6 2 に格納する。

【 0 1 6 2 】

顔パラメータ H M M 6 5 0 の学習

まず、データセット選択部 8 1 0 が、音声 - 顔パラメータ DB 6 6 2 から処理対象のデータセット 7 2 0 を選択する。そして、当該データセット 7 2 0 に含まれる収録音声データ 1 5 0 と動的特徴を含む顔パラメータの系列 7 5 4 とをそれぞれ、視覚素列生成部 2 1 2 とラベリング部 8 1 4 とに与える。

【 0 1 6 3 】

視覚素列生成部 2 1 2 は、第 1 の実施の形態の場合と同様に動作し、音声に対応する音素列を推定し、さらに各音素に対応する視覚素からなる視覚素列 2 0 8 を生成し、ラベリング部 8 1 4 に与える。ラベリング部 8 1 4 は、視覚素列 2 0 8 に基づき、動的特徴を含む顔パラメータの系列 7 5 4 内の各顔パラメータ $1 7 0 A , \dots , 1 7 0 N$ 、速度パラメータ $7 7 2 A , \dots , 7 7 2 N$ 、及び加速度パラメータ $7 7 4 A , \dots , 7 7 4 N$ に対するラベリングを行なう。発話時の視覚素の時間変化と動的特徴を含む顔パラメータとの対応付けを、データセット 7 2 0 の各々について行なうことになる。ラベリング部 8 1 4 は、視覚素によりラベリングされた学習用のデータセット 8 0 0 を生成し、これを学習用 DB 8 0 4 に格納する。

【 0 1 6 4 】

H M M 学習部 8 0 6 は、作成された学習用 DB 8 0 4 に格納された学習用データセット 8 0 0 を用いて、顔パラメータ H M M 6 5 0 の学習を行なう。この際 H M M 学習部 8 0 6 が、三つの視覚素を 1 組とする三つ組視覚素を処理単位として、顔パラメータ H M M 6 5 0 の学習を行なう点は第 1 の実施の形態の場合と同様である。

【 0 1 6 5 】

以上のようにして顔パラメータ H M M 6 5 0 の学習を行なうことにより、顔パラメータ H M M 6 5 0 に基づき、視覚素列から動的特徴を含む顔パラメータの系列を合成することが可能になる。

【 0 1 6 6 】

(顔パラメータの合成)

以下、図 1 6 に示すアニメーション作成システム 6 8 0 の動作について説明する。キャラクターの声を表す音声データ 4 2 が準備され、図 1 6 に示す顔パラメータ合成部 6 8 2 に与えられる。図 2 1 を参照して、顔パラメータ合成部 6 8 2 に音声データ 4 2 が入力されると、視覚素列生成部 2 4 0 が、音素 H M M 6 4 及び視覚素対応表 6 6 を用いて、音声データ 4 2 から視覚素列及び当該視覚素列を構成する各視覚素の継続長を推定する。この動作は、学習システム 6 8 の視覚素生成部 2 1 2 (図 7 参照) の動作と同様である。これに

10

20

30

40

50

より、音声データ 4 2 により表される音声の発話時における口の形の変化が特定される。

【 0 1 6 7 】

HMM マッチング部 8 4 2 は、視覚素列生成部 2 4 0 により生成された視覚素列と顔パラメータ HMM 6 5 0 とのマッチングを行ない、発話全体で最尤の顔パラメータの系列 8 4 4 を合成する。顔パラメータの系列 8 4 4 には、各顔パラメータの出力の際に HMM マッチング部 8 4 2 によるマッチングで選択された HMM の各状態の平均ベクトルと共分散行列とが付され、変換部 8 4 6 に与えられる。

【 0 1 6 8 】

変換部 8 4 6 は、与えられた顔パラメータの系列 8 4 4 に含まれる顔パラメータに対し、その顔パラメータに付随している平均ベクトル及び共分散行列とを用いて、式 (4) による演算を行なって、換算後の顔パラメータの行列 C を算出し、換算後の顔パラメータの系列 6 8 4 を出力する。

10

【 0 1 6 9 】

以上のようにして顔パラメータ合成部 6 8 2 により合成された顔パラメータの系列 6 8 4 は、音声データ 4 2 により表現される音声の発話中における口の形の変化から得られたものである。またこの系列 6 8 4 は、第 1 の実施の形態の場合と異なり、顔の特徴点の位置ベクトルだけでなく、その速度ベクトル及び加速度ベクトルをも用いて学習した HMM から合成されたものである。したがって顔パラメータの系列 6 8 4 によって、第 1 の実施の形態に係る顔アニメーションの作成システム 4 0 により合成されたアニメーションよりも滑らかにアニメーションを作成できると考えられ、現実にもそうした効果が得られることが後述するように確認できた。

20

【 0 1 7 0 】

顔パラメータの系列 6 8 4 が作成されれば、図 1 6 に示すマッピング部 9 0、及び画像化部 9 4 によるアニメーション 6 4 6 の作成は第 1 の実施の形態と同様に行なえる。

【 0 1 7 1 】

< 第 2 の実施の形態による効果 >

図 2 2 は、図 1 5 に、第 2 の実施の形態に係る顔アニメーションの作成システム 6 4 0 によって合成されたアニメーションの口の動きを図 2 2 (C) として追加した図である。図 2 2 (A) (B) はそれぞれ図 1 5 (A) (B) と同一の図である。

【 0 1 7 2 】

図 2 2 (C) と図 2 2 (B) とを比較すると、図 2 2 (C) では図 2 2 (B) に存在していたステップ上の変化が除去されて全体として滑らかなグラフとなっていること、しかもグラフが単になまっているわけではなく、図 2 2 (A) に非常によく似た形のピークを持つグラフが得られていることが分かる。

30

【 0 1 7 3 】

すなわち、本実施の形態のように、発話時の顔の特徴点の位置ベクトルだけでなく、その速度ベクトル及び加速度ベクトルという動的特徴までも含めて学習を行なった顔パラメータ HMM 7 8 0 を使用することにより、音声からその発話者の顔のアニメーションを作成でき、しかもその動きが滑らかで実際の発話者の顔の動きに忠実なアニメーションが作成できることが分かる。

40

【 0 1 7 4 】

この第 2 の実施の形態では、学習時の顔の特徴点の速度ベクトル及び加速度ベクトルを算出する際に、特徴点の位置ベクトルの差分を用いている。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。仮に速度ベクトルを精度よく測定できる装置が利用可能であれば、速度ベクトルを位置ベクトルから算出するのではなく、直接測定するようにしてもよい。この場合、加速度ベクトルは速度ベクトルの差分から算出することができる。

【 0 1 7 5 】

加速度ベクトルも速度ベクトルと同様、直接測定できるような装置があればそれを利用し、直接測定するようにしてもよい。

【 0 1 7 6 】

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及び範囲内のすべての変更を含む。

【図面の簡単な説明】

【0177】

【図1】本発明の実施の形態に係るシステム全体の構成を示すブロック図である。

【図2】収録システム60の構成を示す図である。

【図3】発話者110における特徴点の位置を示す概略図である。

【図4】データセット120の構成を示す概略図である。

【図5】音素HMMの概要を示す図である。

【図6】視覚素対応表66の一例を示す図である。

【図7】学習システム68の構成を示すブロック図である。

【図8】顔パラメータ合成部82の構成を示すブロック図である。

【図9】マッピング部90の構成を示すブロック図である。

【図10】顔オブジェクト44及び仮想マーカの一例を示す図である。

【図11】仮想マーカ配置部270が各ノードの仮想マーカによるラベリングを行なう処理の構造を示すフローチャートである。

【図12】顔オブジェクト44における選択ノードと、境界エッジ及び対応マーカとの関係を模式的に示す図である。

【図13】本発明の実施の形態に係る学習システム68及びアニメーション作成システム80の機能を実現するコンピュータシステムの外観の一例を示す図である。

【図14】図13に示すコンピュータシステムのブロック図である。

【図15】本発明の第1の実施の形態に係るシステムにより合成されたアニメーションによる口の動きを、実際の口の動きと対比して示す図である。

【図16】本発明の第2の実施の形態に係る顔アニメーションの作成システム640の全体の構成を示すブロック図である。

【図17】収録システム660の概略構成を示すブロック図である。

【図18】収録システム660のデータセット作成装置722により作成されるデータセット720の構成を示す図である。

【図19】図16に示す学習システム668の構成を示すブロック図である

【図20】顔パラメータHMM780の概略構成と各状態における出力パラメータとの関係を示す図である。

【図21】図16に示す顔パラメータ合成部682のより詳細な構成を示すブロック図である。

【図22】第2の実施の形態に係る顔アニメーションの作成システム640により合成された顔のアニメーションにおける口の動きを、実際の口の動き、及び第1の実施の形態に係る顔アニメーションの作成システム40により合成されたアニメーションにおける口の動きと対比して示す図である。

【符号の説明】

【0178】

40, 640 アニメーション作成システム

42 音声データ

44 顔オブジェクト

46, 646 アニメーション

50, 650 顔パラメータHMM

60, 660 収録システム

62, 662 音声 - 顔パラメータDB

64 音素HMM

66 視覚素対応表

10

20

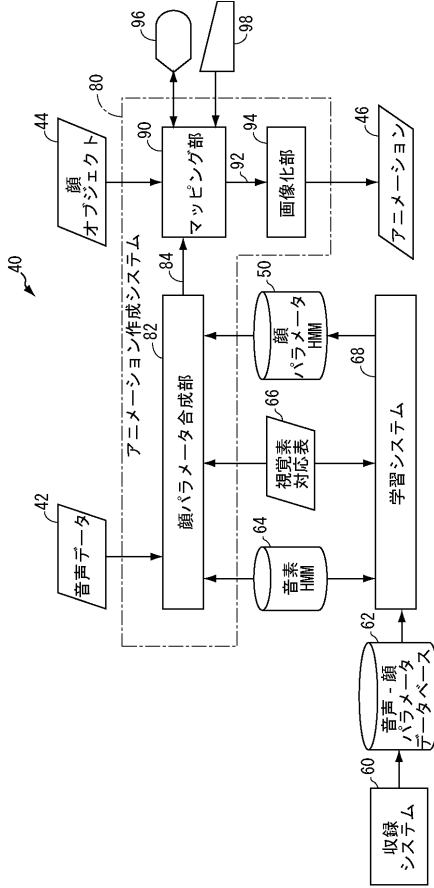
30

40

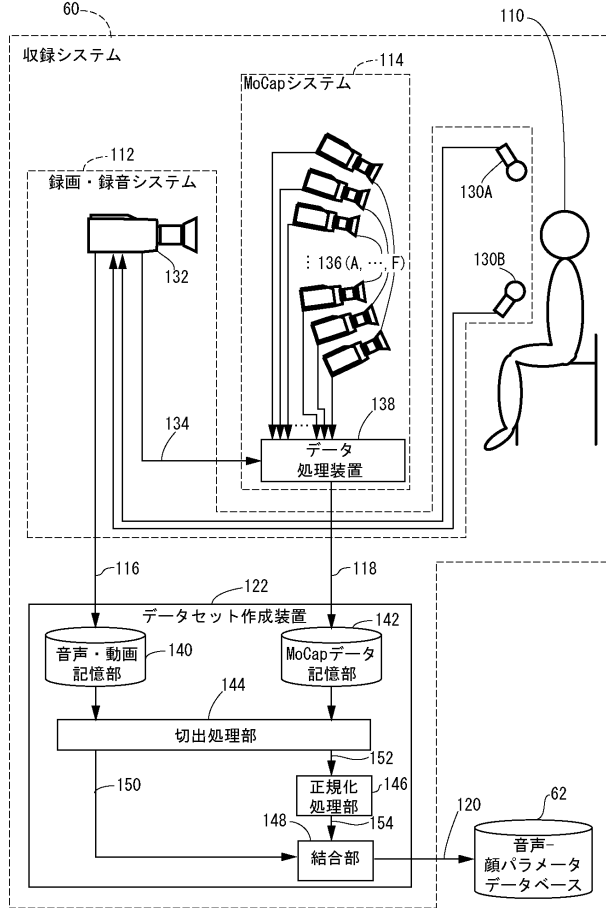
50

6 8 , 6 6 8	学習システム	
8 0 , 6 8 0	アニメーション作成システム	
8 2 , 6 8 2	顔パラメータ合成部	
9 0	マッピング部	
9 4	画像化部	
1 1 0	発話者	
1 1 2	録画・録音システム	
1 1 4	M o C a pシステム	
1 2 2 , 7 2 2	データセット作成装置	
1 3 0 A , 1 3 0 B	マイクロホン	10
1 3 2	カムコーダ	
1 3 6	赤外線カメラ	
1 3 8	データ処理装置	
1 4 0	音声・動画記憶部	
1 4 2	M o C a pデータ記憶部	
1 4 4	切出処理部	
1 4 6	正規化処理部	
1 4 8 , 7 4 8	結合部	
1 5 0	収録音声データ	
1 6 0 A , ... , 1 6 0 M	特徴点	20
1 7 0 A , ... , 1 7 0 N	顔パラメータ	
2 0 2 , 8 0 2	前処理部	
2 0 4 , 8 0 4	学習用 D B	
2 0 6 , 8 0 6	H M M学習部	
2 1 0 , 8 1 0	データセット選択部	
2 1 2 , 2 4 0	視覚素列生成部	
2 1 4 , 8 1 4	ラベリング部	
2 2 0 , 2 5 0	特徴量抽出部	
2 2 2 , 2 5 2	音素列推定部	
2 2 4 , 2 5 4	音素・視覚素変換部	30
2 4 2 , 8 4 2	H M Mマッチング部	
2 7 0	仮想マーカ配置部	
2 7 2	マーカラベリング部	
2 7 4	マーカラベリングデータ記憶部	
2 7 6	顔オブジェクト変形部	
7 4 6	動的特徴算出部	
7 7 2 A , ... , 7 7 2 N	速度パラメータ	
7 7 4 A , ... , 7 7 4 N	加速度パラメータ	
7 8 0	顔パラメータ H M M	
8 4 6	変換部	40

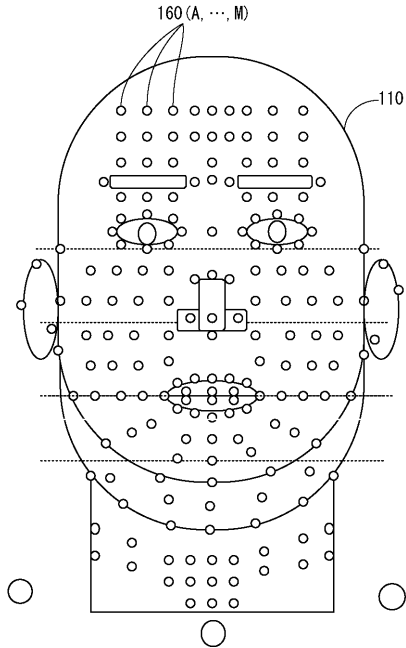
【図1】



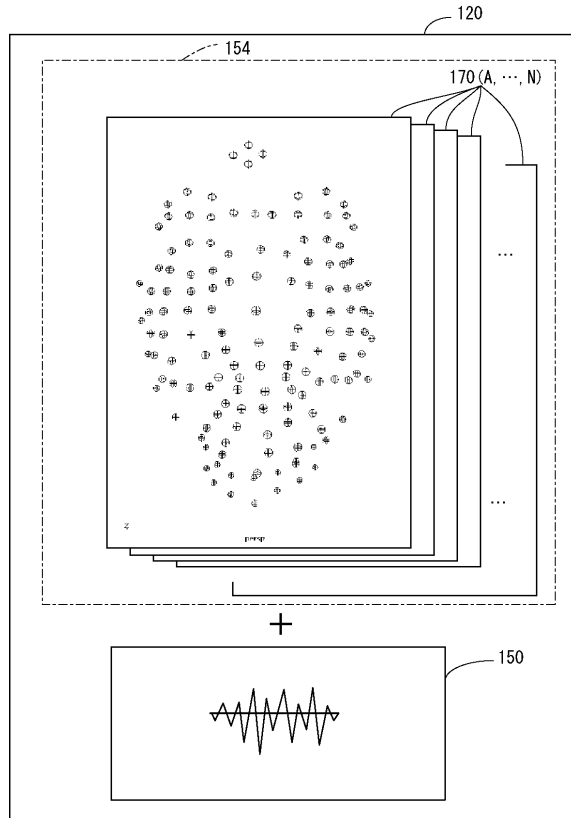
【図2】



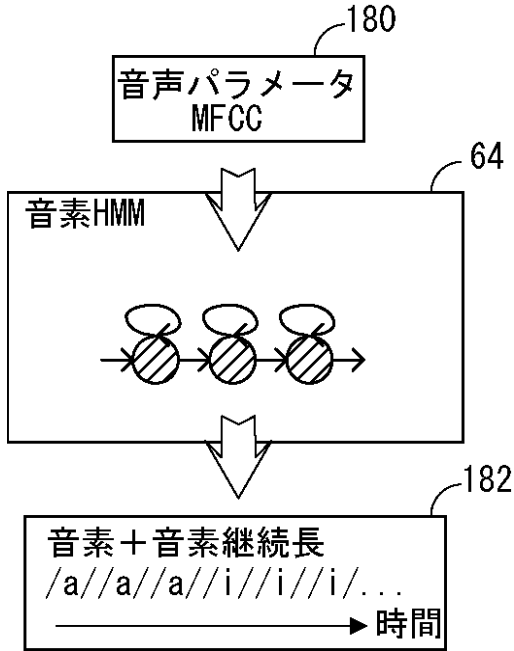
【図3】



【図4】



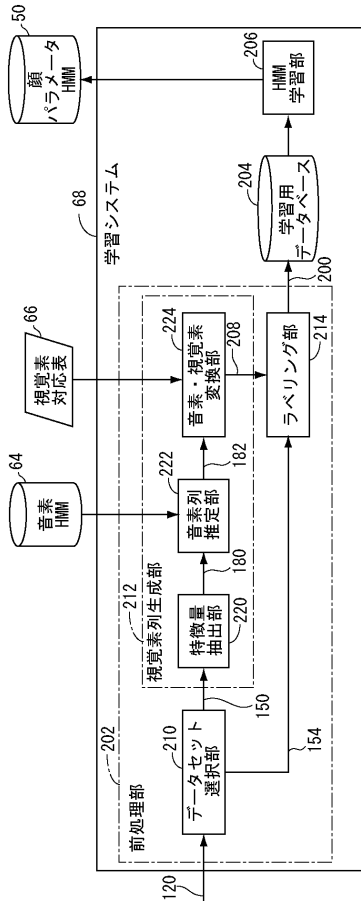
【図5】



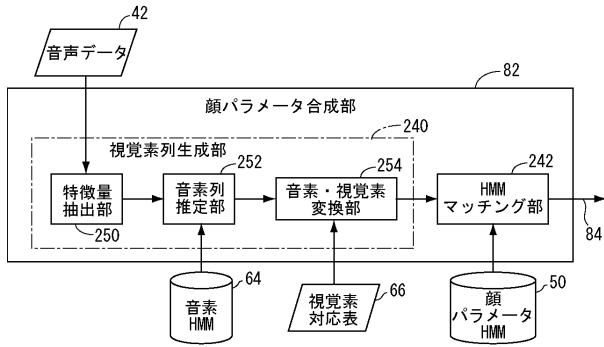
【図6】

視覚素	音素						
A	a	A					
I	i	I					
U	u	U	w				
E	e	E					
O	o	O					
B	m	b	p	my	by	py	f
D	s	sh	t	ch	ts	z	j
	n	ny	hy	y	d	dy	
K	k	r	g	ky	ry	gy	
sp	si	IE	sp	q			
N	N						
***	h						

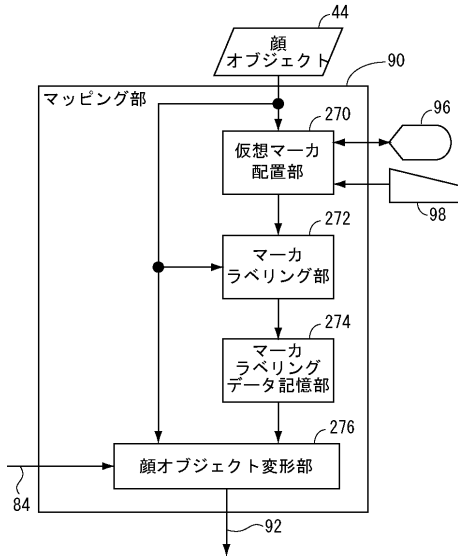
【図7】



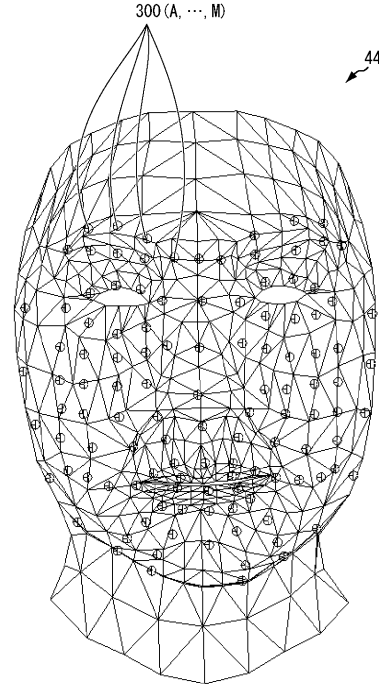
【図8】



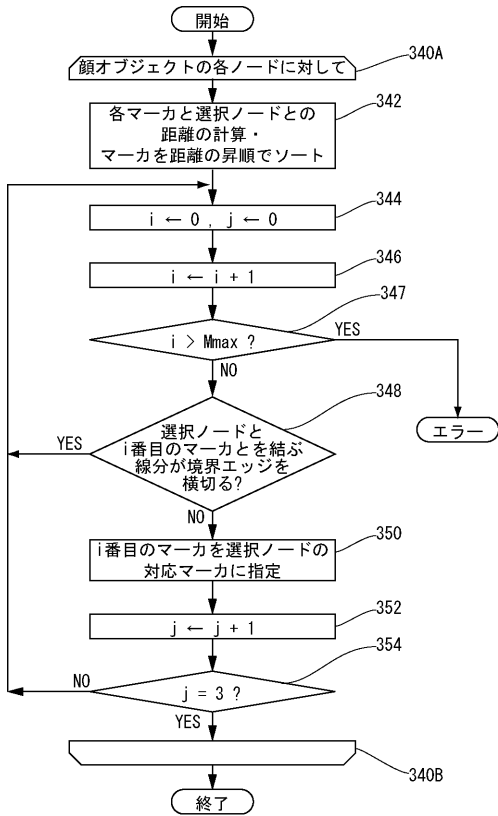
【図9】



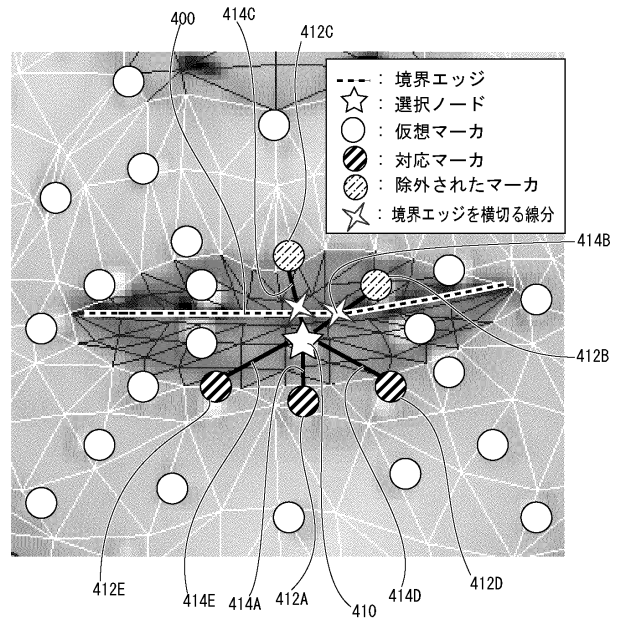
【図10】



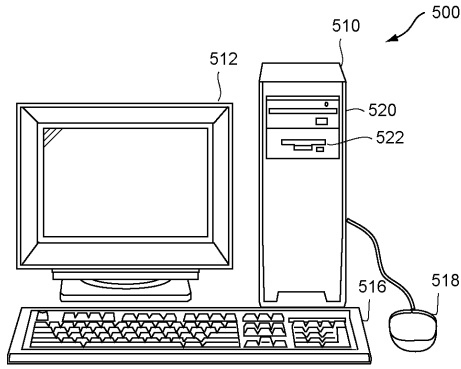
【図11】



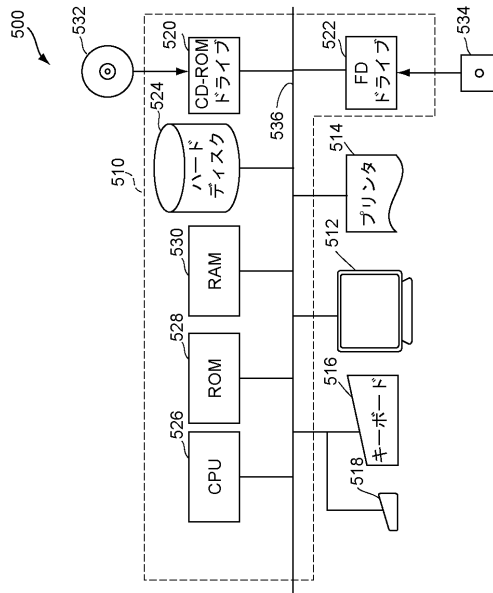
【図12】



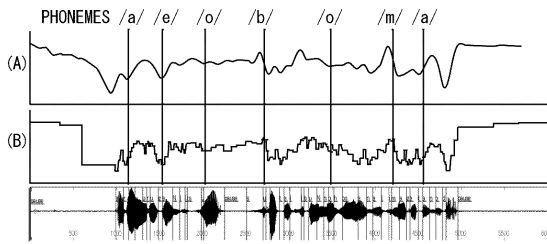
【図13】



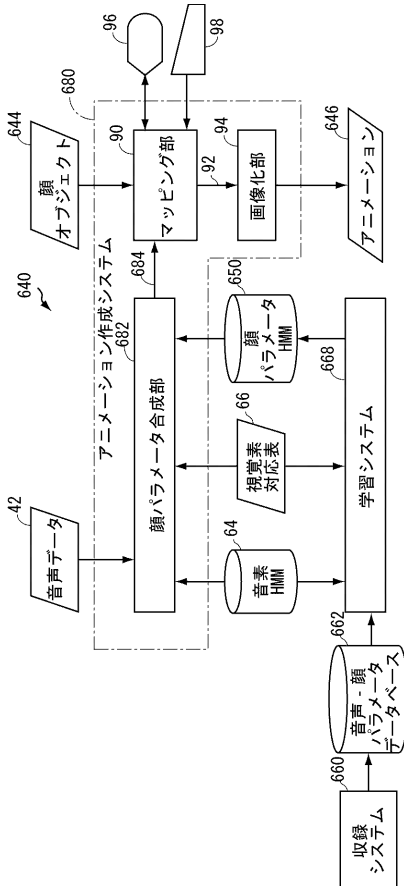
【図14】



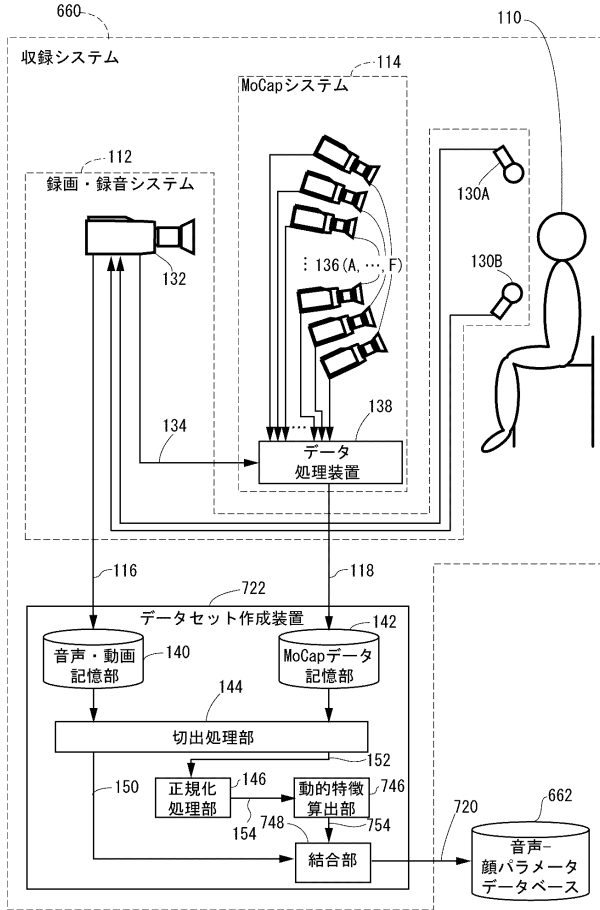
【図15】



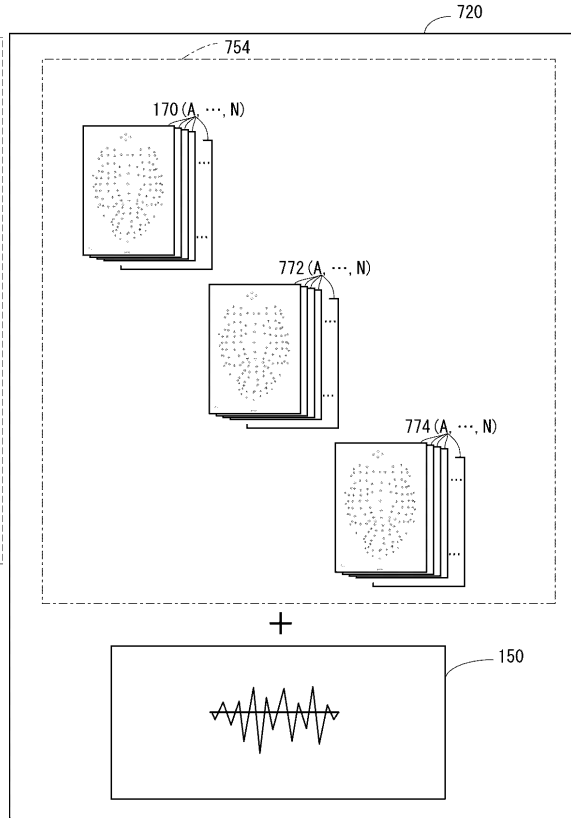
【図16】



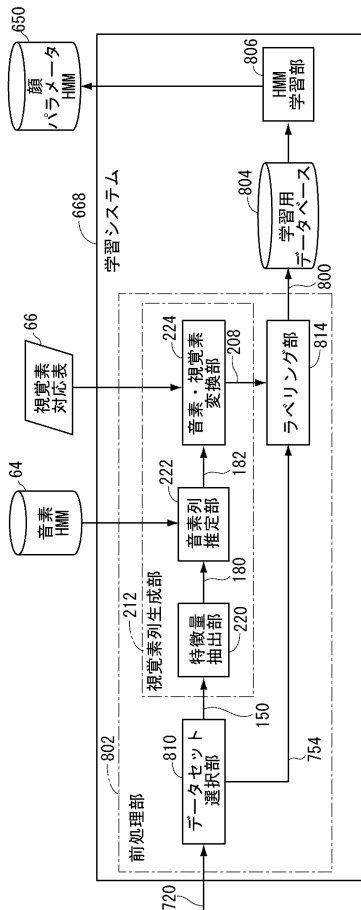
【図17】



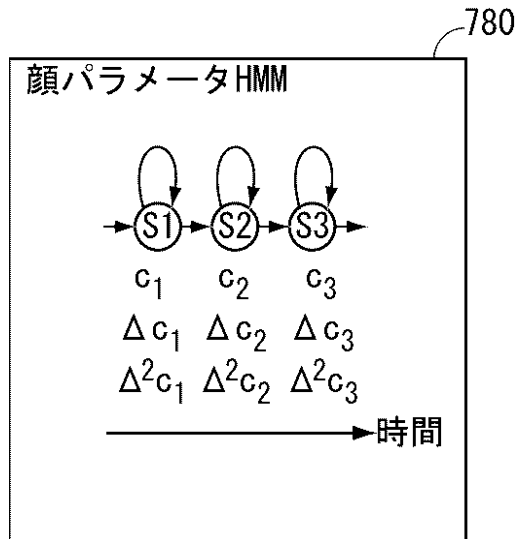
【図18】



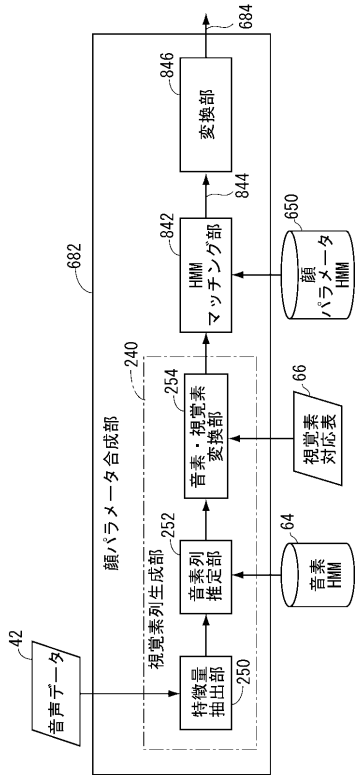
【図19】



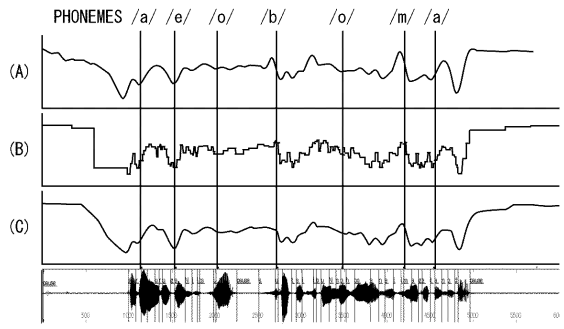
【図20】



【図 2 1】



【図 2 2】



フロントページの続き

審査官 加内 慎也

- (56)参考文献 特開2002-244689(JP,A)
特開2000-123192(JP,A)
HMMを用いた自然な発話動画画像合成 Facial Movement Synthesis by HMM from Audio Sp, 電子情報通信学会論文誌, 2000年11月25日, 2498-2506
唇情報を利用した混合音声の分離 - 方向情報を考慮した Lip Reading - Speech, 情報処理学会第66回全国大会, 2004年 3月 9日, 2-197~2-198

- (58)調査した分野(Int.Cl., DB名)
G06T 15/70