

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4379616号
(P4379616)

(45) 発行日 平成21年12月9日(2009.12.9)

(24) 登録日 平成21年10月2日(2009.10.2)

(51) Int.Cl. F I
G06T 17/40 (2006.01) G O 6 T 17/40 A
G06T 1/00 (2006.01) G O 6 T 1/00 3 4 O A

請求項の数 8 (全 32 頁)

<p>(21) 出願番号 特願2005-56592 (P2005-56592) (22) 出願日 平成17年3月1日(2005.3.1) (65) 公開番号 特開2006-243975 (P2006-243975A) (43) 公開日 平成18年9月14日(2006.9.14) 審査請求日 平成18年12月21日(2006.12.21)</p> <p>特許法第30条第1項適用 2004年9月10日 社団法人情報処理学会発行の「情報処理学会研究報告 情報処研報V o l . 2 0 0 4 N o . 9 0」に発表</p> <p>(出願人による申告)平成16年度独立行政法人情報通信研究機構、研究テーマ「大規模コーパス音声対話翻訳技術の研究開発」に関する委託研究、産業活力再生特別措置法第30条の適用を受ける特許出願</p> <p>特許権者において、実施許諾の用意がある。</p>	<p>(73) 特許権者 393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2 (74) 代理人 100099933 弁理士 清水 敏 (72) 発明者 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 (72) 発明者 森島 繁生 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内 (72) 発明者 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</p> <p style="text-align: right;">最終頁に続く</p>
---	---

(54) 【発明の名称】 モーションキャプチャデータ補正装置、マルチモーダルコーパス作成システム、画像合成装置、及びコンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

発話時の発話者の動画像から得られたモーションキャプチャデータを補正するためのモーションキャプチャデータ補正装置であって、前記モーションキャプチャデータは、複数フレームを含み、前記複数フレームの各々は、当該フレーム撮影時における前記発話者の頭部の複数個の特徴点の位置データを含み、前記複数個の特徴点は、前記発話者の首部より上であってかつ前記発話者の表情変化の影響を受けない所定箇所に配置された第1の種類の特徴点と、その他の特徴点とを含み、

前記複数フレームの各々に対する前記複数個の特徴点の位置データから、前記第1の種類の特徴点の位置データを選択するための選択手段と、

前記複数フレームの各々に対し、前記選択手段により選択された位置データを基準として、前記複数個の特徴点の各々の位置データを補正するための補正手段とを含み、

前記第1の種類の特徴点は、前記発話者の頭部であってかつ前記発話者の表情変化の影響を受けない所定箇所に配置された第1の基準特徴点と、前記発話者の首部であってかつ前記発話者の表情変化の影響を受けない所定箇所に配置された第2の基準特徴点とを含み

前記選択手段は、

前記複数フレームの各々に対し、前記複数個の特徴点の位置データを、前記発話者の頭部の特徴点の位置データと、前記発話者の首部の特徴点の位置データとに分類するための分類手段と、

前記分類手段により分類された前記発話者の頭部の特徴点の位置データから、前記第1の基準特徴点のデータを選択するための頭部基準特徴点選択手段と、

前記頭部基準特徴点選択手段により選択された前記第1の基準特徴点のデータを基準に、同じフレームの前記頭部の特徴点の位置データを補正するための頭部補正式を算出するための頭部補正式算出手段と、

前記複数フレームの各々に対し、前記頭部の特徴点の位置データに前記頭部補正式算出手段により算出された頭部補正式を適用して補正するための頭部補正手段と、

前記分類手段により分類された前記発話者の首部の特徴点の位置データから、前記第2の基準特徴点のデータを選択するための首部基準特徴点選択手段と、

前記首部基準特徴点選択手段により選択された前記第2の基準特徴点のデータを基準に、同じフレームの前記首部の特徴点の位置データを補正するための首部補正式を算出するための首部補正式算出手段と、

前記複数フレームの各々に対し、前記首部の特徴点の位置データに前記首部補正式算出手段により算出された首部補正式を適用して補正するための首部補正手段とを含む、モーションキャプチャデータ補正装置。

【請求項2】

前記第1の種類の特徴点は、前記発話者の額領域、こめかみ領域、及び鼻の先端領域のいずれかに配置される、請求項1に記載のモーションキャプチャデータ補正装置。

【請求項3】

前記補正手段は、

前記複数フレームの各々に対し、前記選択手段により選択された前記第1の種類の特徴点の位置データを基準に、同じフレームの前記複数個の特徴点の位置データを補正するための補正式を算出するための補正式算出手段と、

前記複数フレームの各々に対し、前記複数個の特徴点の位置データに前記補正式算出手段により算出された補正式を適用して補正するための補正式適用手段とを含む、請求項1又は請求項2のいずれかに記載のモーションキャプチャデータ補正装置。

【請求項4】

発話時の発話者の顔画像を含む動画データと、当該発話時の音声の録音データと、発話時における前記発話者の顔の予め定める複数の特徴点に関するモーションキャプチャデータとを発話ごとにそれぞれ分離し、互いに対応付けて保存するための発話分離手段と、

前記発話分離手段により分離された各発話の前記モーションキャプチャデータを補正するための、請求項1～請求項3のいずれかに記載のモーションキャプチャデータ補正装置とを含む、マルチモーダルコーパス作成システム。

【請求項5】

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3のいずれかに記載のモーションキャプチャデータ補正装置として動作させる、コンピュータプログラム。

【請求項6】

顔オブジェクトの形状を第1の座標空間における複数のノードの座標値を用いて定義した形状モデルと、所定の発話を行なっている発話者の顔画像から得られた、前記発話者の頭部の複数の特徴点の所定の第2の座標系における位置情報とを基に、前記所定の発話を行なう前記顔オブジェクトの表情を表す画像を合成するための画像合成装置であって、

前記発話者の頭部の複数の特徴点は、請求項1に記載のモーションキャプチャデータ補正装置により、各々の位置データが補正され、

前記複数の特徴点と、前記形状モデル内の任意の点との対応関係を定義することにより、前記形状モデル内に前記複数の特徴点にそれぞれ対応する複数の仮想特徴点を設定するための仮想特徴点設定手段と、

前記複数のノードの各々に対し、前記複数の仮想特徴点のうちで、当該ノードからの距離が小さいものから順番に、かつ当該ノードと仮想特徴点とを結ぶ線分が前記形状モデルに対し所定の制約条件を充足するものを所定個数だけ選定するための仮想特徴点選定手段

10

20

30

40

50

と、

前記複数のノードの各々に対し、前記仮想特徴点選定手段により選定された所定個数の仮想特徴点の位置情報の間の内挿により算出される座標値を割当てることにより前記形状モデルを変形させるための形状モデル変形手段と、

前記形状モデル変形手段により得られた形状モデルに基づいて前記顔オブジェクトの画像を生成するための画像生成手段とを含む、画像合成装置。

【請求項 7】

前記仮想特徴点選定手段は、前記複数のノードの各々に対し、前記複数の仮想特徴点のうちで、当該ノードからの距離が小さいものから順番に、かつ当該ノードと仮想特徴点とを結ぶ線分が前記形状モデルの境界エッジを横切らないものを所定個数だけ選定するための手段を含む、請求項 6 に記載の画像合成装置。

10

【請求項 8】

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 6 又は請求項 7 に記載の画像合成装置として動作させる、コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声言語処理技術に関し、特に発話時における音声及び表情変化に関する情報を含むマルチモーダルコーパスを作成するためのマルチモーダルコーパス作成装置及びシステム、並びに発話時の表情変化を表現するアニメーションを作成するための画像合成装置に関する。

20

【背景技術】

【0002】

人間にとって容易かつ自然なヒューマンマシンインタフェースを実現するための技術が研究されている。例えば、音声認識及び音声合成等の音声言語処理技術においては、大規模な音声コーパスと統計的な手法とにより、高性能の処理が実現されるようになっている。さらには、音声だけでなく視聴覚情報を用いるヒューマンマシンインタフェースを実現するための技術が盛んに研究されている。そのひとつに、音声合成技術を発展させて、発話時の顔画像を合成する技術がある。

【0003】

30

後掲の非特許文献 1 には、仮想空間上のメッシュで構成された顔の形状モデルを用いて、発話時の顔の表情変化を表現する技術が開示されている。この技術では、メッシュの各点の動きを推定し、推定した動きを基に顔のメッシュ形状を変形させる。この手法では、顔の形状モデル及びその表情変化のバリエーションに制限がなく、表情変化の豊かな顔画像を自在に合成することができる。

【0004】

後掲の非特許文献 2 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化しておく。そして、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。この手法で合成される画像は、撮影された画像を再構成したものであるため、大規模かつ適切なデータベースを用意すれば、自然な顔画像合成を行なうことができる。

40

【0005】

また、視聴覚情報を用いるヒューマンマシンインタフェースを実現するために、音声言語処理技術における音声コーパスに相当するデータベースの整備が進められている。音声と顔の画像情報とを、言語情報に対応付けたマルチモーダルコーパスの整備が進められている。後掲の非特許文献 3 には、マルチモーダルコーパスを構築する種々の試みが紹介されている。

【0006】

マルチモーダルコーパスに収められた画像の特徴量を抽出し分析することにより、当該

50

画像情報に基づく顔画像の合成等が可能となる。非特許文献1に記載の技術では、発話時の顔を撮影した動画像におけるピクセル情報からオプティカルフローを求めることにより、発話時における顔の各部位の動きを推定し、画像の特徴量として用いている。また、後掲の非特許文献4には、唇領域の画像におけるピクセル情報をもとに、唇の変化量を求める技術が開示されている。

【0007】

【非特許文献1】モリシマ, S., イワサワ, S., サカグチ, T., カワカミ, F., アンドウ, M., 「より良い顔のコミュニケーション」、ACMシーグラフ'95、インタラクティブ コミュニティ ビジュアル予稿集、117頁、1995年 (Morishima, S., Iwasawa, S., Sakaguchi, T., Kawakami, F., and Ando, M., "Better Face Communication", Visual Proceedings of ACM SIGGRAPH '95, Interactive Communities, p.117, 1995)

10

【非特許文献2】エザット, T., ガイガー, G., ポッジョ, T. 「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予稿集、2002年 (Ezzat, T., Geiger, G. and Poggio, T. "Trainable Videorealistic Speech Animation", Proceedings of ACM SIGGRAPH 2002).

【非特許文献3】ナカムラ, S., 「最近のマルチモーダルコーパス活動の概要」、COCOSDAワークショップ 2000 (Nakamura, S., "Overview on Recent Activities in Multi-Modal Corpora", COCOSDA Workshop, 2000)

【非特許文献4】タムラ, T., コンドウ, S., マスコ, T., コバヤシ, T., 「HMMからのパラメータ生成に基づくテキスト-発話音声画像合成」 EURO SPEECH '99予稿集、959-962頁、1999年 (Tamura, T., Kondo, S., Masuko, T., and Kobayashi, T., "Text-to-Audio-Visual Speech Synthesis Based on Parameter Generation from HMM", Proceeding of EURO SPEECH, pp.959-962, 1999)

20

【発明の開示】

【発明が解決しようとする課題】

【0008】

非特許文献2及び4のいずれに記載の技術においても、発話時の表情の特徴量を画像情報から得ている。しかし、この場合、次のような問題点が発生する。すなわち、顔及びその表情は立体的であるのに対し、動画像は2次元の情報である。そのため、3次元での形状変化に関する特徴量を得るのは困難である。例えば、発話中には表情を形成する顔の器官だけではなく、頭部及び首部も自由に移動回転する。顔の各器官の動画像上での位置及び形状は、頭部の動きに応じて表情とは無関係に変化する。よって、動画像から顔の器官の変化についての情報を得るのは困難である。また、画像情報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、画像情報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。

30

【0009】

また、非特許文献1に記載の手法をはじめとする、モデルベースの顔画像の合成方法で発話時の顔の画像を作成するには、アニメーションの各フレームにおいて、モデルのメッシュの位置を定義する必要がある。現在のアニメーションに用いられる形状モデルは、膨大な数のメッシュから構成されている。形状モデルの変形によってアニメーションを生成するには、膨大な量のメッシュについて位置をいちいち定義しなければならず、膨大な作業を要する。非特許文献2に記載の手法をはじめとする動画像の再生成手法では、作成可能な顔の表情は、顔画像のコーパスに格納されている画像により限定されてしまう。多様な容貌の顔及び多彩な表情で発話時の表情変化を表現するには、その分膨大な量の顔の画像をコーパス化する必要がある。

40

【0010】

それゆえに、本発明の目的は、顔の表情を形成する各器官の動きについての正確な位置データを得ることができるモーションキャプチャデータ補正装置及びマルチモーダルコー

50

パス作成システムを提供することである。

【 0 0 1 1 】

本発明の別の目的は、多様な表情を持つ顔画像、または実際の発話者の表情を再現する顔画像のアニメーションを高精度かつ容易に合成することを可能にする画像作成装置を提供することである。

【 課題を解決するための手段 】

【 0 0 1 2 】

本発明の第1の局面に係るモーションキャプチャデータ補正装置は、発話時の発話者の動画像から得られたモーションキャプチャデータを補正するためのモーションキャプチャデータ補正装置である。モーションキャプチャデータは、複数フレームを含む。複数フレームの各々は、当該フレーム撮影時における発話者の頭部の複数個の特徴点の位置データを含む。複数個の特徴点は、発話者の首部より上であってかつ発話者の表情変化の影響を受けない所定箇所に配置された第1の種類の特徴点と、その他の特徴点とを含む。このモーションキャプチャデータ補正装置は、複数フレームの各々に対する複数個の特徴点の位置データから、第1の種類の特徴点の位置データを選択するための選択手段と、複数フレームの各々に対し、選択手段により選択された位置データを基準として、複数個の特徴点の各々の位置データを補正するための補正手段とを含む。

10

【 0 0 1 3 】

発話者の表情変化の影響を受けない所定箇所に配置された第1の種類の特徴点の位置データを基準として、発話者の顔の特徴点の位置データが補正される。一般的に発話者の頭部のモーションキャプチャデータには、頭部全体の動き、又は首部の動きによる影響が含まれる。第1の種類の特徴点の位置データは、頭部全体の動き、又は首部の動きのみによる影響を受けると考えられ、これらを基準として特徴点の位置データを補正することにより、表情変化のみに起因する特徴点の位置の変化が得られる。

20

【 0 0 1 4 】

好ましくは、第1の種類の特徴点は、発話者の額領域、こめかみ領域、及び鼻の先端領域のいずれかに配置される。

【 0 0 1 5 】

これら領域は、いずれも表情変化による影響を受けないか、きわめて少ない。したがってこれらの領域に配置された特徴点を基準に特徴点の位置データを補正することで、それら特徴点の、表情変化に起因する位置変化のみが正確に得られる。

30

【 0 0 1 6 】

より好ましくは、補正手段は、複数フレームの各々に対し、選択手段により選択された第1の種類の特徴点の位置データを基準に、同じフレームの複数個の特徴点の位置データを補正するための補正式を算出するための補正式算出手段と、複数フレームの各々に対し、複数個の特徴点の位置データに補正式算出手段により算出された補正式を適用して補正するための補正式適用手段とを含む。

【 0 0 1 7 】

第1の種類の特徴点の位置データを基準に補正式が算出され、この補正式を各特徴点の位置データに適用する。定型化した処理により、対象が別の発話者になっても新たに補正式を算出でき、安定して位置データの補正を行なうことができる。

40

【 0 0 1 8 】

さらに好ましくは、補正式算出手段は、複数フレームの各々に対して、第1の種類の特徴点の位置データに対する特異値分解により、同一フレーム内の複数個の特徴データを変換するためのアフィン変換行列を算出するための手段を含む。

【 0 0 1 9 】

特異値分解により座標変換のためのアフィン変換行列が得られる。その結果、簡単な行列演算で位置データの補正を行なうことができる。

【 0 0 2 0 】

好ましくは、第1の種類の特徴点は、発話者の頭部であってかつ発話者の表情変化の影

50

響を受けない所定箇所に配置された第1の基準特徴点と、発話者の首部であってかつ発話者の表情変化の影響を受けない所定箇所に配置された第2の基準特徴点とを含み、選択手段は、複数フレームの各々に対し、複数個の特徴点の位置データを、発話者の頭部の特徴点の位置データと、発話者の首部の特徴点の位置データとに分類するための分類手段と、分類手段により分類された発話者の頭部の特徴点の位置データから、第1の基準特徴点のデータを選択するための頭部基準特徴点選択手段と、頭部基準特徴点選択手段により選択された第1の基準特徴点のデータを基準に、同じフレームの頭部の特徴点の位置データを補正するための頭部補正式を算出するための頭部補正式算出手段と、複数フレームの各々に対し、頭部の特徴点の位置データに頭部補正式算出手段により算出された頭部補正式を適用して補正するための頭部補正手段と、分類手段により分類された発話者の首部の特徴点の位置データから、第2の基準特徴点のデータを選択するための首部基準特徴点選択手段と、首部基準特徴点選択手段により選択された第2の基準特徴点のデータを基準に、同じフレームの首部の特徴点の位置データを補正するための首部補正式を算出するための首部補正式算出手段と、複数フレームの各々に対し、首部の特徴点の位置データに首部補正式算出手段により算出された首部補正式を適用して補正するための首部補正手段とを含む。

10

【0021】

首部の特徴点は、頭部とは別に首部の動きによる影響を受ける。したがって、頭部とは別に首部に対しても基準となる特徴点を定め、それらに基づいて首部の特徴点の位置データを補正する。こうして、顔面を含む頭部の特徴点と、首部の特徴点との各々について、表情の変化のみに起因する位置変化を算出することができる。

20

【0022】

本発明の第2の局面に係るマルチモーダルコーパス作成システムは、発話時の発話者の顔画像を含む動画データと、当該発話時の音声の録音データと、発話時における発話者の顔の予め定める複数の特徴点に関するモーションキャプチャデータとを発話ごとにそれぞれ分離し、互いに対応付けて保存するための発話分離手段と、発話分離手段により分離された各発話のモーションキャプチャデータを補正するための、上記したいずれかのモーションキャプチャデータ補正装置とを含む。

【0023】

このマルチモーダルコーパス作成システムによれば、発話ごとに、発話者の顔画像の動画データと、音声の録音データと、発話者の顔の特徴点のモーションキャプチャデータが得られる。そのモーションキャプチャデータをモーションキャプチャデータ補正装置を用いて補正することにより、発話者の顔の特徴点の、発話による表情変化のみに起因する位置変化が得られる。その結果、発話に起因するこの発話者の顔の特徴点の位置変化が正確に表され、発話と表情との間の関係を研究するための正確な基礎データが得られる。

30

【0024】

本発明の第3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、上記したいずれかのモーションキャプチャデータ補正装置として動作させる。したがってこのコンピュータプログラムにより、第1の局面に係るモーションキャプチャデータ補正装置と同様の効果を得ることができる。

40

【0025】

本発明の第4の局面に係る画像合成装置は、顔オブジェクトの形状を第1の座標空間における複数のノードの座標値を用いて定義した形状モデルと、所定の発話を行なっている発話者の顔画像から得られた、発話者の頭部の複数の特徴点の所定の第2の座標系における位置情報とを基に、所定の発話を行なう顔オブジェクトの表情を表す画像を合成するための画像合成装置であって、複数の特徴点と、形状モデル内の任意の点との対応関係を定義することにより、形状モデル内に複数の特徴点にそれぞれ対応する複数の仮想特徴点を設定するための仮想特徴点設定手段と、複数のノードの各々に対し、複数の仮想特徴点のうちで、当該ノードからの距離が小さいものから順番に、かつ当該ノードと仮想特徴点とを結ぶ線分が形状モデルに対し所定の制約条件を充足するものを所定個数だけ選定するた

50

めの仮想特徴点選定手段と、複数のノードの各々に対し、仮想特徴点選定手段により選定された所定個数の仮想特徴点の位置情報の間の内挿により算出される座標値を割当てることにより形状モデルを変形させるための形状モデル変形手段と、形状モデル変形手段により得られた形状モデルに基づいて顔オブジェクトの画像を生成するための画像生成手段とを含む。

【0026】

顔オブジェクトの形状モデルに、顔オブジェクトの特徴点と対応する仮想特徴点が設定され、さらに形状モデルを構成する各ノードと、当該ノードとの距離が近く、かつ所定の制約条件を充足する所定個数の仮想特徴点とが対応付けられる。各ノードに、それらに対応付けられた仮想特徴点の位置情報の間の内挿により得られた座標値を割当てることにより、各ノードに割当てられた座標値はもとの発話者の顔においてそのノードに対応する点の位置とほぼ正確に一致する。その結果、こうして得られた座標値を用いて顔オブジェクトの形状を変化させることで、元の発話者の表情変化を顔オブジェクトにより再現できる。

10

【0027】

好ましくは、仮想特徴点選定手段は、複数のノードの各々に対し、複数の仮想特徴点のうちで、当該ノードからの距離が小さいものから順番に、かつ当該ノードと仮想特徴点とを結ぶ線分が形状モデルの境界エッジを横切らないものを所定個数だけ選定するための手段を含む。

【0028】

一般に顔には、目、口、鼻の穴等、顔面を構成しない切れ目があり、形状モデルでは、それらと顔面との間は境界エッジで仕切られている。こうした切れ目を挟んだ両側のノードは互いに別の動きをするため、それらの座標位置を互いに関連付けて計算するのは不相当である。そこで、このように計算対象のノードと仮想特徴点とを結ぶ線分が境界エッジを横切るような仮想特徴点はノードの座標値の計算からは除外する。こうすることで、各ノードの座標値をより正確に、かつ実際の顔と同様に適切な表情が得られるように算出できる。

20

【0029】

より好ましくは、複数フレームの位置情報を元に、仮想特徴点設定手段、仮想特徴点選定手段、形状モデル変形手段、及び画像生成手段により生成された顔オブジェクトの画像を各フレームとして時系列的に保存することにより、所定の発話を行なう顔オブジェクトの表情を表す動画を生成するための手段をさらに含む。

30

【0030】

フレームごとに顔画像を作成し、それらを時系列的に保存することにより、発話時の発話者の顔の表情と同様の表情変化を持つ動画を生成できる。

【0031】

本発明の第5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、上記したいずれかの画像合成装置として動作させる。

【0032】

このコンピュータプログラムによれば、上記した第4の局面に係る画像合成装置と同様の効果を得ることができる。

40

【発明を実施するための最良の形態】

【0033】

以下、図面を参照しつつ、本発明の一実施の形態について説明する。なお、以下の説明に用いる図面では、同一の部品に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。

【0034】

[概要]

本実施の形態では、音声及び顔の動画像に加えて、発話時の表情に関するデータを含むマルチモーダルコーパスを作成する。本実施の形態では、音声及び動画像の収録時に、顔

50

の多数の部位について位置計測を併せて行なう。さらに当該位置の計測データから顔の各器官の変化を表すデータを取得し、表情に関する特徴量データとする。そして、当該顔器官の変化を表すデータ（以下、「顔器官変化量データ」と呼ぶ）と音声及び動画像のデータとを対応付けてデータベース化することにより、マルチモーダルコーパスを作成する。本実施の形態ではさらに、発話時の表情変化を表現するアニメーションを、マルチモーダルコーパスをもとに作成する。この際、顔の形状モデルに顔器官の変化を順次割り当てる。

【 0 0 3 5 】

[図 1 システム全体の構成]

図 1 に、本実施の形態に係るマルチモーダルコーパス作成システム 1 0 0 全体の構成を示す。図 1 を参照して、このマルチモーダルコーパス作成システム 1 0 0 は、発話者 1 0 2 の音声及び顔の動画像を収録すると同時に、発話者 1 0 2 の顔の各部位について位置計測を行なうための収録システム 1 0 4 と、収録システム 1 0 4 による位置の計測結果を基に顔器官変化量データを生成し、収録システム 1 0 4 による収録で得られる発話時の音声のデータ及び動画像のデータ、並びに当該顔器官変化量データを発話内容と対応付けることによりマルチモーダルコーパス 1 0 6 を作成するためのマルチモーダルコーパス作成装置 1 0 8 とを含む。

10

【 0 0 3 6 】

このマルチモーダルコーパス作成システム 1 0 0 はさらに、静止状態における所定の顔の形状を表す初期顔モデル 1 1 0 を記憶するための記憶装置と、入力テキストを受け、マルチモーダルコーパス 1 0 6 内の顔器官変化量データを基に、入力テキストを発話中の各時刻における顔の形状モデルを作成し動画像化することにより、入力テキスト発話時の顔の表情変化を表現するアニメーション 1 1 2 を作成するためのアニメーション作成装置 1 1 4 とを含む。

20

【 0 0 3 7 】

マルチモーダルコーパス作成システム 1 0 0 はさらに、マルチモーダルコーパス作成時のユーザの操作を受け、対応する操作信号をマルチモーダルコーパス作成装置 1 0 8 に与えるための入力装置 1 1 6 A と、マルチモーダルコーパス作成装置 1 0 8 から、操作に用いる情報を受けて出力するための出力装置 1 1 8 A と、アニメーション作成時にユーザの操作を受け、対応する操作信号をアニメーション作成装置 1 1 4 に与えるための入力装置 1 1 6 B と、アニメーション作成装置 1 1 4 からの出力される情報を画像及び音声等に変換して出力するための出力装置 1 1 8 B とを含む。

30

【 0 0 3 8 】

初期顔モデル 1 1 0 は、静止状態における所定の顔の形状を多数の多角形（ポリゴン）によって表現した形状モデルである。図 9 に、初期顔モデル 1 1 0 の一例を示す。図 9 を参照して、この初期顔モデル 1 1 0 は、発話者 1 0 2 の顔の静止画像と所定のワイヤフレームモデルとを整合させることにより準備された形状モデルである。この顔モデルは、約 7 5 0 のポリゴンで構成されている。アニメーション作成装置 1 1 4 は、顔器官変化量データを基に、発話中における顔の各器官の変化を、初期顔モデル 1 1 0 におけるポリゴンの頂点（ノード）の各々に割り当てて発話中の所定の顔の形状モデルを形成する機能を持つ。

40

【 0 0 3 9 】

[収録システム 1 0 4 の構成]

収録システム 1 0 4 は、発話時における発話者 1 0 2 の顔の各部位の位置及びその軌跡を計測しキャプチャデータとして出力するためのモーションキャプチャシステム 1 2 0 と、発話者 1 0 2 の音声を収録するための録音システム 1 2 2 と、発話時における発話者 1 0 2 の動画像を撮影するための撮影システム 1 2 4 と、発話者に発話すべき内容として提示される所定の文章、単語、文字、及び音節の記号等で構成された発話内容を格納する発話リスト 1 2 6 と、発話リスト 1 2 6 の発話内容のいずれかを発話者 1 0 2 に提示するためのテレプロンプタ 1 2 8 と、モーションキャプチャシステム 1 2 0 及び撮影システム 1

50

24に対してタイムコードを供給するためのタイムコードジェネレータ130とを含む。

【0040】

本実施の形態に係るモーションキャプチャシステム120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。モーションキャプチャシステム120は、発話者102の顔面及び首部の予め定める多数の箇所それぞれ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132A, ..., 132F（以下これらをまとめて「赤外線カメラ132」と呼ぶことがある。）と、赤外線カメラ132からの映像信号を基にフレームごとに各マーカの位置を計測し、タイムコードジェネレータ130からのタイムコードを付与して出力するためのデータ処理装置134とを含む。

10

【0041】

[図2 マーカの配置例]

図2(A)及び図2(B)に、発話者102の首部より上へのマーカの装着例を示す。図2(A)は、発話者102の顔面及び首部の右半分の所定位置にマーカを装着した状態での、発話者102の頭部及び首部の外観を示す右側面図であり、図2(B)は、同状態での発話者102の頭部及び首部の外観を示す正面図である。

【0042】

図2(A)及び図2(B)を参照して、発話者102の顔面及び首部の皮膚上には、多数のマーカ170A, ..., 170M（以下これらをまとめて「マーカ170」と呼ぶことがある。）が、図示しない装着材（接着剤）により装着される。マーカ170は、直径3~4mmの半球状又は球状の形状であり、照射光を再帰反射するよう加工されている。

20

【0043】

図2(A)及び図2(B)に示す例では、マーカ170は、眉部の9箇所、目の輪郭部の9箇所、鼻部の5箇所、口唇部の11箇所、頬部の18箇所、顔の輪郭部の8箇所、顎部の6箇所、首部の8箇所、及び額部の4箇所に装着されている。マルチモーダルコーパス作成においては、発話時の顔部位の詳細な変化量を計測すること、及び複数日にわたり又は複数の発話者102について計測を行なうことが想定される。そのため、マーカ170はそれぞれ、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に、予め定めるルールにしたがい装着される。例えば、口唇部のマーカはそれぞれ次の表に示すルールにより定められた装着位置に、定められた装着順序で装着される。なお、こうして定められた装着位置を、本明細書では「特徴点」と呼ぶ。

30

【0044】

【表1】

装着順	識別番号	装着位置
1	MTH1C	顔中央線と上唇上部輪郭との交点
2	MTH2C	顔中央線と上唇下部輪郭との交点から2mm~3mmほど [MTH1C] 寄りの位置
3	MTH4C	顔中央線と下唇下部輪郭との交点
4	MTH3C	顔中央線と下唇上部輪郭との交点から2mm~3mmほど [MTH4C] 寄りの位置
5	MTH1R(L)	右(左)の口角
6	MTH2R(L)	唇輪郭線[MTH1R(L)]-[MTH1C]間を3等分した時の[MTH1R(L)]寄りの点
7	MTH4R(L)	唇輪郭線[MTH1R(L)]-[MTH1C]間を3等分した時の[MTH1C]寄りの点
8	MTH3R(L)	唇輪郭線[MTH1R(L)]-[MTH4C]間を3等分した時の[MTH1R(L)]寄りの点
9	MTH7R(L)	唇輪郭線[MTH1R(L)]-[MTH4C]間を3等分した時の[MTH4C]寄りの点
10	MTH5R(L)	線分[MTH1R(L)]-[MTH2C]の中点から2mm~3mmほど [MTH4R(L)] 寄りの位置
11	MTH6R(L)	線分[MTH1R(L)]-[MTH3C]の中点から2mm~3mmほど [MTH7R(L)] 寄りの位置

40

再び図2を参照して、マーカ170のうち、額部に装着されるマーカ172A, ..., 172Dは、各マーカ170の位置のデータを頭部の動きに応じて補正するための補正用のデータの計測に用いられるマーカである。図2(A)及び(B)に示す例では、額部の皮膚の動きを抑制する拘束部材174を額部に貼付し、マーカ172A, ..., 172Dを、

50

拘束部材 174 を介して間接的に額部に装着している。なお、本実施の形態では、顔全体にマーカを装着する場合、マーカは、合計 137 箇所に装着される。

【0045】

データ処理装置 134 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モーションキャプチャデータ 160 を生成し、マルチモーダルコーパス作成装置 108 に出力する。なお、モーションキャプチャシステム 120 には、市販の光学式モーションキャプチャシステムを利用できる。市販の光学式モーションキャプチャシステムにおける赤外線カメラ 132 及びデータ処理装置 134 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。

【0046】

再び図 1 を参照して、録音システム 122 は、発話者 102 の発する音声を受音して音響信号を発生するためのマイクロホン 140A 及び 140B と、マイクロホン 140A 及び 140B が発生した音響信号を増幅するためのアンプ 142 と、アンプ 142 により増幅された音響信号を所定の形式でデジタル化して図示しない記録媒体に記録するための録音装置 144 とを含む。記録されたデータ 162 はマルチモーダルコーパス作成装置 108 に与えられる。本明細書では、録音装置 144 が記録し出力するデータ 162 を「音声収録データ」と呼ぶ。

【0047】

撮影システム 124 は、マイクロホン 140A 及び 140B と同様の機能を持つマイクロホン 140C と、テレプロンプタ 128 の後方にマイクロホン 140C からの出力を受け取るように配置され、テレプロンプタ 128 を通して発話者 102 の顔面及び首部の動画を撮影し、マイクロホン 140C が発生する音響信号と撮影した動画像とを、タイムコードジェネレータ 130 からのタイムコードを付与して所定の形式でデータ化し、図示しない記録媒体に記録するためのカムコーダ 150 と、動画像の撮影時の光源となる複数の照明装置 152A、152B、及び 152C（以下これらをまとめて「照明装置 152」と呼ぶことがある。）とを含む。カムコーダ 150 により記録されたデータは、マルチモーダルコーパス作成装置 108 に与えられる。本明細書では、撮影システム 124 が記録し出力するデータ 164 を「カムコーダ収録データ」と呼ぶ。

【0048】

図 1 に示す収録システム 104 はさらに、動画像の背景となるクロマキスクリーン 154 と、カムコーダ 150 により撮影される動画像を発話者 102 が確認できるように表示するためのモニタ 156 とを含む。

【0049】

[図 3 マルチモーダルコーパス作成装置の構成]

図 3 に、マルチモーダルコーパス作成装置 108（図 1 参照）の機能的構成をブロック図で示す。図 3 を参照して、マルチモーダルコーパス作成装置 108 は、モーションキャプチャデータ 160 をデータ処理装置 134 から取込むためのモーションキャプチャデータ取込部 180 と、音声収録データ 162 を録音装置 144 から取込むための音声収録データ取込部 182 と、カムコーダ収録データ 164 をカムコーダ 150 から取込むためのカムコーダ収録データ取込部 184 と、取込まれたモーションキャプチャデータ 160、音声収録データ 162、及びカムコーダ収録データ 164（以下、これらのデータをまとめて「収録データ」と呼ぶことがある。）を発話リスト 126 を構成する発話内容ごとに切出して、発話内容ごとの収録データのセット（以下、「発話別収録データセット」と呼ぶ。）200A、…、200L（以下これらをまとめて「発話別収録データセット 200」と呼ぶことがある。）を生成するための切出処理部 186 と、発話別収録データセット 200 を記憶するための発話別収録データセット記憶部 188 とを含む。なお、カムコーダ収録データ 164 のうちの音声データよりも音声収録データ 162 の方が高音質であるため、本実施の形態では音声収録データ 162 を用いる。

【0050】

発話別収録データセット 200 は、発話内容別に収録データをまとめたものである。発

10

20

30

40

50

話別収録データセット200A, ..., 200Lは各々、発話内容を表す言語データ210と、発話者102による当該発話内容の発話時に計測されたフレームのマーカデータ(マーカの測定位置データ)からなる発話別モーションキャプチャデータ212と、当該発話内容の発話時に収録された部分の音声収録データからなる発話別音声データ214と、当該発話内容が発話された区間に収録された動画像のデータからなる発話別動画像データ216とを含む。

【0051】

マルチモーダルコーパス作成装置108はさらに、モーションキャプチャデータの入力を受けて、これを頭部全体の動きをキャンセルするように正規化し、顔の器官の変化を表す顔器官変化量データ220を出力するため正規化処理部190と、発話別収録データセット記憶部188内の発話別収録データセット200A, ..., 200Lのいずれかを読み出し、その中の発話別モーションキャプチャデータ212を正規化処理部190に入力し、これに応答して正規化処理部190により出力される顔器官変化量データ220で、発話別モーションキャプチャデータ212を置換して発話別データセット202A, ..., 202L(以下これらをまとめて「発話別データセット202」と呼ぶことがある。)を生成し、マルチモーダルコーパス106(図1参照)に格納するための発話別データセット生成部192を含む。

【0052】

発話別データセット202は、マルチモーダルコーパス106を構成するデータを発話内容別にまとめたものである。発話別データセット202A, ..., 202Lはそれぞれ、同様のデータ構成を有する。例えば、発話別データセット202Aは、言語データ210と、発話別モーションキャプチャデータ212を正規化することにより得られる顔器官変化量データ220と、発話別音声データ214と、発話別動画像データ216とを含む。

【0053】

図4に、切出処理部186の構成をブロック図で示す。図4を参照して、切出処理部186は、取込まれたモーションキャプチャデータ160、音声収録データ162、及びカムコーダ収録データ164をそれぞれ一時的に記憶しておくための、モーションキャプチャデータ記憶部230、音声収録データ記憶部232、及びカムコーダ収録データ記憶部234と、入力装置116A及び出力装置118Aを用いて行なわれるユーザの操作、並びに発話リスト126に基づき、言語データ210の生成、及びカムコーダ収録データ164からの発話別動画像データ216の切出を行なうための動画像データ切出部240と、モーションキャプチャデータ160のタイムコード及び発話別動画像データ216のタイムコードに基づいて、モーションキャプチャデータ160から発話別モーションキャプチャデータ212を切出すためのモーションキャプチャデータ切出部242と、音声収録データ162を、カムコーダ収録データ164の音声データと同期させることにより音声収録データにタイムコードを付与するための同期処理部244と、この音声収録データ162のタイムコードと発話別動画像データ216のタイムコードとに基づいて、音声収録データ162からの発話別動画像に同期した発話別音声データ214を切出すための音声データ切出部246と、動画像データ切出部240により生成される言語データ210及び発話別動画像データ216、並びに当該データに対応する発話別モーションキャプチャデータ212及び発話別音声データ214をそれぞれ受けて一時的に保持し、発話内容ごとに発話別収録データセット200(A, ..., L)を形成して出力するためのデータセット形成部248とを含む。

【0054】

図1に示す録音装置144は、音声収録データにタイムコードを付与する機能を持たない。しかし音質はカムコーダ150により録音されたものよりも録音装置144により得られた音声収録データの方がよい。そこで、上記したように同期処理部244により音声収録データをカムコーダ収録データ164内の音声データに付与されたタイムコードと同期させる。より具体的には、同期処理部244は、カムコーダ収録データ164における音声のデータと、音声収録データ162との相互相関を計算し、相互相関が最大となるよ

10

20

30

40

50

うに音声収録データとカムコード収録データ164の音声データとのずれを計算し、その結果に基づいて音声収録データにタイムコードを付与する。

【0055】

正規化処理部190は、発話別モーションキャプチャデータ212を構成する各マーカデータに対しアフィン変換を行なうことにより、顔の各器官の変化に起因するマーカ位置の変化量のみからなる（頭部の動きに起因する変化量を除いた）顔器官変化量データを生成する機能を持つ。ここに、発話別モーションキャプチャデータ212におけるマーカデータを同次座標系で $P = P_x, P_y, P_z, 1$ と表現し、当該マーカデータを基に生成される顔器官変化量データを $P' = P'_x, P'_y, P'_z, 1$ と表現すると、アフィン行列 M は、次の式のように表現される。

【0056】

【数1】

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

上記式において、アフィン行列 M は、頭部の動きのみが含まれていると考えられる4箇所以上のマーカに対応するマーカデータから、特異値分解によって算出される。本実施の形態では、正規化用のマーカとして額部に4点、こめかみ部に2点、及び鼻部に2点のマーカを設け、それらを基準として各マーカの変化量の正規化を行なう。

【0057】

なお、首部のマーカの変化量は頭部の動きには影響を受けず、首自身の動きに影響される。そのため、上記の頭部に対するものと同様の考え方にしたいが、別途、首部の動きの補正用マーカ4点を用意し、頭部の動きに対する正規化と同様の処理を首部のマーカに対し行なう。

【0058】

図5に、正規化処理部190の構成をブロック図で示す。図5を参照して、正規化処理部190は、発話別モーションキャプチャデータ212の入力を受け、発話別モーションキャプチャデータ212の各フレームにおいて、マーカデータから、首部以外の、顔を含む頭部に装着されたマーカの位置を表す頭部マーカデータと、首部に装着されたマーカの位置を表す首部マーカデータと分類して出力するためのデータ分類部260とを含む。

【0059】

正規化処理部190はさらに、データ分類部260から頭部マーカデータを受け、当該マーカデータの中から補正用のマーカデータを選択するための頭部補正用マーカデータ選択部262と、頭部補正用マーカデータ選択部262により選択されたマーカデータをもとに特異値分解を行ない、頭部正規化のためのアフィン行列を算出するための頭部アフィン行列算出部264と、頭部アフィン行列算出部264により算出されたアフィン行列を用いて、データ分類部260により出力された頭部マーカデータに対しアフィン変換を行なうことにより、頭部に装着された各マーカの変化量を算出するための頭部マーカデータ変換部266とを含む。

【0060】

正規化処理部190はさらに、データ分類部260から首部マーカデータを受け、当該マーカデータの中から補正用のマーカデータを選択するための首部補正用マーカデータ選択部272と、首部補正用マーカデータ選択部272により選択されたマーカデータを基に特異値分解を行ない、首部正規化のためのアフィン行列を算出するための首部アフィン

10

20

30

40

50

行列算出部 274 と、データ分類部 260 により出力された首部マーカデータに対して、首部アフィン行列算出部 274 により算出されたアフィン行列を用いてアフィン変換を行なうことにより、首部に装着された各マーカの変化量を算出するための首部マーカデータ変換部 276 とを含む。

【0061】

正規化処理部 190 はさらに、頭部マーカデータ変換部 266 から頭部に装着された各マーカの変化量を表すデータを、首部マーカデータ変換部 276 から首部に装着された各マーカの変化量を表すデータを、それぞれ受け、フレームごとに当該データを統合することにより、正規化された顔器官変化量データ 220 を作成し、発話別データセット生成部 192 に出力するためのデータ統合部 278 を含む。

10

【0062】

図 6 に、アニメーション作成装置 114 (図 1 参照) の構成をブロック図で示す。図 6 を参照して、アニメーション作成装置 114 は、入力装置 116B 及び出力装置 118B に接続され、ユーザの操作にしたがい、図 2 に示すマーカ 170 に対応する仮想のマーカ (以下、単に「仮想マーカ」と呼ぶ。) を初期顔モデル 110 上に配置することにより、当該各仮想マーカの、初期顔モデル 110 を規定する座標系上での座標を設定するための仮想マーカ設定部 300 と、初期顔モデル 110 内の各ノードに対して、各ノードに近接する所定数 (本実施の形態では 3 個) の仮想マーカを当該ノードに対応するマーカに選び、その対応関係を付与した顔モデル (以下、「マーカ対応顔モデル」と呼ぶ。) 310 を作成するためのマーカ対応顔モデル作成部 302 とを含む。図 2 に示すマーカ 170 と、仮想マーカとの対応関係がこのようにして定義されることにより、発話時の発話者に装着された各マーカの位置を、顔モデル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、モーションキャプチャデータの座標系と顔モデルの座標系との間の変換も行なわれる。

20

【0063】

アニメーション作成装置 114 はさらに、入力装置 116B 及び出力装置 118B に接続され、ユーザの操作にしたがい、マルチモーダルコーパス 106 内の発話別データセット 202 中のいずれかを、作成予定のアニメーション 112 における発話内容に応じて選択し取得するための発話別データセット取得部 304 と、取得された発話別データセットにおける顔器官変化量データ 220 に基づき、初期顔モデル 110 が表現する顔の形状から、変形した顔モデルを順次作成するための顔モデル変形部 306 と、顔モデル変形部 306 により順次作成される変形した顔モデルに対し、テクスチャ等を付与して画像化することにより、アニメーション 112 を生成するための画像化部 308 とを含む。

30

【0064】

マーカ対応顔モデル作成部 302 は、初期顔モデル 110 のノードの中から、処理の対象となるノードを選択するためのノード選択部 312 と、ノード選択部 312 により選択されたノード (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、仮想マーカの座標の設定値に基づき選択するための仮想マーカ選択部 314 と、仮想マーカ選択部 314 により、各ノードに対し適切な仮想マーカが所定数選択されるように仮想マーカ選択部 314 を制御し、選択された所定数の仮想マーカ (以下これらの仮想マーカを選択ノードに対する「対応マーカ」と呼ぶ。) を特定する情報を処理対象のノードに付与するための選択マーカ検査部 316 とを含む。

40

【0065】

具体的には、選択マーカ検査部 316 は、仮想マーカ選択部 314 により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ために必要な条件を充足するかを検査する。条件が充足されなければ仮想マーカ選択部 314 に対し次にこのノードに近い仮想マーカを選択するように要求する。条件が充足されていればこの仮想マーカを当該ノードの対応マーカに指定する。さらに、対応マーカが 1 個指定されるたびに、対応マーカが 3 個選択されたかを検査し、3 個に満たない場合には新たな仮想マーカを選択するように仮想マーカ選択部 314 に対し要求す

50

る。3個となれば、選択マーカ検査部316は、ノード選択部312に対する次の処理対象のノードの選択要求を発生する。

【0066】

図7に、マーカ対応顔モデル作成部302により実行される、対応マーカの指定処理を実現するコンピュータプログラムの制御構造をフローチャートで示す。図7を参照して、対応マーカの指定処理が開始されると、ステップ340Aとステップ340Bとで囲まれた、ステップ342からステップ354までの処理を、初期顔モデル110における全ノードに対して処理が完了するまで実行する。

【0067】

ステップ342では、初期顔モデル110を構成するノードのうち、未処理のノードを1つ選択する。これを選択ノードとする。ステップ344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する。ステップ345では、以下の繰返しを制御するための変数*i*及び選択されたマーカの数を表す変数*j*に0を代入する。ステップ346では、変数*i*に1を加算する。

10

【0068】

ステップ347では、変数*i*の値が仮想マーカの数*Mmax*を超えているか否かを判定する。変数*i*の値が数*Mmax*を超えていればエラーとし、処理を終了する。このようなことは普通はないが、念のためにこのようなエラー処理を設けておく。変数*i*の値が数*Mmax*以下であれば制御はステップ348に進む。

20

【0069】

ステップ348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、初期顔モデル110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345に戻る。さもなければステップ350に進む。

【0070】

ステップ350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ひとつに指定する。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352に進む。ステップ352では、変数*j*に1を加算する。ステップ354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340Bに進む。さもなければステップ345に進む。

30

【0071】

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モデル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。例えば目の上まぶたと下まぶたとのように、間に境界エッジ(例えば目のふち)が存在する場合がある。この場合、上まぶたに位置するノードと、下まぶた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、例えば上まぶたのノードの変化量を算出する際に、下まぶたに存在するマーカの変化量を用いることは適当ではない。なお、線分がある境界エッジを横切っているか否かは、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基づいて判定する。

40

【0072】

図10に、初期顔モデル110における目輪郭部周辺のポリゴンと、仮想マーカとを示す。図10を参照して、初期顔モデル110の目輪郭部の周囲には、多数の三角形ポリゴンが存在する。このうち例えばポリゴン402は、3つのエッジ404A、404B、及び404Cにより囲まれている。エッジ404A及び404Bは、他のポリゴンと共有されている。しかし、エッジ404Cは、他のポリゴンと共有されていない。エッジ404C等2つのポリゴンにより共有されていないエッジは、初期顔モデル110の切れ目との接線又は外縁にあたる。このようなエッジが境界エッジとなる。

【0073】

50

再び図6を参照して、顔モデル変形部306は、あるフレームにおいて測定された、座標変換済みのマーカ変化量を各仮想マーカに付与する。さらに顔モデル変形部306は、マーカ対応顔モデル310のマーカ・ノード対応情報に基づき、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、顔モデルの変形を行なう。マーカ対応顔モデル310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 、変形後の顔モデルにおけるマーカの座標を M'_i とすると、顔モデル変形部306は、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。

【0074】

【数2】

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

10

なお、本実施の形態においては、 $n = 3$ である。すなわち、1つのノードに対応付ける仮想マーカの数は3である。

【0075】

[動作]

本実施の形態に係るマルチモーダルコーパス作成システム100は以下のように動作する。まず、図1に示すマルチモーダルコーパス作成システム100の収録システム104を用いた、音声、動画像、及びモーションキャプチャデータの収録プロセスについて説明する。

20

【0076】

発話者102の顔面及び首部には、事前に、表1に示すようなルールにしたがい図2に示すように多数のマーカを予め装着しておく。図1を参照して、赤外線カメラ132はそれぞれ、各マーカからの反射光を受光可能な所定の位置に、受光部を発話者102の顔面及び首部に向けて設置される。マイクロホン140A及び140Bはそれぞれ、発話者102の上部及び胸部等、発話者102の発する音声を受音可能な所定の位置に設置される。カムコーダ150は、発話者102の正面等、顔面及び首部の撮影に好適な位置に、受光部を発話者102に向けて設置される。なお、マイクロホン140Cは、発話者102の発する音声を受音可能で、かつカムコーダ150に接続可能な位置に設置される。照明装置152はそれぞれ、発話者102の顔にセルフシャドウが起こることを防止できる位置に設置される。例えば、照明装置152A、152B、及び152Cはそれぞれ、発話者102の左右、及び正面ローアングルから、発話者に向けて光が照射されるように設置される。クロマキスクリーン154は、カムコーダ150から見て発話者102の背後に設置される。

30

【0077】

テレプロンプタ128は、発話者102とカムコーダ150との間に、発話者102側からの光がカムコーダ150側に透過するよう設置される。モニタ156は、テレプロンプタ128の上部に画面を発話者102に向けて設置される。カムコーダ150は、テレプロンプタ128越しに発話者102を撮影することになる。そのため発話者102がテレプロンプタ128及びモニタの表示を見ると、発話者102の視線はテレプロンプタ128越しにカムコーダ150に向けられることになる。

40

【0078】

収録時には、発話リスト126を構成する文章等を発話リスト126にしたがいテレプロンプタ128が表示する。発話者102は、テレプロンプタ128及びモニタ156の表示を確認しながら、発話リスト126により指定された内容の文章等を順次発話する。

【0079】

発話時における顔の各部位の位置は、モーションキャプチャシステム120により次の

50

ようにして計測される。マーカはそれぞれ、発話時における顔の各器官の変化並びに頭部及び首部の動きに追従して移動する。赤外線カメラ132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134に出力する。データ処理装置134は、赤外線カメラ132からの映像信号の各フレームにタイムコードジェネレータ130からのタイムコードを付与し、当該映像信号を基に各マーカの位置をフレームごとに算出する。データ処理装置134は、各マーカの位置のデータをフレームごとにまとめてモーションキャプチャデータ160として蓄積する。

【0080】

発話時における発話者102の音声は、録音システム122により、次のようにして収録される。すなわち、マイクロホン140A及び140Bは、発話者102の音声を受音して、音響信号を発生する。アンプ142は、発生した音響信号の入力を受け、当該音響信号の各々を増幅して録音装置144に出力する。録音装置144は、増幅された音響信号をアンプ142から受け音声収録データ162として記録する。

【0081】

発話時における発話者102の顔の動画は、撮影システム124により、次のようにして収録される。すなわち、マイクロホン140Cは、140A及び140Bと同様に発話者102の音声を受音して音響信号を発生する。この音響信号は、カムコーダ150に与えられる。同時にカムコーダ150は、テレプロンプタ128越しに、発話中の発話者102のバストアップの動画を正面から撮影する。カムコーダ150は、動画とマイクロホン140Cからの音響信号とから所定の形式のカムコーダ収録データを形成し記録する。この際カムコーダ150は、タイムコードジェネレータ130のタイムコードをカムコーダ収録データ164の各フレームに付与する。

【0082】

以上の収録プロセスにより、タイムコードジェネレータ130のタイムコードが付与されたモーションキャプチャデータ160と、同じタイムコードが付与された音声及び動画のデータからなるカムコーダ収録データ164と、音声収録データ162とが同時に収録される。これらのデータは、マルチモーダルコーパス作成装置108に与えられる。

【0083】

[マルチモーダルコーパス作成装置108の動作]

図3を参照して、マルチモーダルコーパス作成装置108のモーションキャプチャデータ取込部180は、図1に示す収録システム104のデータ処理装置134よりモーションキャプチャデータ160を取込む。この際モーションキャプチャデータ取込部180は、モーションキャプチャデータ160を、3次元コンピュータグラフィックスを扱うソフトウェアで利用可能な形式で取込む。

【0084】

図8は、1フレーム分のモーションキャプチャデータ160に含まれるマーカデータを基に、各マーカの位置をコンピュータグラフィックスで表現した図である。図8を参照して、円形の目印はそれぞれ、当該フレームにおけるマーカの位置を表す。1フレーム分のモーションキャプチャデータは、マーカと同数のマーカデータを含む。

【0085】

再び図3を参照して、音声収録データ取込部182は、録音装置144より音声収録データ162を取込む。カムコーダ収録データ取込部184は、カムコーダ150よりカムコーダ収録データ164を取込む。取込まれたモーションキャプチャデータ160と、音声収録データ162と、カムコーダ収録データ164とはそれぞれ、切出処理部186に与えられる。

【0086】

切出処理部186は、発話リスト126を構成する文章、単語等の発話内容ごとに以下の動作により、発話別収録データセットを作成する。

【0087】

10

20

30

40

50

図4を参照して、モーションキャプチャデータ160と、音声収録データ162と、カムコード収録データ164はそれぞれ、モーションキャプチャデータ記憶部230、音声収録データ記憶部232、及びカムコード収録データ記憶部234に格納される。音声収録データ162とカムコード収録データ164が格納されると、同期処理部244は、カムコード収録データ164における音声のデータと音声収録データ162との相互相関を計算し、最大の相関が得られるように音声収録データをフレームに分割し、各フレームに対応するカムコード収録データ164の音声データのフレームに付与されていたものと同じタイムコードを付与する。同期処理部244は、処理後の音声収録データ162を音声収録データ記憶部232に格納する。

【0088】

ユーザが発話別のデータセット形成を指示するために入力装置116Aを用いて所定の操作を行なうと、動画像データ切出部240は、カムコード収録データ記憶部234からカムコード収録データ164を読み出す。動画像データ切出部240はさらに、発話リスト126を取得する。動画像データ切出部240は、ユーザの操作に応じて、カムコード収録データ164の動画像及び音声、並びに発話リストを出力装置118Aを介して出力する。出力装置118Aによる出力を参考にユーザが入力装置116Aを用いて、1つの発話内容に対応する動画像の収録された区間の開始位置及び終了位置を指定すると、動画像データ切出部240は、この入力にしたがい、指定された区間の動画像のデータをカムコード収録データ164から抽出し、発話別動画像データ216を生成してデータセット形成部248に与える。発話別動画像データ216のうち、その開始と終了とを表すタイムコードが、モーションキャプチャデータ切出部242、音声データ切出部246、及びデータセット形成部248に与えられる。動画像データ切出部240はさらに、抽出した部分の動画像に対応する言語データ210を、ユーザによる入力及び発話リスト126を基に生成する。生成された言語データ210は、データセット形成部248に与えられる。

【0089】

発話別動画像データ216の開始と終了とを表すタイムコードに回答して、モーションキャプチャデータ切出部242は、指定された区間を特定する。モーションキャプチャデータ切出部242は、モーションキャプチャデータ記憶部230内のモーションキャプチャデータ160から、当該区間に対応するデータを抽出して発話別モーションキャプチャデータ212を生成し、データセット形成部248に与える。

【0090】

音声データ切出部246は、発話別動画像データ216の開始と終了とを表すタイムコードに回答して、音声収録データからタイムコードにより指定された区間に対応するデータを抽出して発話別音声データ214を生成する。生成された発話別音声データ214は、データセット形成部248に与えられる。

【0091】

データセット形成部248は、言語データ210、発話別動画像データ216、発話別モーションキャプチャデータ212、及び発話別音声データ214が与えられたことに回答して、これら与えられたデータをまとめて発話別収録データセット200を生成し、図3に示す発話別収録データセット記憶部188に格納する。以上の動作により、発話内容ごとの発話別収録データセット200が形成され、発話別データセット記憶部188に格納される。

【0092】

発話別データセット生成部192及び正規化処理部190は、発話別収録データセット200の各々について以下の処理を行ない、発話別データセット202を生成する。すなわち、発話別データセット生成部192は発話別収録データセット記憶部188から発話別収録データセットを1セット分読み出す。発話別データセット生成部192はさらに、発話別モーションキャプチャデータ212から、1フレーム分のデータを正規化処理部190に与える。

【0093】

10

20

30

40

50

図5を参照して、正規化処理部190のデータ分類部260は、1フレーム分のデータが与えられたことに応答して、当該フレームにおけるマーカデータを、頭部マーカデータと、首部マーカデータとに分類する。データ分類部260は、頭部マーカデータを頭部補正用マーカデータ選択部262及び頭部マーカデータ変換部266に与え、首部マーカデータを首部補正用マーカデータ選択部272及び首部マーカデータ変換部276に与える。

【0094】

頭部補正用マーカデータ選択部262は、与えられたマーカデータの中から、予め定められた8箇所のマーカデータを補正用のマーカとして選択し、それぞれ頭部アフィン行列算出部264に与える。頭部アフィン行列算出部264は、与えられたマーカデータからの特異値分解によってアフィン行列Mを算出し、頭部マーカデータ変換部266に与える。頭部マーカデータ変換部266は、与えられた頭部マーカデータを、このアフィン行列Mによって変換する。この変換により、マーカデータはそれぞれ、頭部の動きを除いた正規化した変化量に変換される。頭部マーカデータ変換部266は、各マーカの正規化後の変化量をデータ統合部278に与える。

10

【0095】

首部補正用マーカデータ選択部272は、与えられたマーカデータの中から予め定められた4箇所の首部補正用のマーカデータを選択し、それぞれ首部アフィン行列算出部274に与える。首部アフィン行列算出部274は与えられたマーカデータを用いて首部補正用のアフィン行列を算出し首部マーカデータ変換部276に与える。首部マーカデータ変換部276は、データ分類部260から与えられた首部マーカデータを首部アフィン行列算出部274から与えられたアフィン行列で変換する。この変換により、マーカデータはそれぞれ、首部の動きを除いた正規化した変化量に変換される。首部マーカデータ変換部276は、各マーカの変化量をデータ統合部278に与える。

20

【0096】

データ統合部278は、頭部マーカデータ変換部266と首部マーカデータ変換部276とからそれぞれ与えられるマーカの変化量のデータを統合して、1フレーム分の顔器官変化量データ220を生成する。データ統合部278は、生成した顔器官位置変化量データ220を発話別データセット生成部192(図3参照)に返す。

【0097】

図3を参照して、発話別データセット生成部192は、正規化処理部190から1フレーム分の顔器官変化量データ220が返されると、発話別モーションキャプチャデータ212における当該フレームのデータを、そのフレームの顔器官変化量データ220で置換し、言語データ210、発話別動画データ216、及び発話別音声データ214とともに発話別データセット202に出力する。発話別データセット生成部192はこの後、新たに1フレーム分のマーカデータを正規化処理部190に与え、上記と同様の処理を繰り返す。

30

【0098】

正規化処理部190及び発話別データセット生成部192は、以上の動作を発話別収録データセット200の各々の全フレームについて繰り返すことにより、発話別データセット202を形成する。形成された発話別データセット202は、図1に示すマルチモーダルコーパス106に格納される。

40

【0099】

[アニメーションの作成]

次に、アニメーション作成装置114がアニメーション112を作成する動作について説明する。図6を参照して、アニメーション作成装置114に初期顔モデル110が与えられると、アニメーション作成装置114は、動作を開始する。図9に、初期顔モデル110の一例を示す。図9を参照して、この初期顔モデル110は、発話者102の顔の静止画像と所定のワイヤフレームモデルとを整合させることにより準備された形状モデルである。この顔モデルは、約750のポリゴンで構成されている。初期顔モデル110は、

50

仮想マーカ設定部 300 と、マーカ対応顔モデル作成部 302 のノード選択部 312 及び選択マーカ検査部 316 とに与えられる。

【0100】

仮想マーカ設定部 300 は、初期顔モデル 110 を画像化して出力装置 118B に出力する等して、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 116B を介して受ける。初期顔モデル 110 上での仮想マーカの位置は、既に述べた表 1 と同様のルールにしたがって指定される。そのため、初期顔モデル 110 における顔器官と仮想マーカとの位置関係は、発話者 102 の顔器官と当該発話者 102 に装着されたマーカとの位置関係に対応する。

【0101】

仮想マーカ設定部 300 は、ユーザによる指定を基に、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。仮想マーカ設定部 300 は、当該各仮想マーカの識別子と当該仮想マーカの座標とを、マーカ対応顔モデル作成部 302 の仮想マーカ選択部 314 に与える。

【0102】

マーカ対応顔モデル作成部 302 は、初期顔モデル 110 と仮想マーカの識別子及び座標とが与えられたことに応答して、初期顔モデル 110 の各ノードに対して、当該ノードの対応マーカを次のようにして特定する。まず、ノード選択部 312 が、初期顔モデル 110 を構成するノードの中からノードを 1 つ選択する。このノードが選択ノードである。選択ノードと全ての仮想マーカとの距離を算出し、仮想マーカを距離の昇順にソートしてリスト化する。このリストの先頭の 1 つを選び、その仮想マーカと選択ノードとを結ぶ線が顔モデルの境界エッジを横切るか否かを判定する。横切らなければこの仮想マーカを選択ノードの対応ノードの 1 つに選択する。横切っていればリストの次の仮想マーカを選択し、同じ処理を繰り返す。

【0103】

こうして、選択ノードに対し 3 つの仮想マーカが当該選択ノードの対応ノードとして特定される。対応ノードと選択ノードとを結ぶ線分のいずれも、顔モデルの境界エッジを横切らない。

【0104】

例えば、図 10 を参照して、ノード 412 が選択ノードであるときを考える。なお、初期顔モデル 110 の目輪郭部周囲において、仮想マーカ 410A, ..., 410L が設定されているものとする。仮想マーカ選択部 314 は、ノード 412 の座標と、仮想マーカ 410A, ..., 410L の座標データとを基に、選択ノード 412 と仮想マーカ 410A, ..., 410L との間の距離をそれぞれ算出する。仮想マーカ選択部 314 は、仮想マーカ 410A, ..., 410Lの中から、ノード 412 に最も近い位置にある仮想マーカ 410K を選択する。

【0105】

選択マーカ検査部 316 は、ノード 412 と選択された仮想マーカ 410K とを結ぶ線分が境界エッジを横切るか否かを検査する。ノード 412 と選択された仮想マーカ 410K とを結ぶ線分は、いずれの境界エッジも横切らない。そのため、選択マーカ検査部 316 は、当該仮想マーカ 410K をノード 412 の対応マーカに指定する。選択マーカ検査部 316 はさらに、新たな仮想マーカの選択要求を仮想マーカ選択部 314 に与える。

【0106】

仮想マーカ選択部 314 は、選択マーカ検査部 316 からの通知及び要求に応答して、仮想マーカ 410K の次にノード 412 に近い位置にある仮想マーカを選択する。図 10 に示す例では、この選択により、ノード 412 に 2 番目に近接する仮想マーカ 410B が選択される。

【0107】

選択マーカ検査部 316 は、選択された仮想マーカ 410B についての検査を上記の動

10

20

30

40

50

作と同様の動作で行なう。この場合、ノード412と仮想マーカ410Bとを結ぶ線分は、境界エッジを横切る。そのため、選択マーカ検査部316は、当該仮想マーカ410Bをノード412の対応マーカに指定せず対象から除外する。選択マーカ検査部316はさらに、新たな仮想マーカの選択を仮想マーカ選択部314に要求する。

【0108】

仮想マーカ選択部314及び選択マーカ検査部316が以上の動作を繰返し、ノード412の対応マーカとして3個の仮想マーカ(図10に示す例では仮想マーカ410J, 410K, 及び410L)が指定されると、ノード412に対する仮想マーカの対応付けが完了する。選択マーカ検査部316はノード412とその対応マーカに関するマーカ・ノード対応情報をマーカ対応顔モデルの一部として出力し、ノード選択部312に対し新たなノードの選択要求を与える。

10

【0109】

ノード選択部312は、選択マーカ検査部316からの要求に応答して、初期顔モデル110を構成するノードのうち、対応付けが未完了のノードから1つのノードを選択する。以下、上記したノード選択部312、仮想マーカ選択部314、及び選択マーカ検査部316の動作が、全てのノードに対して対応マーカが決定されるまで繰返される。

【0110】

こうして、マーカ対応顔モデル作成部302により、各ノードに対し3個の仮想マーカを対応付けるマーカ対応顔モデル310(図6参照)が生成される。マーカ対応顔モデル310は、図6に示す顔モデル変形部306に与えられる。

20

【0111】

次に、アニメーション作成装置114が、マーカ対応顔モデル310を用いてアニメーションを作成する動作について説明する。図6を参照して、ユーザが入力装置116Bを用いて、発話内容等を入力すると、当該入力は、発話別データセット取得部304に与えられる。発話別データセット取得部304は、マルチモーダルコーパス106から、入力された発話内容等に対応する発話別データセット202A, ..., 202L(図3参照)を讀出し、当該発話別データセット内の顔器官変化量データ220(図3参照)を顔モデル変形部306に与える。

【0112】

この時点で顔モデル変形部306には、顔器官変化量データ220と、マーカ対応顔モデル310とが与えられている。マーカ対応顔モデル310の各ノードには、当該顔モデル上の仮想マーカが3個指定されている。顔モデル変形部306は、顔器官変化量データ220をもとに、各マーカの位置の変化量に基づき、マーカ対応顔モデル310中の各ノードの変化量を次のようにして算出する。

30

【0113】

すなわち、顔モデル変形部306はまず、マーカ対応顔モデル310上における仮想マーカの座標を取得する。仮想マーカはそれぞれ、顔器官変化量データ220におけるマーカと対応関係にある。そこで、顔モデル変形部306は、顔器官変化量データ220における1フレーム分のデータを基に、仮想マーカの各々に、当該仮想マーカに対応するマーカの変化量を付与し、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。

40

【0114】

さらに顔モデル変形部306は、1つのノードの変化量を、ノードに対し指定された3個の対応マーカの座標を基に決定する。ここに、あるノード座標をNとする。また当該ノードの対応マーカの変化前の座標をそれぞれ M_i ($1 \leq i \leq n = 3$)とする。さらに、当該対応マーカについて、1フレーム分の変化量が付与された後の座標を M'_i とする。顔モデル変形部306は、ノードの変化量ベクトル v を次の式により算出する。

【0115】

【数 3】

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

上記の式でノードの変化量ベクトル v を、変形前の当該ノードに対し付与することにより、変化後のノードの座標が算出される。顔モデル変形部 306 は、フレームごとに、マーカ対応顔モデル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、変形した顔モデルがフレームごとに生成される。顔モデル変形部 306 は、変形した顔モデルの各々を、画像化部 308 に与える。

10

【0116】

画像化部 308 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してそれらを画像化することにより、アニメーション 112 における各コマの画像を生成する。さらに、必要に応じて、コマの間引き等の処理を行ない、一連の動画像を形成する。形成した動画像が、アニメーション 112 となる。

【0117】

図 11 に、発話中における発話者 102 の顔画像と、顔器官変化量データ 220 及び図 9 に示す初期顔モデル 110 をもとに作成したアニメーション 112 における顔画像とを対比して示す。図 11 を参照して、1 段目には、マルチモーダルコーパス 106 に格納された動画像データのうち、異なる 5 つの発話内容の発話中にそれぞれ撮影された動画像中のフレームの画像を示す。発話内容に応じて、発話者 102 の口及び目等顔の各器官の形状が変化している。これらの画像のいずれにおいても、発話者の頭部の向き、大きさ、及び傾きは他の画像におけるそれらとは僅かながら異なる。この相違は、各画像における額部のマーカ及び拘束部材の位置に顕著に現れている。

20

【0118】

2 段目の画像は、1 段目の各画像の収録と同じ時点での顔器官変化量データ 220 に基づき変形した顔モデルにおけるポリゴンの形状を表す画像である。3 段目は、アニメーション 112 において、2 段目に示すポリゴン形状の顔モデルをもとに画像化されたフレームの画像である。対応する 1 段目の画像と比較すると、2 段目及び 3 段目の画像における口及び目等、顔の各器官の形状は、1 段目の動画像と同様に変化している。また、2 段目及び 3 段目の画像においては、顔の各器官の変化量に応じて顔モデルを変形させているため、頭部の向き、大きさは一定に保たれている。3 段目の各画像における額のマーカの位置は、一定している。

30

【0119】

以上のように、本実施の形態では、マルチモーダルコーパスを、発話者による発話中の音声、動画像、及び顔部位の位置の計測データを基に作成する。顔部位の位置の計測には光学式モーションキャプチャシステムを用いるため、顔部位の位置を動画像から推定しなくてもよく、高速度で 3 次元の位置計測が行なえる。その結果、顔部位の特徴量の算出が容易になる。また、発話の収録時には、顔部位の同定に用いる多数のマーカを、事前に定めたルールにしたがい発話者 102 に装着する。したがって、高精度かつ詳細に顔部位の変化量を得ることができる。また、複数の発話者から、又は同一の発話者から複数回にわたってそれぞれ収録を行なう場合であっても、計測条件を安定させることが容易で再現性の高い計測をすることが可能となる。その結果、大規模なマルチモーダルコーパスを作成することが可能になる。

40

【0120】

本実施の形態のマルチモーダルコーパス作成装置は、モーションキャプチャデータを基に、顔器官の変化量の算出を行なう。そのため、動画像の光学的な誤差に影響を受けることなく発話中の各器官の変化を正確にコーパス化できる。マルチモーダルコーパス作成装置は、モーションキャプチャデータを正規化して、発話中の顔器官の変化量を算出するた

50

め、発話者の頭部全体の回転及び移動等に影響されることなく、顔器官の変化量を得ることができる。よって、顔器官の変化量をより高精度にコーパス化できる。

【 0 1 2 1 】

また、音声、動画、及び顔器官の変化量のデータを同期させてコーパス化するため、音声と顔器官の変化量との対応関係を詳細に得ることができる。そのため、音声言語処理技術において確立している種々の手法を、発話中の顔器官の変化量に関する処理に適用することができる。

【 0 1 2 2 】

さらに、コーパスを構成するデータに対応する顔器官の位置は一定に保たれる。よって、当該コーパスの利用が容易になる。

【 0 1 2 3 】

本実施の形態のアニメーション作成装置は、発話者の顔部位の計測データを基に構築されたマルチモーダルコーパスに基づき、発話中の顔器官の変化量をモデルに割当てることにより、アニメーションを作成する。よって、動画像を用いた手法と同様に自然なアニメーションを作成することができる。

【 0 1 2 4 】

また、本実施の形態のアニメーション作成装置は、モデルベースでアニメーションを作成するため、バリエーションに富んだアニメーションの作成が可能となる。

【 0 1 2 5 】

さらに、本実施の形態のアニメーション作成装置は、マルチモーダルコーパスを基に、発話中の顔器官の変化量をモデルに割当てることによりアニメーションを作成する。顔器官の特徴点は事前にルールとして定められている。したがって、どのようなモデルに対しても、当該モデルにおける特徴点をルールにしたがい指定するだけで、モデルを発話時の音声及び各器官の動きに適切に同期した自然なアニメーションを作成できる。よって、手軽に高度なアニメーションを作成することができる。

[コンピュータによる実現及び動作]

なお、本実施の形態のマルチモーダルコーパス作成装置 1 0 8 及びアニメーション作成装置 1 1 4 は、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 2 はこのコンピュータシステム 5 0 0 の外観を示し、図 1 3 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 6 】

図 1 2 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 7 】

図 1 3 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。コンピュータ 5 1 0 はさらに、データ処理装置 1 3 4 (図 1 参照) 及びバス 5 3 6 に接続されたデータインタフェース 5 4 0 と、録音装置 1 4 4 (図 1 参照) 及びバス 5 3 6 に接続されたメディアコンバータ 5 4 2 と、カムコーダ 1 5 0 (図 1 参照) 及びバス 5 3 6 に接続されたキャプチャカード 5 4 4 とを含む。

【 0 1 2 8 】

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

10

20

30

40

50

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。

【0129】

コンピュータシステム500にマルチモーダルコーパス作成装置108又はアニメーション作成装置114の機能を実現させるためのコンピュータプログラムは、CD-ROMドライブ520又はFDドライブ522に挿入されるCD-ROM532又はFD534に記憶され、さらにハードディスク524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510に送信されハードディスク524に記憶されてもよい。プログラムは実行の際にRAM530にロードされる。CD-ROM532から、FD534から、又はネットワークを介して、直接にRAM530にプログラムをロードしてもよい。

10

【0130】

このプログラムは、コンピュータ510にこの実施の形態のマルチモーダルコーパス作成装置108又はアニメーション作成装置114の機能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510上で動作するオペレーティングシステム(OS)又はサードパーティのプログラム、若しくはコンピュータ510にインストールされる各種ツールキットのモジュールにより提供される。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、上記したマルチモーダルコーパス作成装置108又はアニメーション作成装置114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500の動作は周知であるので、ここでは繰返さない。

20

【0131】

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

【図面の簡単な説明】

【0132】

【図1】マルチモーダルコーパス作成システム100全体の構成を示す図である。

30

【図2】マーカが設置された状態での、発話者102の顔面及び首部の外観の一例を示す正面図及び側面図である。

【図3】マルチモーダルコーパス作成装置108の構成を示すブロック図である。

【図4】切出処理部186の構成を示すブロック図である。

【図5】正規化処理部190の構成を示すブロック図である。

【図6】アニメーション作成装置114の構成を示すブロック図である。

【図7】対応マーカの指定処理の制御構造を示すフローチャートである。

【図8】1フレーム分のモーションキャプチャデータ160により表現されるマーカの位置を模式的に示す図である。

【図9】初期顔モデル110の一例を示す図である。

40

【図10】初期顔モデル110の目輪郭部周辺におけるポリゴン、仮想マーカの概要を示す図である。

【図11】動画像における発話者102の顔の画像と、アニメーション112における顔の画像との変化を示す図である。

【図12】本発明の実施の形態に係るマルチモーダルコーパス作成装置108及びアニメーション作成装置114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13】図12に示すコンピュータシステムのブロック図である。

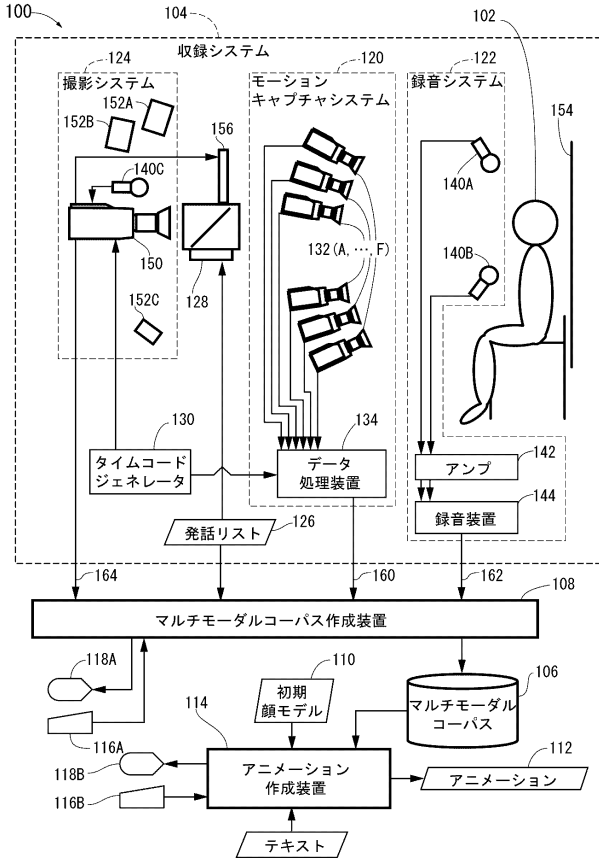
【符号の説明】

【0133】

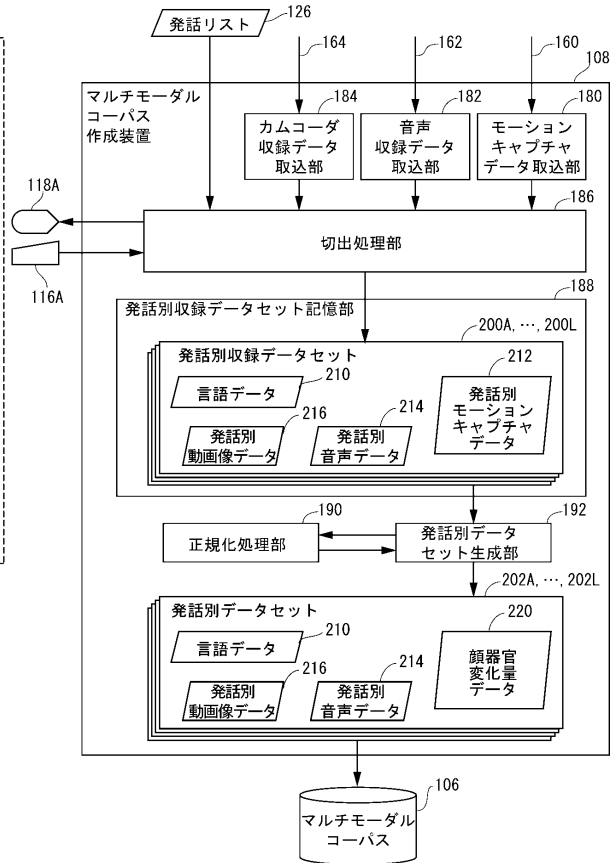
50

1 0 0	システム	
1 0 4	収録システム	
1 0 6	マルチモーダルコーパス	
1 0 8	マルチモーダルコーパス作成装置	
1 1 0	初期顔モデル	
1 1 2	アニメーション	
1 1 4	アニメーション作成装置	
1 2 2	録音システム	
1 2 4	撮影システム	
1 3 4	データ処理装置	10
1 7 0 A , ... , 1 7 0 M	マーカ	
1 8 0	モーションキャプチャデータ取込部	
1 8 2	音声収録データ取込部	
1 8 4	カムコーダ収録データ取込部	
1 8 6	切出処理部	
1 8 8	発話別収録データセット記憶部	
1 9 0	正規化処理部	
1 9 2	発話別データセット生成部	
2 0 0 A , ... , 2 0 0 L	発話別収録データセット	
2 0 2 A , ... , 2 0 2 L	発話別データセット	20
2 1 2	発話別モーションキャプチャデータ	
2 1 4	発話別音声データ	
2 1 6	発話別動画像データ	
2 2 0	顔器官変化量データ	
2 3 0	モーションキャプチャデータ記憶部	
2 3 2	音声収録データ記憶部	
2 3 4	カムコーダ収録データ記憶部	
2 4 0	動画像データ切出部	
2 4 2	モーションキャプチャデータ切出部	
2 4 4	同期処理部	30
2 4 6	音声データ切出部	
2 4 8	データセット形成部	
2 6 0	データ分類部	
2 6 2	頭部補正用マーカデータ選択部	
2 6 4	頭部アフィン行列算出部	
2 6 6	頭部マーカデータ変換部	
2 7 2	首部補正用マーカデータ選択部	
2 7 4	首部アフィン行列算出部	
2 7 6	首部マーカデータ変換部	
2 7 8	データ統合部	40
3 0 0	仮想マーカ設定部	
3 0 2	マーカ対応顔モデル作成部	
3 0 4	発話別データセット取得部	
3 0 6	顔モデル変形部	
3 0 8	画像化部	

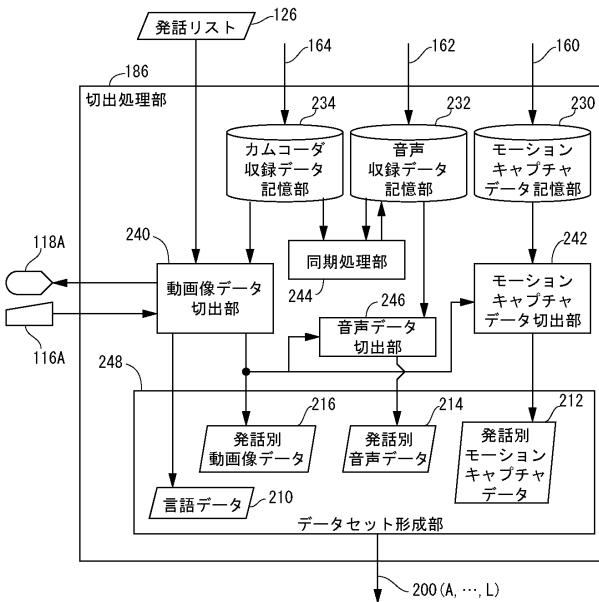
【図1】



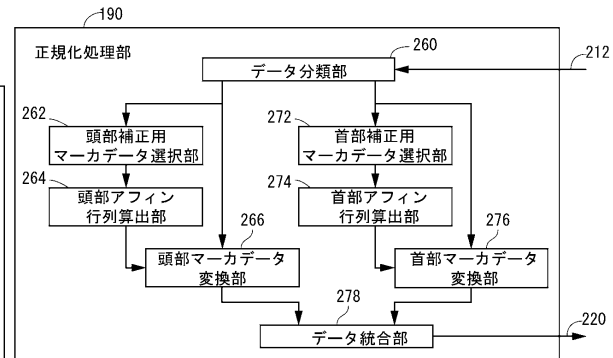
【図3】



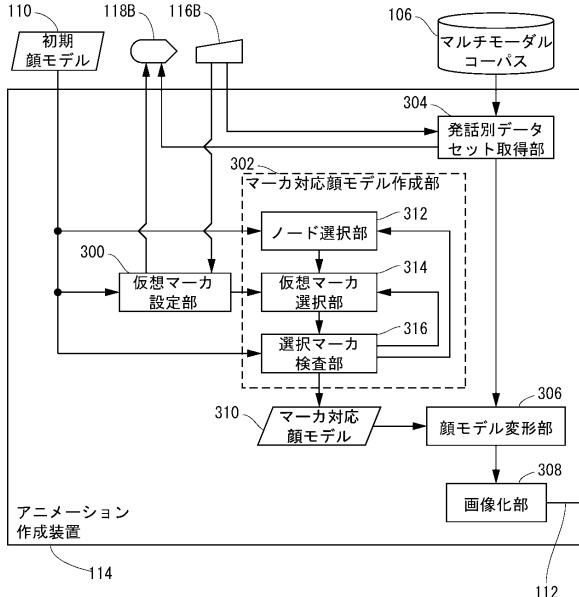
【図4】



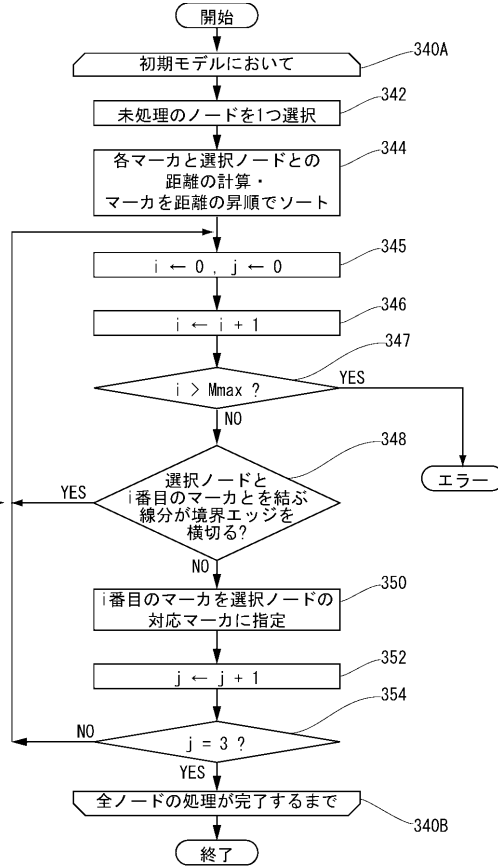
【図5】



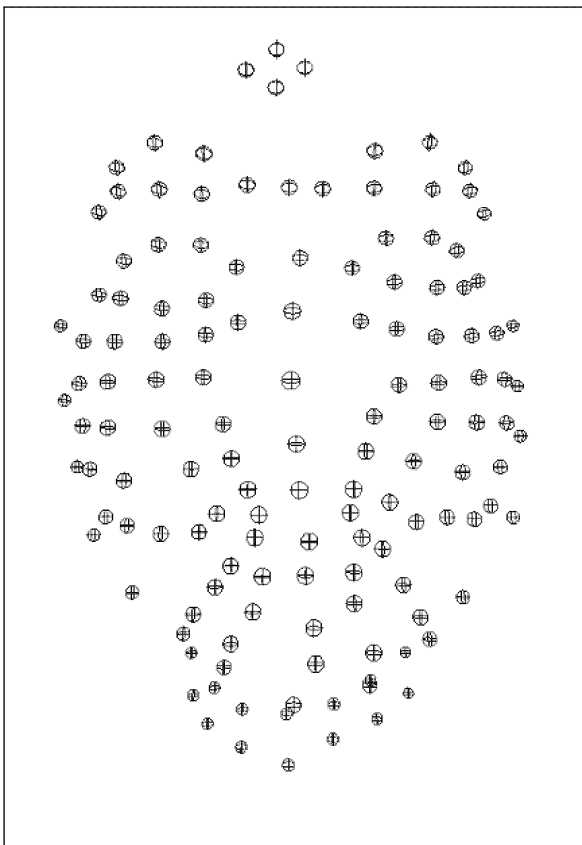
【図6】



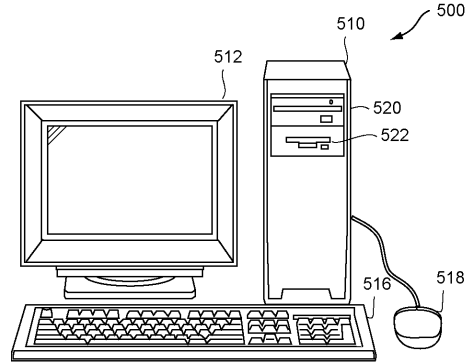
【図7】



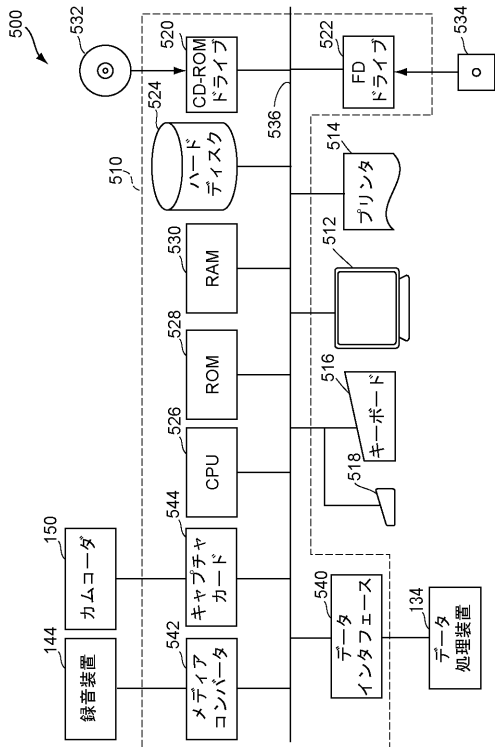
【図8】



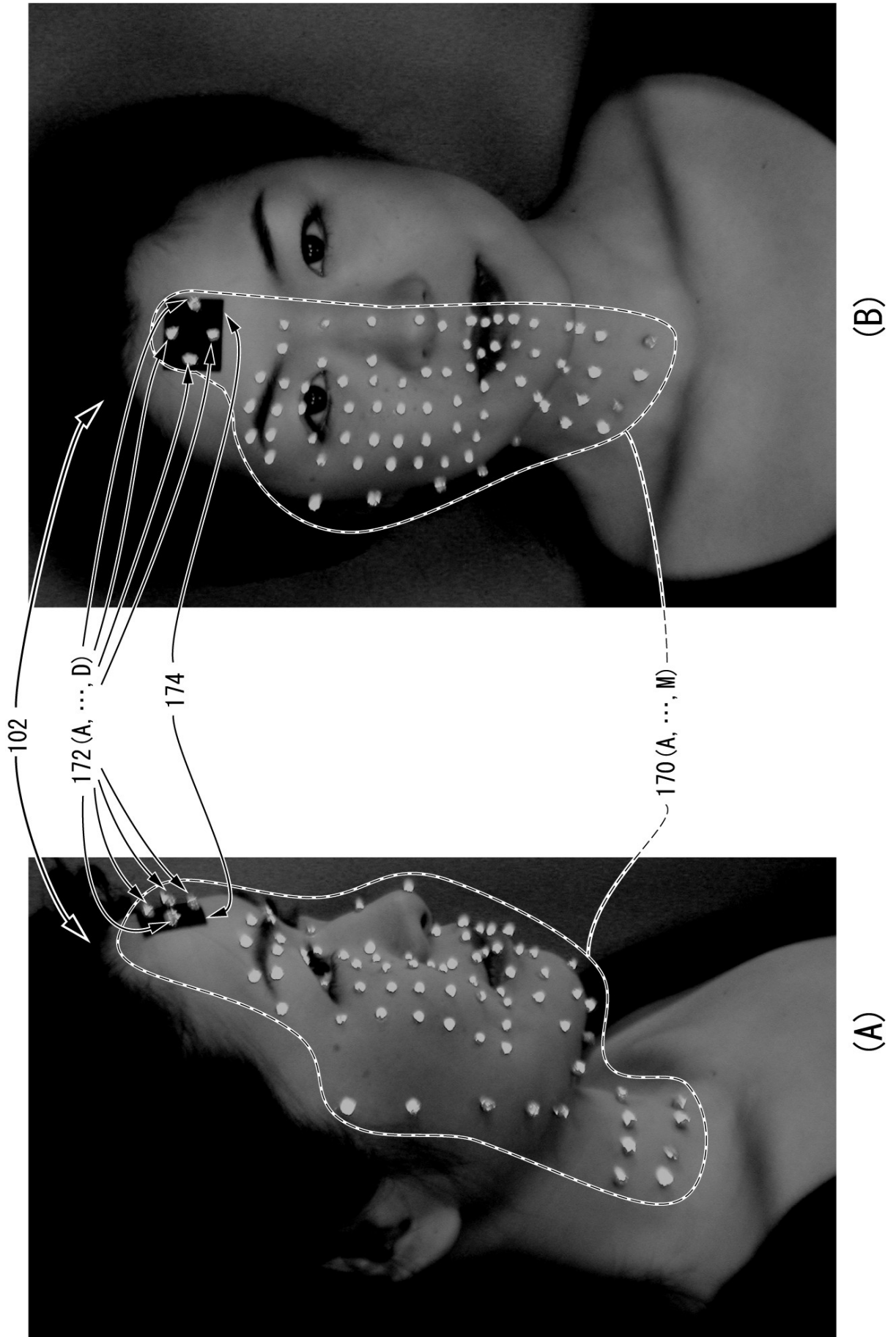
【図12】



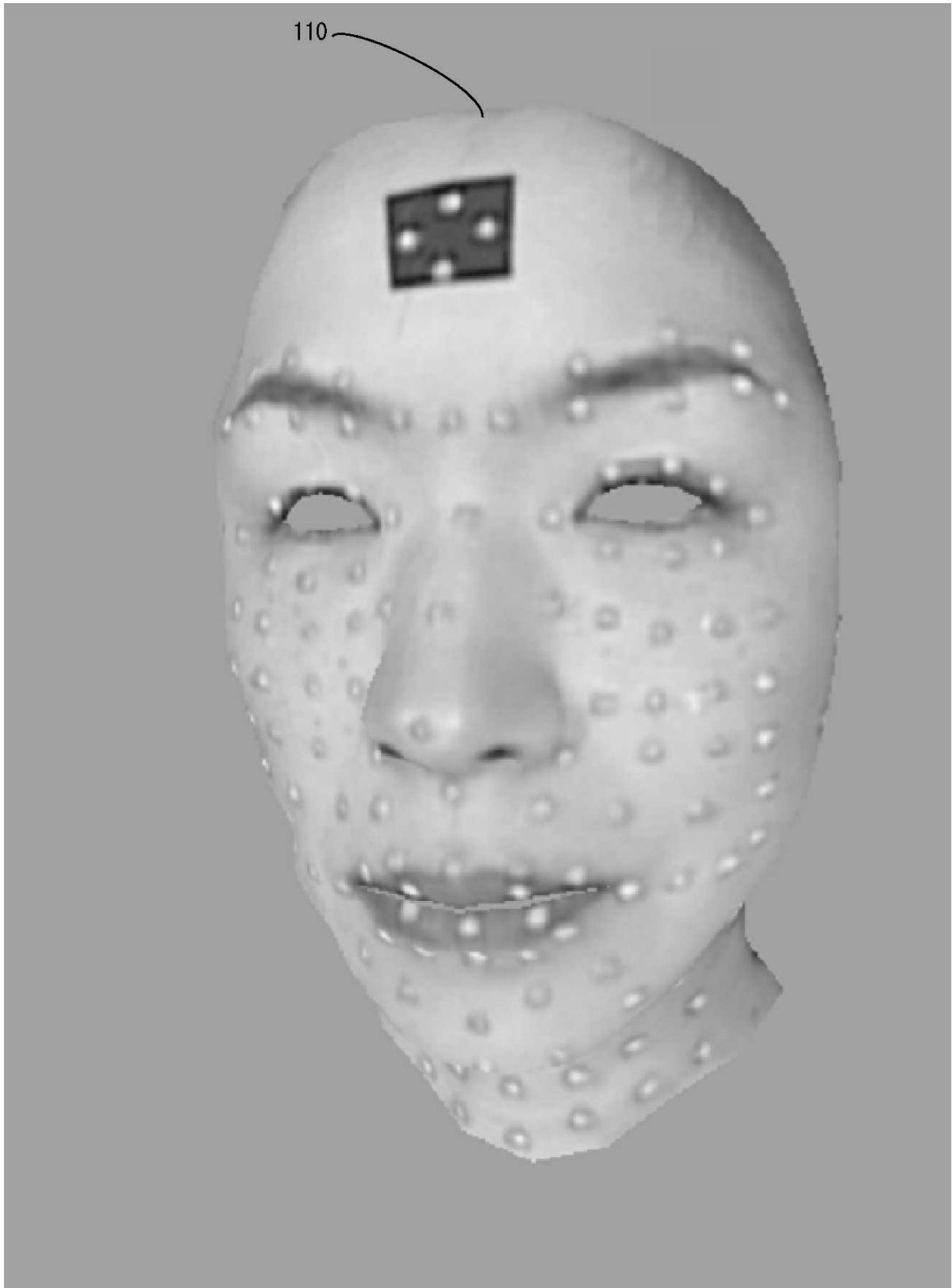
【 図 13 】



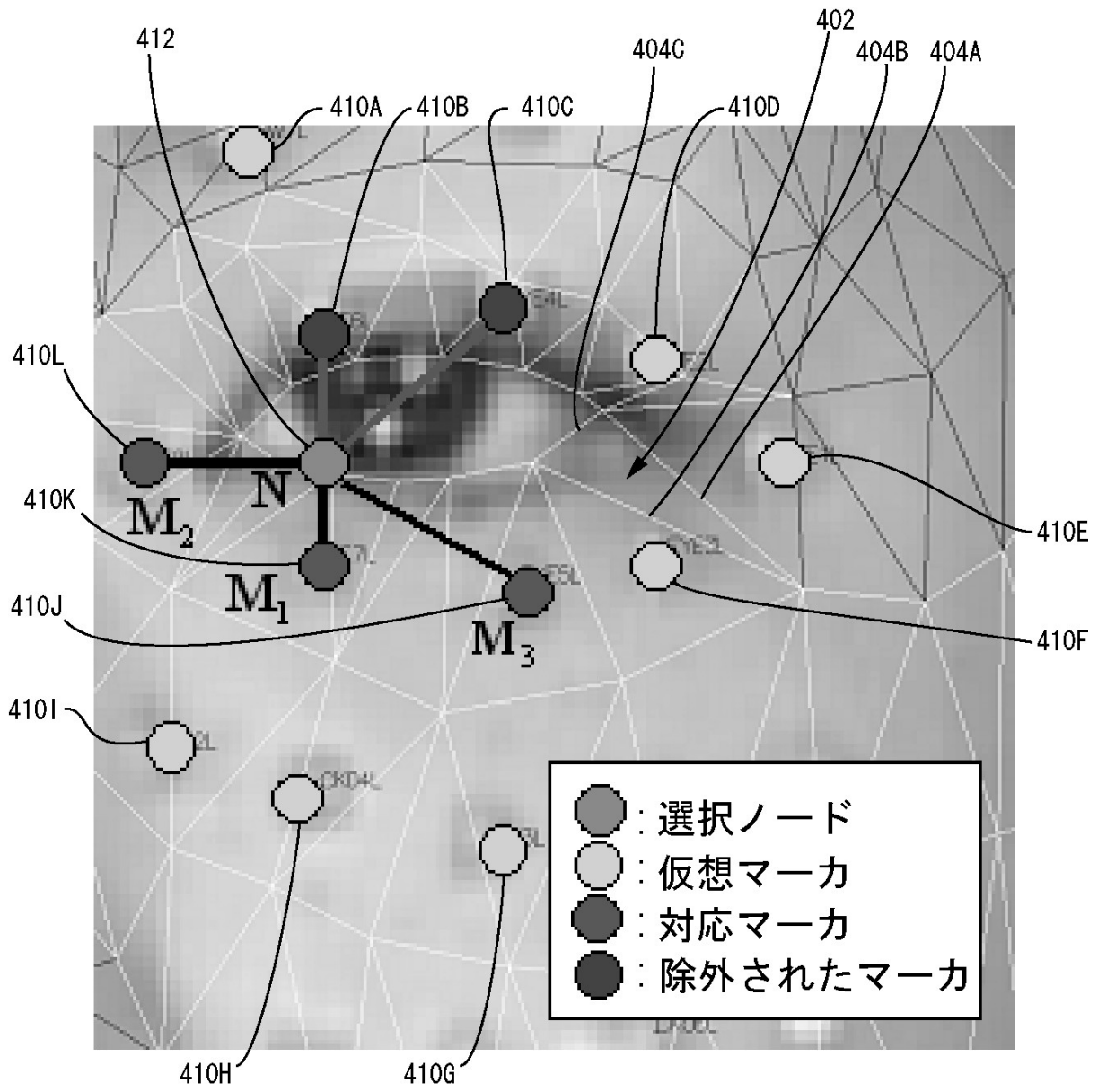
【 図 2 】



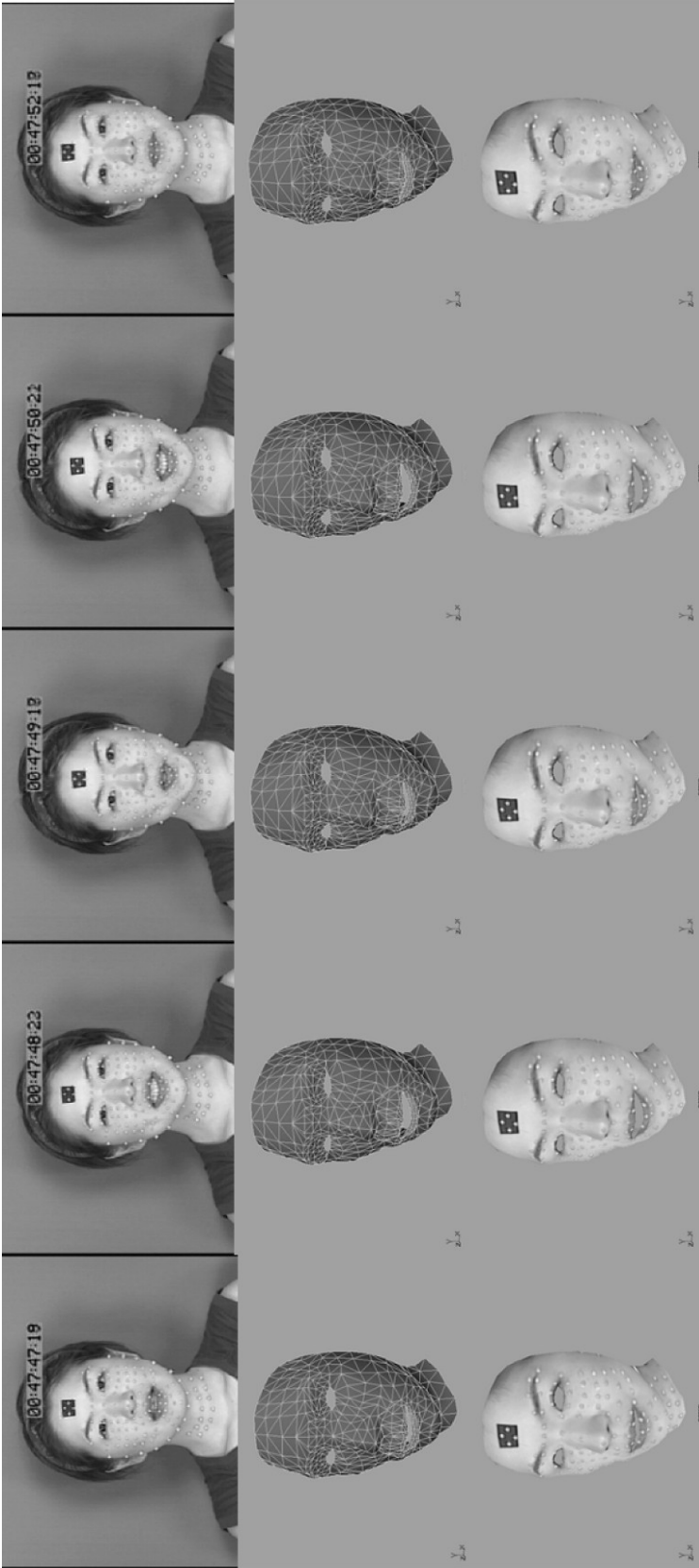
【図9】



【図10】



【 1 1 】



動画像

顔モデル

アニメーション

フロントページの続き

審査官 松永 隆志

(56)参考文献 特開2000-315259(JP,A)

楊大昭、外5名、人体の行動計測による自然な人体モーションと顔表情の合成手法の検討、映像情報メディア学会技術報告 ヒューマンインフォメーション ネットワーク映像メディア、日本、社団法人映像情報メディア学会、1998年 6月 1日、Vol.22 No.28、P65~P72

(58)調査した分野(Int.Cl., DB名)

G06T 17/40

G06T 1/00