

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4631077号
(P4631077)

(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
G O 6 T 13/40 (2011.01) G O 6 T 15/70 B

請求項の数 6 (全 28 頁)

(21) 出願番号	特願2006-128110 (P2006-128110)	(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5月2日(2006.5.2)	(74) 代理人	100099933 弁理士 清水 敏
(65) 公開番号	特開2007-299300 (P2007-299300A)	(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11月15日(2007.11.15)	(72) 発明者	川本 真一 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
審査請求日	平成19年12月18日(2007.12.18)	(72) 発明者	中村 哲 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

最終頁に続く

(54) 【発明の名称】 アニメーション作成装置

(57) 【特許請求の範囲】

【請求項1】

音声信号を受け、当該音声信号の表す音素列内の各音素の継続時間中の、所定のキーフレーム時刻における画像により構成されるキーフレーム画像を表すキーフレームデータを作成するためのキーフレームデータ作成手段と、

前記キーフレームデータ作成手段により作成される前記キーフレームデータに基づき、前記音声信号に同期して変化する一連の画像からなる画像のアニメーションを生成するためのアニメーション生成手段と、前記所定のキーフレーム時刻は、前記音素列内の各音素の継続時間の開始時刻であり、

予め定められた複数種類のテキストをユーザに選択させるためのテキスト選択手段と、
前記テキスト選択手段によりテキストが選択されたことに基づき、ユーザの音声録音して前記音声信号に変換し、選択された前記テキストとともに前記キーフレームデータ作成手段に与えるための手段とを含む、アニメーション作成装置であって、

前記キーフレームデータ作成手段は、
音素を、所定の基準画像を含む所定の複数個の画像のいずれかにマッピングするマッピングデータを記憶するためのマッピングデータ記憶手段と、

前記音声信号及び前記テキストを受け、前記テキストに基づいて、前記音声信号に対する音素セグメンテーションを行ない、得られる音素列と、各音素の継続時間長を表す時間情報とを含む音素列データを出力するための音素セグメンテーション手段と、

前記音素セグメンテーション手段より出力される前記音素列データに含まれる各音素に

10

20

対し、当該音素の前記時間情報と、前記マッピングデータとを参照することにより、当該音素がマッピングされる画像を特定する識別子と、当該音素に対する所定の特徴量に対応して定められるブレンド率とを付すことによりキーフレームデータを作成して出力するための手段と、

前記音素セグメンテーション手段より出力される前記音素列データに含まれる各音素に対し、前記マッピングデータを参照して得られる画像の識別子と、所定の定数からなるブレンド率とを付し、画像マッピング済の音素列データを出力するためのマッピング処理手段と、

前記マッピング処理手段の出力する前記画像マッピング済の音素列データの各音素に対し、当該音素の継続長の単調増加関数として、前記ブレンド率を調整するための第1のブレンド率調整手段とを含む、アニメーション作成装置。

10

【請求項2】

前記キーフレームデータ作成手段はさらに、前記第1のブレンド率調整手段の出力する、ブレンド率が調整された音素列データの各音素に対し、当該音素の継続期間内のパワーの大きさの単調増加関数として、前記ブレンド率を調整するための第2のブレンド率調整手段を含む、請求項1に記載のアニメーション作成装置。

【請求項3】

前記アニメーション生成手段は、

アニメーションの画像を生成するための生成時刻を、前記音声の録音時間と関係付けて決定するための時刻決定手段と、

20

前記時刻決定手段により決定された前記生成時刻におけるフレームの画像を、当該生成時刻をはさむ複数のキーフレームの画像の間の補間により算出するための補間手段とを含む、請求項1又は請求項2に記載のアニメーション作成装置。

【請求項4】

前記補間手段は、前記時刻決定手段により決定された前記生成時刻におけるフレームの画像を、当該生成時刻をはさんで互いに隣接する二つのキーフレームの画像の間の補間により算出するための手段を含む、請求項3に記載のアニメーション作成装置。

【請求項5】

前記算出するための手段は、

前記生成時刻をはさんで互いに隣接する第1及び第2のキーフレームのうち、第1のキーフレームにおいて100%、第2のキーフレームにおいて0%となる第1の補間関数により、前記生成時刻における第1のブレンド率を前記第1のキーフレームにおけるブレンド率から補間するための第1のブレンド率補間手段と、

30

前記第1のキーフレームにおいて0%、前記第2のキーフレームにおいて100%となる第2の補間関数により、前記生成時刻における第2のブレンド率を前記第2のキーフレームにおけるブレンド率から補間するための第2のブレンド率補間手段と、

前記第1のブレンド率及び前記第2のブレンド率を用いた、前記第1のキーフレームにマッピングされた画像のデータ及び前記第2のキーフレームにマッピングされた画像のデータの間の加重和により、前記生成時刻における画像のデータを算出するための手段とを含む、請求項4に記載のアニメーション作成装置。

40

【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アニメーション作成装置を構成する各手段として機能させる、コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は音声からアニメーションを作成するアニメーション作成装置に関し、特に、発話音声にあわせて口等の形が変わる顔画像等のアニメーションを自動的に生成する装置に関する。

50

【背景技術】

【0002】

コンピュータ技術の発達により、以前は大部分が手作業で行なわれていた仕事がコンピュータによる作業に置き換えられるケースが多くなっている。その代表的なものに、アニメーションの作成がある。

【0003】

以前は、アニメーションといえば次のような手法で作成されることが一般的であった。登場するキャラクタをアニメーションの演出家が決め、絵コンテと呼ばれる、主要なシーンのラフな原画を作成する。これら絵コンテに基づき、アニメーションの各フレームの絵をアニメータと呼ばれる作業者が作成する。それら絵を仕上げ担当者がセル画に仕上げる。セル画を順にフィルムに写し、所定のフレームレートで再生すればアニメーションの画像の部分が出来上がる。

10

【0004】

このアニメーションの画像を再生しながら、声優がアニメーションの台本に基づいて台詞をつけていく。いわゆる「アフレコ」である。

【0005】

このような作業で最も人手がかかるのはセル画の作成である。一方、原画をCG（コンピュータ・グラフィックス）で作成する場合、原画を加工してセル画を作成するのは比較的単純な作業である。一枚一枚撮影する必要もない。そのため、この部分については原画のCG化とあわせてかなりコンピュータ化されている。

20

【0006】

一方、残りの作業のうちで比較的むずかしいのは、アフレコの作業である。アニメーションの動きにあわせて、なおかつ状況にあわせた声で台詞をしゃべる必要があるため、アフレコの作業にはそれなりの時間がかかり、習熟も必要である。

【0007】

そこで、アフレコの逆に、先に音声を収録し、その音声にあわせてアニメーションを作成する手法が考えられた。これは「プレスコ」又は「プレレコ」（以下「プレスコ等」と呼ぶ。）と呼ばれる。これはもともと米国等で手作業でアニメーションを作成する際に採用されていた手法である。この手法でアニメーションを作成する場合には、次のような作業手順となる。

30

【0008】

まず、アニメーションに登場するキャラクタを決める。絵コンテも従来と同様に作成する。声優が、絵コンテと台本に基づいて発話し、それを音声として収録する。この音声にあわせて、アニメーションを作成する。

【0009】

このプレスコ等の手法によるアニメーション作成をコンピュータで実現する場合には、音声からアニメーションをいかにして自動的に作成するか、という点が問題となる。特に、人物等のアニメーションの口の動きを、予め録音した声優の音声にあわせて自然な形で生成するのは難しく、これを自動的に行なう手法が求められている。

40

【0010】

このための一手法として提案されたものに、特許文献1に記載された手法がある。特許文献1に記載された手法では、口形状の基本パターンを予め複数個用意しておく。そして、任意の音声に対応する口形状を、これら基本パターンの加重和により求める。そのために、声優の音声の所定の特徴量から、各基本パターンの加重パラメータに変換するための変換関数を、重回帰分析によって予め求めておく。台本に沿って録音された声優の音声の所定の特徴量をこの変換関数で加重パラメータに変換し、その加重パラメータを用いて口形状の基本パターンの加重和を算出することで、声優の音声に対応する口形状及び顔画像を作成する。こうした処理をアニメーションの各フレームに相当する時刻に行なうことで、アニメーションのフレームシーケンスを作成する。

【特許文献1】特開平7-44727号公報

50

【発明の開示】

【発明が解決しようとする課題】

【0011】

現代では、例えば遠隔会議とか、テレビ電話等、動画像を伴う通信量が増大している。そのため、いかにして動画像のデータ量を削減するかが問題となっている。そのための一つの方策は、通信では音声のみを送信するが、受信側ではその音声から顔画像を合成する、というものである。こうした技術を一般化させるためには、不特定多数の人間の音声であっても、それらに対応する口画像を適切に生成する技術が必要である。

【0012】

また、上記したアニメーションの作成を用いるサービスとして、例えば、不特定多数の話者の音声にあわせ、特定のキャラクタの顔画像を用いたアニメーションを作成するようなサービスが考えられる。そうしたサービスでは、不特定の話者の音声から適切に顔画像の口の動きを生成する必要がある。

10

【0013】

しかし、上記した特許文献1に開示の技術では、予め変換関数を求める必要がある。そのため、特定の話者に対しては有効であっても、不特定多数の話者に対しては適用できない。なぜなら、話者により、発声する音素が同一でもその音声から得られる音響特徴量は様々だからである。

【0014】

それ故に本発明の目的は、話者に依存せず、音声に応じて適切に一部の形状を変化させる動画像を生成できるアニメーション作成装置を提供することである。

20

【課題を解決するための手段】

【0015】

本発明の第1の局面に係るアニメーション作成装置は、音声信号を受け、当該音声信号の表す音素列内の各音素の継続時間中の、所定のキーフレーム時刻における画像により構成されるキーフレーム画像を表すキーフレームデータを作成するための手段と、キーフレームデータ作成手段により作成されるキーフレームデータに基づき、音声信号に同期して変化する一連の画像からなる画像のアニメーションを生成するためのアニメーション生成手段とを含む。

【0016】

キーフレームデータ作成手段は、音声信号を受け、当該音声信号の表す音素列内の各音素の継続時間中の、所定のキーフレーム時刻における画像により構成されるキーフレーム画像を表すキーフレームデータを作成する。アニメーション生成手段は、キーフレームデータ作成手段により作成されるキーフレームデータに基づき、音声信号に同期して変化する一連の画像からなる画像のアニメーションを生成する。音声信号からキーフレームデータを作成し、そのキーフレームデータから画像のアニメーションが生成される。キーフレームデータは、音声信号の発話者に依存せずに定められる。従って、話者に依存せず、音声から適切に動画像を生成できるアニメーション作成装置を提供することができる。

30

【0017】

好ましくは、所定のキーフレーム時刻は、音素列内の各音素の継続時間の開始時刻である。

40

【0018】

ある音素を発音するときの口の形状の特徴は、その音素を発音する最初の時期のときに最もよく現われている。従って、所定のキーフレーム時刻を、音素の継続時間の開始時刻とすることにより、得られる動画像は、音声の変化をよく反映した、適切なものとなる。

【0019】

より好ましくは、アニメーション作成装置は、予め定められた複数種類のテキストをユーザに選択させるためのテキスト選択手段と、テキスト選択手段によりテキストが選択されたことに基づき、ユーザの音声を録音して音声信号に変換し、選択されたテキストとともにキーフレームデータ作成手段に与えるための手段とをさらに含む。キーフレームデー

50

タを作成するための手段は、音素を、所定の基準画像を含む所定の複数個の画像のいずれかにマッピングするマッピングデータを記憶するためのマッピングデータ記憶手段と、音声信号及びテキストを受け、テキストに基づいて、音声信号に対する音素セグメンテーションを行ない、得られる音素列と、各音素の継続時間長を表す時間情報とを含む音素列データを出力するための音素セグメンテーション手段と、音素セグメンテーション手段より出力される音素列データに含まれる各音素に対し、当該音素の時間情報と、マッピングデータとを参照することにより、当該音素がマッピングされる画像を特定する識別子と、当該音素に対する所定の特徴量に対応して定められるブレンド率とを付すことによりキーフレームデータを作成して出力するためのキーフレームデータ作成手段とを含む。

【 0 0 2 0 】

10

音素セグメンテーションは、テキスト選択手段により選択されたテキストに基づいて行なわれる。音声信号を構成する各音素が予め判明しているため、音素セグメンテーションを正しく行なうことができる。

【 0 0 2 1 】

さらに好ましくは、キーフレームデータ作成手段は、音素セグメンテーション手段より出力される音素列データに含まれる各音素に対し、マッピングデータを参照して得られる画像の識別子と、所定の定数からなるブレンド率とを付し、画像マッピング済の音素列データを出力するためのマッピング処理手段と、マッピング処理手段の出力する画像マッピング済の音素列データの各音素に対し、当該音素の継続長の単調増加関数として、ブレンド率を調整するための第1のブレンド率調整手段とを含む。

20

【 0 0 2 2 】

音素の継続長は、音素を発音するときの口等の形状の変化の割合を反映している。従って、ブレンド率を音素の継続時間長に対する単調増加関数として調整することにより、口等の形状の実際の変化を適切に反映したアニメーションを得ることができる。

【 0 0 2 3 】

キーフレームデータ作成手段はさらに、第1のブレンド率調整手段の出力する、ブレンド率が調整された音素列データの各音素に対し、当該音素の継続期間内のパワーの大きさの単調増加関数として、ブレンド率を調整するための第2のブレンド率調整手段を含んでもよい。

【 0 0 2 4 】

30

音素の継続期間中のパワーは、音素を発音するときの強さ、従ってそのときの口等の形状の変化の割合を反映している。従って、ブレンド率を音素の継続期間中におけるパワーに対する単調増加関数として調整することにより、口等の形状の実際の変化を適切に反映したアニメーションを得ることができる。

【 0 0 2 5 】

好ましくは、アニメーション生成手段は、アニメーションの画像を生成するための生成時刻を、音声の録音時間と関係付けて決定するための時刻決定手段と、時刻決定手段により決定された生成時刻におけるフレームの画像を、当該生成時刻をはさむ複数のキーフレームの画像の間の補間により算出するための補間手段とを含む。

【 0 0 2 6 】

40

補間手段が、ある生成時刻におけるフレームの画像を、その時刻を含む複数のキーフレームの画像の間の補間により生成する。ある時刻における口等の形状は、その前の音素から次の音素への遷移の途中の形状となる。このように補間によりある生成時刻の口等の形状を算出することにより、音素の遷移に対応した適切な画像のアニメーションを作成できる。

【 0 0 2 7 】

より好ましくは、補間手段は、時刻決定手段により決定された生成時刻におけるフレームの画像を、当該生成時刻をはさんで互いに隣接する二つのキーフレームの画像の間の補間により算出するための手段を含む。

【 0 0 2 8 】

50

補間を、生成時刻をはさんで隣接する二つのキーフレームの間で行なって、生成時刻におけるフレームの画像を生成する。計算量を少なくしながら、適切な補間ができ、滑らかに変化するアニメーションを得ることができる。

【0029】

さらに好ましくは、算出するための手段は、生成時刻をはさんで互いに隣接する第1及び第2のキーフレームのうち、第1のキーフレームにおいて100%、第2のキーフレームにおいて0%となる第1の補関数により、生成時刻における第1のブレンド率を第1のキーフレームにおけるブレンド率から補間するための第1の補間手段と、第1のキーフレームにおいて0%、第2のキーフレームにおいて100%となる第2の補関数により、生成時刻における第2のブレンド率を第2のキーフレームにおけるブレンド率から補間するための第2の補間手段と、第1のブレンド率及び第2のブレンド率を用いた、第1のキーフレームにマッピングされた画像のデータ及び第2のキーフレームにマッピングされた画像のデータの間の加重和により、生成時刻における画像のデータを算出するための手段とを含む。

10

【0030】

第1のキーフレームにおけるブレンド率と、第2のキーフレームにおけるブレンド率とを第1及び第2の補関数により別個に算出し、次に、これらを用い、第1のキーフレームにマッピングされた画像のデータ及び第2のキーフレームにマッピングされた画像のデータの間の加重和を算出する。単純な計算を組み合わせることにより、二つのキーフレームの間の画像の滑らかなアニメーションを算出することができる。

20

【0031】

時刻決定手段は、補間手段によりあるフレームの画像が得られた時刻を、次のフレームの画像を生成するための生成時刻として決定するための手段を含んでもよい。

【0032】

補間手段による画像の生成が終了すると、その時刻が生成時刻として決定される。生成時刻が決定されると、その生成時刻におけるフレームの画像が、当該生成時刻をはさむ複数のキーフレームの画像の間の補間により補間手段により算出される。従って補間手段は休むことなく常に画像の生成のために動作していることになり、補間手段を有効に利用することができる。

30

【0033】

なお、画像は発話時の口の形状の変化を反映した顔画像でもよい。

【0034】

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、上記したいずれかのアニメーション作成装置を構成する各手段として機能させる。

【発明を実施するための最良の形態】

【0035】

以下、本発明について、実施の形態に基づいて説明する。以下の説明では、基本となる顔画像を6種類使用しているが、顔画像の数はこれには限定されない。6種類よりも少なくてもよいし、6種類よりも多くてもよい。

40

【0036】

[第1の実施の形態]

<構成>

図1に、本発明の第1の実施の形態に係るアニメーション作成装置によるアニメーション作成過程30の概略を示す。図1を参照して、アニメーション作成過程30においては、話者40が台本44に基づき台詞を発話すると、その音声信号42に対し、音声認識装置による音素セグメンテーション(発話から、発話を構成する音素列を生成すること)が行なわれる。

【0037】

予め、主要な音素については、その音素を発音するときの口の形状を含む顔画像60～

50

68が準備されており、音声認識の結果得られる各音素50～58に対し、これら顔画像を割当ててアニメーション化する。

【0038】

なお、個々の音素に対して発話画像を一つずつ割当てても滑らかな画像が得られないため、本実施の形態では、後述するように、主要な顔画像として「あ（/a/）」「い（/i/）」「う（/u/）」「え（/e/）」「お（/o/）」という5つの音素に対する5つの顔画像、及び無表情の顔画像の、合計6つの顔画像を準備する。「あ」～「お」の5つの音素はそれぞれ対応の顔画像に割当て、残りの音素についてはそれぞれ上記した6つの顔画像のいずれかに割当てる。これを以下、音素から顔画像へのマッピングと呼ぶ。

【0039】

さらに、音素ごとに、このようにマッピングされた顔画像を割当ててそれらを単純につないでアニメーションを作成すると、画像の動きが過大になって、いわゆる「うるさい」アニメーションとなる。そのため、本実施の形態では、音素の継続時間長及びそのパワーによって、各画像の「強さ」を調整し、調整後の画像を用い、音素間の遷移過程での顔画像を補間により生成する。また、所定のしきい値より小さな継続時間長又はパワーしか持たない音素については、あえてその音素に対応する画像を挿入せず、その直前の音素の画像に統合してしまう。こうすることで、滑らかに変化する、自然なアニメーションを音声にあわせて生成することができる。

【0040】

図2に、本実施の形態に係るアニメーション生成システム80の概略の機能的構成を示す。このアニメーション生成システム80は、予め複数の書起しテキストを準備しておき、それらのいずれかを話者に選択させて発話させ、その発話音声に合致して変化する顔画像のアニメーションを、予め準備した6つの顔画像から補間により生成するものである。

【0041】

図2を参照して、アニメーション生成システム80は、発話者が書起しテキストを選択する際に使用するテキスト選択インターフェイス90と、発話者の音声を音声信号に変換するマイクロフォン92と、予め複数種類のテキストを記憶しておき、話者にそのうちの一つをテキスト選択インターフェイス90を用いて選択させた上で、マイクロフォン92の出力する音声信号を録音しデジタル化した音声データファイルを作成するための入力指示ユニット94と、入力指示ユニット94から与えられる音声データファイルに対する音素セグメンテーションを、入力指示ユニット94から与えられる対応する書起しテキストを用いて行ない、その結果と、入力指示ユニット94からの音声データファイルとに基づき、アニメーションのキーフレームを規定するキーフレームデータを作成するためのキーフレームデータ作成ユニット96とを含む。

【0042】

アニメーション生成システム80はさらに、入力指示ユニット94の出力する音声データファイルと、キーフレームデータ作成ユニット96により出力されるキーフレームデータとを用い、音声データファイルの音声に同期して口形状が変化する顔画像のアニメーションを作成し、音声とともに出力するためのアニメーション再生ユニット98と、いずれもアニメーション再生ユニット98に接続された、アニメーションを表示するためのモニタ102及び音声を再生するためのスピーカ100とを含む。

【0043】

入力指示ユニット94は、予め複数種類の書起しテキストを記憶しておくためのテキスト記憶部110と、テキスト記憶部110に記憶されたテキストをテキスト選択インターフェイス90により話者40に提示し、いずれか一つを選択させてそのテキストをキーフレームデータ作成ユニット96に対し与えるとともに、テキスト選択インターフェイス90を用いて、話者に対し当該テキストを発話するように指示を与えるためのテキスト選択部112と、話者がテキスト選択部112の指示に対して発話するテキストの音声についてマイクロフォン92から出力される音声信号を、所定のフレーム長及びフレームシフト長でフレーム化し音声データとして保存し、キーフレームデータ作成ユニット96及びア

10

20

30

40

50

ニメーション再生ユニット 98 に与えるための音声収録部 114 とを含む。

【0044】

キーフレームデータ作成ユニット 96 は、テキスト選択部 112 から与えられるテキストに基づいて音声収録部 114 からの音声データに対する音素セグメンテーションを行い、音素列と、その継続時間長とを含む音素列データを出力するための音声認識装置 120 と、日本語を構成する全ての音素を、前述した 6 つの顔画像の識別子にマッピングするマッピングテーブルを記憶したマッピングテーブル記憶部 130 と、音声認識装置 120 から出力される音素列ファイル、テキスト選択部 112 から与えられるテキスト、及び音声収録部 114 から与えられる音声データに基づき、アニメーションのうち主要時点でのフレームの顔画像を、前述した 6 つの顔画像から作成するためのパラメータを生成してキーフレームデータとして出力するためのキーフレームデータ作成部 136 とを含む。

10

【0045】

音声認識装置 120 は、音素セグメンテーションをし、音素列と、それぞれの継続時間長が分かる時間データとを出力できるものであればどのようなものでもよい。発話内容が予め分かっているので、音声認識装置 120 は音声データを確実に音素列に変換できる。

【0046】

図 5 に、音声認識装置 120 の出力する音素列ファイル 160 の構成例を示す。図 5 を参照して、音素列ファイル 160 は、音声認識の結果得られた音素列と、各音素列の継続時間長が分かる時間情報との組を複数個含んでいる。図 5 において、継続時間長はミリ秒単位で示してある。

20

【0047】

アニメーション再生ユニット 98 は、5 つの音素 (/ a / , / i / , / u / , / e / , / o /) に対応する顔画像と、無表情の顔画像との 6 つの顔画像を、ワイヤフレームモデルとして保持する顔データファイルを記憶した顔データファイル記憶部 132 と、キーフレームデータ作成部 136 によって作成されたキーフレームにおける顔画像を作成するためのパラメータを用い、アニメーションを構成する所定時点のフレームの顔画像を顔データファイル記憶部 132 に記憶された 6 つの顔画像から作成するために使用される補間関数を記憶するための補間関数記憶部 134 と、キーフレームデータ作成部 136 から与えられるキーフレームデータと、顔データファイル記憶部 132 に記憶された顔データファイルと、補間関数記憶部 134 に記憶された補間関数とを用い、アニメーションでの所定の時点でのフレームの顔画像を補間により生成するためのアニメーション生成部 138 とを含む。

30

【0048】

顔データファイル記憶部 132 に記憶される顔画像の例を図 3 に示す。図 3 (A) ~ (E) は、それぞれ音素 / a / , / i / , / u / , / e / , / o / に対応する顔画像であり、図 3 (F) は、無表情に対応する顔画像である。本明細書では、これら画像をそれぞれ顔画像 / A / , / I / , / U / , / E / , / O / , 及び / / と表記することにする。

【0049】

なお、本実施の形態では、顔画像 / A / , / I / , / U / , / E / , / O / は、いずれも顔画像 / / を基準とし、各特徴点が、顔画像の定義されている 3 次元空間において、顔画像 / / の対応する特徴点からどの程度移動しているかを示す 3 次元ベクトル情報によって定義されている。従って、例えば顔画像 / A / と顔画像 / / との間で、その中間の顔画像を定義することもできる。本実施の形態では、特定の顔画像と顔画像 / / との間で中間の顔画像を定義するために、「ブレンド率」という概念を導入する。ブレンド率とは、特定の顔画像を 100 %、顔画像 / / を 0 % として、顔画像 / / から特定の顔画像に至るまでの特徴点の移動量の割合で中間の顔画像を表すものである。従って、顔画像 / A / , / I / , / U / , / E / , / O / をそのまま音素に割当てた場合、そのブレンド率はいずれも 100 % となる。ブレンド率 50 % の顔画像 / A / とは、顔画像 / / からの特徴点の移動量の割合が、顔画像 / A / の特徴点の移動量の 50 % となっているような顔画像のことをいう。顔画像 / / での位置を始点とするベクトルで顔画像の特徴点の

40

50

移動量を表せば、ブレンド率 B % の顔画像とは、各特徴点を表すベクトルが、方向はブレンド率 100 % の顔画像のベクトルと等しく、長さがブレンド率 B % に相当するだけ縮小されたものとなっている顔画像に相当する。

【0050】

二つの顔画像の間の補間については後述する。

【0051】

図4に、マッピングテーブル記憶部130に記憶されたマッピングテーブルの例の一部を示す。図4を参照して、本実施の形態では、マッピングテーブル記憶部130は、音素 / a / を顔画像 / A / に、音素 / b / を顔画像 / / に、音素 / d / を顔画像 / U / に、音素 / e / を顔画像 / E / に、それぞれ対応付けている。マッピングテーブルでは、図3に示す顔画像 / A / , / I / , / U / , / E / , / O / のように、予めある音素に対して準備された顔画像には、その音素を必ず対応付けるようにする。さもないと得られる顔の動画が発話内容とちぐはぐになってしまう。また / b / , / m / 等、唇を閉じるような音素は無表情の顔画像 / / に対応付ける。それ以外の音素は、前述した6つの顔画像のうち、口の形状が最も近いと思われるものに適宜割当てるようにする。

【0052】

再び図2を参照して、アニメーション再生ユニット98はさらに、音声収録部114の出力する音声データを格納した音声ファイルを記憶するための音声ファイル記憶部140と、アニメーション生成部138から順次与えられる顔画像と、音声ファイル記憶部140に記憶された音声ファイルからの音声とを、互いに同期させてモニタ102及びスピーカ100にそれぞれ与えるための出力部142と、入力指示ユニット94のテキスト選択部112及び音声収録部114、キーフレームデータ作成ユニット96のキーフレームデータ作成部136及び音声認識装置120、並びにアニメーション再生ユニット98のアニメーション生成部138及び出力部142を所定のシーケンスで動作させ、それらの協働によってアニメーション生成システムを実現するようこれらを制御するためのシーケンス制御部144とを含む。

【0053】

図6に、図2のキーフレームデータ作成部136の構成の詳細を示す。図6を参照して、キーフレームデータ作成部136は、音声認識装置120からの音素列データ内の各音素に対し、マッピングテーブル記憶部130を参照して顔画像をマッピングし、マッピングされた顔画像の識別子と、ブレンド率「100%」とを付して出力するためのマッピング処理部180と、マッピング処理部180により出力された、継続時間長、対応顔画像の識別子及びそのブレンド率が付された音素列の各ブレンド率を、各音素の継続時間長に基づいて調整するための、継続時間長によるブレンド率調整部182と、継続時間長によるブレンド率調整部182の出力する、継続時間長、対応顔画像及びその継続時間長により調整されたブレンド率が付された音素列のブレンド率を、各音素の継続期間におけるパワーの大きさに基づいて調整するための、パワーによるブレンド率調整部184とを含む。パワーによるブレンド率調整部184の出力は、各音素に、その継続時間長と、対応顔画像と、継続時間長及びパワーにより調整されたブレンド率とが付された音素列となる。この音素列がキーフレームデータである。なお、本実施の形態では、キーフレームとは、各音素の継続期間の先頭時刻でフレームが作成される場合のそのフレームのことをいう。

【0054】

図7に、アニメーション生成部138のより詳細なブロック図を示す。図7を参照して、アニメーション生成部138は、二つのキーフレームにおける、それぞれ所定のブレンド率が割当てられた顔画像と、それら二つのキーフレームに対応する時刻と、その二つのキーフレームの間で、アニメーションを生成すべき時刻（ここでは便宜のため、「生成時刻」と呼ぶ。生成時刻は、二つのキーフレームの時刻を基準とする相対時刻で表される。）とが与えられると、その生成時刻の顔画像を、二つのキーフレームの顔画像から補間関数記憶部134に記憶された補間関数を用いた補間処理により生成して出力部142に対して出力するための補間処理部204と、所定の生成時刻が決まると、その生成時刻をは

10

20

30

40

50

さむ二つのキーフレームを定め、それらのキーフレームにおける顔画像のデータ及びブレンド率、ならびにそれら二つのキーフレームの時刻の間における生成時刻の相対的位置を示す情報を補間処理部 204 に与え、生成時刻における顔画像を作成させるとともに、補間処理部 204 による顔画像の生成が終わると、そのときの時刻を次の生成時刻として次の顔画像を作成する処理を繰り返す機能を持つアニメーション生成制御部 200 と、アニメーション生成制御部 200 が時刻を定めるために参照するタイマ 202 とを含む。この補間処理とアニメーションの生成処理とについては後述する。

【0055】

図 8 から図 14 を参照して、本実施の形態に係るキーフレームデータ作成部 136 及びアニメーション生成部 138 による顔のアニメーションの作成処理についてより詳細に説明する。

10

【0056】

例えば図 5 に示すような音素列ファイル 160 が与えられたとする。この場合、図 6 に示すマッピング処理部 180 の出力を図示すると図 8 のようになる。図 8 を参照して、時間軸上で、各音素 / a / , / i / , / u / , / e / , / o / の発話期間がそれぞれ継続時間 300、300、100、30 及び 250 (いずれもミリ秒) で割当てられる。各期間の先頭時刻がキーフレームとなる。各キーフレームでのブレンド率は、いずれも 100% である。

【0057】

図 8 に示されるような音素列が継続時間長によるブレンド率調整部 182 により処理される途中の結果の一例を図 9 に示す。図 9 を参照して、まず、各音素のうちで、所定のしきい値よりも小さな継続時間長しか持たない音素については、その直前の音素の期間に統合してしまう。図 8 に示す例では音素 / e / の継続時間長が 30 ミリ秒であり、しきい値が 50 ミリ秒であったものとする。音素 / e / は削除され、その継続時間長はその直前の音素 / u / に統合される。従って図 9 に示されるように、音素列は / a / , / i / , / u / , / o / となり、その継続時間長はそれぞれ 300、300、130、及び 250 (ミリ秒) となる。音素が一つ削除されるので、キーフレームの数も 5 つから 4 つに減少する。また、以下の説明では、これらのキーフレームに対応する時刻をそれぞれ T_0 , T_1 , T_2 及び T_3 とし、最後の音素 / o / の直後のキーフレームの時刻を T_4 とする。なお、以下、一般的に、 k 番目の音素の開始時刻 T_k により規定されるキーフレームを「キーフレーム T_k 」と呼ぶ。

20

30

【0058】

継続時間長によるブレンド率調整部 182 は、さらに、各キーフレームに割当てられたブレンド率を、その継続時間長に応じて調整する。具体的には、継続時間長によるブレンド率調整部 182 は、音素列の中の継続時間長の最大値 L_{MAX} を探し出し、各音素の継続時間長のブレンド率を次の式 (1) により調整する。

【0059】

【数 1】

$$B(n) = B(n) \times \frac{L(n) - C_1}{L_{MAX} - C_1} \quad (1)$$

40

ただし、 $B(n)$ は n 番目の音素のブレンド率を、 $L(n)$ は n 番目の音素の継続時間長を、それぞれ表す。 C_1 は所定の定数で、例えば短い継続時間長の音素を削除したときに使用されるしきい値と同程度の大きさに選ばれる。こうしてブレンド率が継続時間長により調整された音素列の、時間軸上での配置とそのブレンド率とを図 10 に模式的に示す。図 10 において、例えば、調整後の音素 / a / , / i / , / u / , / o / のブレンド率はそれぞれ a 、 a 、 b 、及び c (%) である。

【0060】

パワーによるブレンド率調整部 184 は、継続時間長によるブレンド率調整部 182 と同様にして、各音素のブレンド率を、各音素の継続期間における音声のパワーによって調

50

整する。具体的には、まず、所定のしきい値以下のパワーしかない音素については削除し、その継続期間を直前の音素の継続期間に統合する。こうして得られた各継続期間の先頭がキーフレームである。各キーフレームには、ブレンド率が割当てられている。パワーによるブレンド率調整部184はこのブレンド率（ n 番目の音素のブレンド率を前と同様 $B(n)$ とする。）を以下の式（2）により調整する。

【0061】

【数2】

$$B(n) = B(n) \times \frac{P(n) - C_2}{P_{MAX} - C_2} \quad (2)$$

10

ただし P_{MAX} は全体でのパワーの最大値であり、 $P(n)$ は n 番目の音素の継続期間のパワーであり、 C_2 は所定のしきい値である。このしきい値も、前述した音素の削除のときに使用されたしきい値と同程度の大きさに選ばれる。

【0062】

こうして最終的に得られた音素列と、その継続時間長と、各キーフレームにおける調整後のブレンド率とを模式的に示したものが図11である。図11を参照して、キーフレーム T_0, T_1, T_2, T_3 の音素はそれぞれ / a /, / i /, / u /, / o / であり、ブレンド率はそれぞれ a', a'', b' 及び c' (ただし $a' a, a'' a, b' b$ 及び $c' c$) であり、継続時間長はそれぞれ 300、300、130 及び 250 (ミリ秒) である。

20

【0063】

このようにしてキーフレームデータが作成される。

【0064】

次に、図7に示す補間処理部204による補間処理について説明する。補関数記憶部134に記憶される補関数としては、様々なものが考えられるが、本実施の形態では計算処理の容易さと高速さとに重点をおき、線形補間を与えるものを使用する。線形補間の概念について図12を参照して説明する。

【0065】

図12を参照して、時間軸を横軸、各キーフレームの時間 $T_0, T_1, T_2, T_3 \dots$ におけるブレンド率を縦軸のグラフで表すものとする。本実施の形態での補関数は、線分 220, 222, 224 及び 226 で表されるように、各キーフレームでのブレンド率と、隣接するキーフレームの時刻でのブレンド率「0」の点とを結んだ線分に沿って、各時間でのブレンド率を線形補間する関数である。すなわち、一方における率が 100%、他方における率が 0% となるように線形補間を行なう関数である。

30

【0066】

例えば、時刻 T_0 と時刻 T_1 との中間の時刻 t が生成時刻であるものとする。キーフレーム T_0 及び T_1 での音素はそれぞれ / a /、/ i / である。各キーフレームでのブレンド率はパワーによるブレンド率調整部184により算出されている。時刻 T_0 でのブレンド率 a' と、時刻 T_1 でのブレンド率「0」の点とを結んだ線分 220 によって、時刻 t におけるキーフレーム T_0 のブレンド率 が線形補間される。同様に、時刻 T_0 でのブレンド率「0」の点と、時刻 T_1 でのブレンド率 a'' の点とを結んだ線分 222 に沿って、時刻 t におけるキーフレーム T_1 のブレンド率 が算出される。

40

【0067】

時刻 T_0 におけるキーフレーム T_0 のブレンド率を $B(T_0)$ 、時刻 T_1 におけるキーフレーム T_1 のブレンド率を $B(T_1)$ 、補間により求められた、時刻 t におけるキーフレーム T_0, T_1 のブレンド率をそれぞれ B_0 及び B_1 とする。すると、 B_0 及び B_1 は次の式(3)により求められる。

【0068】

【数 3】

$$\alpha = \frac{T_1 - t}{T_1 - T_0} B(T_0)$$

$$\beta = \frac{t - T_0}{T_1 - T_0} B(T_1)$$
(3)

本実施の形態では、このようにして算出された二つのブレンド率（例えば α 及び β ）を用い、図 13 に示されるようにして時刻 t における顔画像を作成する。

【0069】

今、キーフレーム T_0 での顔画像の各特徴点の、顔画像 A / I / での対応特徴点を基準とした移動量を要素とする 3 次元ベクトルを $X(T_0)$ 、同様にキーフレーム T_1 での 3 次元ベクトルを $X(T_1)$ とする。すると、 T_0 から T_1 における顔画像の各特徴点の、顔画像 A / I / での対応特徴点を基準とした移動量を要素とする 3 次元ベクトル $X(t)$ は、次の式 (4) で表されるベクトル加重和で算出される。

【0070】

【数 4】

$$X(t) = \alpha X(T_0) + \beta X(T_1)$$
(4)

補間処理部 204 は、こうした計算を顔画像の各特徴点に対して実行する。後述するようにこうした演算はグラフィックプロセッサユニット (GPU) が得意とするところである。従ってアニメーション生成システム 80 は、GPU を備えていることが望ましい。

【0071】

アニメーション生成制御部 200 によるアニメーションの生成制御について説明する。図 14 を参照して、アニメーション生成制御部 200 が、最初のキーフレームの時刻 T_0 に等しい時刻 t からアニメーションの作成を開始するものとする。すなわち、時刻 t_0 ($= T_0$) においてアニメーション生成制御部 200 は、補間処理部 204 に対して顔画像の生成の指示 240 を与える。すなわち、生成時刻 $t = t_0$ である。

【0072】

この場合には、まず $T_{k-1} < t < T_k$ となるような整数 k を探す。ここでは $t = t_0 = T_0$ であるから、 $k = 1$ となる。補間処理部 204 は、時刻 T_0 における音素 a / の顔画像 A / を構成する各特徴点の 3 次元ベクトル $X(T_0)$ に、このときのブレンド率 α' (図 11 参照) を乗算する。さらに、時刻 T_1 における音素 i / の顔画像 I / を構成する各特徴点の 3 次元ベクトル $X(T_1)$ に、このときのブレンド率 α'' (図 11 参照) を乗算する。補間処理部 204 は次に式 (3) を用いて α 、 β を算出する。ここでは $t = T_0$ なので、 $\alpha = B(T_0)$ 、 $\beta = 0$ である。これらの結果を用い、式 (4) によって時刻 t における顔画像 242 を生成し出力する。

【0073】

この生成処理に時間 s_1 を要したものとする。顔画像 242 を生成し出力すると補間処理部 204 は、アニメーション生成制御部 200 に対して処理を終了したことを示す完了通知 244 を与える。この時刻 t_2 を新たな生成時刻 t とする。

【0074】

アニメーション生成制御部 200 は、新たな生成時刻 $t = t_2$ において完了通知 244 を受けたことに応答し、この生成時刻 t をはさむ二つのキーフレーム、図 14 に示す例では時刻 T_0 及び T_1 におけるキーフレームを特定し、これら二つのキーフレームにおける顔画像 A / 及び I / と、二つの時刻 T_0 及び T_1 と、時間 $T_0 \sim T_1$ の中における時刻 $t = t_2$ の相対時間とを補間処理部 204 に与え、顔画像の生成の指示 246 を与える。

【0075】

補間処理部 204 は、この指示に応答し、時間 s_2 をかけて生成時刻 $t = t_2$ における顔画像 248 を生成し、出力する。このとき補間処理部 204 は、アニメーション生成制

10

20

30

40

50

御部 200 に対して完了通知 250 を与える。このときの時刻 t_3 を新たな生成時刻 t とする。

【0076】

すると、アニメーション生成制御部 200 は、この新たな生成時刻 $t = t_3$ に対し、直前の生成時刻 t_2 で行なったものと同様の処理を行ない、顔画像の生成指示 252 を補間処理部 204 に対し与える。以下同様に、時間 s_3 後に顔画像 254 が出力され、完了通知 256 が時刻 t_3 でアニメーション生成制御部 200 に与えられ、これに応答してアニメーション生成制御部 200 から時刻 t_3 における顔画像生成の指示 258 が補間処理部 204 に与えられる。以下同様である。

【0077】

すなわち本実施の形態では、補間処理部 204 が常にその能力をフルに発揮するように、アニメーション生成制御部 200 がアニメーション生成のためのタイミングを制御する。

【0078】

<動作>

図 2 ~ 図 14 を参照して、上記したアニメーション生成システム 80 は以下のように動作する。なお、以下の各部の動作は、図 2 に示すシーケンス制御部 144 による制御によって所定のシーケンスで行なわれるが、説明を分かりやすくするため、以下ではシーケンス制御部 144 の制御については言及しない。

【0079】

予め、アニメーションのキャラクタの顔画像を、上記した 6 種類の音素について準備し、顔データファイル記憶部 132 に記憶させておく。各音素に対して顔画像をマッピングするマッピングテーブルも予め準備し、マッピングテーブル記憶部 130 に記憶させておく。補間関数を実現するプログラムも予め準備し、補間関数記憶部 134 に記憶させておく。さらに、ユーザ発話のための書起しテキストも予め何種類か準備し、テキスト記憶部 110 に記憶させておく。

【0080】

テキスト選択部 112 は、テキスト記憶部 110 に記憶されている書起しテキストを全て読み出し、テキスト選択インターフェイス 90 に表示して、いずれかを選択するように促すメッセージを表示する。

【0081】

ユーザがいずれかのテキストを選択すると、テキスト選択部 112 はそのテキストをキーフレームデータ作成部 136 に与えるとともに、テキスト選択インターフェイス 90 上に、そのテキストを発話することを促すメッセージを表示する。同時にテキスト選択部 112 は、音声収録部 114 を起動し、マイクロフォン 92 からの音声信号の収録を開始する。

【0082】

音声収録部 114 は、入力される音声を所定フレーム長、所定シフト長でフレーム化した音声データを作成し、ハードディスク内に音声データファイルとして記憶する。

【0083】

音声信号の収録が終了すると、テキスト選択部 112 及び音声収録部 114 は、それぞれ、書起しテキストと音声データファイルとを、キーフレームデータ作成部 136 及び音声認識装置 120 の各々に与える。

【0084】

キーフレームデータ作成部 136 及び音声認識装置 120 は、このデータに対し以下のように動作する。まず音声認識装置 120 が、音声収録部 114 から与えられた音声データファイルに対し、書起しデータを参照して音素セグメンテーションを行ない、図 5 に示すような音素列ファイル 160 (継続時間長を特定できる時間情報を含む)を作成する。音声認識装置 120 は、この音素列ファイル 160 のデータをキーフレームデータ作成部 136 のマッピング処理部 180 に与える。

10

20

30

40

50

【 0 0 8 5 】

図 6 を参照して、キーフレームデータ作成部 1 3 6 のマッピング処理部 1 8 0 は、音素列ファイル 1 6 0 内の音素の各々に対し、マッピングテーブル記憶部 1 3 0 を参照してそれぞれ顔画像の識別子を付与し、継続時間長によるブレンド率調整部 1 8 2 に与える。

【 0 0 8 6 】

継続時間長によるブレンド率調整部 1 8 2 は、各音素の継続時間長としきい値とを比較し、しきい値未満の継続時間長しか持たない音素を削除し、その継続期間を直前の音素の継続期間に統合する。継続時間長によるブレンド率調整部 1 8 2 はさらに、各音素のブレンド率を、音素継続時間長の最大値と、その音素の継続時間長とに基づき、式 (1) に従って調整する。継続時間長によるブレンド率調整部 1 8 2 は、このようにして作成された、継続時間長、顔画像の識別子、及びブレンド率の付された音素列をパワーによるブレンド率調整部 1 8 4 に与える。

10

【 0 0 8 7 】

パワーによるブレンド率調整部 1 8 4 は、継続時間長によるブレンド率調整部 1 8 2 から与えられた音素列の各音素のうち、その期間中のパワーの値が所定のしきい値未満のものがあれば、その音素を削除する。そしてその音素の継続期間を直前の音素の継続期間と統合する。

【 0 0 8 8 】

パワーによるブレンド率調整部 1 8 4 はさらに、各音素のブレンド率を、パワーの最大値と、各音素のパワーとに基づき、式 (2) に従って調整する。パワーによるブレンド率調整部 1 8 4 は、このようにしてブレンド率が調整された音素列からなるキーフレームデータを図 7 に示すアニメーション生成制御部 2 0 0 に与える。

20

【 0 0 8 9 】

図 7 を参照して、アニメーション生成制御部 2 0 0 は、まず、与えられたキーフレームデータのうちの最初のキーフレーム (先頭の音素の開始時刻) の時刻を生成時刻 t とし、 $T_{k-1} < t < T_k$ となる整数 k を探す。この場合 $t = T_0 = T_{k-1}$ なので、 $k = 1$ となる。アニメーション生成制御部 2 0 0 は、キーフレーム T_0 及び T_1 に対応する顔画像のデータを顔データファイル記憶部 1 3 2 から読み出し、時刻 T_0 、 T_1 、生成時刻 t 、及びキーフレーム T_0 及び T_1 に対応する顔画像のデータを補間処理部 2 0 4 に与える。

30

【 0 0 9 0 】

補間処理部 2 0 4 は、与えられた時刻 T_0 、 T_1 と、それらキーフレーム T_0 、 T_1 のブレンド率と、生成時刻 t の値とに基づき、キーフレーム T_0 及び T_1 のブレンド率に対する、生成時刻 t におけるブレンド率 α 、 β を、補間関数記憶部 1 3 4 に記憶された補間関数 (式 (3)) を用いてそれぞれ算出する。補間処理部 2 0 4 はさらに、算出されたブレンド率 α 、 β と、キーフレーム T_0 及び T_1 における顔画像データの各特徴点ベクトル $X(T_0)$ 、 $X(T_1)$ とを用い、前述した式 (4) を用いて生成時刻 t における顔画像の各特徴点ベクトル $x(t)$ を算出し、出力部 1 4 2 に与える。出力部 1 4 2 はこの画像をディスプレイ 1 0 2 上に表示する。出力部 1 4 2 は、この顔画像の表示と同期して音声ファイル記憶部 1 4 0 に記憶された音声ファイルの再生を開始する。

40

【 0 0 9 1 】

アニメーション生成部 1 3 8 は、生成時刻 t における顔画像の算出が終了すると、処理の終了を示す信号をアニメーション生成制御部 2 0 0 に与える。アニメーション生成制御部 2 0 0 は、この信号を受信すると、そのときの時刻をタイマ 2 0 2 を参照して定める。アニメーション生成制御部 2 0 0 は、この時刻を新たな生成時刻 t に設定し、 $T_{k-1} < t < T_k$ となる整数 k を定める。そして、時刻 T_{k-1} 、 T_k 、生成時刻 t 、キーフレーム T_{k-1} 及び T_k のデータを補間処理部 2 0 4 に与え、時刻 t における顔画像の生成を実行させる。

【 0 0 9 2 】

補間処理部 2 0 4 は、前のサイクルと同様にして、与えられた時刻 T_{k-1} 、 T_k と、それらのキーフレーム T_{k-1} 、 T_k のブレンド率と、生成時刻 t の値とに基づき、キー

50

フレーム T_{k-1} 及び T_k のブレンド率に対する、生成時刻 t におけるブレンド率 α を補間関数記憶部 134 に記憶された補間関数 (式 (3)) を用いてそれぞれ算出する。補間処理部 204 はさらに、算出されたブレンド率 α と、キーフレーム T_{k-1} 及び T_k における顔画像データの各特徴点ベクトル $X(T_{k-1})$ 、 $X(T_k)$ とを用い、式 (4) を用いて生成時刻 t における顔画像の各特徴点ベクトル $x(t)$ を算出し、出力部 142 に与える。出力部 142 はこの画像をディスプレイ 102 上に表示する。音声ファイル記憶部 140 の再生は、画像の出力と同期して継続される。

【0093】

補間処理部 204 は、生成時刻 t における顔画像の算出が終了すると、それを示す信号をアニメーション生成制御部 200 に与える。アニメーション生成制御部 200 は、この信号を受信すると、そのときの時刻をタイマ 202 を参照して求める。そしてその時刻を新たな生成時刻 t に定める。

10

【0094】

以下、同様の処理を繰り返す。そして、生成時刻 t が音声の収録時間を上回ると、アニメーション生成システム 80 はアニメーション生成の処理を終了し、その後の状態は最初の書起しテキスト選択時の表示時の状態に戻る。

【0095】

このようにして、ある発話テキストをユーザが選択して読み上げると、その音声データに基づき、顔データファイル記憶部 132 に記憶された顔画像データを用い、口の形状が音声データに同期して変形する顔画像のアニメーションが得られる。最初に収録された音声も顔画像に同期して再生されるため、アニメーションのキャラクタが発話しているように見える。その結果、ユーザの声でキャラクタが発話するアニメーションを得ることができる。

20

【0096】

本実施の形態では、顔画像は、限定された音素に対応するものしか準備されていないが、マッピングテーブルを用いて各音素に対し、適切な顔画像をマッピングすることにより、十分に自然なアニメーションを得ることができる。音素の継続時間長が極端に短かったり、パワーが極端に小さかったりした場合、その音素については、画像の生成を省略している。通常は、このような音素を発音する際の実際の顔の動きも非常に小さい。そのため、この省略により、得られる顔画像のアニメーションは自然な動きに近く感じられる効果がある。さらに、ブレンド率という概念を用いて、各音素の発話の強さに応じて顔の変形量 (各特徴点の、基準画像 (無表情という画像) の各特徴点位置からの 3 次元的な移動量) を調整している。そのため、音素の発話の強さに応じて自然な動きの顔画像のアニメーションを得ることができる。また、隣り合うキーフレームの間の顔画像は、隣り合うキーフレームの顔画像を、キーフレームにおけるブレンド率と、キーフレームの時間と、画像の生成時間とに応じた加重和により内挿して得ている。従って、音素から音素への変化の際の口の形状変化が滑らかなものとなり、得られた顔画像のアニメーションも自然なものに感じられる。

30

【0097】

[第2の実施の形態]

40

上記した第1の実施の形態に係るアニメーション生成システム 80 は、十分な性能のコンピュータがあれば、そのコンピュータ一台でも実現可能である。しかし、ある程度短い時間で作業を完了させるためには、複数のコンピュータを用いることが実際的である。

【0098】

図 15 に、本発明の第2の実施の形態に係るアニメーション生成システム 280 の概略構成を示す。図 15 を参照して、このアニメーション生成システム 280 は、不特定のユーザによる音声入力を受け、アニメーション生成システム 280 でのアニメーションの生成を開始させる処理を行なう音声入力用のコンピュータ 292 と、音声入力用のコンピュータ 292 によって入力された音声に対する音素セグメンテーションを行なってキーフレームデータを作成するための音声認識サーバ 294 と、音声入力用のコンピュータ 292

50

による音声入力を受け、音声認識サーバ294が出力するキーフレームデータを利用して、入力された音声と同期して口の形状が変化する、所定のキャラクタの顔画像のアニメーションを作成し表示するためのアニメーション表示用コンピュータ296とを含む。音声入力用のコンピュータ292、音声認識サーバ294、及びアニメーション表示用コンピュータ296はいずれもネットワーク290を介して互いに所定のプロトコルで通信可能となっている。

【0099】

図2と比較すると、音声入力用のコンピュータ292が図2の入力指示ユニット94に、音声認識サーバ294が図2の音声認識装置120に、アニメーション表示用コンピュータ296が図2のアニメーション再生ユニット98に、それぞれ相当する。音声入力用のコンピュータ292、音声認識サーバ294、及びアニメーション表示用コンピュータ296の機能構成は、それぞれ図2の入力指示ユニット94、音声認識装置120、及びアニメーション再生ユニット98の構成と同様であるので、ここではその詳細は繰返さない。

10

【0100】

音声入力用のコンピュータ292は、タッチパネル300と、マイク302とを有する。アニメーション表示用コンピュータ296は、スピーカ312を有するモニタ310と、モニタ310の下に配置されたコンピュータ筐体314とを含む。

【0101】

図16に、アニメーション表示用コンピュータ296のハードウェア構成を示す。図16を参照して、アニメーション表示用コンピュータ296は、図15に示すモニタ310及びスピーカ312に加え、いずれもコンピュータ本体314内に配置された、CPU(中央演算処理装置)350と、読出専用メモリ(ROM)352と、随時読出書込可能メモリ(RAM)354と、ハードディスクドライブ356と、DVD(Digital Versatile Disk)330を装着可能なDVDドライブ358と、顔画像の演算処理を実行するためのGPU360とを含む。これらはいずれもバス362によりCPU350に接続されている。

20

【0102】

アニメーション表示用コンピュータ296はさらに、いずれもコンピュータ本体314内に配置され、バス362に接続された、ネットワークインターフェイス(I/F)368、フラッシュメモリからなる持運び可能なメモリ332を装着可能なメモリポート366、及びスピーカ312が接続されるサウンドボード364を含む。

30

【0103】

なお、アニメーション表示用コンピュータ296においては、アニメーション作成処理においてキーボードを使用する必要がないため、キーボードを備えていない。もちろん、アニメーション表示用コンピュータ296を通常のコンピュータとして使用する際には、コンピュータ本体314にキーボード及びマウス等の入力装置を接続することが可能である。

【0104】

図には示していないが、音声入力用のコンピュータ292及び音声認識サーバ294のハードウェア構成もアニメーション表示用コンピュータ296とほぼ同様である。相違点といえば、音声入力用のコンピュータ292において、モニタ310と入力装置とが一体となってタッチパネル300を構成していること、音声入力用のコンピュータ292がさらにマイクロフォン302を備えていること、音声入力用のコンピュータ292及び音声認識サーバ294ではGPU360が不要であること等である。

40

【0105】

図17に、音声入力用のコンピュータ292で実行されることにより、音声入力用のコンピュータ292を図2に示す入力指示ユニット94として動作させるためのコンピュータプログラムの制御構造をフローチャート形式で示す。

【0106】

50

図17を参照して、音声入力用のコンピュータ292の電源が投入され、このプログラムが起動されると、ステップ400で初期化処理が実行される。この処理では、音声入力用のコンピュータ292内で処理に必要な資源の確保及び初期化、通信機能の確認、発話テキストファイルからの発話テキストの読み込み等が行なわれる。

【0107】

初期化処理が終了すると、ステップ402において、音声入力用のコンピュータ292の準備が完了したことをアニメーション表示用コンピュータ296に通知する。続いてステップ404において、アニメーション表示用コンピュータ296より、音声入力用のコンピュータ292、音声認識サーバ294及びアニメーション表示用コンピュータ296がともに準備完了状態となったことを示す準備完了通知を受信したか否かを判定する。準備完了通知を受けたらステップ406に進む。準備完了通知を受取るまで、ステップ404の判定処理を繰り返す。

10

【0108】

このようにアニメーション表示用コンピュータ296からの準備完了通知を待つのは、同時期に音声入力用のコンピュータ292、音声認識サーバ294及びアニメーション表示用コンピュータ296が起動されたとして、全てにおいて準備が完了しないと、アニメーション生成システム280全体として機能することができないためである。

【0109】

続いてステップ406において、いくつかの発話テキストをタッチパネル300に表示し、「テキストを一つ選択してください」という、入力待ちメッセージを表示する。そしてステップ408で入力待ちの状態となる。入力があると、すなわちテキストがユーザにより選択されるとステップ410に進む。

20

【0110】

ステップ410では、選択されたテキストを発話するようにユーザに促すメッセージを表示し、録音を開始する。録音が終了するとステップ412に進む。

【0111】

ステップ412では、アニメーション表示用コンピュータ296に対して録音完了を通知する。続くステップ414において、アニメーション表示用コンピュータ296から処理開始通知を受信したか否かを判定する。処理開始通知とは、音声認識サーバ294における音素セグメンテーション処理と、アニメーション表示用コンピュータ296におけるアニメーション生成処理との実行を開始したことを示すメッセージである。

30

【0112】

ステップ416では、処理中を示す表示をタッチパネル300上に表示する。ステップ418で、ステップ410において録音した音声データと、対応するテキストデータ(書起しデータ)とをアニメーション表示用コンピュータ296に送信する。そして、ステップ420で、アニメーション表示用コンピュータ296から処理完了通知を受信するまで待機する。処理完了通知とは、ステップ418でアニメーション表示用コンピュータ296に対し送信した音声データに対して、音素セグメンテーション処理とその後のキーフレームデータ作成処理までが完了したことを示す通知である。

【0113】

40

処理完了通知を受信すると、ステップ422において、アニメーション表示用コンピュータ296に対し、アニメーションの出力命令を送信する。後述するように、アニメーション表示用コンピュータ296は、この出力命令に対してアニメーションの生成処理及び出力処理を開始する。ステップ424ではアニメーション表示用コンピュータ296から出力中通知を受信するまで待機し、出力中表示を受けるとステップ426に進む。ステップ426では、タッチパネル300上に、アニメーションをアニメーション表示用コンピュータ296のモニタ310上に出力中であることを示すメッセージを表示する。そしてステップ428で、アニメーション表示用コンピュータ296からアニメーションの出力処理が完了したことを示す出力完了通知を待つ。出力完了通知を受信すると、ステップ410で録音した音声に対するアニメーションの生成及び表示が全て完了したということで

50

ある。従って制御はステップ 4 0 6 に戻り、次のユーザ入力を待つ。

【 0 1 1 4 】

音声入力用のコンピュータ 2 9 2 は、上記した処理を繰り返す。

【 0 1 1 5 】

図 1 8 は、音声認識サーバ 2 9 4 が実行する処理のフローチャートである。このプログラムが起動されると、ステップ 4 4 0 において初期化処理が実行される。初期化処理が完了すると、ステップ 4 4 2 においてアニメーション表示用コンピュータ 2 9 6 に対し音声認識サーバ 2 9 4 の準備が完了したことを示す通知を送信する。

【 0 1 1 6 】

ステップ 4 4 4 では、アニメーション表示用コンピュータ 2 9 6 から音素セグメンテーションの依頼を受信したか否かを判定する。音素セグメンテーションとは音声認識処理と同様の処理であって、入力された音声を、音響モデルを用いて音素に分割する処理のことをいう。依頼を受信すると、ステップ 4 4 8 に進む。

【 0 1 1 7 】

ステップ 4 4 8 で、アニメーション表示用コンピュータ 2 9 6 に対し、音声認識サーバ 2 9 4 が音素セグメンテーション処理を開始したことを通知する。

【 0 1 1 8 】

続くステップ 4 5 0 において、依頼に従い、音素セグメンテーションを行なうべき音声データと書起しデータとをアニメーション表示用コンピュータ 2 9 6 より取得する。この取得が完了したら、ステップ 4 5 2 において対象データの受信が完了したことを示す通知をアニメーション表示用コンピュータ 2 9 6 に送信する。ステップ 4 5 4 では、受信した音声データに対し、図示しない音響モデルと、受信した書起しデータとを用いた音素セグメンテーション処理を実行する。この処理では、書起しデータが存在するので、正確な音素セグメンテーションをすることが可能である。

【 0 1 1 9 】

音素セグメンテーションが終了し、音素列ファイルの生成が完了したら、ステップ 4 5 6 において、音素列ファイルの生成が完了したことをアニメーション表示用コンピュータ 2 9 6 に通知する。

【 0 1 2 0 】

さらに、この音素列ファイルに基づき、ステップ 4 5 8 において、キーフレームデータの生成処理を実行する。キーフレームデータの生成処理の詳細については図 1 9 を参照して後述する。キーフレームデータの生成処理が完了すると、ステップ 4 6 0 においてキーフレームデータ生成が完了したことをアニメーション表示用コンピュータ 2 9 6 に対して通知する。さらに、ステップ 4 6 2 において、音素列ファイルと、キーフレームデータとをアニメーション表示用コンピュータ 2 9 6 に対して送信する。ステップ 4 6 4 では、アニメーション表示用コンピュータ 2 9 6 に対して音声認識サーバ 2 9 4 における処理が全て完了したことを通知し、ステップ 4 4 4 に戻る。

【 0 1 2 1 】

図 1 9 は、図 1 8 のステップ 4 5 8 で実行されるキーフレームデータの作成処理のフローチャートである。図 1 9 を参照して、ステップ 4 8 0 において、与えられた音素列の中で、継続時間長が所定のしきい値より小さい音素、又はパワーが所定のしきい値より小さい音素があるか否かを判定する。もしあれば、ステップ 4 8 2 において、その音素を削除し、その音素の継続時間長を直前の音素の継続時間長に統合する処理を行ない、ステップ 4 8 0 に戻る。上記したような音素が存在しなくなると、ステップ 4 8 4 に進む。

【 0 1 2 2 】

ステップ 4 8 4 では、与えられた音素列を構成する全ての音素に対して、ブレンド率の初期値として 1 0 0 % を設定する。続くステップ 4 8 6 では、図 2 に示すマッピングテーブル記憶部 1 3 0 に記憶されたマッピングテーブルを用い、音素列中の各音素に対し、図 3 に示す顔画像 / A / ~ / / の中のいずれかを割当て、その顔画像の識別子を音素に付す。こうして割当てられた顔画像が、その音素の開始時点フレーム時刻とするキーレ

10

20

30

40

50

ームとなる。

【0123】

続いてステップ488において、与えられた全ての音素列を調べ、音素の最大継続時間長と最大パワーとを探索する。探索された最大継続時間長を L_{MAX} 、最大パワーを P_{MAX} とする。

【0124】

ステップ490では、各音素のブレンド率を、前述した式(1)により更新する。なお、式(1)で $B(n)$ は n 番目の音素のブレンド率を表す。同様に、ステップ492では、各音素のブレンド率を、前述した式(2)により更新する。

【0125】

ステップ494では、上記したように算出されたブレンド率と、対応の顔画像の識別子と、時間情報とが付された音素列を、キーフレームデータとしてファイルに出力し、キーフレームデータの作成処理を終了する。

【0126】

図20は、アニメーション表示用コンピュータ296により実行されるアニメーション生成制御処理を実現するコンピュータプログラムの制御構造を示すフローチャートである。図20を参照して、アニメーション生成制御処理が起動されると、ステップ500において初期化処理を行ない、ステップ502において音声入力用のコンピュータ292及び音声認識サーバ294からの準備完了通知を待つ。

【0127】

音声入力用のコンピュータ292及び音声認識サーバ294から準備完了通知を受信すると、ステップ504において音声入力用のコンピュータ292に対しアニメーション生成システム280の全体が準備完了していることを示す準備完了通知を送信する。続いてステップ506で、音声入力用のコンピュータ292から録音完了通知を受信するまで待機する。

【0128】

録音完了通知を受信すると、ステップ508において、音声入力用のコンピュータ292に対し音声認識サーバ294及びアニメーション表示用コンピュータ296がアニメーションを作成するための一連の処理を実行開始することを示す処理開始通知を送信する。続いてステップ510で、音声入力用のコンピュータ292から書起しテキストデータ及び音声データを受信するまで待機し、これらデータを受信するとステップ512に進む。

【0129】

ステップ512では、音声認識サーバ294に対し、ステップ510で受信した書起しテキストデータ及び音声データを送信し、音素セグメンテーションを依頼する。そしてステップ514では、音素セグメンテーションの結果得られるキーフレームデータを音声認識サーバ294から受信するまで待機する。キーフレームデータを受信すると、ステップ520以下のアニメーション生成のための処理を実行する。

【0130】

ステップ520において、本実施の形態では、顔画像のアニメーションの先頭フレームの時刻(生成時刻) t として、音素列の最初の音素の時刻 T_0 を選択する。

【0131】

続いてステップ522において、直前のステップで決定されたフレームの生成時刻 t に対し、 $T_{k-1} < t < T_k$ となる k を決定する。ただし T_k は音素列中の k 番目の音素の期間の開始時刻を指す。例えば $t = T_0$ であれば $T_0 < t < T_1$ であるから、 $k = 1$ となる。

【0132】

続いてステップ524において、時刻 T_k 及び T_{k-1} と、キーフレーム T_{k-1} 、 T_k のブレンド率と、生成時刻 t と、時刻 T_k 及び T_{k-1} での音素に対応する顔画像データとをGPU360に渡し、生成時刻 t における顔画像を補間により生成することを依頼する。

10

20

30

40

50

【 0 1 3 3 】

これに応答し、GPU360が実行するプログラムは、生成時刻 t における、キーフレーム T_{k-1} のブレンド率から補間演算されるブレンド率、及びキーフレーム T_k のブレンド率から補間演算されるブレンド率をそれぞれ前述した補間式(3)により算出し、さらに生成時刻 t における顔画像の各特徴点ベクトル $X(t)$ を、キーフレーム T_{k-1} 及び T_k における顔画像の各特徴点ベクトル $X(T_{k-1})$ 及び $X(T_k)$ と、 α 、 β とを用い、前述の式(4)によるベクトル加重和によって算出する。GPU360は、この計算が顔画像の全ての特徴点に対し終了すると、生成された時刻 t における顔画像を出力し、さらに処理終了通知をCPU350に対して送信する。

【 0 1 3 4 】

図20を参照して、アニメーション生成制御処理のプログラムは、ステップ526でGPU360からの終了通知を受信するまで待ち状態となる。終了通知を受信するとステップ528に進む。

【 0 1 3 5 】

ステップ528では、タイマ202の時刻を読む。この時刻を新たな生成時刻 t とする。続いてステップ530では、生成時刻 t が、録音の最終時刻よりも後か否かを判定する。生成時刻 t が録音時刻より後であれば、処理を終了する。さもなければ、この新たな生成時刻 t における顔画像データを生成すべく、ステップ522に戻る。

【 0 1 3 6 】

以下、ステップ522～ステップ530の処理を、生成時刻 t が録音時間より大きくなるまで繰り返す。生成時刻 t が録音時間より大きくなると、ステップ532に進む。

【 0 1 3 7 】

ステップ532では、音声入力用のコンピュータ292に対し、音声認識サーバ294及びアニメーション表示用コンピュータ296における処理が完了したことを示す通知を送信する。音声入力用のコンピュータ292はこの通知を図17のステップ428で受信し、これに回答してステップ406に戻り、上記した一連の処理が音声入力から繰返される。

【 0 1 3 8 】

図17にフローチャートで示す制御構造を有するプログラムを音声入力用のコンピュータ292で、図18及び図19にフローチャートで示す制御構造を有するプログラムを音声認識サーバ294で、図20にフローチャートで示す制御構造を有するプログラムをアニメーション表示用コンピュータ296で、それぞれ実行することにより、第1の実施の形態に係るアニメーション生成システム80と同様の機能を持つアニメーション生成システム280を実現することができる。

【 0 1 3 9 】

なお、第1の実施の形態に係るアニメーション生成システム80をコンピュータで実現する際にも、上記した図17～図20に示す制御構造を有するコンピュータプログラムと同様のプログラムを利用することができる。

【 0 1 4 0 】

< 動作 >

第2の実施の形態に係るアニメーション生成システム280の動作は、第1の実施の形態に係るアニメーション生成システム80と同様である。従って、ここではその詳細については繰返さない。

【 0 1 4 1 】

本実施の形態では、各コンピュータに処理を分散させている。そのため、各コンピュータの性能はそれほど高くなくてもよい。また、音声認識サーバ294としては高性能なものを準備しておき、複数の音声入力用のコンピュータ292とアニメーション表示用コンピュータ296との組からのキーフレームデータ作成要求を単一の音声認識サーバ294で処理することも可能である。

【 0 1 4 2 】

10

20

30

40

50

さらに、本実施の形態では、音声入力用のコンピュータ 292 とアニメーション表示用コンピュータ 296 とは別のコンピュータとしたが、これらをまとめて一つのコンピュータによって実現するようにしてもよい。

【0143】

どのような音素に対応する顔画像を準備するか、及びどれだけの数の顔画像を準備するかは、アニメーション製作時の設計事項である。また、どの音素に対しどの顔画像をマッピングするかもアニメーション製作時の設計事項である。また、音素の組は対象とする言語によっても異なり、従ってマッピングも異なってくることは当然である。

【0144】

上記した実施の形態では、ある音素に対しては必ず一つの顔画像が対応するように音素と顔画像とのマッピングがされているが、そうでなくてもよい。すなわち、同一の音素でも、その前後の音素によって異なる顔画像を割当てるようにしてもよい。

10

【0145】

上記した実施の形態では、ブレンド率の算出に式(1)及び(2)を使用している。しかし本発明は、式(1)及び(2)を用いるものには限定されない。継続時間長又はパワーが短くなればブレンド率が低くなるようなものであれば、すなわち継続時間長及びパワーに対する単調関数であれば、どのような関数を用いてブレンド率を算出するようにしてもよい。また、継続時間長及びパワーに限らず、それ以外の音声特徴量を考慮してブレンド率を決定してもよい。

【0146】

20

上記した実施の形態では、補間関数として図12に示されるような直線式に対応するものを用いた。しかし本発明はそのような実施の形態には限定されない。補間関数として、時間に対する2次以上の多項式を用いたり、非線形関数を用いたりしてもよい。本実施の形態では、キーフレームに相当する時刻においてブレンド率が最も高くなり、そこから遠ざかるにつれてブレンド率が低くなるような補間関数であれば、どのようなものを用いてもよい。補間関数として複数のものを用意しておき、ユーザが切替えて使用できるようにしておいてもよい。

【0147】

また上記実施の形態では、キーフレームの位置は、各音素の継続期間の先頭位置としたが、本発明はそのような実施の形態には限定されない。キーフレームの位置を、各音素の継続期間の途中にしてもよい。キーフレームの位置についても、ユーザが任意に変更可能としてもよい。

30

【0148】

なお、上記した実施の形態では、音素列中のある音素の継続時間長又はパワーがしきい値より小さい場合には、その音素を削除し、その継続時間長を直前の音素の継続時間長に統合した。こうすることにより、口形状の変化が滑らかで自然なものとなる効果が得られる。

【0149】

しかし本発明はそのような実施の形態には限定されない。例えば、ある音素の継続時間長のみを考慮したり、パワーのみを考慮するようにしてもよい。又は、継続時間長及びパワーの双方がそれぞれしきい値より小さいときに、その音素を削除するようにしてもよい。これらの間で、切換を行なうようにしてもよい。

40

【0150】

さらに、上記した実施の形態では、最終的にアニメーションとともに再生される音声は、最初に収録されたユーザの音声そのままである。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。口形状は主として音素との関係で決定されるので、音素の位置にさえ大きな変更を加えないのであれば、ユーザの音声に何らかの加工を加えるようにしてもよい。この場合でも、最終的に再生される音声にはユーザの発話の特徴が生かされることが多く、より多彩なアニメーションを作成できる。

【0151】

50

上記した実施の形態では、ユーザによる書起しテキストの発話の録音後、キーデータファイルを生成し、キーデータファイルを生成した後はGPU360による顔画像の作成処理の終了時に次の顔画像の生成を開始するようにしている。従って、顔画像の生成は一定のサイクルで行なわれているわけではない。こうすることにより、GPU360はその性能をフルに発揮できる。しかし、本発明はこのようにして顔画像を作成するものには限定されない。

【0152】

例えば、図12に示されるように、各キーフレームの補間によるブレンド率の分布を求めた後、一定のフレーム間隔で顔画像を生成すべき時刻の系列を求め、各時刻での顔画像を生成し、全ての顔画像が生成された後にそれらをアニメーションとして再生するようにしてもよい。この場合には、フレーム間隔が短くなると処理に長時間を要するようになり、逆にフレーム間隔が長くなるとアニメーションの動きがぎこちなくなる可能性がある。

【0153】

なお、上記した実施の形態は、ユーザの音声にあわせて顔画像のアニメーションを作成し、再生する。音声の書起しテキストは決まっているため、ユーザが不特定であっても音素セグメンテーションを精度高く行なえ、滑らかなアニメーションを作成できる。

【0154】

上記実施の形態では、音声が入力されると、それに基づいて作成したアニメーションを一回だけ再生し、次の音声の入力を待つ。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。音声を入力し、キーフレームデータを作成した後は、そのキーフレームデータに基づいて、何回でもアニメーションの再生を行なうことができる。特に、この再生においては、使用される顔画像を変えたり、補間の関数を変えたり、音素を間引く際のしきい値を変えたりして、同じ音声から様々なアニメーションを生成できる。そのため、いわゆるプレスコ(プレレコ)方式によってアニメーションを作成するためのツールとして利用することが可能である。

【0155】

さらに、上記した実施の形態は音声に基づいて顔画像のアニメーションを生成するものであった。しかし本発明はそのような実施の形態に限定されるわけではない。音声に伴って形状が変化するものであり、その形状とある音素とのマッピングが可能なものであれば、どのようなものにも適用可能である。例えば、音声にあわせて声道形状のアニメーションを作成したり、調音機構のアニメーションを作成したりすることも考えられる。

【0156】

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

【図面の簡単な説明】

【0157】

【図1】本発明の第1の実施の形態に係るアニメーション作成装置によるアニメーション作成過程30の概略を示す図である。

【図2】第1の実施の形態に係るアニメーション生成システム80の概略の機能的構成を示すブロック図である。

【図3】顔データファイル記憶部132に記憶される顔画像の例を示す図である。

【図4】マッピングテーブル記憶部130に記憶されたマッピングテーブルの例の一部を示す図である。

【図5】音声認識装置120の出力する音素列ファイル160の構成例を示す図である。

【図6】図2のキーフレームデータ作成部136の構成の詳細を示すブロック図である。

【図7】アニメーション生成部138のより詳細なブロック図である。

【図8】音素列とブレンド率との関係を示す図である。

【図9】音素列とブレンド率との関係を示す図である。

10

20

30

40

50

- 【図10】音素列とブレンド率との関係を示す図である。
- 【図11】音素列とブレンド率との関係を示す図である。
- 【図12】ブレンド率の補間の概略を示す図である。
- 【図13】キーフレームにおける顔画像のベクトル加重和を説明するための図である。
- 【図14】アニメーション生成制御部200によるアニメーションの生成制御処理を説明するための図である。
- 【図15】本発明の第2の実施の形態にかかるアニメーション生成システム280の概略構成を示すブロック図である。
- 【図16】アニメーション表示用コンピュータ296のハードウェア構成を示すブロック図である。
- 【図17】音声入力用のコンピュータ292を図2に示す入力指示ユニット94として動作させるためのコンピュータプログラムの制御構造を示すフローチャートである。
- 【図18】音声認識サーバ294が実行するコンピュータプログラムの制御構造を示すフローチャートである。
- 【図19】キーフレームデータの作成処理を実現するコンピュータプログラムの制御構造を示すフローチャートである
- 【図20】アニメーション生成制御処理を実現するコンピュータプログラムの制御構造を示すフローチャートである。

10

【符号の説明】

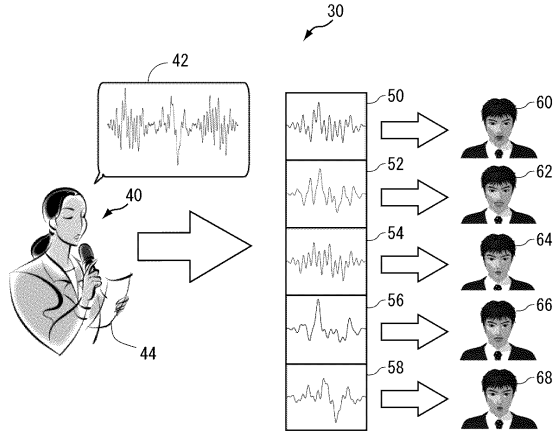
【0158】

20

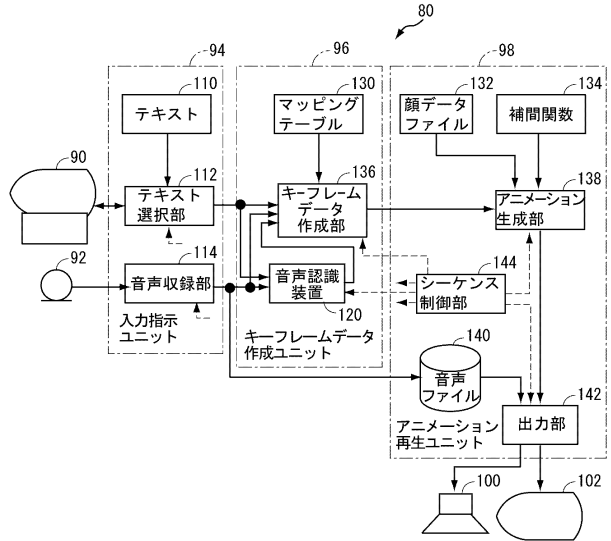
40 話者、42 音声信号、44 台本、60～68 顔画像、80, 280 アニメーション生成システム、90 テキスト選択インターフェイス、92 マイクロフォン、94 入力指示ユニット、96 キーフレームデータ作成ユニット、98 アニメーション再生ユニット、100 スピーカ、102, 310 モニタ、110 テキスト記憶部、112 テキスト選択部、114 音声収録部、120 音声認識装置、130 マッピングテーブル記憶部、132 顔データファイル記憶部、134 補間関数記憶部、136 キーフレームデータ作成部、138 アニメーション生成部、140 音声ファイル記憶部、142 出力部、160 音素列ファイル、180 マッピング処理部、182 継続時間長によるブレンド率調整部、184 パワーによるブレンド率調整部、200 アニメーション生成制御部、202 タイマ、204 補間処理部、290 ネットワーク、292 音声入力用のコンピュータ、294 音声認識サーバ、296 アニメーション表示用コンピュータ、300 タッチパネル、302 マイクロフォン

30

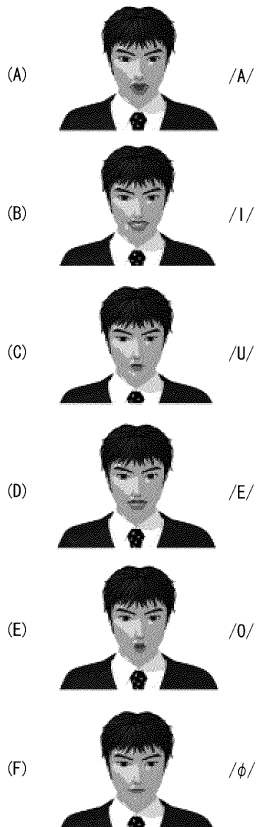
【図1】



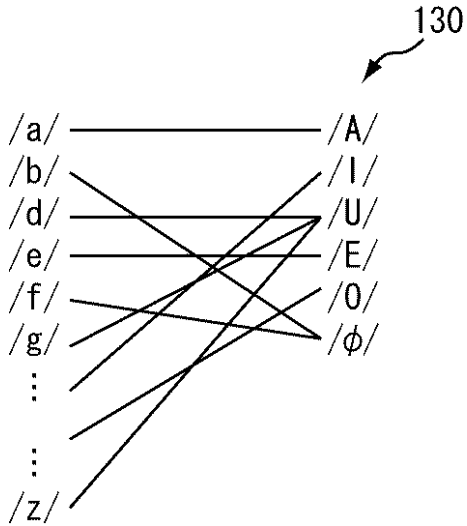
【図2】



【図3】



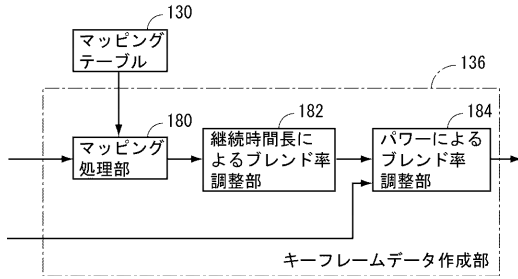
【図4】



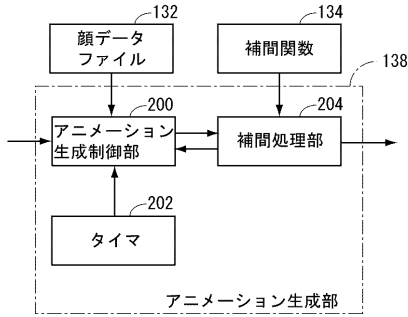
【図5】

160	
音素列	/a/, /i/, /u/, /e/, /o/, ...
継続時間長	300, 300, 100, 30, 250, ...

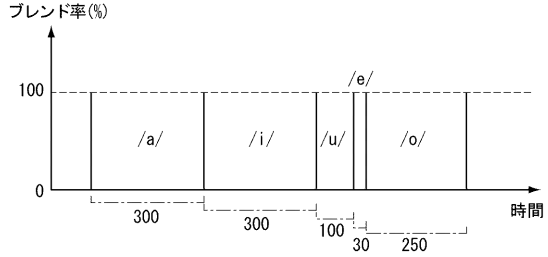
【図 6】



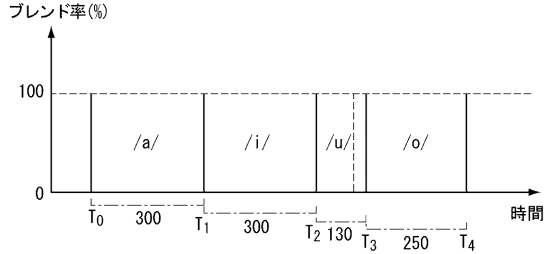
【図 7】



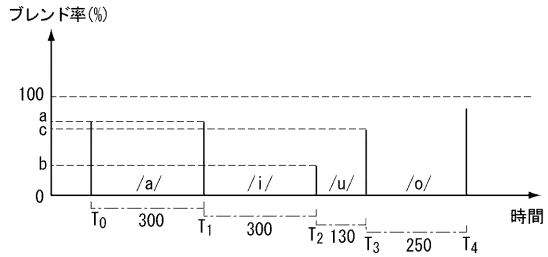
【図 8】



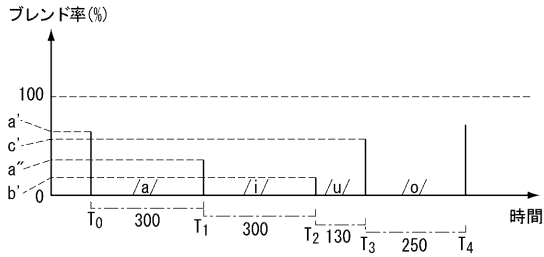
【図 9】



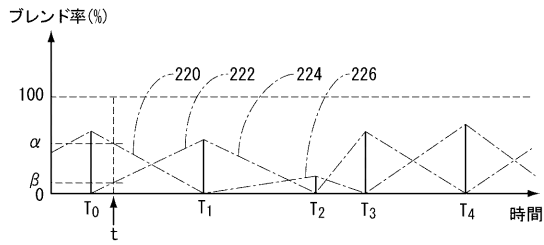
【図 10】



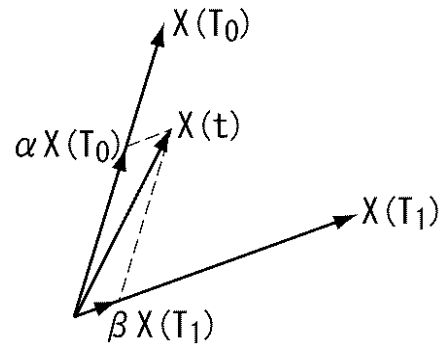
【図 11】



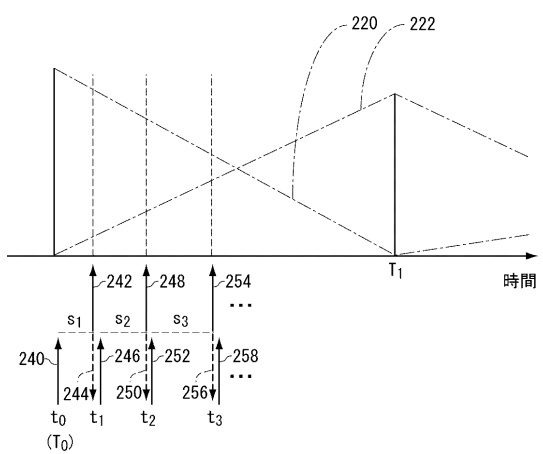
【図 12】



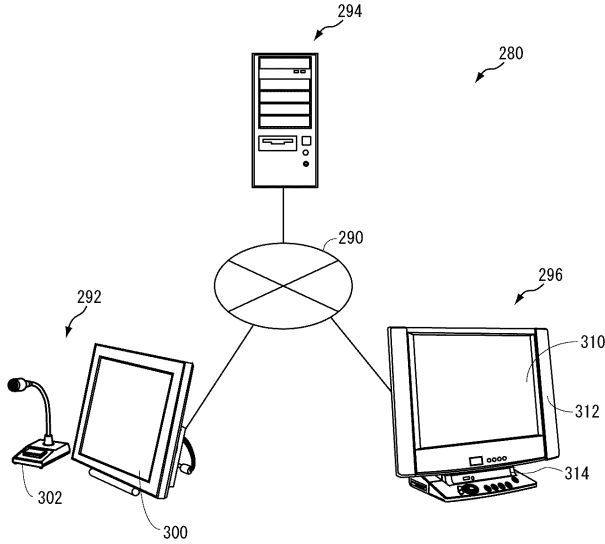
【図 13】



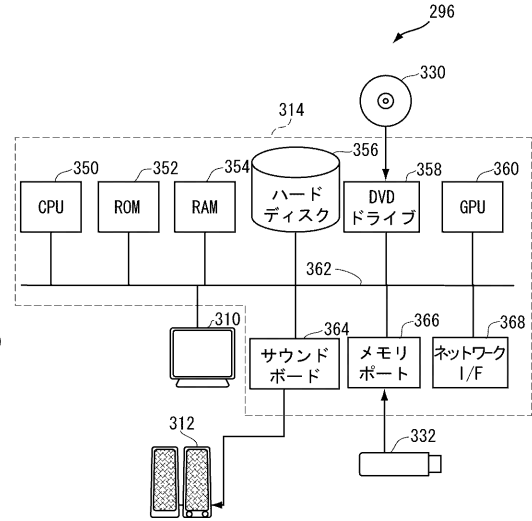
【図 14】



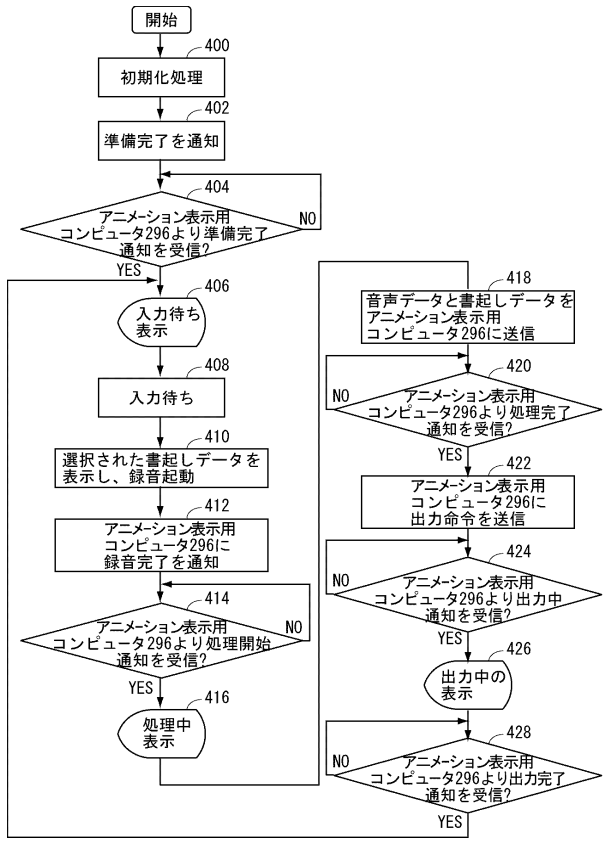
【図15】



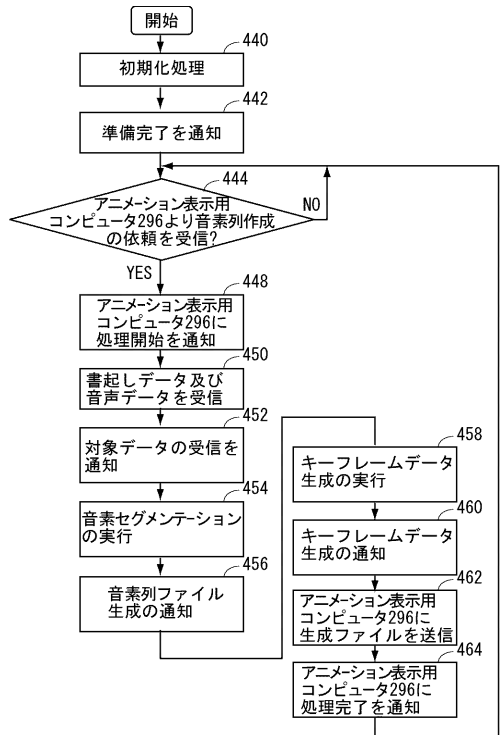
【図16】



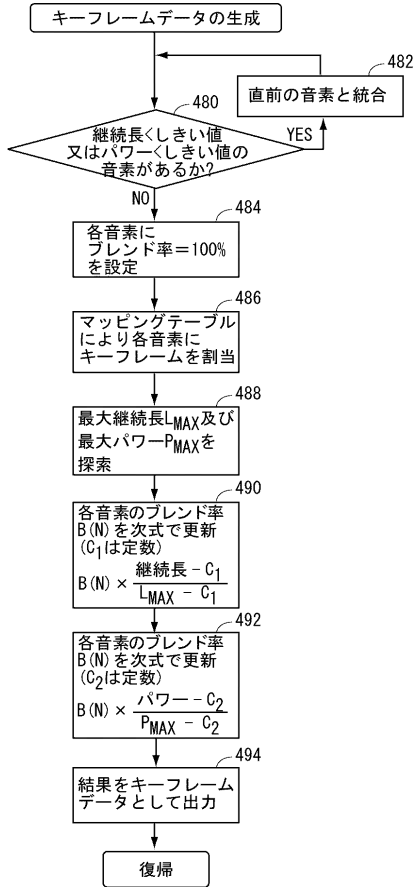
【図17】



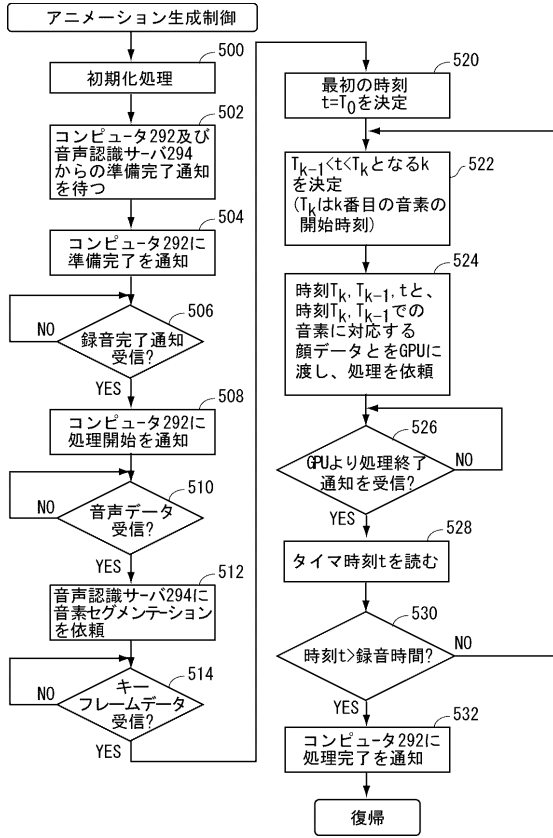
【図18】



【図19】



【図20】



フロントページの続き

審査官 加内 慎也

- (56)参考文献 特開平08 - 123977 (JP, A)
特開平06 - 052290 (JP, A)
特開2003 - 337956 (JP, A)
特開2003 - 132363 (JP, A)
特開2002 - 133445 (JP, A)

- (58)調査した分野(Int.Cl., DB名)
G06T 15/70