

(19)日本国特許庁 (J P)

(12) 特 許 公 報 (B 2)

(11)特許番号

特許第3485508号
(P3485508)

(45)発行日 平成16年1月13日(2004.1.13)

(24)登録日 平成15年10月24日(2003.10.24)

(51)Int.Cl. ⁷	識別記号	F I
G 0 6 T 11/20	1 0 0	G 0 6 T 11/20 1 0 0
G 0 6 N 3/00	5 6 0	G 0 6 N 3/00 5 6 0 C
G 1 0 L 15/00		H 0 4 M 11/06
		G 1 0 L 3/00 5 3 9
H 0 4 M 11/06		5 5 1 A

請求項の数17(全 15 頁) 最終頁に続く

(21)出願番号	特願平11-304295	(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11年10月26日(1999.10.26)	(72)発明者	へ二・ヤヒヤ 京都府相楽郡精華町大字乾谷小字三平谷 5番地 株式会社エイ・ティ・アール人 間情報通信研究所内
(65)公開番号	特開2001-126077(P2001-126077A)	(72)発明者	倉立 尚明 京都府相楽郡精華町大字乾谷小字三平谷 5番地 株式会社エイ・ティ・アール人 間情報通信研究所内
(43)公開日	平成13年5月11日(2001.5.11)	(74)代理人	100064746 弁理士 深見 久郎 (外5名)
審査請求日	平成13年3月12日(2001.3.12)	審査官	伊知地 和之

最終頁に続く

(54)【発明の名称】 顔画像伝送方法およびシステムならびに当該システムで用いられる顔画像送信装置および顔画像再生装置

1

(57)【特許請求の範囲】

【請求項1】 話者の発する音声を受け、当該音声を発声するときの当該話者の表情を推定する信号を出力する表情推定手段を、受信側において準備するステップを備え、

前記準備するステップは、
前記送信側と前記受信側とにおいて同一のニューラルネットワークを準備するステップと、
前記送信側において、前記話者の発する音声を受けて前記話者の顔画像を特定する情報を出力するように前記ニューラルネットワークを学習させるステップと、
前記送信側の学習済みの前記ニューラルネットワークの特性パラメータを前記受信側に送信し、前記受信側のニューラルネットワークを前記送信された特性パラメータにより設定するステップとを含む、

2

前記話者の発する音声を送信側から前記受信側に送信し、前記表情推定手段に与えて前記話者の表情を推定する信号を得るステップと、
前記表情推定手段の出力する前記話者の表情を推定する信号に基づいて、前記話者の表情の動画像を生成するステップとを含む、顔画像伝送方法。

【請求項2】 前記学習させるステップは、
前記送信側において、予め定められた学習のための文章を朗読する前記話者の顔の、予め定められた箇所の座標を時系列的に測定するステップと、
前記学習のための文章を朗読する際の前記話者の発する音声について、時系列的に音声の特徴量を求めるステップと、
前記音声の特徴量を入力とし、測定された前記座標を教師信号として前記送信側のニューラルネットワークの特

10

性パラメータを調整するステップとを含む、請求項 1 に記載の顔画像伝送方法。

【請求項 3】 前記送信側の前記ニューラルネットワークの学習後において、話者の発声時に、前記話者の顔の予め定められた箇所の座標を時系列的に測定するステップと、

前記話者の発声する音声を前記送信側の学習済みの前記ニューラルネットワークに与えて前記話者の顔の前記予め定められた箇所の座標の推定値を学習済みの前記ニューラルネットワークの出力として得るステップと、前記測定された座標と、前記座標の推定値とを比較し、前記座標の推定値の有する誤差を求めるステップとをさらに含む、請求項 2 に記載の顔画像伝送方法。

【請求項 4】 前記誤差を前記受信側に送信するステップと、前記受信側において前記誤差を受信し、前記受信側のニューラルネットワークの出力を受信された前記誤差に基づいて補正するステップとをさらに含む、請求項 3 に記載の顔画像伝送方法。

【請求項 5】 前記誤差の大きさと所定のしきい値とを比較するステップと、前記誤差の大きさが前記所定のしきい値を超えたことに応答して、前記誤差を前記受信側に送信するステップと、

前記受信側において前記誤差を受信し、前記受信側のニューラルネットワークの出力を受信された前記誤差に基づいて補正するステップとをさらに含む、請求項 3 に記載の顔画像伝送方法。

【請求項 6】 送信側の話者の発声時の顔の動画像を受信側に送信するための顔画像伝送システムであって、話者の発する音声信号を前記受信側に送信する送信装置を備え、

前記送信装置は、
前記話者の発する音声の所定の特徴量を時系列的に抽出するための特徴量抽出手段と、

話者の発する音声に基づいて前記特徴量抽出手段が出力する特徴パラメータを入力とし、前記話者の顔画像を特定する情報を出力するように学習可能なニューラルネットワークと、

前記ニューラルネットワークの特性パラメータを前記受信側に送信する手段とを含み、

前記送信装置からの信号を受ける受信装置をさらに備え、

前記受信装置は、
前記送信装置から送信されてくる前記話者の発する音声信号を受け、当該音声を発声するときの当該話者の表情を推定する信号を出力する表情推定手段を含み、

前記表情推定手段は、
前記送信装置の前記ニューラルネットワークと同一構成のニューラルネットワークと、

前記受信装置の前記ニューラルネットワークを前記送信装置から送信された前記特性パラメータにより設定するための手段とを有し、

前記表情推定手段が前記送信側から前記話者の発する音声信号を受信して出力する前記話者の表情を推定する信号に基づいて、前記話者の表情の動画像を生成する顔画像生成手段をさらに含む、顔画像伝送システム。

【請求項 7】 前記送信装置は、
話者の発声時に、前記話者の顔の予め定められた箇所の座標を時系列的に測定するための測定手段と、
10 予め定められた学習のための文章を話者が朗読する際の、前記話者の顔の予め定められた箇所に関して前記測定手段によって得られた座標データを教師信号とし、前記特徴量抽出手段によって得られる音声の特徴量を入力として前記送信装置の前記ニューラルネットワークを学習させるための学習手段とをさらに含む、請求項 6 に記載の顔画像伝送システム。

【請求項 8】 前記送信装置は、前記送信側の前記ニューラルネットワークの学習後において、前記話者の発する音声を前記送信側の学習済みの前記ニューラルネットワークに与えて得られた前記話者の顔の前記予め定められた箇所の座標の推定値を前記測定された座標と比較し、前記座標の推定値の有する誤差を求めるための手段をさらに含む、請求項 7 に記載の顔画像伝送システム。

【請求項 9】 前記送信装置は、前記誤差を前記受信側に送信するための手段をさらに含み、
前記受信装置は、前記誤差を受信して、前記受信側のニューラルネットワークの出力を受信された前記誤差に基づいて補正するための誤差補正手段をさらに含む、請求項 8 に記載の顔画像伝送システム。

【請求項 10】 前記送信装置は、前記誤差の大きさと所定のしきい値とを比較し、前記誤差の大きさが前記所定のしきい値を超えたことに応答して、前記誤差を前記受信側に送信するための手段をさらに含み、
前記受信装置は、前記誤差を受信して、前記受信側のニューラルネットワークの出力を受信された前記誤差に基づいて補正するための誤差補正手段をさらに含む、請求項 8 に記載の顔画像伝送システム。

【請求項 11】 送信側の話者の発声時の顔の動画像を受信側に送信するための顔画像伝送システムで使用される顔画像送信装置であって、
話者の音声を前記受信装置に送信するための手段と、
話者の発する音声の特徴量を入力とし、前記話者の顔画像を特定する情報を出力するように学習可能なニューラルネットワークと、
前記ニューラルネットワークの特性パラメータを前記受信側に送信するための手段とを含む、顔画像送信装置。

【請求項 12】 前記話者の顔の予め定められた箇所の座標を時系列的に測定するための測定手段と、
50 前記話者の発する音声の所定の特徴量を時系列的に抽出

するための特徴量抽出手段と、前記話者が予め定められた文章を朗読する際に前記特徴量抽出手段により求められた前記特徴量を入力とし、予め定められた文章を朗読する際に前記測定手段により測定された前記座標を教師信号として前記ニューラルネットワークを学習させるための手段とをさらに含む、請求項 1 1 に記載の顔画像送信装置。

【請求項 1 3】 前記ニューラルネットワークの学習後に、前記話者の発声する音声に対して前記特徴量抽出手段により出力された特徴量を、学習済みの前記ニューラルネットワークに与えて得られた前記話者の顔の前記予め定められた箇所の座標の推定値を、前記測定手段によって測定された座標と比較し、前記座標の推定値の有する誤差を求めるための手段をさらに含む、請求項 1 2 に記載の顔画像送信装置。

【請求項 1 4】 前記送信装置は、前記誤差を前記受信側に送信するための手段をさらに含む、請求項 1 3 に記載の顔画像送信装置。

【請求項 1 5】 前記送信装置は、前記誤差と所定のしきい値とを比較し、前記誤差が前記所定のしきい値を超えたことに応答して、前記誤差を前記受信側に送信するための手段をさらに含む、請求項 1 3 に記載の顔画像送信装置。

【請求項 1 6】 話者の発する音声から得られた所定の特徴量を受け、当該音声を発声するときの当該話者の表情を推定する信号を出力する表情推定手段を備え、前記表情推定手段は、予め定められた構成のニューラルネットワークと、顔画像の再生に先立って前記ニューラルネットワークを設定するための特性パラメータを所定の信号源から受け、前記特性パラメータによって前記ニューラルネットワークを設定するための手段とを含み、前記表情推定手段が前記特徴量を受けて出力する前記話者の表情を推定する信号に基づいて、前記話者の表情の動画像を生成する顔画像生成手段をさらに備える、顔画像再生装置。

【請求項 1 7】 前記受信顔画像受信装置の前記ニューラルネットワークの設定後に、前記所定の信号源から送信される誤差を受信して、前記顔画像再生装置の前記ニューラルネットワークの出力を受信された前記誤差に基づいて補正するための誤差補正手段をさらに含む、請求項 1 6 に記載の顔画像再生装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は電気通信によるコミュニケーション技術に関し、特に、話者の表情を少ない情報量で再生できるようにする顔画像受信方法、システム、ならびにそのための顔画像送信装置および顔画像再生装置に関する。

【0002】

【従来の技術】ネットワークの発達およびそのための電気通信技術の発達により、ネットワークを介して遠隔地にいる人同士が会議を行う TV 会議システムが実用化されつつある。また、簡易なカメラを装着した電話機同士での、いわゆる TV 電話も実用化されつつある。

【0003】一般的に、こうしたシステムで動画像信号を送信する場合、音声のみを送信する場合と比較して送信されるべき情報量ははるかに大きい。画像信号は電話音声の 1000 倍程度の情報量を持つことが知られており、そのため、たとえば電話回線を用いたリアルタイムの動画像送信はほぼ不可能である。前述したシステムでも、画面の動きは極めて少なく、1 秒間に数フレームの送信しか行われない。

【0004】動画像を送信するための媒体として非常に伝送容量の大きい回線を用いればこうした問題は解決できる可能性があるが、そうした回線を敷設するための経費の問題から現実的でなく、また将来、リアルタイムの動画像伝送を用いた通信が非常に増大することが予想されることから、回線容量の増大にあわせて伝送されるべき情報量の削減を図る必要がある。そのための一つの技術が、いわゆる画像符号化の技術である。

【0005】画像符号化技術の中でも、画像信号の情報量を圧縮する高効率符号化の技術が重要である。

【0006】そのために、たとえば画像信号の持つ冗長性を削減したり、人間の視覚の特性を利用して、画像信号のうちで画質に与える影響の少ない部分で情報を削減したりする方法が用いられる。

【0007】

【発明が解決しようとする課題】しかしながら、上記したような画像符号化技術を用いた場合、圧縮率を高くすると画質が落ちるという問題がある。そのため、限られた伝送容量の経路におけるリアルタイムの動画像通信ではこうした技術を用いることができない。

【0008】それ故にこの発明の目的は、限られた容量の媒体でも高画質で送信側の人物の顔の動画像を伝送または再生できる顔画像伝送方法、システム、ならびにそのための顔画像送信装置および顔画像再生装置を提供することである。

【0009】

【課題を解決するための手段】請求項 1 に記載の発明によれば、顔画像伝送方法は、話者の発する音声を受け、当該音声を発声するときの当該話者の表情を推定する信号を出力する表情推定手段を、受信側において準備するステップを備え、準備するステップは、送信側と受信側とにおいて同一のニューラルネットワークを準備するステップと、送信側において、話者の発する音声を受けて話者の顔画像を特定する情報を出力するようにニューラルネットワークを学習させるステップと、送信側の学習済みのニューラルネットワークの特性パラメータを受信側に送信し、受信側のニューラルネットワークを送信さ

れた特性パラメータにより設定するステップとを含み、話者の発する音声を送信側から受信側に送信し、表情推定手段に与えて話者の表情を推定する信号を得るステップと、表情推定手段の出力する話者の表情を推定する信号に基づいて、話者の表情の動画像を生成するステップとを含む。

【0010】音声信号のみを送信側から受信側に送信することにより、受信側では送信側にいる話者の表情の動画像を生成することができる。画像情報を伝送する場合と比較してはるかに小さい伝送容量の媒体を用いて、効率的に顔画像の伝送を行なうことができる。さらに、通信に先立って送信側のニューラルネットワークを学習させてその特性パラメータを受信側に送信することにより、受信側のニューラルネットワークを送信側のニューラルネットワークと同一に設定できる。少ない送信データ量で、顔の動画像のリアルタイム送信の準備を完了することができる。

【0011】

【0012】

【0013】請求項2に記載の発明によれば、請求項1に記載の発明の構成に加えて、学習させるステップは、送信側において、予め定められた学習のための文章を朗読する話者の顔の予め定められた箇所の座標を時系列的に測定するステップと、学習のための文章を朗読する際の話者の発する音声について、時系列的に音声の特徴量を求めるステップと、音声の特徴量を入力とし、測定された座標を教師信号として送信側のニューラルネットワークの特性パラメータを調整するステップとを含む。

【0014】請求項2に記載の発明によれば、請求項1に記載の発明の作用に加えて、送信側において、学習の文章を話者が朗読することにより得られる音声の特徴量と話者の顔の所定位置の座標とに基づいてニューラルネットワークの学習が行なわれる。話者ごとにニューラルネットワークの学習が行なわれるので、任意の話者について、リアルタイムの顔画像の送信を効率的に行なうことが可能となる。

【0015】請求項3に記載の発明によれば、顔画像伝送方法は、請求項2に記載の発明の構成に加えて、送信側のニューラルネットワークの学習後において、話者の発声時に、話者の顔の予め定められた箇所の座標を時系列的に測定するステップと、話者の発声する音声を送信側の学習済みのニューラルネットワークに与えて話者の顔の予め定められた箇所の座標の推定値を学習済みのニューラルネットワークの出力として得るステップと、測定された座標と、座標の推定値とを比較し、座標の推定値の有する誤差を求めるステップとをさらに含む。

【0016】請求項3に記載の発明によれば、請求項2に記載の発明の作用に加えて、送信側においてニューラルネットワークによって推定された話者の顔の所定位置の座標の推定値と、実際の話者の顔の所定位置の座標と

の誤差が求められるので、通信により受信側で生成される話者の顔画像が音声と合致したものとなるか否かの判定が可能となる。

【0017】請求項4に記載の発明によれば、請求項3に記載の発明の構成に加えて、顔画像伝送方法は、誤差を受信側に送信するステップと、受信側において誤差を受信し、受信側のニューラルネットワークの出力を受信された誤差に基づいて補正するステップとをさらに含む。

10 【0018】請求項4に記載の発明によれば、請求項3に記載の発明の作用に加えて、送信側で得られた誤差信号を受信側に送信することで、受信側において得られる画像をこの誤差信号を用いて補正することができる。その結果、音声信号とよく合致した話者の顔の動画像を生成することができる。

20 【0019】請求項5に記載の発明によれば、請求項3に記載の発明の構成に加えて、顔画像伝送方法は、誤差の大きさと所定のしきい値とを比較するステップと、誤差の大きさが所定のしきい値を超えたことに応答して、誤差を受信側に送信するステップと、受信側において誤差を受信し、受信側のニューラルネットワークの出力を受信された誤差に基づいて補正するステップとをさらに含む。

30 【0020】請求項5に記載の発明によれば、請求項3に記載の発明の作用に加えて、送信側で得られた誤差信号の大きさが所定のしきい値より大きくなったときに送信側で得られた誤差信号を受信側に送信することで、受信側において得られる画像をこの誤差信号を用いて補正することができる。その結果、音声信号とよく合致した話者の顔の動画像を生成することができ、また誤差信号を常に送信するわけではないので、音声信号の送信に支障が生ずるおそれも少ない。

40 【0021】請求項6に記載の発明によれば、送信側の話者の発声時の顔の動画像を受信側に送信するための顔画像伝送システムは、話者の発する音声信号を受信側に送信する送信装置を備え、送信装置は、話者の発する音声の所定の特徴量を時系列的に抽出するための特徴量抽出手段と、話者の発する音声に基づいて特徴量抽出手段が出力する特徴パラメータを入力とし、話者の顔画像を特定する情報を出力するように学習可能なニューラルネットワークと、ニューラルネットワークの特性パラメータを受信側に送信する手段とを含み、送信装置からの信号を受ける受信装置をさらに備え、受信装置は、送信装置から送信されてくる話者の発する音声信号を受け、当該音声が発声するときの当該話者の表情を推定する信号を出力する表情推定手段を含み、表情推定手段は、送信装置のニューラルネットワークと同一構成のニューラルネットワークと、受信装置のニューラルネットワークを送信装置から送信された特性パラメータにより設定する
50 ための手段とを有し、表情推定手段が送信側から話者の

発する音声信号を受信して出力する話者の表情を推定する信号に基づいて、話者の表情の動画像を生成する顔画像生成手段をさらに含む。

【0022】請求項6に記載の発明によれば、音声信号のみを送信側から受信側に送信することにより、受信側では送信側にいる話者の表情の動画像を生成することができる。画像情報を伝送する場合と比較してはるかに少ない伝送容量の媒体を用いて、効率的に顔画像の伝送を行なうことができる。さらに、通信に先立って送信側のニューラルネットワークを学習させてその特性パラメータを受信側に送信することにより、受信側のニューラルネットワークを送信側のニューラルネットワークと同一に設定できる。少ない送信データ量で、顔の動画像のリアルタイム送信の準備を完了することができる。

【0023】

【0024】

【0025】請求項7に記載の発明によれば、請求項6に記載の発明の構成に加えて、送信装置は、話者の発声時に、話者の顔の、予め定められた箇所の座標を時系列的に測定するための測定手段と、予め定められた学習のための文章を話者が朗読する際の、話者の顔の予め定められた箇所に関して測定手段によって得られた座標データを教師信号とし、特徴量抽出手段によって得られる音声パラメータを入力として送信装置のニューラルネットワークを学習させるための学習手段とをさらに含む。

【0026】請求項7に記載の発明によれば、請求項6に記載の発明の作用に加えて、学習のための文章を話者が朗読することにより、送信側において話者ごとにニューラルネットワークの学習を行なうことができるので、任意の話者について、リアルタイムの顔画像の効率的な送信が可能となる。

【0027】請求項8に記載の発明によれば、請求項7に記載の発明の構成に加えて、送信装置は、送信側のニューラルネットワークの学習後において、話者の発声する音声を送信側の学習済みのニューラルネットワークに与えて得られた話者の顔の予め定められた箇所の座標の推定値を測定された座標と比較し、座標の推定値の有する誤差を求めめるための手段をさらに含む。

【0028】請求項8に記載の発明によれば、請求項7に記載の発明の作用に加えて、送信側においてニューラルネットワークによって推定された話者の顔の所定位置の座標の推定値と、実際の話者の顔の所定の位置の座標位置との誤差が求められるので、通信により受信側で生成される話者の顔画像が音声と合致したものとなるか否かの判定が可能となる。

【0029】請求項9に記載の発明によれば、請求項8に記載の発明の構成に加えて、送信装置は、誤差を受信側に送信するための手段をさらに含み、受信装置は、誤差を受信して、受信側のニューラルネットワークの出力を受信された誤差に基づいて補正するための誤差補正手

段をさらに含む。

【0030】請求項9に記載の発明によれば、請求項8に記載の発明の作用に加えて、送信側で得られた誤差信号を受信側に送信することで、受信側において得られる画像をこの誤差信号を用いて補正することができる。その結果、音声信号とよく合致した話者の顔の動画像を生成することができる。

【0031】請求項10に記載の発明によれば、請求項8に記載の発明の構成に加えて、送信装置は、誤差と所定のしきい値とを比較し、誤差が所定のしきい値を超えたことに応答して、誤差を受信側に送信するための手段をさらに含み、受信装置は、誤差を受信して、受信側のニューラルネットワークの出力を受信された誤差に基づいて補正するための誤差補正手段をさらに含む。

【0032】請求項10に記載の発明によれば、請求項8に記載の発明の作用に加えて、送信側で得られた誤差信号の大きさが所定のしきい値より大きくなったときに送信側で得られた誤差信号を受信側に送信することで、受信側において得られる画像をこの誤差信号を用いて補正することができる。その結果、音声信号とよく合致した話者の顔の動画像を生成することができ、また誤差信号を常に送信するわけではないので、音声信号の送信に支障が生ずるおそれも少ない。

【0033】請求項11に記載の発明によれば、送信側の話者の発声時の顔の動画像を受信側に送信するための顔画像伝送システムで使用される顔画像送信装置は、話者の音声を受信装置に送信するための手段と、話者の発する音声の特徴量を入力とし、話者の顔画像を特定する情報を出力するように学習可能なニューラルネットワークと、このニューラルネットワークの特性パラメータを受信側に送信するための手段とを含む。

【0034】請求項11に記載の発明によれば、通信に先立って送信側のニューラルネットワークを学習させてその特性パラメータを受信側に送信することにより、受信側のニューラルネットワークを送信側のニューラルネットワークと同一に設定できる。そのため受信装置側でも音声信号を受けるだけでニューラルネットワークを用いて話者の表情を推定し顔の動画像を生成することができ、少ない送信データ量で、顔の動画像のリアルタイム送信をすることができる。

【0035】請求項12に記載の発明によれば、請求項11に記載の発明の構成に加えて、顔画像送信装置は、話者の顔の予め定められた箇所の座標を時系列的に測定するための測定手段と、話者の発する音声の所定の特徴量を時系列的に抽出するための特徴量抽出手段と、話者が予め定められた文章を朗読する際に特徴量抽出手段により求められた特徴量を入力とし、予め定められた文章を朗読する際に測定手段により測定された座標を教師信号としてニューラルネットワークを学習させるための手段とをさらに含む。

【0036】請求項12に記載の発明によれば、請求項11に記載の発明の作用に加えて、学習のための文章を話者が朗読することにより、話者ごとにニューラルネットワークの学習を行なうことができるので、任意の話者について、リアルタイムの顔画像の効率的な送信が可能となる。

【0037】請求項13に記載の発明によれば、請求項12に記載の発明の構成に加えて、ニューラルネットワークの学習後に、話者の発声する音声に対して特徴量抽出手段により出力された特徴量を、学習済みのニューラルネットワークに与えて得られた話者の顔の予め定められた箇所の座標の推定値を、測定手段によって測定された座標と比較し、座標の推定値の有する誤差を求めめるための手段をさらに含む。

【0038】請求項13に記載の発明によれば、請求項12に記載の発明の作用に加えて、送信側においてニューラルネットワークによって推定された話者の顔の所定位置の座標の推定値と、実際の話者の顔の所定位置の座標との誤差が求められるので、通信により受信側で生成される話者の顔画像が音声と合致したもとなるか否かの判定が可能となる。

【0039】請求項14に記載の発明によれば、請求項13に記載の発明の構成に加えて、送信装置は、誤差を受信側に送信するための手段をさらに含む。

【0040】請求項14に記載の発明によれば、請求項13に記載の発明の作用に加えて、送信側で得られた誤差信号を受信側に送信することで、受信側において得られる画像をこの誤差信号を用いて補正することができる。その結果、音声信号とよく合致した話者の顔の動画画像を生成することができる。

【0041】請求項15に記載の発明によれば、請求項13に記載の発明の構成に加えて、送信装置は、誤差の大きさと所定のしきい値とを比較し、誤差の大きさが所定のしきい値を超えたことに応答して、誤差を受信側に送信するための手段をさらに含む。

【0042】請求項15に記載の発明によれば、請求項13に記載の発明の作用に加えて、送信側で得られた誤差の大きさが所定のしきい値より大きくなったときに送信側で得られた誤差を受信側に送信することで、受信側において得られる画像をこの誤差を用いて補正することができる。その結果、音声信号とよく合致した話者の顔の動画画像を生成することができ、また誤差を常に送信するわけではないので、音声信号の送信に支障が生ずるおそれも少ない。

【0043】請求項16に記載の発明によれば、顔画像再生装置は、話者の発する音声から得られた所定の特徴

量を受け、当該音声を発声するときの当該話者の表情を推定する信号を出力する表情推定手段を備え、表情推定手段は、予め定められた構成のニューラルネットワークと、顔画像の再生に先立ってニューラルネットワークを設定するための特性パラメータを所定の信号源から受け、特性パラメータによってニューラルネットワークを設定するための手段とを含み、表情推定手段が特徴量を受けて出力する話者の表情を推定する信号に基づいて、話者の表情の動画画像を生成する顔画像生成手段をさらに備える。

【0044】請求項16に記載の発明によれば、音声信号のみを受けることにより、顔画像再生装置でその音声を発声した話者の表情の動画画像を生成することができる。画像情報を伝送したり記憶し再生する場合と比較してはるかに少ない容量の媒体を用いて、効率的に顔画像の再生を行なうことができる。さらに、予め話者の画像を再生するために必要な特性パラメータを顔画像再生装置に与えることにより、この特性パラメータによりニューラルネットワークをその音声に対応する顔画像を再生するように適切に設定できる。少ないデータ量で、顔の動画画像の再生の準備を完了することができる。

【0045】

【0046】

【0047】請求項17に記載の発明によれば、請求項16に記載の発明の構成に加えて、顔画像再生装置は、受信顔画像受信装置のニューラルネットワークの設定後に、所定の信号源から送信される誤差を受信して、顔画像再生装置のニューラルネットワークの出力を受信された誤差に基づいて補正するための誤差補正手段をさらに含む。

【0048】請求項17に記載の発明によれば、請求項16に記載の発明の作用に加えて、実際に測定された顔座標データとの誤差を顔画像再生装置に与えることで、顔画像再生装置において得られる画像をこの誤差信号を用いて補正することができる。その結果、音声信号とよく合致した話者の顔の動画画像を生成することができる。

【0049】

【発明の実施の形態】以下の実施の形態の説明において用いられる式を一括して下に掲げる。以後の説明では、各式に付された式番号を用いて各式を参照する。なお、以下の記載中において、電子出願で使用可能な記法の制約により、変数の上に付される記号(「~」など)を変数の前に付して表すことがある。

【0050】

【数1】

(7)

$$\mathbf{f}_m = [f_{1m} \ f_{2m} \ \dots \ f_{pm}]^t \quad (1)$$

$$\mathbf{F}_{tr} = [f_1 \ f_2 \ \dots \ f_{M_{tr}}] \quad (2)$$

$$\mathbf{x}_m = [x_{1m} \ x_{2m} \ \dots \ x_{3Nm}]^t \quad (3)$$

$$\mathbf{X}_{tr} = [x_1 \ x_2 \ \dots \ x_{M_{tr}}] \quad (4)$$

$$\mathbf{C}_{xx} = \frac{1}{M_{tr}} [\mathbf{X}_{tr} - \boldsymbol{\mu}_x] [\mathbf{X}_{tr} - \boldsymbol{\mu}_x]^t \quad (5)$$

$$\mathbf{C}_{xx} = \mathbf{U} \mathbf{S}_{xx} \mathbf{U}^t \quad (6)$$

$$\mathbf{x} = \mathbf{U}_x \mathbf{p}_x + \boldsymbol{\mu}_x \quad (7)$$

$$\mathbf{p}_x = \mathbf{U}_x^t (\mathbf{x} - \boldsymbol{\mu}_x) \quad (8)$$

$$p_k \approx \tilde{p}_k = w_k^2 [\tanh(\mathbf{W}_k^1 \mathbf{f} + b_k^1) + b_k^2], \quad k = 1, \dots, K \quad (9)$$

$$\mathbf{W}_k^1 = \begin{bmatrix} w_{11k}^1 & \dots & w_{1pk}^1 \\ \vdots & & \vdots \\ w_{Q1k}^1 & \dots & w_{Qpk}^1 \end{bmatrix} \quad (10)$$

$$b_k^1 = [b_{1k}^1 \ \dots \ b_{Qk}^1] \quad (11)$$

$$w_k^2 = [w_{1k}^2 \ \dots \ w_{Qk}^2] \quad (12)$$

【0051】図1を参照して、本願発明の実施の形態にかかる顔画像伝送システムの原理について説明する。図1に示すように、送信側の装置は、話者の音声を音声信号に変換するマイク38と、マイク38の音声から抽出された特徴量を入力信号とし、音声から逆に推定される話者の顔の所定位置の座標の推定値を出力するための、ニューラルネットワークからなる表情評価部58と、ニューラルネットワークの学習後、話者の発声する音声に基づいて表情評価部58によって推定された話者の顔の所定位置の座標の推定値と、発声時の話者の顔の所定位置について実際に測定された座標値とを比較して表情評価部58の出力の座標の推定値に含まれる誤差を評価するための誤差評価部60とを含む。誤差評価部60によって評価された誤差の大きさが所定のしきい値よりも大きな場合には、その誤差は受信側に送信される。なお表情評価部58に含まれるニューラルネットワークは、実際の処理に先立って、図示しないビデオカメラなどにより取得されるテスト文章を朗読する話者の表情を測定した測定顔画像26と、マイク38が出力する話者の音声信号とに基づいて予め学習が行なわれる。

【0052】実際の通信に先立って、表情評価部58を構成するニューラルネットワークを学習させておき、話

者の発する音声から、話者の顔の所定位置の座標の推定値を出力するように学習させておく必要がある。この学習の結果、ニューラルネットワークの各ノードの重み付けなどの特性パラメータが定められるが、この特性パラメータは実際の通信に先立って一度だけ受信側に送信される。

【0053】受信側の装置は、表情評価部58から送信されてきた特性パラメータ22にしたがって設定される、表情評価部58を構成するニューラルネットワークと同一構成のニューラルネットワークを含む表情評価部74を含む。表情評価部74は、一旦設定されると、送信側から送信されてくる音声信号を受けると、話者の顔の所定位置の座標の推定値を出力することになる。

【0054】受信側の装置はさらに、表情評価部74が出力する話者の顔の所定位置の座標の推定値に基づいて、話者の顔の動画像を生成しモニタ44に表示するための顔画像生成部78を含む。送信側から送信されてきた音声は受信側のスピーカ42によって再生され、このときに音声の時系列的な変換に対応して表情の変化する話者の顔画像がモニタ44に表示されることになる。

【0055】なお、顔画像生成部78が誤差評価部60から送信されてきた誤差信号24を表情評価部74の出

力に加算することにより、話者の顔の動画像がより忠実に再現される。

【0056】図2を参照して、このシステムについてより詳細に説明する。このシステムは、公衆回線網36を介して互いに接続される顔画像送信装置32および顔画像受信装置34を含む。以下の実施の形態のシステムは、顔画像送信装置32が音声信号を送信し、顔画像受信装置34がその音声信号に基づいて話者の顔画像を生成し表示するという形態であるが、顔画像受信装置34側に顔画像送信装置32と同様の構成を設け、顔画像送信装置32側には顔画像受信装置34と同様の構成を設けることにより、双方向の音声（および顔画像）の通信が行えることはいうまでもない。

【0057】顔画像送信装置32は、話者の音声を音声電気信号に変換するためのマイク38からアナログの音声電気信号を受けこれを標本化および量子化することによりデジタル信号に変換するためのA/D変換部50と、A/D変換部50によりデジタル化された音声信号から、音声の時系列的な特徴パラメータを抽出し行列化したデータに変換するための音声パラメータ行列生成部52と、音声パラメータ行列生成部52によって特徴パラメータの抽出されたデジタル化された音声信号に対して圧縮などの処理を行って顔画像受信装置34に対して送信処理するための音声信号送信部54と、ビデオカメラ40により撮像された話者の顔の画像から、顔の特定位置の座標を測定し出力するための表情の動き測定部64と、表情の動き測定部64が出力した顔の特定位置の座標データを行列化して出力するための表情パラメータ行列生成部62と、音声パラメータ行列生成部52の出力を入力とする、ニューラルネットワークからなる表情評価部58と、音声パラメータ行列生成部52の出力および表情パラメータ行列生成部62の出力を受け、音声パラメータ行列生成部52の出力に対する表情評価部58の出力と表情パラメータ行列生成部62から与えられる顔の特定位置の座標データとの誤差が最小となるように、誤差逆伝搬法により表情評価部58を学習させ、学習した結果設定された表情評価部58のニューラルネットワークの特性パラメータを顔画像受信装置34に対して送信する処理を行うためのパラメータ学習部56と、表情評価部58のニューラルネットワークの設定が済んだ後、表情評価部58から出力される話者の顔の所定位置の座標の推定値と、表情の動き測定部64によって測定された話者の顔の所定位置の座標の測定値との間の誤差を評価し、誤差の大きさが所定のしきい値より大きくなったときに当該誤差信号を顔画像受信装置34に対して送信する処理を行うための誤差評価部60と、音声信号送信部54、表情評価部58および誤差評価部60に接続され、これらから出力されるデジタルの信号を変調して公衆回線網36上に送出することにより顔画像受信装置34に送信するための送信部66とを含む。

【0058】顔画像受信装置34は、公衆回線網36から受信したデータを復調し、音声信号20と、特性パラメータ22と、誤差信号24とに分離して出力するための受信部68と、受信部68が受信した音声信号20を受け、圧縮された音声信号を再生する処理などを行うための音声信号受信・再生部70と、音声信号受信・再生部70によって再生されたデジタルの音声信号をアナログ信号に変換してスピーカ42に与えるためのD/A変換部72と、受信部68を介して顔画像送信装置32から受信された特性パラメータによって設定されるニューラルネットワークからなり、音声信号受信・再生部70から与えられる音声信号を入力として、話者の顔の特定位置の座標の推定値を出力するための表情評価部74と、表情評価部74の出力に、受信部68を介して誤差評価部60から受けた誤差信号を加算して話者の顔の所定位置の座標の推定値を補正するための誤差加算部76と、誤差加算部76によって補正された話者の顔の所定位置の座標の推定値に基づいて、話者の顔画像を生成しモニタ44に出力して表示させるための顔画像生成部78とを含む。

【0059】図3を参照して、表情評価部74は、音声信号受信・再生部70の出力するオーディオ信号102に対してLPC（線形予測）処理を行ないLSP（線スペクトル対）パラメータを出力するためのLPC分析部104と、LPC分析部104の出力するLSPパラメータを入力として、顔座標ベクトルのPCA（主成分分析）表現を出力するための、ニューラルネットワークからなる非線形推定部106と、非線形推定部106の出力するPCA表現に対して逆PCA変換を行なって顔位置情報110を出力するための逆PCA変換部108とを含む。

【0060】図4を参照して、顔画像生成部78は、表情評価部74から出力され誤差加算部76によって補正が行なわれた顔位置情報110に基づいて線形予測を行ない顔画像を形成するメッシュ画像のPCA表現を出力するための線形推定部112と、線形推定部112の出力する顔のメッシュ画像のPCA表現に対して逆PCA処理を行ない顔の所定位置の座標表現を得るための逆PCA変換部114と、逆PCA変換部114の出力に基づいて形成された顔画像の表面に人の肌のテクスチャをマッピングして顔画像118を出力するためのテクスチャマッピング部116とを含む。

【0061】図5を参照して、顔画像送信装置32において行われる、表情評価部58のニューラルネットワークの学習時の処理の流れについて説明する。この実施の形態のシステムでは、通信処理を開始する前に、予め準備された定型の文章（テスト文）を話者が朗読し、そのときの音声信号と話者の顔の所定位置の座標値とからニューラルネットワークの学習を行う。このとき、音声についての処理と話者の顔画像についての処理とが平行し

て行われるが、以下では説明の都合上、まず音声の処理について述べ、次に顔画像の処理について述べる。

【0062】最初に、テスト音声をもとのサンプリングレートでサンプリング（標準化および量子化）する（80）。サンプリング値について、所定の長さのフレームを単位としてハミングウィンドウによる処理を施し、各フレームごとにLPC分析を行なう（82）。LPC分析によって得られた係数をLSP表現に変換する（84）。こうしてLSP表現に変換されたパラメータを行列表現Fに変換する（86）。

【0063】なお、ここではLSP表現を用いているが、これはLSPパラメータが人間の発声のホルマントと密接な関係があり、かつホルマントは人間が母音を発声する際の声道の形の共鳴周波数であり、声道の形が人間の顔の表情をかなりの程度支配しているという事実に基づいている。

【0064】この行列化処理についてより具体的に説明すると以下のとおりである。デジタル化された音声信号のm番目のフレームについて、LPC分析（P次とする）により抽出されたパラメータは、式（1）に示されるP次元ベクトル f_m として表現される。なお式（1）および他の式において、ベクトルまたは行列の右肩に付された「t」は転置ベクトルまたは転置行列を意味する。

【0065】テスト文（フレーム数 = Mtr）についてこうして得られたベクトル f_m （ $m = 1 \sim Mtr$ ）をマトリクス形式にまとめ、式（2）に示される行列Ftrを得る。この行列Ftrが、ステップ86の出力である。この行列Ftrはニューラルネットワークの学習処理（98）への入力として用いられる。

【0066】一方、上記した音声の処理に平行して、テスト文を朗読する際の話者の顔の画像から、顔の所定位置の座標データが測定されサンプリングされる（90）。ここでは顔の画像データから画像認識技術を用いて測定を行なうが、話者の顔の所定位置に直接位置センサを固定してその位置の座標データを測定するようにしてもよい。ここでも、データはフレーム化されて処理される。話者の顔の所定位置としてどのような箇所を選択したかについて、例を図6および図7に示す。図6に示す例では話者の顔のうち、あご、口周辺およびほほを中心とした12箇所が選択された。図7に示す例では同じく18箇所が選択された。

【0067】こうして得られた顔の所定位置のm番目のフレームの座標値のベクトル（顔座標ベクトルとよぶ） x_m は式（3）のように表される。式（3）でNは、測定対象となる顔の所定位置の数を示す。ベクトル x_m は3N次元となる。

【0068】こうして得られた座標位置のベクトルを行列化すると式（4）に示されるXtrのようになる。式（4）においても、Mtrはテスト文の朗読の際のサン

リングデータのフレーム数である。

【0069】こうして得られた行列Xtrに対して、主成分分析（PCA）のために固有値分解処理（94）を行ないさらに主成分を決定する（96）。主成分分析とは、統計的パターン認識において、特徴抽出後のパターンが多次元の特徴ベクトルで表現される場合に、これらの多次元の特徴ベクトルに一般的に含まれる冗長性を削減するために、統計的な手法を用いて次元を低減する手法の一つである。主成分分析では、パターンの分布の分散の大きな低次元部分空間でパターンを表現する。主成分分析は少ないデータ量でなるべく二乗誤差が少なくなるように元の情報を近似しようとする手法であり、主成分として抽出された固有ベクトルの一次結合で各特徴ベクトルが表される。

【0070】主成分の決定手法は以下のとおりである。まず、行列Xtrの共分散行列Cxxを式（5）にしたがって計算する。なお式（5）において μ_x は平均顔座標ベクトルである。平均顔座標ベクトルを行列Xtrから減算することによりベクトルの正規化が行なわれる。

【0071】共分散行列Cxxを式（6）のように表現すべく、固有値分解が行なわれる。式（6）でUはCxxの固有ベクトルを各列に有するユニタリー行列である。式（6）のSxxは、Cxxの固有値を対角線に有する対角行列である。

【0072】ここで、全ての固有値の和は、Cxxに関連して観測される分散の総計に等しい。そこで、固有値のうち最初のK個の和が、全固有値の和のうち所定の割合（たとえば99%）以上となる場合、Cxxのその最初のK個の固有値（Cxxの最初のK列に含まれる）の和は、テスト文に基づいて得られた学習用データセットの全分散のうち、前記した割合（たとえば99%）に等しくなる。したがって、任意のベクトルxは、Cxxの最初のK個の固有ベクトルの一次結合として非常に良く近似できる。この最初のK個の固有ベクトルを行列XtrのK主成分と呼ぶ。

【0073】本願発明者の研究によれば、上記した割合として99%を選べば十分であり、8種の顔形状を用いた場合、これに対応するK = 7である。このUの固有ベクトルのうち最初の7個を最初の7列に有する行列をUxと表せば、所与の顔座標ベクトルxと主成分とは式（7）によって関係付けられる。ただしPxは、式（8）で表される主成分係数ベクトルである。また、特に話者の表情の細かな動きまで再現する必要がなければ、K = 2または3でも十分な動きを表現することが可能である。

【0074】こうして得られた主成分は、ステップ86で得られたLSPパラメータの行列Fとともにニューラルネットワークの学習（98）に用いられる。

【0075】上記したKに対応して、表情評価部58は7個の別個のニューラルネットワークを含み、これら全

てに対してステップ 9 8 で学習処理が行なわれる。この実施の形態のシステムで用いられるニューラルネットワークはフィードフォワード型であり、いずれも非線形の一つの隠れ層と、一つの出力層とを含み、L S P パラメータと顔座標位置との関係を示す。7 つのニューラルネットワークは、ステップ 8 4 で得られた L S P ベクトル f を入力として、前記した主成分係数ベクトル Px の 7 つの成分 p_k ($k = 1 \sim 7$) をそれぞれ導出するためのものである。なお、ニューラルネットワークの層の数は 2 に限定されず、3 以上の層状構造を有するニューラルネットワークを用いても良い。また、フィードフォワード型に限らず、自己の出力を入力にフィードバックする機構を持つことにより時系列的なリアルタイム処理に適したいわゆるリカレントニューラルネットワークを用いてもよい。

【0076】各ニューラルネットワーク ($1 \sim K$) は、式 (9) によって前記したベクトル f の関数として成分 p_k を導出する。なお以下の説明では、一般化するために主成分の数を K として説明する。ここで、各ニューラルネットワークの隠れ層は式 (10) に示す重み行列 W_k^1 と、バイアスベクトル b_k^1 とによって定義される。また出力層は式 (12) に示される重みベクトル w_k^2 とバイアス係数 b_k^2 とによって定義される。また Q は各ニューラルネットワークの隠れ層に含まれるニューロンの数を示し、本実施の形態の装置では $Q = 10$ である。

【0077】この 7 個のニューラルネットワークの学習 (9 8) は、本実施の形態ではレヴェンバーク - マーカート (Levenberg-Marquardt) 法を用いて行なわれた。もちろん、レヴェンバーク - マーカート法以外にも種々の変数問題の最適化法を用いることができる。たとえばニュートン法 (Newton's method)、準ニュートン法 (quasi Newton's method)、共役方向法 (conjugate direction method) などがある。

【0078】ステップ 9 8 の結果、ニューラルネットワークの振舞いを決定する特性パラメータ (重み係数など) が各ニューラルネットワークごとに得られる (100)。これら特性パラメータを表情評価部 7 4 の各ニューラルネットワークに与えて各ニューラルネットワークを設定することにより、表情評価部 7 4 は表情評価部 5 8 と同じ動作を行なうことになる。すなわち、顔画像受信装置 3 4 の表情評価部 7 4 は、顔画像送信装置 3 2 においてテスト文を用いてトレーニングされた表情評価部 5 8 のニューラルネットワークと同一の振舞いを示すことになる。

【0079】原理的には、こうして表情評価部 7 4 のニューラルネットワークを設定することにより、音声信号受信・再生部 7 0 から話者の音声の L S P パラメータを表情評価部 7 4 に与えると、対応の主成分係数ベクトルの推定値 $\sim px$ が得られ、一旦推定値ベクトル $= px$ が得られると式 (7) によって顔の所定位置の座標の推定値

ベクトル px が得られる。

【0080】以上説明したこの実施の形態のシステムは以下のように動作する。図 2 を参照して、通信を開始するに先立って、まず所定のテスト文を話者が読み上げる。マイク 3 8 によってオーディオ信号に変換された音声は、A / D 変換部 5 0 によってデジタル化され、音声パラメータ行列生成部 5 2 によってフレームごとに L S P パラメータとしてベクトル化された上で行列表現 F に変換される。音声パラメータ行列生成部 5 2 で生成された音声信号の行列表現データはパラメータ学習部 5 6 に与えられる。

【0081】一方、上記した音声の処理と平行して、テスト文を朗読したときの話者の顔の表情がビデオカメラ 4 0 により撮影され表情の動き測定部 6 4 に与えられる。表情の動き測定部 6 4 は、画像認識により顔の所定位置の座標を測定し音声パラメータ行列生成部 5 2 に与える。音声パラメータ行列生成部 5 2 は、この座標データをこれもフレームごとにベクトル化し、行列表現 X に変換して P C A 分析をして結果をパラメータ学習部 5 6 に与える。

【0082】パラメータ学習部 5 6 は、音声パラメータ行列生成部 5 2 から与えられた音声の L S P 表現を入力とし、表情パラメータ行列生成部 6 2 から与えられた P C A 成分を教師信号として表情評価部 5 8 中の 7 個のニューラルネットワークの学習を行なう。学習の後、パラメータ学習部 5 6 は表情評価部 5 8 内の 7 個のニューラルネットワークの重みなどの特性パラメータを送信部 6 6 および公衆回線網 3 6 を介して顔画像受信装置 3 4 に送る。

【0083】顔画像受信装置 3 4 では、表情評価部 7 4 は受信部 6 8 が受信したこの特性パラメータによって、表情評価部 7 4 中の、表情評価部 5 8 に含まれる 7 個のニューラルネットワークと同じ構成の 7 個のニューラルネットワークを設定する。この処理により、表情評価部 5 8 と表情評価部 7 4 とは同じ動作を行なうこととなる。

【0084】さて、このようにして顔画像送信装置 3 2 の表情評価部 5 8 内のニューラルネットワークの学習および顔画像受信装置 3 4 の表情評価部 7 4 内のニューラルネットワークの設定が終了すると、次のようにして実際の通信が行なわれる。

【0085】話者は、顔画像受信装置 3 4 に向けて伝達すべき事項をマイク 3 8 に向けて発声する。ここでも音声パラメータ行列生成部 5 2 は学習時と同様に動作し、音声から得られた L S P パラメータを生成して表情評価部 5 8 に与える。表情評価部 5 8 は、入力された L S P パラメータに基づいて話者の顔座標を推定し顔座標データを出力する。この出力は誤差の評価に用いられる。一方、音声パラメータ行列生成部 5 2 が出力する L S P パラメータおよびもとのデジタル音声信号は音声信号送信

部 5 4 を介して顔画像受信装置 3 4 に送信される。

【0086】この L S P パラメータおよびもとのデジタル音声信号は受信部 6 8 を介して音声信号受信・再生部 7 0 に与えられる。音声信号受信・再生部 7 0 は、デジタル音声信号を D / A 変換部 7 2 に与え、D / A 変換部 7 2 によってデジタルからアナログ信号に変換された音声信号はスピーカ 4 2 に与えられてもとの音声再生される。

【0087】一方、音声信号受信・再生部 7 0 が受けた L S P パラメータは表情評価部 7 4 に与えられる。表情評価部 7 4 内のニューラルネットワークは、既に学習済みの表情評価部 5 8 内のニューラルネットワークと同じように設定されているため、入力された L S P パラメータから顔座標データを推定し出力する。誤差が小さい場合であれば、このようにして推定され表情評価部 7 4 から出力された顔座標データに基づいて顔画像生成部 7 8 で顔画像を生成しモニタ 4 4 上に話者の顔画像を表示する。話者の顔画像は、音声信号から得られているため、音声の再生と同期して、リアルタイムで話者の顔の動画画像を合成表示することができる。

【0088】一方、顔画像送信装置 3 2 の表情の動き測定部 6 4 はテスト時と同じように話者の顔の所定位置の座標を測定している。誤差評価部 6 0 は、こうして実際に測定された話者の顔座標データから、表情評価部 5 8 が出力する顔座標データを減算し誤差信号を出力する。この誤差信号の絶対値は誤差評価部 6 0 によって所定のしきい値と比較され、誤差信号の絶対値がしきい値より大きければ誤差評価部 6 0 は誤差信号を顔画像受信装置 3 4 に送信する。

【0089】顔画像受信装置 3 4 の受信部 6 8 は、誤差信号を受信した場合にこの誤差信号を誤差加算部 7 6 に与える。誤差加算部 7 6 は、表情評価部 7 4 の出力する顔座標データにこの誤差信号を加算して顔画像生成部 7 8 に与える。

【0090】この補正により、顔画像生成部 7 8 に与えられる顔座標データと、実際に顔画像送信装置 3 2 で測定された顔座標データとの間の誤差が前記したしきい値以上となることが防止できる。

【0091】図 8 に、この発明を適用した実験システムにおいて、米国人の話者についてある文章 ("When the sunlight strikes raindrops in the air, ...") を朗読したときに実際に測定された話者の顔面各部の動き (濃い細い実線) と、このときの音声信号からニューラルネットワークを用いて推定された話者の顔面の対応する各部の動き (濃い太い実線) とを時系列的に示す。比較のために、各グラフにはニューラルネットワークではなく、アフィン変換などの線形推定法により推定された結果 (色合いの薄い実線) をも示す。

【0092】図 8 において、最上段には測定された音声信号が、2 段目 ~ 1 1 段目にはそれぞれほぼ、口元 (唇

の端部)、上唇、下唇、アゴの、それぞれの垂直方向と水平方向とについて得られた結果が、それぞれ示されている。図 8 において、横軸は時間軸を表し、縦軸は移動距離 (cm) を表している。

【0093】図 8 の各グラフの各段に付した数字のうち、下段の数字は、本発明によりニューラルネットワークを用いて得られた顔座標データと、実際に測定された顔座標データとの相関係数を示す。上段の数字は、線形推定法により得られた顔座標データと、実際に測定された顔座標データとの相関係数を示す。

【0094】同じく、図 9 に、この発明を適用した実験システムにおいて、日本人の話者についてある文章 (「桃を割ってみると中から男の子が出てきました。」) を朗読したときの結果を示す。

【0095】図 8 および図 9 からよく分かるように、ニューラルネットワークを用いると、実際に測定された話者の顔の動きとよく一致した結果が得られる。また、送信側でこのように誤差を得た場合、前述のようにその誤差を表す誤差信号を受信側に送信し、受信側でこの誤差信号を用いてニューラルネットワークの出力を補正することができる。その結果、実際に測定された話者の顔の動きと同じ顔の動画画像を受信側で生成することが可能となる。この誤差信号の送信は、上述した実施の形態におけるように所定のしきい値を超えたときに行なうようにすれば、音声信号の送信に支障が生じないので好ましい。しかし、回線状況が許せば常に誤差信号を送信し常に受信側で補正を行なうようにしてもよい。

【0096】以上のようにこの実施の形態のシステムでは、通信に先立ってニューラルネットワークの学習を行ない、ニューラルネットワークの特性パラメータを受信側に送信すれば、以後はほとんど音声信号のみの送信をすることで、話者の顔のリアルタイムの動画画像を受信側で合成表示することができる。最初の特性パラメータのデータ量も多くはないので、電話回線など、伝送容量の非常に少ない回線でも快適に顔の動画画像のリアルタイム送信を行なうことができる。また、誤差信号を送るようにすると、合成された顔の動画画像と音声との間のずれも最小限にすることができるという効果を奏する。ただし、誤差信号の送信は必ずしも必要ではなく、最初のパラメータの送信およびそのあとの音声信号の送信のみでも十分な品質の通信を行なうことができる。

【0097】なおまた、上記の説明では各通信の前に必ず表情評価部 5 8 内のニューラルネットワークの学習が必要であったが、話者が特定の一人に限定されている場合には表情評価部 5 8 の学習を毎回行なう必要は無く、前回の通信時の設定をそのまま用いることができる。また、顔画像受信装置 3 4 側が特定の顔画像送信装置 3 2 の特定の話者としか通信を行なわない場合には、表情評価部 7 4 の設定として前回の設定を利用できるので、表情評価部 7 4 のニューラルネットワークの特性パラメー

タを顔画像送信装置 3 2 から顔画像受信装置 3 4 側に送信する必要もない。

【0098】なお、上記した顔画像送信装置および顔画像再生装置がいずれも汎用のコンピュータと、その上で実行されるプログラム、およびマイク、スピーカ、モニタ、モデム、ターミナルアダプタなどの一般的な周辺機器とによって実現可能であることはいうまでもない。

【0099】今回開示された実施の形態はすべての点で例示であって制限的なものではないと考えられるべきである。本発明の範囲は上記した説明ではなくて特許請求の範囲によって示され、特許請求の範囲と均等の意味および範囲内でのすべての変更が含まれることが意図される。

【図面の簡単な説明】

【図1】 本願発明の実施の形態にかかる顔画像伝送方法およびシステムの原理を示すブロック図である。

【図2】 図1に示される顔画像伝送システムのより詳細なブロック図である。

【図3】 受信装置側における表情評価部のブロック図である。

*【図4】 受信装置側における顔画像生成部のブロック図である。

【図5】 顔画像送信装置で行われる処理の流れを示すフローチャートである。

【図6】 第1の被験者における顔の座標の測定位置を示す図である。

【図7】 第2の被験者における顔の座標の測定位置を示す図である。

【図8】 実験システムで米国人の話者について得られた結果を示す図である。

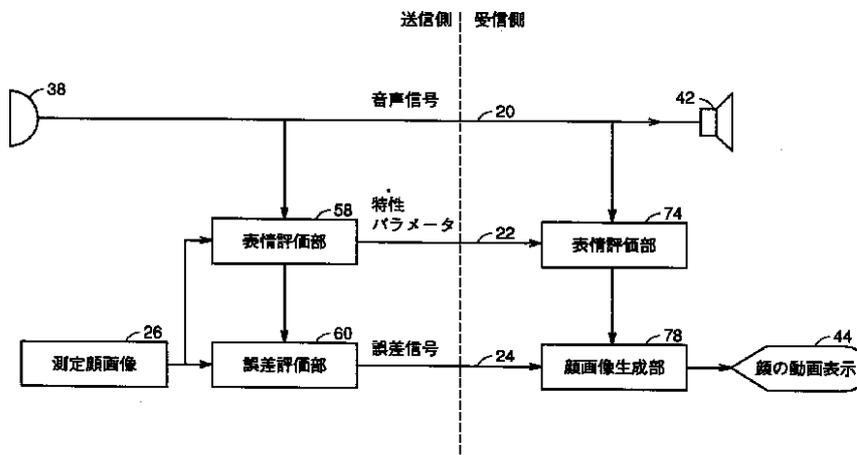
【図9】 実験システムで日本人の話者について得られた結果を示す図である。

【符号の説明】

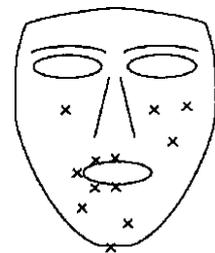
20 音声信号、22 特性パラメータ、24 誤差信号、32 顔画像送信装置、34 顔画像受信装置、38 マイク、52 行列生成部、56 パラメータ学習部、58 表情評価部、60 誤差評価部、62 行列生成部、64 表情の動き測定部、74 表情評価部、76 誤差加算部、78 顔画像生成部。

* 20

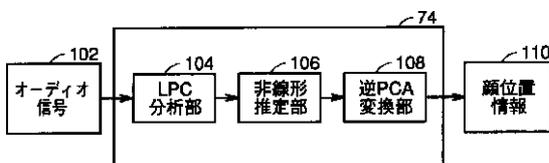
【図1】



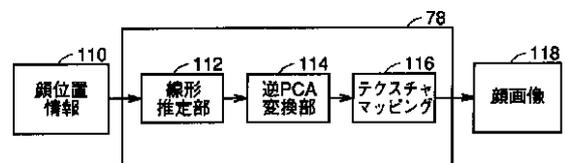
【図6】



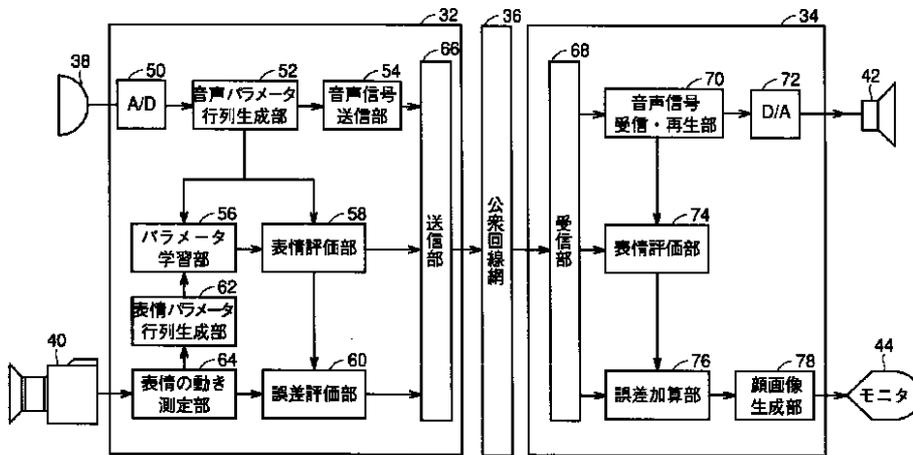
【図3】



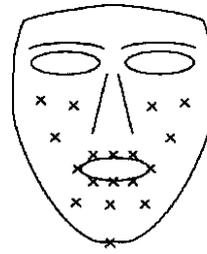
【図4】



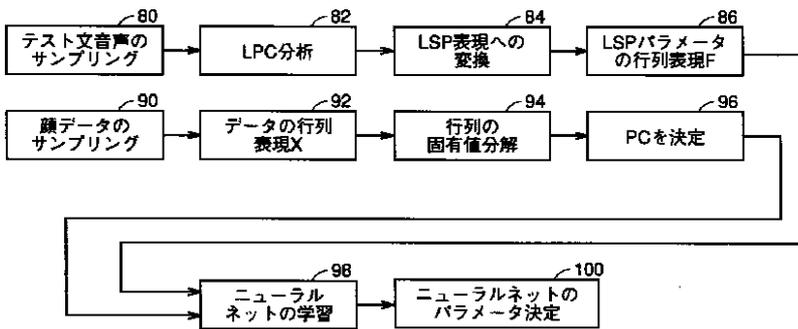
【図2】



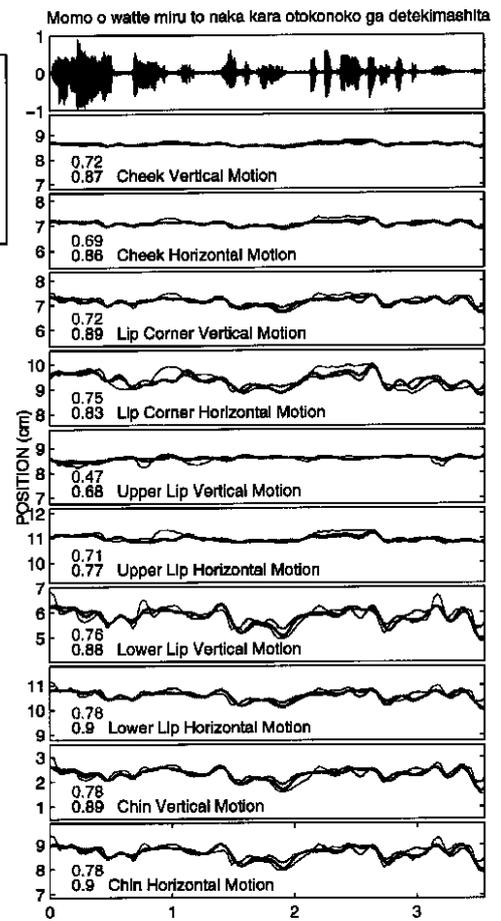
【図7】



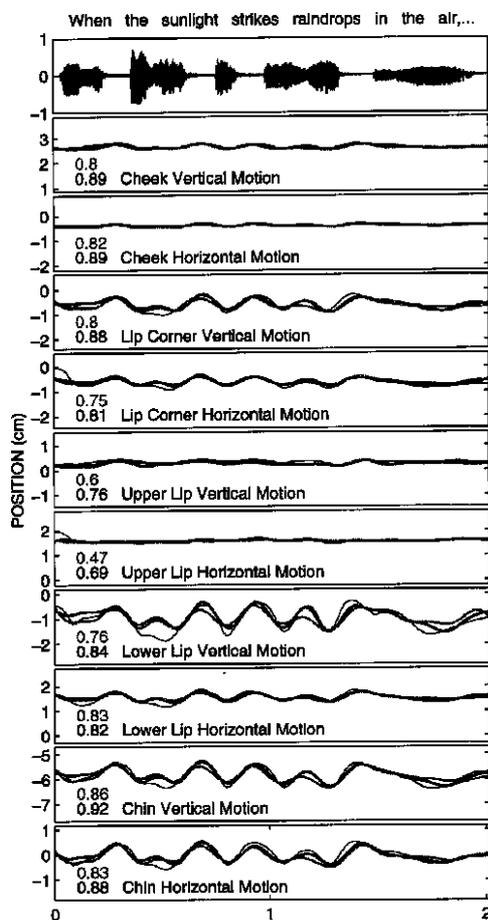
【図5】



【図9】



【図 8】



フロントページの続き

(51) Int. Cl.⁷

H 0 4 N 7/24

識別記号

F I

H 0 4 N 7/13

Z

(72) 発明者

エリック・パティキオティス・ベイツン
京都府相楽郡精華町大字乾谷小字三平谷
5 番地 株式会社エイ・ティ・アール人
間情報通信研究所内

(56) 参考文献

特開 平11 - 21942 (J P , A)

特開2001 - 34787 (J P , A)

特開 平10 - 49707 (J P , A)

特開2000 - 11200 (J P , A)

特開 平 5 - 128261 (J P , A)

今村達也、外 2 名、音声からの感情推定と実時間メディア変換システム、1999 年電子情報通信学会総合大会講演論文集基礎・境界、日本、電子情報通信学会、1999年 3 月 8 日、A14 - 4、305

四倉達夫 外 2 名、" サイバースペース上の仮想人物による実時間対話システムの構築 "、情報処理学会論文誌、社団法人情報処理学会、1999年 2 月15日、第40巻、第 2 号、p . 677 - 686

(58)調査した分野(Int.Cl.⁷, D B名)

G06T 11/20

G06N 3/00

G10L 15/00

G10L 15/16

H04M 11/06

H04N 7/24

C S D B (日本国特許庁)