

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5055486号
(P5055486)

(45) 発行日 平成24年10月24日(2012.10.24)

(24) 登録日 平成24年8月10日(2012.8.10)

(51) Int. Cl. F I
G 1 0 L 15/00 (2006.01) G 1 0 L 15/00 2 0 0 H
G 1 0 L 15/28 (2006.01) G 1 0 L 15/28 5 0 0

請求項の数 4 (全 15 頁)

(21) 出願番号	特願2006-197112 (P2006-197112)	(73) 特許権者	393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2
(22) 出願日	平成18年7月19日(2006.7.19)	(74) 代理人	100090181 弁理士 山田 義人
(65) 公開番号	特開2008-26485 (P2008-26485A)	(72) 発明者	石井 カルロス 寿憲 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
(43) 公開日	平成20年2月7日(2008.2.7)	(72) 発明者	西尾 修一 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21年6月17日(2009.6.17)	(72) 発明者	石黒 浩 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】 遠隔操作アンドロイドの発話動作制御システム

(57) 【特許請求の範囲】

【請求項1】

操作者によって遠隔操作されるアンドロイドの発話動作を前記操作者の発話音声に応じて制御するためのシステムであって、

前記操作者の発する音声の音声データを取得する音声取得手段、

取得した音声データの音響特徴から当該音声データを発音するための口唇形状を非線形モデルを用いて推定する口唇形状推定手段、

前記口唇形状を形成するための動作指令を発行してから当該口唇形状が形成されるまでにかかる時間情報を示す動作遅延を推定する動作遅延推定手段、

音声データの再生を、前記音声データを取得してから、動作遅延推定手段が推定する最大動作遅延の値よりも大きい値の再生遅延時間経過後に開始する音声再生手段、

前記音声再生手段による再生開始タイミングを基準として、当該再生開始タイミングより前にかつ前記動作遅延に基づいて前記動作指令の発行タイミングを設定する動作指令設定手段、および

前記動作指令設定手段によって設定された発行タイミングに従って各動作指令を発行する動作指令発行手段を備える、システム。

【請求項2】

前記口唇形状推定手段は、前記音響特徴の変動量が閾値を超えたかどうか判断する判断手段、および前記判断手段によって前記音響特徴の変動量が閾値を超えたと判断したときその時点の前後の音声データを抽出する抽出手段を含み、抽出した音声データに基づいて

口唇形状を推定する、請求項 1 記載のシステム。

【請求項 3】

前記口唇形状推定手段によって推定された所定時間の区間の口唇形状の時系列に基づいて当該区間を通じた動作の簡略化を含む最適化を行う最適化手段をさらに備える、請求項 1 または 2 記載のシステム。

【請求項 4】

前記最適化手段による最適化を行った後、最適化した動作に基づいて前記動作遅延を再度推定する再推定手段をさらに備える、請求項 3 記載のシステム。

【発明の詳細な説明】

【技術分野】

10

【0001】

この発明は遠隔操作アンドロイドの発話動作制御システムに関し、特にたとえば、操作者の発話音声に基づいてアンドロイドの発話時の口唇動作を制御するシステムに関する。

【背景技術】

【0002】

一般に、音声合成に用いる音素情報を利用して口唇形状を出力する手法は有効であるが、音素情報が無い場合には、発話音声の音響特徴から口唇形状を推定する必要が生じる。発話音声のみを元に人の口唇の動きを再現する技術は、たとえば非特許文献 1 - 4 で提案されている。音声のみから口唇の形状を推定する手法は、音響特徴抽出と、音響特徴から口唇形状へのマッピングの 2 ステップに分けられる。音声信号から抽出される音響特徴としては、LPC (linear predictive coding) - Cepstral 係数、MFCC (Mel-Frequency Cepstral Coefficients) 係数、LSP (Line Spectrum Pair) 係数などが挙げられる。音響特徴から口唇形状へのマッピング手法としては、線形回帰分析、ニューラル・ネットワーク、HMM (Hidden Markov Model)、KNN (K-Nearest Neighbor) などが挙げられる。上記の手法では、音声と口唇形状の画像的な情報との 1 対 1 のマッピングが定期のフレームごとに行われる。

20

【非特許文献 1】Lavagetto, F., "Converting speech into lip movements: A multimedia telephone for hard of hearing people", IEEE Trans. on Rehabilitation Engineering, Vol. 3, No. 1, pp. 90-102, March 1995.

【非特許文献 2】Yehia, H., Rubin, P., Vatikiotis-Bateson, E., "Quantitative association of vocal-tract and facial behavior", Speech Communication, vol. 26, no. 1-2, pp. 23-43, 1998.

30

【非特許文献 3】Hong, P., Wen, Z., Huang, T., "Real-time speech-driven face animation with expressions using neural networks", IEEE Trans. on Neural Networks, Vol. 13, No. 4, pp. 916-927, July 2002.

【非特許文献 4】Gutierrez-Osuna, R., Kakumanu, P., et al., "Speech-driven facial animation with realistic dynamics", IEEE Trans. on Multimedia, Vol. 7, No. 1, pp. 33-41, Feb. 2005.

【発明の開示】

【発明が解決しようとする課題】

40

【0003】

しかしながら、音響特徴と口唇の形状との間には非線形な関係が存在することから、的確に動作するものは未だ存在しない。実物のロボットまたはアンドロイド (人間に酷似した姿形を有し、また、人間に酷似した動作を行うロボット) の制御においては、アクチュエータの応答は動作によって変わってしまう。たとえば、口を開く動作は、閉じる動作に比べて遅い。そのため、上記関連技術のように 1 フレームごとに音響特徴とターゲットとした口唇形状との 1 対 1 のマッピングをするのは困難となる。

【0004】

また、以下のような問題も存在する。発声においては、口唇や舌等の形が定まってから空気が通ることによって音が発せられる。無声破裂音 (/p/, /t/, /k/) などでは、音

50

に先立って口唇や舌の形が定まっている。したがって、音を元に口唇等を動かしたのでは、自然でない動作になってしまう。つまり、音を元に口唇を主とした顔面動作の制御を行う過程には、必然的に遅延が伴う。

【 0 0 0 5 】

また、アンドロイドの顔面は、人の顔面とは動き方が異なる。これは、人の顔面筋肉と同様にアクチュエータを内在させることが現状の技術では不可能であるためである。また、反応速度等も異なるため、人とまったく同一の動きをさせることはできない。一般には、人では可能であってもアンドロイドでは不可能な動きが数多く存在する。

【 0 0 0 6 】

それゆえに、この発明の主たる目的は、アンドロイドにおいて操作者の音声に適合した口唇動作を行うことができる、遠隔操作アンドロイドの発話動作制御システムを提供することである。

【課題を解決するための手段】

【 0 0 0 7 】

請求項 1 の発明は、操作者によって遠隔操作されるアンドロイドの発話動作を操作者の発話音声に応じて制御するためのシステムであって、操作者の発する音声の音声データを取得する音声取得手段、取得した音声データの音響特徴から当該音声データを発音するための口唇形状を非線形モデルを用いて推定する口唇形状推定手段、口唇形状を形成するための動作指令を発行してから当該口唇形状が形成されるまでにかかる時間情報を示す動作遅延を推定する動作遅延推定手段、音声データの再生を、音声データを取得してから、動作遅延推定手段が推定する最大動作遅延の値よりも大きい値の再生遅延時間経過後に開始する音声再生手段、音声再生手段による再生開始タイミングを基準として、当該再生開始タイミングより前であつ動作遅延に基づいて動作指令の発行タイミングを設定する動作指令設定手段、および動作指令設定手段によって設定された発行タイミングに従って各動作指令を発行する動作指令発行手段を備える、システムである。

【 0 0 0 8 】

請求項 1 の発明では、システム（ 1 0 : 後述する実施例で相当する参照符号。以下同じ。）は、遠隔操作されるアンドロイド（ 1 2 ）の発話動作、すなわち発話時の口唇を主とした顔面動作を、操作者の発話音声に応じて制御するためのものである。システムは、たとえば、遠隔操作端末（ 2 0 ）とアンドロイドの制御装置（ 1 4 ）を含み、以下のような各手段を備えている。操作者の発話した音声は、音声取得手段（ 2 4 、 7 0 、 7 4 、 7 8 、 5 3 ）によって取得される。音声再生手段（ 7 0 、 8 0 、 4 8 、 5 7 、 5 9 ）は、取得された音声を一定時間の遅延のもとに再生する。口唇形状推定手段（ 7 0 、 5 5 、 5 2 1 - 5 2 7 ）は、取得された音声の音響特徴から口唇形状を非線形モデルを用いて推定する。たとえば、音響特徴としては M F C C 係数が適用され、当該音響特徴が高い変動量を示す時点の前後所定時間の音響特徴を用いて、口唇形状の推定が行われる。また、音声信号と口唇形状との間には非線形な関係が存在するため、ニューラル・ネットワークのような非線形モデルが使用される。動作遅延推定手段（ 7 0 、 5 3 1 ）は、口唇形状を形成するための動作指令を発行してから実際にアンドロイドのアクチュエータ（ 4 0 ）が駆動して顔面において当該口唇形状が形成されるまでにかかる時間情報を示す動作遅延を推定する。動作指令設定手段（ 7 0 、 5 3 7 ）は、推定された動作遅延に基づいて、一定遅延で再生される発話音声の再生開始タイミングを基準として、当該動作指令の発行タイミングを設定する。具体的には、或る音声を発する際に形成される口唇形状のための動作指令の発行タイミングは、動作遅延を考慮して、当該形状に対応する音声が出力されるタイミングよりも早められる。動作指令発行手段（ 7 0 、 5 3 9 ）は、設定された発行タイミングに従って、各動作指令を発行する。

【 0 0 0 9 】

したがって、請求項 1 の発明によれば、アンドロイドのアクチュエータが駆動されてその顔面において実際に口唇形状が形成されてから、当該口唇形状に対応する音声が出力されることとなる。このように、遠隔操作アンドロイドにおいて、操作者の音声に適合した

10

20

30

40

50

口唇動作を実現することができる。

【0010】

請求項2の発明は、口唇形状推定手段は、音響特徴の変動量が閾値を超えたかどうか判断する判断手段、および判断手段によって音響特徴の変動量が閾値を超えたと判断したときその時点の前後の音声データを抽出する抽出手段を含み、抽出した音声データに基づいて口唇形状を推定する、請求項1記載のシステムである。

請求項3の発明は、口唇形状推定手段によって推定された所定時間の区間の口唇形状の時系列に基づいて当該区間を通じた動作の簡略化を含む最適化を行う最適化手段をさらに備える、請求項1または2記載のシステムである。

請求項4の発明は、最適化手段による最適化を行った後、最適化した動作に基づいて動作遅延を再度推定する再推定手段をさらに備える、請求項3記載のシステムである。

10

【0011】

請求項3の発明では、最適化手段(70、535)は、所定区間の口唇形状の時系列に基づいて、当該区間における動作の再構成が行われて、当該区間を通じた動作の最適化が行われる。所定区間は、たとえば、複数の音素や単語単位が含まれる程度の時間に設定される。具体的には、当該区間を通して、動作の簡略化、重要音素のみでの提示、動作量の相対化、動作速度の変換などのような変換が試みられる。そして、最適化された口唇動作に含まれる各動作の指令の発行タイミングが設定されて、当該最適化された口唇動作が提示されることとなる。したがって、口唇動作をより自然に見せることや、動作をより素早く行うことなどが可能になる。

20

【発明の効果】

【0012】

この発明によれば、発話音声の再生開始タイミングを基準として、動作指令発行から当該口唇形状を形成するまでにかかる動作遅延を考慮してアクチュエータの動作指令を発行するようにしたので、操作者の発話音声に適合した口唇動作を行うことができる。

【0013】

この発明の上述の目的、その他の目的、特徴および利点は、図面を参照して行う以下の実施例の詳細な説明から一層明らかとなろう。

【発明を実施するための最良の形態】

【0014】

図1を参照して、この実施例のアンドロイド制御システム(以下、単に「システム」という。)10は、アンドロイド12を含む。アンドロイド12は、人間に酷似した姿形(外観など)を有する人型ロボットであり、人間に酷似した動作(振り、振る舞い、発話)を行う。アンドロイド12は、制御装置14に接続され、この制御装置14には、環境センサ16が接続される。

30

【0015】

また、制御装置14は、インターネットや電話通信回線のようなネットワーク18を介して遠隔操作端末20に接続される。遠隔操作端末20は、PC或いはPDAのような汎用のコンピュータであり、この遠隔操作端末20には、スピーカ22、マイク24およびモニタ26が接続される。図示は省略するが、当然のことながら、遠隔操作端末20には、キーボードおよびコンピュータマウスのような入力装置が含まれる。また、遠隔操作端末20の動作を制御するためのプログラムおよびデータは、その図示しないメモリに記憶されており、図示しないCPUによって遠隔操作端末20の全体的な動作が制御される。

40

【0016】

図2は、アンドロイド12、制御装置14および環境センサ16の電気的な構成を示すブロック図である。この図2を参照して、アンドロイド12は複数のアクチュエータ(たとえば、エアアクチュエータ)40を含み、各アクチュエータ40は、制御装置14(14b)の制御ボード(アクチュエータ制御ボード)98に接続される。アクチュエータ40は、アンドロイド12の身体動作を提示するために設けられる。いくつかのアクチュエータ40はアンドロイド12の口唇を主とした顔面動作を行うために用いられる。たとえ

50

ば、上顎、下顎、上唇、下唇、口の左右の側面（頬）などの動きのためのアクチュエータ 40 がそれぞれ設けられる。また、他のアクチュエータ 40 は、アンドロイド 12 の各関節や眼および腹部などその他の部位または部分を動かすために用いられる。ただし、簡単のため、この実施例では、コンプレッサは省略してある。

【0017】

また、アンドロイド 12 は、触覚センサ 42、眼カメラ 44、衝突センサ 46、スピーカ 48 およびマイク 50 を含む。この触覚センサ 42、眼カメラ 44 および衝突センサ 46 は、制御装置 14（14b）のセンサ入出力ボード 100 に接続される。

【0018】

触覚センサ 42 ないし皮膚センサは、たとえばタッチセンサであり、アンドロイド 12 の触覚の一部を構成する。つまり、触覚センサ 42 は、人間や他の物体等がアンドロイド 12 に触れたか否かを検出するために用いられる。触覚センサ 42 からの出力（検出データ）は、センサ入出力ボード 100 を介して MPU 90 に与えられる。MPU 90 は、触覚センサ 42 からの検出データを第 1 コンピュータ 14a の CPU 70 に送信する。したがって、CPU 70 は、人間や他の物体等がアンドロイド 12 に触れたことを検出することができる。ただし、触覚センサ 42 としては、圧力センサを用いることもできる。かかる場合には、人間や他の物体等がアンドロイド 12 の肌に触れたか否かのみならず、その触れ方（強弱）を知ることができる。

【0019】

眼カメラ 44 は、イメージセンサであり、アンドロイド 12 の視覚の一部を構成する。つまり、眼カメラ 44 は、アンドロイド 12 の眼から見た映像ないし画像を検出するために用いられる。この実施例では、眼カメラ 44 の撮影映像（動画ないし静止画）に対応するデータ（画像データ）は、センサ入出力ボード 100 を介して MPU 90 に与えられる。MPU 90 は、画像データを第 1 コンピュータ 14a の CPU 70 に送信し、CPU 70 は、撮影映像の変化を検出するのみならず、その画像データを、ネットワーク 18 を介して遠隔操作端末 20 に送信する。そして、遠隔操作端末 20 は、受信した画像データをモニタ 26 に出力する。したがって、眼カメラ 44 の撮影映像がモニタ 26 に表示される。

【0020】

衝突センサ 46 は、人間や他の物体等がアンドロイド 12 に衝突したか否かを判断する。衝突センサ 46 の出力（検出データ）は、センサ入出力ボード 100 を介して MPU 90 に与えられる。MPU 90 は、衝突センサ 46 からの検出データを第 1 コンピュータ 14a の CPU 70 に送信する。したがって、CPU 70 は、人間や他の物体等がアンドロイド 12 に衝突したことを検出することができる。

【0021】

また、スピーカ 48 およびマイク 50 は、制御装置 14（14a）の音声入出力ボード 80 に接続される。スピーカ 48 は、アンドロイド 12 が発話を行う際に音声を出力する。遠隔操作端末 20 の操作者ないしオペレータ（以下、「遠隔オペレータ」という。）が直接発話を行う場合、当該音声出力される。具体的には、遠隔オペレータがマイク 24 を通して発話すると、対応する音声データが遠隔操作端末 20 からネットワーク 18 を介して制御装置 14a（CPU 70）に与えられる。そして、CPU 70 は、その音声データを、音声入出力ボード 80 を介してスピーカ 48 から出力する。なお、予めプログラミングされた所定の動作を行う場合には、スピーカ 48 からは合成音声出力される。

【0022】

マイク 50 は、音センサであり、アンドロイド 12 の聴覚の一部を構成する。このマイク 50 は、指向性を有し、主として、アンドロイド 12 と対話（コミュニケーション）する人間（ユーザ）の音声を検出するために用いられる。

【0023】

制御装置 14 は、第 1 コンピュータ 14a および第 2 コンピュータ 14b によって構成される。たとえば、第 1 コンピュータ 14a がメインのコンピュータであり、第 2 コンピ

10

20

30

40

50

ユーザ 14 b がサブのコンピュータである。なお、この実施例では、制御装置 14 を 2 台のコンピュータ (14 a, 14 b) で構成するようにしてあるが、処理能力が高ければ、1 台のコンピュータで構成することもできる。

【0024】

第 1 コンピュータ 14 a は、CPU 70 を含み、CPU 70 には内部バス 72 を介してメモリ 74、通信ボード 76、LAN ボード 78、音声入出力ボード 80 およびセンサ入出力ボード 82 が接続される。メモリ 74 は、たとえば、ハードディスク装置 (HDD) のような主記憶装置、ROM および RAM を含む。詳細な説明は省略するが、このメモリ 74 には、制御装置 14 の全体の動作を制御するためのプログラムおよびデータが記憶されており、特にたとえば、アンドロイド 12 の動作についてのコマンド名に対応して、そのコマンド名が示す動作を実行するための制御情報が記憶されている。

10

【0025】

ここで、動作とは、振り、振る舞いのような身体動作および発話動作をいう。したがって、この実施例では、制御情報は、アクチュエータ 40 を駆動制御するための制御データのみならず、必要に応じて、発話内容についての合成音声の音声データを含む。ただし、身体動作には、自然の動作 (無意識動作) も含まれる。無意識動作の代表的な例としては、瞬きや呼吸が該当する。また、このような生理的な動作のみならず、人間の癖による動作も無意識動作に含まれる。たとえば、癖による動作としては、たとえば、髪の毛を触る動作、顔を触る動作や爪を噛む動作などが該当する。

【0026】

20

このような動作は、アンドロイド 12 が外部 (環境) からの刺激に対応して実行されたり、遠隔オペレータからの命令 (遠隔操作命令) に従って実行されたりする。

【0027】

図 2 に戻って、通信ボード 76 は、他のコンピュータ (この実施例では、第 2 コンピュータ 14 b) とデータ通信するためのインターフェイスである。たとえば、通信ボード 76 は、後述する第 2 コンピュータ 14 b の通信ボード 96 と、RS 232C のようなケーブル (図示せず) を用いて接続される。LAN ボード 78 は、ネットワーク 18 を介して他のコンピュータ (この実施例では、遠隔操作端末 20) とデータ通信するためのインターフェイスである。この実施例では、LAN ボード 78 は、LAN ケーブル (図示せず) を用いて接続される。

30

【0028】

なお、この実施例では、各コンピュータがケーブルを用いた有線のネットワークを構成するように説明してあるが、これに限定される必要はなく、無線のネットワークを構成するようにしてもよく、有線と無線とが混在するネットワークを構成するようにしてもよい。

【0029】

音声入出力ボード 80 は、音声を入力および出力するためのインターフェイスであり、上述したように、アンドロイド 12 のスピーカ 48 およびマイク 50 が接続される。この音声入出力ボード 80 は、CPU 70 によって与えられた音声データを音声信号に変換して、スピーカ 48 に出力する。また、音声入出力ボード 80 は、マイク 50 を通して入力された音声信号を音声データに変換して、CPU 70 に与える。

40

【0030】

なお、詳細な説明は省略するが、制御装置 14 a は、音声認識機能を備える。したがって、たとえば、人間がアンドロイド 12 に対して発話した内容はマイク 50 を通して入力されると、CPU 70 は、DP マッチングや隠れマルコフ法により、人間の発話内容を音声認識するのである。ただし、音声認識用の辞書データはメモリ 74 に記憶されているものとする。

【0031】

センサ入出力ボード 82 は、各種センサからの出力を CPU 70 に与え、制御データを各種センサに出力するためのインターフェイスである。この実施例では、センサ入出力ボ

50

ード 82 には、環境センサ 16 が接続され、環境センサ 16 は、全方位カメラ 60、PTZカメラ 62 およびフロアセンサ 64 を含む。

【0032】

全方位カメラ 60 は、アンドロイド 12 が配置される部屋ないし場所（区画）に設置され、当該部屋ないし場所の全体（360度）を撮影することができる。全方位カメラ 60 の撮影映像に対応する画像データは、CPU70 に与えられる。CPU70 は、画像データに基づいて撮影映像の変化を検出するのみならず、その画像データを遠隔操作端末 20 に送信する。遠隔操作端末 20 は、受信した画像データをモニタ 26 に出力する。したがって、全方位カメラ 60 の撮影映像がモニタ 26 に表示される。

【0033】

PTZカメラ 62 は、オペレータの指示に従って、パン（Pan）、チルト（Tilt）およびズーム（Zoom）の各々を制御（調整）することができるカメラである。たとえば、遠隔オペレータが、パン、チルト、ズームの指示を入力すると、対応する制御信号（コマンド）が遠隔操作端末 20 からネットワーク 18 を介して第 1 コンピュータ 14a に与えられる。すると、第 1 コンピュータ 14a の CPU70 は、そのコマンドに従って PTZカメラ 62 を駆動制御する。PTZカメラ 62 の撮影映像に対応する画像データもまた、CPU70 に与えられる。CPU70 は、画像データに基づいて撮影映像の変化を検出するのみならず、その画像データを遠隔操作端末 20 に送信する。遠隔操作端末 20 は、受信した画像データをモニタ 26 に出力する。したがって、PTZカメラ 62 の撮影映像もモニタ 26 に表示される。

【0034】

なお、この実施例では、眼カメラ 44、全方位カメラ 60 および PTZカメラ 62 の撮影画像が、遠隔操作端末 20 に接続されるモニタ 26 に画面を分割されて表示される。したがって、遠隔オペレータはモニタ 26 を見て、アンドロイド 12 の視線の映像やアンドロイド 12 の周囲の状況を知ることができる。

【0035】

フロアセンサ 64 ないし床圧力センサは、図示は省略するが、多数の検出素子（感圧センサ）を含み、この多数の検出素子はアンドロイド 12 が配置される部屋ないし場所の床に埋め込まれる（敷き詰められる）。したがって、フロアセンサ 64 からの出力に基づいて、アンドロイド 12 の周囲に人間が存在するか否か、存在する人間の人数、アンドロイド 12 から見た人間の方向、アンドロイド 12 と人間との距離などを知ることができる。

【0036】

第 2 コンピュータ 14b は、MPU90 を含み、MPU90 には、内部バス 92 を介して、メモリ 94、通信ボード 96、制御ボード 98 およびセンサ入出力ボード 100 が接続される。

【0037】

メモリ 94 は、HDD、ROM および RAM を含み、メモリ 94 には、制御装置 14b の動作を制御するためのプログラムおよびデータが記憶されている。通信ボード 96 は、他のコンピュータ（この実施例では、第 1 コンピュータ 14a）とデータ通信するためのインターフェイスである。制御ボード 98 は、制御対象としての複数のアクチュエータ 40 を制御するための制御データを出力するとともに、各アクチュエータ 40 からの角度情報（回転角度、曲げ角度）を入力するためのインターフェイスである。したがって、MPU90 は複数のアクチュエータ 40 をフィードバック制御することができる。ただし、MPU90 は、第 1 コンピュータ 14a の CPU70 からの動作指令（制御データ）に従って各アクチュエータ 40 を駆動制御する。

【0038】

センサ入出力ボード 100 は、各種センサからの出力を MPU90 に与え、MPU90 からの制御データを各種センサに出力するためのインターフェイスである。このセンサ入出力ボード 100 には、上述したように、触覚センサ 42、眼カメラ 44 および衝突センサ 46 が接続される。

10

20

30

40

50

【 0 0 3 9 】

なお、この実施例では、アンドロイド 1 2 とは別に制御装置 1 4 を設けるようにしてあるが、制御装置 1 4 を含めてアンドロイドと呼んでもよい。さらには、環境センサ 1 6 も含めてアンドロイドと呼んでもよい。

【 0 0 4 0 】

また、アンドロイド 1 2 の遠隔地にオペレータが存在することを想定して制御装置 1 4 と遠隔操作端末 2 0 とをネットワーク 1 8 を介して接続しているが、オペレータがアンドロイド 1 2 の存在する場所やその近傍に存在するような場合には、遠隔操作端末 2 0 をネットワーク 1 8 を介さずに制御装置 1 4 に直接接続することも可能である。

【 0 0 4 1 】

アンドロイド 1 2 は、或る会社や或るイベント会場などに配置され、人間の代役として働くことができる。たとえば、アンドロイド 1 2 は、会社やイベント会場の受付や案内役として機能し、訪問者に対応する。アンドロイド 1 2 は、アンドロイド 1 2 や環境センサ 1 6 のセンサ群 (4 2、4 4、4 6、5 0、6 0、6 2、6 4) によって検出される外部刺激に応じて所定の対応動作を実行する。状況に応じてアンドロイド 1 2 が実行すべき (実行可能な) 動作、すなわち CPU 7 0 が指示するべき制御情報は予め決定されており、その内容はメモリ 7 4 に記憶されている。

【 0 0 4 2 】

ただし、センサ群の反応に応じた対応動作がメモリ 7 4 に記憶されていない場合や外部からの刺激を認識できない場合などのように、アンドロイド 1 2 が刺激に対して対応動作を実行することができない場合には、アンドロイド 1 2 は遠隔オペレータを呼び出す。つまり、制御装置 1 4 (CPU 7 0) は、アンドロイド 1 2 が対応できない旨の情報を遠隔操作端末 2 0 に通知する。たとえば、アンドロイド 1 2 が人間の質問 (発話内容) を理解 (音声認識) できない場合には、呼び出された遠隔オペレータが人間の質問を理解し、遠隔操作によって、アンドロイド 1 2 の動作を制御する。

【 0 0 4 3 】

また、上述したように、眼カメラ 4 4、全方位カメラ 6 0 および PTZ カメラ 6 2 の撮影映像が遠隔操作端末 2 0 のモニタ 2 6 に表示されるため、遠隔オペレータは、その撮影映像等によりアンドロイド 1 2 が存在する近傍、周囲およびアンドロイド 1 2 と対話する人間の様子を知ることができる。したがって、遠隔オペレータは、アンドロイド 1 2 (制御装置 1 4) からの呼び出しが無くても、必要に応じて遠隔操作し、アンドロイド 1 2 に命令動作を実行させることもできる。

【 0 0 4 4 】

遠隔オペレータは、所定の合成音声の出力を指示する代わりに、マイク 2 4 に発話することによって、自身の発話音声をアンドロイド 1 2 のスピーカ 4 8 から出力して、人間と直接対話することができる。つまり、遠隔オペレータが発話すると、マイク 2 4 で検出された音声入力に対応する音声データが取得され、当該音声データが遠隔操作端末 2 0 からネットワーク 1 8 を介して制御装置 1 4 に与えられ、アンドロイド 1 2 のスピーカ 4 8 から出力される。

【 0 0 4 5 】

なお、アンドロイド 1 2 と対話する人間が発話したときには、その音声はアンドロイド 1 2 のマイク 5 0 を通して入力され、対応する音声データが制御装置 1 4 からネットワーク 1 8 を介して遠隔操作端末 2 0 に送信され、スピーカ 2 2 から出力される。

【 0 0 4 6 】

アンドロイド 1 2 は、人間に酷似した姿形を有して人間の動作に酷似した動作を行うロボットであるから、遠隔オペレータの発話音声を出力する際に、たとえば口を動かさなかったり単に音声に関係なく口を動かしたりするだけでは人間に強い違和感を与えてしまう。したがって、このアンドロイド 1 2 では、出力される遠隔オペレータの発話音声に合わせてその口唇を主とした顔面を動作させる。

【 0 0 4 7 】

このシステム10の動作を図3および図4に示すフロー図を参照しながら説明する。図3には、制御装置14の発話処理の動作の一例が示される。制御装置14aのCPU70は、この発話処理を一定時間ごとに繰り返し実行する。

【0048】

図3のステップS1では、音声データを受信したか否かを判断する。遠隔オペレータが発話したとき、遠隔操作端末20からマイク24で取得された発話音声の音声データが送信されてくるので、この音声データをネットワーク18を介して受信したか否かが判断される。なお、遠隔操作端末20は、発話音声を所定のサンプリングレート（たとえば8kHz）で音声データとして取得し、取得した音声データを所定の packets 長（たとえば20ms）で一定時間ごとに送信する。

10

【0049】

ステップS1で“YES”であれば、ステップS3で、音声記憶処理を開始する。音声記憶処理はCPU70によって他の処理と並行的に実行される。この音声記憶処理によって、受信される音声データが順次メモリ74に記憶される。音声記憶処理は、発話音声を検出されなくなって音声データが受信されなくなったときに終了される。

【0050】

続いて、ステップS5で、口唇動作制御処理を開始する。口唇動作制御処理はCPU70によって他の処理と並行的に実行される。この口唇動作制御処理では、取得された発話音声の解析が行われて、当該音声に基づいて口唇動作が制御される。口唇動作制御処理の動作の一例は後述する図4に示される。

20

【0051】

ステップS7では、音声取得から一定時間経過したか否かを判断する。この実施例では、取得した発話音声を一定量の遅延のもとに再生するようにしているので、この判定によって、音声データの取得（受信）から一定時間の経過を待つ。

【0052】

ステップS7で“YES”であれば、ステップS9で、音声再生処理を開始する。音声再生処理はCPU70によって他の処理と並行的に実行される。この音声再生処理では、取得された音声データがメモリ74から読み出されて音声入出力ボード80に与えられ、これによって、アンドロイド12のスピーカ48から当該音声が出力される。音声再生処理は、取得した音声データをすべて再生し終わったときに終了される。

30

【0053】

なお、ステップS1で“NO”の場合、つまり、発話が行われていないときには、そのまま図3の発話処理を終了する。

【0054】

ステップS5で開始される口唇動作制御処理の動作の一例を図4を参照して説明する。まず、ステップS21で、音響特徴の変動量を抽出する。

【0055】

アンドロイド12のような物体の場合、画像のようにフレームごとに口唇形状を制御することは困難である。従って、まず、遠隔オペレータの音声の周波数やケプストラムの解析を行い、音響特徴の変動が高い位置を検出する。音響特徴の変動量は、たとえば、ある時刻における前後所定時間（たとえば20ms程度）のフレームのパラメータ（たとえばMFCC）の平均二乗誤差として算出される。なお、音声信号から取得する音響特徴としては、LPC-Cepstral係数、MFCC係数、LSP係数、フォルマント周波数、およびF0（基本周波数）、RMS（Root Mean Square）などが挙げられる。フォルマント周波数と口唇形状には、母音の場合、直感的な関係があるが、精度のよい自動検出は難しいので、この実施例では、音声認識で多く使用されるMFCC係数を用いる。

40

【0056】

次に、ステップS23で、この変動量（MFCC平均二乗誤差など）が閾値を超えたか否かを判断する。実験によって、音素の変化を表す程度に、この変動量に閾値を設定しておく。閾値を超えた変動量のピーク位置がアンドロイド12の動作指令発行時点を決める

50

際の基礎となる。

【 0 0 5 7 】

ステップ S 2 3 で “ N O ” の場合、処理はステップ S 2 1 へ戻り、次の時刻を基点とする音声データについて処理を繰り返す。

【 0 0 5 8 】

一方、ステップ S 2 3 で “ Y E S ” の場合、ステップ S 2 5 で、音響特徴の高い変動量が検出された時点の前後所定時間（たとえば 1 0 0 m s 程度）の音声から音響特徴（たとえば M F C C ）を抽出し、ステップ S 2 7 で、非線形モデルを用いて口唇形状の推定を行う。推定の手法として、線形回帰分析、ニューラル・ネットワーク、HMM、KNNなどが挙げられる。また、一度音素情報を認識して、その音素における口唇形状をマッピングする手法があるが、音素認識の精度はあまり高くないので、音響特徴から口唇形状の直接マッピングの方が効率がよいと考えられる。また、音響特徴と口唇形状の間には非線形な関係があるので、ニューラル・ネットワークのような非線形なモデルを用いる。なお、そのためには、予め収録したビデオデータまたはモーションキャプチャによる口唇形状のデータベースによってモデル学習を行っておき、メモリ 7 4 にモデル学習による非線形マッピングのための情報を記憶しておく。

10

【 0 0 5 9 】

続いて、ステップ S 2 9 で、推定された口唇形状を形成するための制御情報を設定し、ステップ S 3 1 で動作遅延を推定する。具体的には、アンドロイド 1 2 のアクチュエータ 4 0 の制御情報に関しては、アクチュエータ制御の静的特徴と動的特徴を考慮する。つまり、静的特徴としては、特定の口唇形状に近づけるためのアンドロイド 1 2 のアクチュエータ 4 0 の制御情報を予め手動的に取得しておき、口唇形状と制御情報とを対応付けたデータベースをメモリ 7 4 に記憶しておく。また、動的特徴としては、特定の形状をターゲットとして口唇を動かした際に、指令を発行した時点からアンドロイド 1 2 が実際にターゲットの形状に辿りつくまでにかかる時間（これを動作遅延と呼ぶ。）を実験により取得しておき、制御情報（口唇形状）と動作遅延とを対応付けたデータベースをメモリ 7 4 に記憶しておく。後述のステップ S 3 7 では、この動作遅延の情報を基に、音声と同期を取るために、動作指令を送る時点が早められたり遅くされたりする。

20

【 0 0 6 0 】

なお、この動作遅延の情報を基にして、音声再生開始までの遅延時間が決められる。つまり、上述のように、この実施例では、音声は常に一定遅延で再生されるようにするので、この再生遅延時間を、最大の動作遅延の値よりも大きい値に設定しておく。

30

【 0 0 6 1 】

ステップ S 3 3 では、所定時間の推定を行ったか否かを判断する。この実施例では、音響特徴を抽出した範囲よりも広い範囲、たとえば複数の音素や単語単位で、口唇動作の再構成をすることを想定しているので、このステップ S 3 3 の判定を行う。ステップ S 3 3 で “ N O ” の場合、ステップ S 2 1 に戻って処理を繰り返す。

【 0 0 6 2 】

ステップ S 3 3 で “ Y E S ” であれば、ステップ S 3 5 で、上述の所定時間の区間を通じて口唇動作の最適化処理を行う。つまり、比較的短い期間の音声に関して、ステップ S 2 1 やステップ S 2 5 の処理を行い、これらの音声を束ねたより長い区間を通じて動作の最適化を試みる。推定された口唇形状は完全にはアンドロイド 1 2 では再現できない場合もあるため、推定された口唇形状の時系列を元に、この口唇動作の変換を行う。たとえば、以下のような変換が考えられる。

40

（ a ）動作の簡略化：一部の部位のみを動かす。動作に応じて使用される部位が異なり、また、音や動作によって特に強調されるべき部位が異なるので、省略可能な部位の動作は省略する。たとえば、音や動作に応じて、顎の関節、唇の前後方向の動き、または口の側面の動きなど、最も特徴的な動きのみになるように変換する。

（ b ）重要音素のみでの提示：イントネーションなどを考慮し、音素系列のうち強調され目立つ動作のみを行う。たとえば、韻律特徴（ピッチや強さなど）を用いて、ピッチが高

50

く、パワーが強く、長めに発声された音節を強調されたものとみなす。

(c) 動作量の相対化：動きの存在だけを提示し、動きの量は追従させないようにする。たとえば、指定されたアクチュエータ40の動きの大きさを元々の指定通りに設定するのではなく、より少ない動きで同等な効果が得られると判断される場合には、少しだけの動きに留めるように動作量を変更する。これによって、より自然な動作を見せることができ、素早く動作を行うことができる。また、動きを目立たせるため、時間的に連続する同じ部位の動作を抑制または増幅する。たとえば、動作の大きさを1 - 5段階で表す(5 = 大きい)場合、単一のアクチュエータ40について、3 5 3という動作系列が実行されるが「3」の部分はそれほど重要ではないとき、1 3 1や0 3 0などに置き換える。相対的な動きの度合いのみが、見る者には重要であることが多いため、このような置換を実行しても、同じように見せることが可能になる。このように、前後の動作も含めて一定期間の動作に抑制または増幅を施すことによって、動作を素早く行うことが可能になるとともに、場合によっては、より強調したい部分を見せるなど、より強い効果を見る者に与えることができる。

10

(d) 動作速度の変換：特定の音声が発せられる時点で、口唇形状はその最大変化を示している方がより自然に見えるため、素早くその目的形状を形成する。

【0063】

なお、この最適化処理では、遅延なども考慮して、可能な限り簡素な動作に変換するようにするのが望ましいが、より自然に見せるためには多からず少なからぬ口唇動作が必要であり、実験によってパラメータを適宜に設定する。なお、発話された音声によっては最適化が行われない場合もあり得る。

20

【0064】

また、必要に応じて、ステップS35では、最適化を行った後に、ステップS31の各動作の動作遅延の推定をやり直して、最適化された動作に適切な遅延を取得するようにしてよい。特に動作の増幅が行われる場合、動作遅延の再推定は必要である。また、動作の抑制や速度変換が行われる場合には、動作遅延の再推定によって、特定形状への動作完了から当該音声の発話までをより素早くまたはより適切なタイミングで行うことができる。

【0065】

なお、他の実施例では、ステップS31の処理は、ステップS35の後に実行するようにしてもよい。また、その他の実施例では、ステップS33およびステップS35の処理は省略されてもよい。つまり、推定された口唇形状をそのまま提示するようにしてもよい。

30

【0066】

続いて、ステップS37で、動作遅延に基づいて、音声再生開始タイミングを基準として、各動作指令の発行タイミングを設定する。つまり、特定の口唇形状を形成するための動作指令の発行タイミングは、当該音声との同期をとるために、当該推定遅延に基づいて音声再生開始タイミングを基準として設定される。基本的には、動作指令の発行タイミングは、当該口唇形状に対応する音声の再生時点よりも早められる。発話を終えて口を閉じる動作の場合はこの限りではなく、遅くされる場合もあり得る。

【0067】

40

そして、ステップS39で、動作指令発行処理を開始する。動作指令発行処理はCPU70によって他の処理と並行的に実行される。この動作指令発行処理では、各動作指令の発行タイミングになったと判断されたときに、当該動作指令が発行される。具体的には、CPU70は、アクチュエータ制御情報を含む動作指令を通信ボード76に出力して制御装置14bのMPU90に与える。これに応じて、MPU90は制御ボード98に制御情報を与え、これによって、制御情報において指定されたアクチュエータ40が駆動され、アンドロイド12の顔面において、制御情報に対応する口唇形状が形成されることとなる。上述のように、発行タイミングには各口唇形状を形成する際の動作遅延が考慮されているので、当該口唇形状が実際にアンドロイド12の顔面において形成されてから、当該口唇形状に対応する音声出力される。また、音声の出力が終わってから当該音声に対応す

50

る口唇形状が変えられる。なお、動作指令発行処理は、すべての動作指令の発行が完了したと判断されたときに終了される。

【0068】

ステップS41では、未処理の音声データが残っているか否かを判断し、“YES”であれば、ステップS21に戻って処理を繰り返す。このようにして、アンドロイド12においては、遠隔オペレータの発話音声に適合した口唇動作を伴って出力される。一方、ステップS41で“NO”であれば、この口唇動作制御処理を終了する。

【0069】

この実施例によれば、遠隔オペレータの発話音声の音響特徴から非線形モデルを用いて口唇形状を推定し、発話音声の再生開始タイミングを基準として当該口唇形状を形成するまでにかかる動作遅延を考慮してアクチュエータ40の動作指令の発行タイミングを設定するようにしたので、アンドロイド12において遠隔オペレータの発話音声に適合させた口唇動作を実現することができる。したがって、アンドロイド12の対応する人間に対して違和感を与えることなく、自然な対話を行うことができる。

【0070】

なお、上述の実施例では、制御装置14で発話処理を実行する場合を説明したが、他の実施例では、遠隔操作端末20が発話処理を実行して、制御装置14に対して動作指令を送信するようにしてもよい。つまり、たとえば、図3において、ステップS1でマイク24で音声を検出したか否かを判定する。さらにステップS7で音声の取得(検出)から一定時間経過したことが判定されたときには、ステップS9で音声データを含む再生指示を制御装置14に対して送信する。制御装置14では、この再生指示の受信に応じて、当該音声を出力する。また、遠隔操作端末20では、図4のステップS39で開始される動作指令発行処理で、制御装置14に対して制御情報を含む各動作指令を送信する。制御装置14では、この動作指令の受信に応じて、当該動作指令に含まれる制御情報に従ってアクチュエータ40を制御する。

【0071】

また、上述の各実施例では、遠隔オペレータの発話音声のみに基づいて、口唇形状を推定するようにしているが、他の実施例では、発話音声と画像処理とを組み合わせる口唇形状を推定するようにしてもよい。つまり、遠隔操作端末20にカメラを設けて、発話する遠隔オペレータの顔を撮影し、当該撮影画像に画像処理を施すことによって、口唇形状を検出する。そして、発話音声から推定した口唇形状と撮影画像から推定した口唇形状とに基づいて、最終的な口唇形状の推定を行う。これによって、口唇形状の推定の精度を高めることができる。あるいは、その他の実施例では、発話音声とモーションキャプチャとを組み合わせる口唇形状を推定するようにしてもよい。つまり、遠隔操作端末20にモーションキャプチャシステムを組み合わせる。具体的には、遠隔オペレータの顔面の適宜な位置にマーカを取り付けるとともに、複数のカメラを設ける。そして、発話する遠隔オペレータの顔を撮影し、当該撮影画像に基づいて当該マーカの3次元位置を計測して、各マーカの位置に基づいて口唇形状を検出する。そして、発話音声から推定した口唇形状と3次元位置から推定した口唇形状とに基づいて、最終的な口唇形状の推定を行う。これによっても、口唇形状の推定の精度をよくすることができる。

【図面の簡単な説明】

【0072】

【図1】この発明のアンドロイド制御システムの一例を示す図解図である。

【図2】図1に示すアンドロイド、制御装置および環境センサの電氣的な構成を示すブロック図である。

【図3】図2に示すCPU70の発話処理における動作の一例を示すフロー図である。

【図4】図3に示す口唇動作制御処理の動作の一例を示すフロー図である。

【符号の説明】

【0073】

10 ... アンドロイド制御システム

10

20

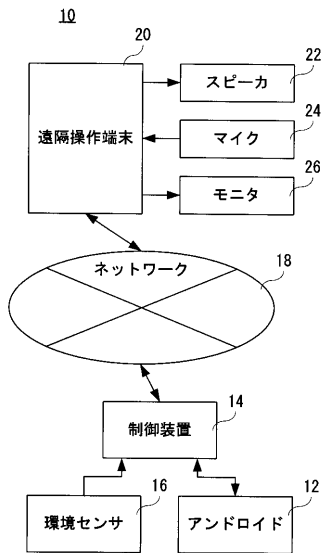
30

40

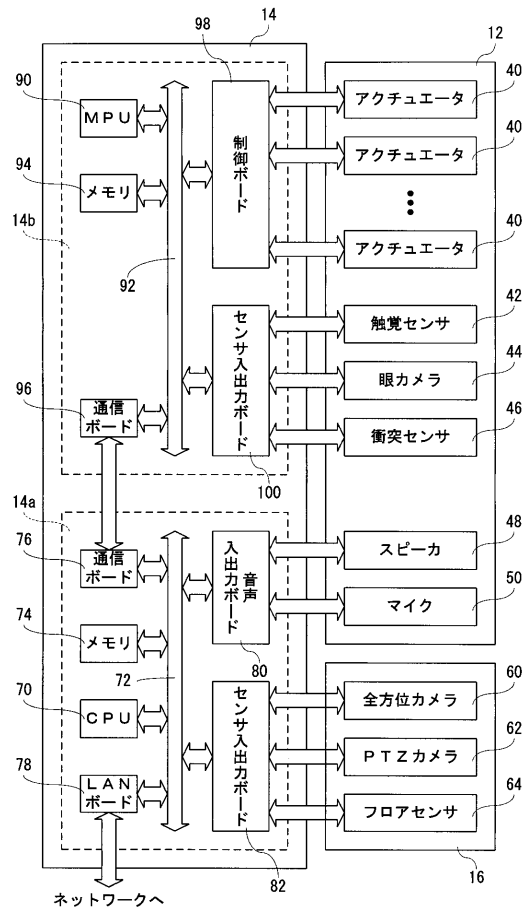
50

- 1 2 ... アンドロイド
- 1 4 ... 制御装置
- 2 0 ... 遠隔操作端末
- 2 2 , 4 8 ... スピーカ
- 2 4 , 5 0 ... マイク
- 4 0 ... アクチュエータ
- 7 0 ... C P U
- 7 4 , 9 4 ... メモリ
- 8 0 ... 音声入出力ボード
- 9 0 ... M P U
- 9 8 ... 制御ボード

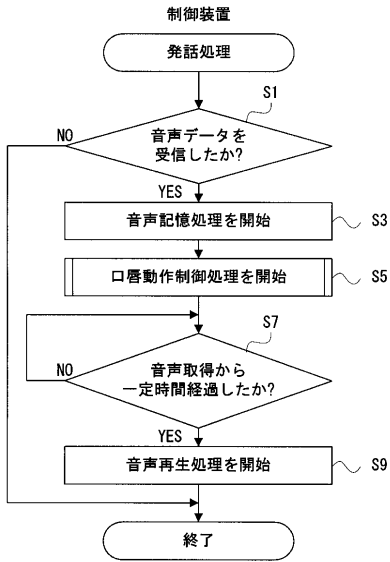
【 図 1 】



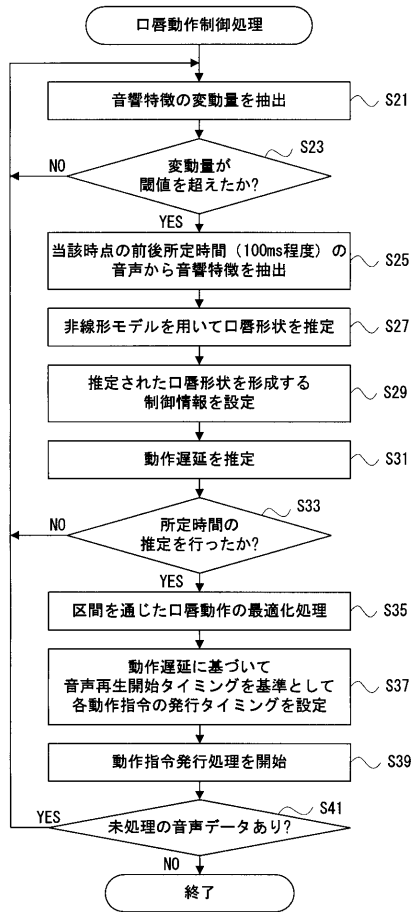
【 図 2 】



【図3】



【図4】



フロントページの続き

(72)発明者 萩田 紀博

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

審査官 前田 祐希

(56)参考文献 特開2003-248841(JP,A)

特開2005-266671(JP,A)

特開2005-012819(JP,A)

特開2004-034273(JP,A)

(58)調査した分野(Int.Cl., DB名)

G10L 11/00-21/06