

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4735965号
(P4735965)

(45) 発行日 平成23年7月27日(2011.7.27)

(24) 登録日 平成23年5月13日(2011.5.13)

(51) Int.Cl. F 1
HO 4M 11/00 (2006.01) HO 4M 11/00 3 0 1

請求項の数 2 (全 24 頁)

<p>(21) 出願番号 特願2005-322141 (P2005-322141)</p> <p>(22) 出願日 平成17年11月7日(2005.11.7)</p> <p>(65) 公開番号 特開2007-129626 (P2007-129626A)</p> <p>(43) 公開日 平成19年5月24日(2007.5.24)</p> <p>審査請求日 平成20年3月31日(2008.3.31)</p> <p>(出願人による申告)平成17年4月1日付け、支出負担行為担当官 総務省大臣官房会計課企画官、研究テーマ「ネットワーク・ヒューマン・インターフェースの総合的な研究開発(ネットワークロボットの技術)」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>前置審査</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</p> <p>(74) 代理人 100090181 弁理士 山田 義人</p> <p>(72) 発明者 小泉 智史 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>(72) 発明者 塩見 昌裕 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>(72) 発明者 神田 崇行 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】 遠隔地間対話システム

(57) 【特許請求の範囲】

【請求項1】

ネットワークを介して接続される2つの対話装置を含む遠隔地間で対話を行うためのシステムであって、

各前記対話装置は、

音声を取得する取得手段、

前記取得手段で取得した前記音声を相手側の前記対話装置へ送信する送信手段、

相手側の前記対話装置から送信された音声を受信する受信手段、および

前記受信手段で受信した前記音声を出力する出力手段を含んでいて、

少なくとも空白時間と発話音声の特徴に関する情報を含む複数の間パターンを記憶する間パターン記憶手段、

各前記対話装置における少なくとも音声取得状態および音声出力状態を含む対話状態の履歴を記録する履歴記録手段、

前記履歴記録手段によって記録された前記履歴に基づいて両方の前記対話装置で無発話状態であると判定されるとき、少なくとも空白時間および当該空白前の発話音声の特徴を含む間の状況と前記複数の間パターンとの照合を行う照合手段、

前記照合手段による照合の結果マッチする前記間パターンがあるとき、当該間パターンに対応する所定の音声を、当該所定の音声の発話者の相手側に存在する前記対話装置の前記出力手段から出力する間制御手段、および

前記履歴記録手段によって記録された前記履歴に基づいて両方の前記対話装置で発話が

10

20

重複したと判定されるとき、一方の音声を録音して、その後発話が終了したときに当該録音音声を他方の前記対話装置の前記出力手段から出力する遅延再生手段を備える、遠隔地間対話システム。

【請求項 2】

前記対話装置の少なくとも一方が身振りを実行可能なロボットであるとき、前記間制御手段は、前記音声の出力とともに、前記間パターンに対応する所定の身振りを当該対話装置に実行させる、請求項 1 記載の遠隔地間対話システム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は遠隔地間対話システムに関し、特にたとえば、遠隔地に離れた対話者同士の音声をネットワークを介して通信する、遠隔地間対話システムに関する。

【背景技術】

【0002】

遠隔地間で対話を行う場合には、遅延により対話の空白時間が長くなり発話のタイミングが取りづらくなって、両対話者の発話が重複する事態が生じ易い。従来、このような発話の重複を防止して、遠隔地間の円滑な対話を生成しようとする技術は存在しなかった。

【0003】

なお、たとえば特許文献 1 には、音声入出力装置の出力する音声と利用者の発話によって入力される音声との重畳を検出する技術の一例が開示されている。

【0004】

また、特許文献 2 および特許文献 3 には、単に対話音声を出力するだけでなく、対話者の前にロボットを設置してゼスチャを行わせる技術の一例が開示されている。特許文献 2 の技術では、話し手の音声に基づいて当該話し手側のロボットが身振りを実行し、一方、聞き手側で受信した音声に基づいて当該聞き手側のロボットが身振りを実行することで、会話の実感を高めている。また、特許文献 3 の技術では、話し手側の身振り情報の送信に応じて相手側ロボットで当該身振りが再現される。

【特許文献 1】特開平 7 - 264103 号公報

【特許文献 2】特開 2000 - 349920 号公報

【特許文献 3】特開 2001 - 156930 号公報

【発明の開示】

【発明が解決しようとする課題】

【0005】

特許文献 1 の技術では、出力音声と入力音声との重畳を検出して、エコーキャンセラの動作が変化されるが、話者の発話タイミングを制御することはできない。また、特許文献 2 および 3 の技術では、対話者の前に設置したロボットに身振りをさせることによって、円滑な対話の実現を図っているが、話者の発話タイミングを制御することはできない。このように、従来技術では、発話のタイミングを制御することができなかつたので、遅延により対話に異常な空白時間が生じても対応できなかつた。したがって、両対話者の発話が重なることを防止することができず、円滑な対話を実現することができなかつた。

【0006】

それゆえに、この発明の主たる目的は、遠隔地間対話に適切な間を与えることができ、円滑な対話を実現できる、遠隔地間対話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】

【0007】

請求項 1 の発明は、ネットワークを介して接続される 2 つの対話装置を含む遠隔地間で対話を行うためのシステムである。各対話装置は、音声を取得する取得手段、取得手段で取得した音声を相手側の対話装置へ送信する送信手段、相手側の対話装置から送信された音声を受信する受信手段、および受信手段で受信した音声を出力する出力手段を含んでいる。当該システムは、少なくとも空白時間と発話音声の特徴に関する情報を含む複数の間

10

20

30

40

50

パターンを記憶する間パターン記憶手段、各対話装置における少なくとも音声取得状態および音声出力状態を含む対話状態の履歴を記録する履歴記録手段、履歴記録手段によって記録された履歴に基づいて両方の対話装置で無発話状態であると判定されるとき、少なくとも空白時間および当該空白前の発話音声の特徴を含む間の状況と複数の間パターンとの照合を行う照合手段、照合手段による照合の結果マッチする間パターンがあるとき、当該間パターンに対応する所定の音声を、当該所定の音声の発話者の相手側に存在する対話装置の出力手段から出力する間制御手段、および履歴記録手段によって記録された履歴に基づいて両方の対話装置で発話が重複したと判定されるとき、一方の音声を録音して、その後発話が終了したときに当該録音音声を他方の対話装置の出力手段から出力する遅延再生手段を備える。

10

【0008】

請求項1の発明では、遠隔地間対話システムは2つの対話装置を含み、各対話装置が取得した音声を通信して相手側で出力することによって、遠隔地間での対話者同士の対話が行われる。間パターン記憶手段には複数の間パターンが記憶されている。間パターンは、対話における適切な間の取り方を示し、少なくとも空白時間と発話音声の特徴に関する情報を含む。たとえば、発話音声の特徴は、基本周波数（ピッチ）、振幅および音節の平均持続時間等を含んでよい。履歴記録手段は、各対話装置における少なくとも音声取得状態および音声出力状態を含む対話状態の履歴を記録する。対話状態の履歴は、後述される実施例では発話フラグテーブルであり、各時刻の発話の有無状態（SPEAKINGフラグ、SILENTフラグ）および処理の状態（RECORDINGフラグ、INTERPOLATINGフラグ）などが記録される。照合手段は、両方の対話装置で無発話状態であると判定されるとき、少なくとも空白時間および当該空白前の発話音声の特徴を含む間の状況と、複数の間パターンとの照合を行う。つまり、対話が無音状態である場合に、現在の間の状況が複数の間パターンのいずれかにマッチしているかが確認される。間制御手段は、照合の結果マッチする間パターンがあるとき、当該間パターンに対応する所定の音声を、当該所定の音声の発話者の相手側に存在する対話装置の出力手段から出力する。ただし、発話が重複してしまった場合には、遅延再生手段によって、一方の音声が録音され、その後発話が終了したときに、当該録音した音声が他方の対話装置から出力される。

20

なお、上記間パターン記憶手段、履歴記録手段、照合手段、間制御手段および遅延再生手段は、2つの対話装置のいずれか一方に、または、このシステムに含まれる別の装置（実施例では発話タイミング制御サーバ）に設けられてよい。あるいは、これらの手段は、2つの対話装置に分散して設けられてもよい。

30

【0009】

請求項1の発明によれば、会話における無音状態が検出されたときに、適切な間を与える間パターンに対応する音声を出力することができる。したがって、発話が重なってしまうのを防止することができ、円滑な対話を成立させることができる。また、万一発話が重複してしまっても、両発話が同時に相手側で出力されるのを回避することができる。

【0010】

請求項2の発明は、請求項1の発明に従属し、対話装置の少なくとも一方が身振りを実行可能なロボットであるとき、間制御手段は、音声の出力とともに、さらに間パターンに対応する所定の身振りを当該対話装置に実行させる。

40

【0011】

請求項2の発明では、対話装置の少なくとも一方は、身振りを実行可能なロボットであってよい。間制御手段は、当該ロボットに、間パターンに対応する音声を出力させるとともに、間パターンに対応する身振りを実行させる。したがって、音声と身振りを使用して、対話に適切な間を与えることができるので、より円滑な遠隔地間対話を成立させることができる。

【発明の効果】

【0014】

この発明によれば、対話に無音状態が検出されたときに適切な間を取るように言葉を挿

50

入するようにしたので、対話の空白時間を適切な長さにすることができる。このため、遅延によって空白時間が長くなって対話者に違和感を与えてしまうようなことを回避できる。したがって、対話者は発話のタイミングを計りやすくなるので、両者の発話の重複を防止することができ、円滑な対話を成立させることができる。

【0015】

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る。

【発明を実施するための最良の形態】

【0016】

図1を参照して、この実施例の遠隔地間対話システム(以下、単に「システム」とも言う。)10は、遠隔地に離れた対話者同士が対話を行うためのものである。システム10は少なくとも2つの対話装置12(12a, 12b)を含む。2つの対話装置12は、ネットワーク、たとえば公衆インターネット網を介して接続されており、対話装置12a側の対話者Aおよび対話装置12b側の対話者Bの発話した音声の音声データを互いに通信する。また、この実施例のシステム10は、発話タイミング制御サーバ(以下、単に「サーバ」とも言う。)14を含み、サーバ14はネットワークを介して少なくとも2つの対話装置12と通信可能に接続される。

10

【0017】

この実施例では、一方の対話装置12としてコンピュータ12aが適用され、他方の対話装置12としてコミュニケーションロボット(以下、単に「ロボット」とも言う。)12bが適用された場合を説明する。

20

【0018】

対話装置12aはマイク16およびスピーカ18を備える。また、対話装置12aはたとえばパーソナルコンピュータであり、CPU、メインメモリ、通信装置および入力装置等を備えている。メインメモリには、この発明の対話装置12として機能するために必要なプログラムおよびデータが記憶される。プログラムおよびデータは、メインメモリに予め固定的に記憶されてもよいし、または、情報記憶媒体やネットワークから取得されてよい。CPUは、当該プログラムに従って、メインメモリのうちのワーキングメモリに一時的なデータを生成または取得しつつ対話のための処理を実行する。

【0019】

マイク16は対話者の発話した音声を取得するためのものであり、当該音声は音声入出力ボードでデータに変換されて、音声データとしてメインメモリに記憶される。スピーカ18は、対話相手の音声およびシステム10の備える音声を出力するためのものである。CPUは受信した音声データを音声入出力ボードに与えて当該音声をスピーカ18から出力する。通信装置は、ネットワークを介して他方の対話装置12やサーバ14にデータを送受信する。また、入力装置は、キーボードまたはポインティングデバイス等である。

30

【0020】

この実施例では、相手側の対話装置12としてロボット12bが使用されるので、ユーザが入力装置を用いてロボット12bの身振りを指示可能になっている。ロボット12bの身振りは、表示装置の画面に表示されたリストから選択されてよいし、あるいは入力装置の各キーに割り当てられてもよい。対話者Aは発話しながら入力装置を用いて動作を指示することによって、相手側のロボット12から自分の音声出力することができ、しかも当該ロボット12に所望の身振りを行わせることができる。

40

【0021】

なお、対話装置12aは、音声入出力可能かつ通信可能なコンピュータであればよく、PCに限られず、ゲーム機、携帯電話、携帯ゲーム機などの他のコンピュータであってよい。

【0022】

他方の対話装置12bは人間のような身体部位を有するロボットであり、身体部位を動かすことによって所定の身振りを対話者Bに提示することができる。このロボット12b

50

は、マイク 20 およびスピーカ 22 を備えている。詳しくは、図 2 にロボット 12 b の外觀の一例が示され、図 3 には当該ロボット 12 の電氣的な構成の一例が示される。

【 0023 】

図 2 を参照して、ロボット 12 b は台車 24 を含み、この台車 24 の下面には、このロボット 12 b を自律移動させる車輪 26 が設けられる。この車輪 26 は、車輪モータ (図 3 において参照番号「 28 」で示す。) によって駆動され、台車 24 すなわちロボット 12 b を前後左右任意の方向に動かすことができる。なお、図 2 では示さないが、この台車 24 の前面には、衝突センサ (図 3 において参照番号「 30 」で示す。) が取り付けられ、この衝突センサ 30 は、台車 24 への人や他の障害物との接触を検知する。ロボット 12 b の移動中に接触を検知すると、直ちに車輪 26 の駆動を停止することができる。

10

【 0024 】

台車 24 の上には、多角形柱のセンサ取付パネル 32 が設けられ、このセンサ取付パネル 32 の各面には、超音波距離センサ 34 が取り付けられる。この超音波距離センサ 34 は、取付パネル 32 すなわちロボット 12 b の周囲の主として人との間の距離を計測するためのものである。

【 0025 】

台車 24 の上には、さらに、ロボット 12 b の胴体が、その下部が上述の取付パネル 32 に囲まれて、直立するように取り付けられる。この胴体は下部胴体 36 と上部胴体 38 とから構成され、これら下部胴体 36 および上部胴体 38 は、連結部 40 によって連結される。連結部 40 には、図示しないが、昇降機構が内蔵されていて、この昇降機構を用いることによって、上部胴体 38 の高さすなわちロボット 12 b の高さを変化させることができる。昇降機構は、腰モータ (図 3 において参照番号「 42 」で示す。) によって駆動される。

20

【 0026 】

上部胴体 38 のほぼ中央には、1つの全方位カメラ 44 と、1つのマイク 20 とが設けられる。全方位カメラ 44 は、ロボット 12 b の周囲を撮影するもので、後述の眼カメラ 46 と区別される。マイク 20 は、上述のように、周囲の音、とりわけ人の声を取り込む。

【 0027 】

上部胴体 38 の両肩には、それぞれ、肩関節 48 R および 48 L によって、上腕 50 R および 50 L が取り付けられる。肩関節 48 R および 48 L は、それぞれ 3 軸の自由度を有する。すなわち、右肩関節 48 R は、X 軸、Y 軸および Z 軸の各軸廻りにおいて上腕 50 R の角度を制御できる。Y 軸は、上腕 50 R の長手方向 (または軸) に平行な軸であり、X 軸および Z 軸は、その Y 軸に、それぞれ異なる方向から直交する軸である。左肩関節 48 L は、A 軸、B 軸および C 軸の各軸廻りにおいて上腕 50 L の角度を制御できる。B 軸は、上腕 50 L の長手方向 (または軸) に平行な軸であり、A 軸および C 軸は、その B 軸に、それぞれ異なる方向から直交する軸である。

30

【 0028 】

上腕 50 R および 50 L のそれぞれの先端には、肘関節 52 R および 52 L を介して、前腕 54 R および 54 L が取り付けられる。肘関節 52 R および 52 L は、それぞれ、W 軸および D 軸の軸廻りにおいて、前腕 54 R および 54 L の角度を制御できる。

40

【 0029 】

なお、上腕 50 R および 50 L ならびに前腕 54 R および 54 L の変位を制御する X、Y、Z、W 軸および A、B、C、D 軸では、「 0 度」がホームポジションであり、このホームポジションでは、上腕 50 R および 50 L ならびに前腕 54 R および 54 L は下方向に向けられる。

【 0030 】

また、図 2 では示さないが、上部胴体 38 の肩関節 48 R および 48 L を含む肩の部分や上述の上腕 50 R および 50 L ならびに前腕 54 R および 54 L を含む腕の部分には、それぞれ、タッチセンサ (図 3 において参照番号「 56 」で包括的に示す。) が設けられ

50

ていて、これらのタッチセンサ56は、人がロボット12bのこれらの部位に接触したかどうかを検知する。

【0031】

前腕54Rおよび54Lのそれぞれの先端には、手に相当する球体58Rおよび58Lがそれぞれ固定的に取り付けられる。なお、この球体58Rおよび58Lに代えて、この実施例のロボット12bと異なり指の機能が必要な場合には、人の手の形をした「手」を用いることも可能である。

【0032】

上部胴体38の中央上方には、首関節60を介して、頭部62が取り付けられる。この首関節60は、3軸の自由度を有し、S軸、T軸およびU軸の各軸廻りに角度制御可能である。S軸は首から真上に向かう軸であり、T軸およびU軸は、それぞれ、このS軸に対して異なる方向で直交する軸である。頭部62には、人の口に相当する位置に、上述のスピーカ22が設けられる。なお、スピーカ22は、ロボット12bが、その周囲の人に対して音声または声によってコミュニケーションを図るために用いられてよい。また、スピーカ22は、ロボット12の他の部位たとえば胴体に設けられてもよい。

10

【0033】

また、頭部62には、目に相当する位置に眼球部64Rおよび64Lが設けられる。眼球部64Rおよび64Lは、それぞれ眼カメラ46Rおよび46Lを含む。なお、左右の眼球部64Rおよび64Lをまとめて参照符号「64」で示し、左右の眼カメラ46Rおよび46Lをまとめて参照符号「46」で示すこともある。眼カメラ46は、ロボット12bに接近した人の顔や他の部分ないし物体等を撮影してその映像信号を取り込む。

20

【0034】

なお、上述の全方位カメラ44および眼カメラ46のいずれも、たとえばCCDやCMOSのように固体撮像素子を用いるカメラであってよい。

【0035】

たとえば、眼カメラ46は眼球部64内に固定され、眼球部64は眼球支持部(図示せず)を介して頭部62内の所定位置に取り付けられる。眼球支持部は、2軸の自由度を有し、 θ 軸および ϕ 軸の各軸廻りに角度制御可能である。 θ 軸および ϕ 軸は頭部62に対して設定される軸であり、 θ 軸は頭部62の上へ向かう方向の軸であり、 ϕ 軸は θ 軸に直交しかつ頭部62の正面側(顔)が向く方向に直交する方向の軸である。この実施例では、頭部62がホームポジションにあるとき、 θ 軸はS軸に平行し、 ϕ 軸はU軸に平行するように設定されている。このような頭部62において、眼球支持部が θ 軸および ϕ 軸の各軸廻りに回転されることによって、眼球部64ないし眼カメラ46の先端(正面)側が変位され、カメラ軸すなわち視線方向が移動される。

30

【0036】

なお、眼カメラ46の変位を制御する θ 軸および ϕ 軸では、「0度」がホームポジションであり、このホームポジションでは、図2に示すように、眼カメラ46のカメラ軸は頭部62の正面側(顔)が向く方向に向けられ、視線は正視状態となる。

【0037】

図3を参照して、このロボット12bは、全体の制御のためにマイクロコンピュータまたはCPU66を含み、このCPU66には、バス68を通して、メモリ70、モータ制御ボード72、センサ入力/出力ボード74および音声入力/出力ボード76が接続される。

40

【0038】

メモリ70は、図示しないが、ROMやHDDおよびRAM等を含み、ROMまたはHDDには、このロボット12bをこの発明の対話装置12として機能させるためのプログラムおよびデータが予め格納されている。CPU66は、このプログラムに従って処理を実行する。また、RAMは、バッファメモリやワーキングメモリとして使用される。

【0039】

モータ制御ボード72は、たとえばDSP(Digital Signal Processor)で構成され、右

50

腕、左腕、頭および眼等の身体部位を駆動するためのモータを制御する。すなわち、モータ制御ボード72は、CPU66からの制御データを受け、右肩関節48RのX、YおよびZ軸のそれぞれの角度を制御する3つのモータと右肘関節52Rの軸Wの角度を制御する1つのモータを含む計4つのモータ(図3ではまとめて「右腕モータ」として示す。)78の回転角度を調節する。また、モータ制御ボード72は、左肩関節48LのA、BおよびC軸のそれぞれの角度を制御する3つのモータと左肘関節52LのD軸の角度を制御する1つのモータとを含む計4つのモータ(図3ではまとめて「左腕モータ」として示す。)80の回転角度を調節する。モータ制御ボード72は、また、首関節60のS、TおよびU軸のそれぞれの角度を制御する3つのモータ(図3ではまとめて「頭部モータ」として示す。)82の回転角度を調節する。モータ制御ボード72は、また、腰モータ42および車輪26を駆動する2つのモータ(図3ではまとめて「車輪モータ」として示す。)28を制御する。さらに、モータ制御ボード72は、右眼球部64Rの軸および軸のそれぞれの角度を制御する2つのモータ(図3ではまとめて「右眼球モータ」として示す。)84の回転角度を調節し、また、左眼球部64Lの軸および軸のそれぞれの角度を制御する2つのモータ(図3ではまとめて「左眼球モータ」として示す。)86の回転角度を調節する。

【0040】

なお、この実施例の上述のモータは、車輪モータ28を除いて、制御を簡単化するためにそれぞれステップモータまたはパルスモータであるが、車輪モータ28と同様に、直流モータであってよい。また、この実施例では、ロボット12bの腕、頭、眼などの身体部位を駆動するアクチュエータとして電力を駆動源とするモータを用いた。しかしながら、このロボット12bとしては、たとえば空気圧(または負圧)、油圧、圧電素子あるいは形状記憶合金などによる他のアクチュエータによって身体部位を駆動するロボットが適用されてもよい。

【0041】

センサ入力/出力ボード74も、同様に、DSPで構成され、各センサやカメラからの信号を取り込んでCPU66に与える。すなわち、超音波距離センサ34の各々からの反射時間に関するデータがこのセンサ入力/出力ボード74を通して、CPU66に入力される。また、全方位カメラ44からの映像信号が、必要に応じてこのセンサ入力/出力ボード74で所定の処理が施された後、CPU66に入力される。眼カメラ46からの映像信号も、同様にして、CPU66に与えられる。また、タッチセンサ56からの信号がセンサ入力/出力ボード74を介してCPU66に与えられる。

【0042】

スピーカ22には音声入力/出力ボード76を介して、CPU66から音声データが与えられ、それに従って、スピーカ22からはそのデータに従った音声または声が出力される。また、マイク20からの音声入力が、音声入力/出力ボード76を介して音声データとしてCPU66に取り込まれる。

【0043】

通信LANボード88も、同様に、DSPで構成され、CPU66から与えられた送信データを無線通信装置90に与えて、当該データを無線通信装置90から送信させる。また、通信LANボード88は無線通信装置90を介してデータを受信し、受信データをCPU66に与える。

【0044】

図1に戻って、サーバ14は、両対話者の発話のタイミングを制御するために設けられる。サーバ14は、CPU、メインメモリ、通信装置等を備える。メインメモリにはこのサーバ14を制御するためのプログラムおよびデータが記憶される。CPUは当該プログラムに従って処理を実行する。

【0045】

また、サーバ14は音声解析履歴データベース(DB)92および間パターンDB94を含む。音声解析履歴DB92には、対話装置12で取得された対話者の音声の解析デー

10

20

30

40

50

タの履歴が記憶される。

【0046】

間パターンDB94には、後述するように、対話に適切な間を与えるための間パターンデータ(図4参照)が記憶されている。間パターンデータは、予め発話の計測を行って得た発話データからパターン認識によって抽出される。計測を実際の使用者を対象として行うと、間の取り方の個人的特徴を抽出できる。ただし、標準的なまたは一般的な間の取り方も存在すると考えられるので、任意の人を被験者としてその発話を計測して間パターンデータを抽出してよい。

【0047】

このシステム10では、各対話装置12が、一定時間Tごとにマイク16または20で音声を検出する。Tはたとえば1フレームまたは所定のフレーム数であってよい。1フレームはたとえば1/30秒である。対話装置12は、検出した結果すなわち発話の有無に応じた処理を行う。対話装置12は、検出時刻における発話状態(音声取得状態)および実行した処理(音声出力状態)など、当該装置12における対話状態に関する情報をサーバ14に送信する。サーバ14は、当該対話状態に関する情報を受信して、当該対話装置12における状態を逐一記憶する。このような対話状態の履歴は、発話フラグテーブル(図5参照)としてメモリに記憶される。発話フラグテーブルでは、後述するように、検出時刻ごとの対話装置12における少なくとも音声取得状態および音声出力状態を含む対話状態を示すフラグが記憶されている。

【0048】

なお、便宜上、ここでは、対話装置12a側からみた動作を説明する。しかし、対話装置12aと対話装置12bとの相違は、主に身体動作の提示に関する機能のみであるから、対話装置12bの動作も、対話装置12aの場合と同様である。

【0049】

対話装置12aは、マイク16から音声を検出した場合には、サーバ14の発話フラグテーブルを参照する。そして、対話装置12aは、1つ前の検出時刻における(つまりT前の)相手側の状態に応じた処理を実行する。具体的には、前検出時刻における相手の状態フラグがSPEAKINGフラグでない場合には、つまり、前検出時刻において相手が自分に音声を送信している状態ではない場合には、対話装置12aは、マイク16で検出した音声データをメモリにローカルファイルとして記録しつつ、当該音声データを相手側の対話装置12bに送信する。これに応じて、対話装置12bは、当該音声データを受信して当該音声をスピーカ22から出力する。このように、一方の対話者Aが発話し、かつ、前検出時刻で他方の対話者Bの音声が送信されていない場合には、対話者Aの音声データが直ちに送信され、当該音声データが相手側の対話装置12bで出力されて、相手Bに聞かされる。

【0050】

また、この場合、対話装置12aは、SPEAKINGフラグをサーバ14に送信し、これに応じて、サーバ14は発話フラグテーブルに当該対話装置12aの当該検出時刻tにおける状態として当該SPEAKINGフラグを記憶する。SPEAKINGフラグは、一方の対話者が発話している状態、すなわち、発話音声データが直接相手側に送信されて再生されている状態を意味する。さらに、対話装置12aは、メモリのSENDフラグをオンにして、自分の処理状態として、音声を相手に送信中であることを記憶する。

【0051】

なお、当該音声を記録したローカルファイルには、当該発話が終わったときに音声解析が実行され、当該発話音声の特徴ないし状態はサーバ14の音声解析履歴DB92に記憶される。

【0052】

一方、対話装置12aはマイク16から音声を検出しなかった場合には、SILENTフラグをサーバ14に送信し、これに応じて、サーバ14は発話フラグテーブルに当該SILENTフラグを記憶する。SILENTフラグは、当該検出時刻tにおいて対話者が

10

20

30

40

50

発話していない状態を意味する。このように、音声を検出されない場合には、発話フラグテーブルに S I L E N T フラグが記録される。

【 0 0 5 3 】

サーバ 1 4 では、両対話者とも発話していない状態 (S I L E N T フラグ) が検出されたとき、間パターンと対話における現在までの発話状況 (間の状況) との照合が行われる。間の状況は、少なくとも対話の空白時間 (無音時間) 、および当該空白前の対話における発話音声の特徴 (音声解析結果) を含む。

【 0 0 5 4 】

間パターン D B 9 4 に記憶される間パターンデータの一例が図 4 に示される。間パターン D B 9 4 には、会話に適切な間を与えることができる複数の間パターンデータが記憶されている。間パターンは、少なくとも空白時間とその前の発話音声の特徴に関する情報を含む。この実施例では、間パターンデータは、会話の空白時間、最終発話者、条件 (I) 、間機能言葉、発話者および動作コマンド等の情報を含む。空白時間 (t) は、両者無音状態が継続している時間の条件である。最終発話者は、当該空白前の対話での最後の発話者 A または B の条件である。条件 (I) は、当該空白前の対話の音声の解析結果の条件であり、たとえば基本周波数 (ピッチ) 、振幅および音節の平均持続時間等の要素を含む。

【 0 0 5 5 】

このような空白時間とその前の発話音声の特徴によって規定される間パターンに見合う間の状況が検出されたとき、当該間パターンに基づく言葉および身振りが対話に挿入される。具体的には、この実施例では、最終発話者、条件 (I) および空白時間条件 (t) に合う対話後の空白時間が生じた場合に、当該パターンで指定された発話者側から間機能言葉が発せられる。また、発話者の代わりに当該間機能言葉を発する対話装置 1 2 がロボット 1 2 b である場合には、当該出力する言葉に対応する身振りも動作コマンドに従って再現される。

【 0 0 5 6 】

間機能言葉は、無音時間に挿入されることによって無音を間として機能させ会話に適切な間を与えるための言葉である。たとえば、応答、合の手、間投詞などの言葉であってよい。図 4 では、「うんうん」、「うーん」、「はいはい」および「えーと」などが間機能言葉として示される。また、各間機能言葉には、当該言葉とともに提示される身体動作を実行するための動作コマンドが対応付けられている。図 4 では、「うんうん」には「うなづく」のコマンド、「うーん」には「首傾げる」のコマンド、「はいはい」には「うなづく」のコマンド、「えーと」には「視線を上方に向ける」コマンドがそれぞれ対応付けられる。

【 0 0 5 7 】

なお、図 4 において、発話者は、人間である対話者を意味しており、当該間機能言葉を出力する対話装置 1 2 は、この発話者の相手側の場所に存在する対話装置 1 2 である。たとえば、図 4 の一番上のパターンの場合、最後の発話者が対話者 A であり、間機能言葉を出力する対話装置は、発話者 B の相手である対話者 A 側の対話装置 1 2 a となる。一番上のパターンは「うんうん」という言葉と「うなづく」行動に対応付けられており、最終発話者の相手側である発話者が応答する動作を、最終発話者側に存在するロボット 1 2 b が表現することで、間が与えられる。一方、上から 2 番目のパターンは「うーん」という言葉と「首を傾げる」行動に対応付けられており、最終発話者の相手側に存在するロボット 1 2 b が発話および身振りをさらに続けることによって、間が与えられる。

【 0 0 5 8 】

間パターンを用いた照合の結果、現在の間の状況にマッチする間パターンデータが間パターン D B 9 4 に存在する場合には、つまり、間パターンに従って適切な間を取る必要がある状況であると判断される場合には、サーバ 1 4 は、必要な対話装置 1 2 に間を取るための言葉の再生を指示する。この実施例では、当該言葉の音声データと再生指示とが送信される。これに応じて、当該対話装置 1 2 側で当該音声出力される。さらに、対話装置 1 2 が身体表現可能な対話装置 1 2 b である場合には、サーバ 1 4 は当該間を取るための

10

20

30

40

50

言葉に対応する身振りの実行を指示する。この実施例では、当該身振りに動作コマンドと再生指示とが送信される。これに応じて、当該対話装置 1 2 b では、対応する身体部位が動かされて当該身振りが実行される。

【 0 0 5 9 】

このように、対話に無発話状態（無音状態）が検出されたときに、現在の間の状況と間パターンとの照合を行うようにした。そして、必要があれば適切な間を取るよう言葉や身振りを空白時間に挿入するようにしたので、対話における空白時間を適切な時間に維持することができる。このため、遅延によって対話における空白時間が長くなって対話者に違和感を与えてしまうようなことを回避できる。したがって、対話者は発話のタイミングを計りやすくなり、両対話者の発話が重複する事態が生ずるのを防止することができる。このように、言葉や身振りの挿入によって対話者の発話タイミングを制御することができるので、対話を継続させたり発話を促進したりすることができるし、自然な会話の流れを作り出すことができる。したがって、円滑な対話を成立させることができる。

10

【 0 0 6 0 】

また、このシステム 1 0 では、万一両対話者の発話が重複した場合には、一方の発話の出力を遅らせることによって、両発話が完全に重なってしまうのを回避する機能を備えるようにしている。

【 0 0 6 1 】

具体的には、対話装置 1 2 a で音声を検出された場合において、前検出時刻の相手の状態フラグが S P E A K I N G フラグであるときには、つまり、両者の発話が重複している場合には、対話装置 1 2 a は、マイク 1 6 で検出した音声の録音を開始し、当該音声データを音声ファイルとしてメモリに記憶する。さらに、対話装置 1 2 a は、R E C O R D I N G フラグをサーバ 1 4 に送信し、これに応じて、サーバ 1 4 は発話フラグテーブルに当該対話装置 1 2 a の当該検出時刻 t における状態として当該 R E C O R D I N G フラグを記憶する。R E D O R D I N G フラグは、音声データを録音中であり、当該音声は相手側に送信されていない状態を意味する。

20

【 0 0 6 2 】

また、対話装置 1 2 a は、メモリの R E C O R D フラグをオンにして、自分の処理状態として、音声を録音中であることを記憶する。

【 0 0 6 3 】

また、サーバ 1 4 の発話フラグテーブルでは、録音した音声ファイルの再生を制御するための情報として P L A Y フラグが記憶される。この実施例では、対話装置 1 2 a は、録音しているときは、P L A Y フラグの値に 1 を加算するようにサーバ 1 4 に指示する。P L A Y フラグの初期値は 0 であり、録音が行われているときは毎検出時刻ごとに前の検出時刻の値に 1 だけ加算され、録音が行われていないときには前の検出時刻の値が維持される。

30

【 0 0 6 4 】

その後、対話装置 1 2 a で音声を検出されなくなったときには、録音した音声ファイルがサーバ 1 4 に送信される。これに応じて、サーバ 1 4 は、受信した音声ファイルを、当該録音が行われた検出時刻 t に対応付けてメモリに記憶する。発話フラグテーブルでは、当該音声ファイルを格納した記憶位置が記憶される。なお、音声ファイルにはサーバ 1 4 に送信される前に音声解析処理が施され、当該解析データがサーバ 1 4 に送信されて音声解析履歴 D B 9 2 に記憶される。

40

【 0 0 6 5 】

サーバ 1 4 は、両対話装置 1 2 とも音声を出力していないことが検出された場合、つまり、両対話者の状態のいずれにも S P E A K I N G フラグが記憶されていないことが検出された場合、いずれかの対話者の音声ファイルが再生されずに記憶されているか否かを判定する。未再生の録音ファイルが残っている場合、つまり、P L A Y フラグの値が 1 以上である場合には、録音の開始された時刻の早い方の音声ファイルの再生が実行される。具体的には、サーバ 1 4 は、当該ファイルの再生が終了するまで、音声データと再生指示と

50

を相手側の対話装置 1 2 に送信する。これに応じて、対話装置 1 2 は、受信した音声データに基づいて、当該音声を出力する。

【 0 0 6 6 】

このようにして、両対話者の発話が重複した場合には、後から発話し始めた側の音声を録音し、その後両方の発話が終了したときに、当該録音音声を相手側で出力することができる。なお、両発話が同時に始まった場合には優先順位に従って音声を遅延再生できる。したがって、重複したときの発話の出力を遅らせることができるので、円滑な遠隔地間対話を成立させることができる。

【 0 0 6 7 】

図 5 には、サーバ 1 4 に記憶される発話フラグテーブルの一例が示される。発話フラグテーブルでは、検出時刻 t ごとに、ユーザ、状態フラグ、対象、保存音声ファイルの記憶位置、保存コマンドファイルの記憶位置、および P L A Y フラグ等の情報が記憶される。ユーザ情報は、当該データの主体であり、たとえば A は当該データが対話装置 1 2 a の状態であることを意味し、B は当該データが対話装置 1 2 b の状態であることを意味する。また、対象は、ユーザの発話対象を示す。

【 0 0 6 8 】

状態フラグは、対話装置 1 2 での音声取得状態および音声出力状態を示し、上述のように、S P E A K I N G フラグ、S I L E N T フラグ、R E C O R D I N G フラグが記憶される。なお、図 5 の時刻 $t = T + 2$ T では、状態フラグは I N T E R P O L A T I N G フラグである。上述のように、両対話者の状態フラグが S I L E N T フラグであった場合において、間パターンに従って間機能言葉が挿入されたときには、当該時刻の状態フラグとして、この I N T E R P O L A T I N G フラグが上書きされるようになっている。これによって、当該検出時刻が、対話における空白時間としては計測されなくなる。

【 0 0 6 9 】

保存音声ファイルは、録音された音声ファイルの記憶位置を示している。たとえば、図 5 では、時刻 $T + 4$ T および時刻 $T + 5$ T において、ユーザ A の状態フラグとして R E C O R D I N G フラグが記憶されており、当該時刻の録音に対応する音声ファイルの保存場所が示されている。また、時刻 $T + 4$ T での P L A Y フラグは 1 であり、録音が開始されたことを意味し、次の時刻 $T + 5$ T での P L A Y フラグは 2 であり、録音が継続されていることを意味し、その次の時刻 $T + 6$ T での P L A Y フラグは 2 のままであり、録音が終了されていることを意味する。

【 0 0 7 0 】

なお、保存コマンドファイルは、録音が行われている間に、ユーザによって当該対話装置 1 2 a で入力された動作コマンドを記録したファイルの記憶位置を示している。このコマンドファイルは音声ファイルと一緒に相手側対話装置 1 2 b に送信され、したがって、対話装置 1 2 b では、録音した音声とともに入力指示された身振りが実行される。

【 0 0 7 1 】

図 6 から図 9 には、対話装置 1 2 の C P U の入力処理における動作の一例が示される。入力処理を開始すると、図 6 の最初のステップ S 1 では、初期化が行われる。たとえば、S E N D フラグがオフされ、R E C O R D フラグがオフされ、また、時刻（またはフレーム番号） t に現在の時刻 T（または初期値 T）が代入される。続くステップ S 3 から図 9 のステップ S 6 9 までの処理は一定時間 T ごとに、たとえば 1 フレームごとに繰り返し実行される。

【 0 0 7 2 】

ステップ S 3 では、マイク 1 6 または 2 0 の入力をチェックし、ステップ S 5 で当該入力データに基づいて、音声入力があるか否かを判断する。ステップ S 5 で “ Y E S ” であれば、つまり、対話者が発話している場合には、ステップ S 7 で、サーバ 1 4 の発話フラグテーブルを参照する。たとえば、対話装置 1 2 は発話フラグの要求をサーバ 1 4 に送信する。サーバ 1 4 はこれに応じて発話フラグテーブルのデータを当該対話装置 1 2 に送信する。対話装置 1 2 は発話フラグテーブルデータを受信してメモリに記憶する。なお、開

10

20

30

40

50

始後には音声入力の無い状態が続くので、最初の発話の前には発話フラグテーブルには両対話者の状態としてSILENTフラグが記憶されている。

【0073】

続いて、ステップS9で、間計測処理を実行する。この間計測処理の動作の一例は図10に詳細に示される。間計測処理を開始すると、図10の最初のステップS81では、発話フラグテーブルに基づいて、現時刻 t の T 前の時刻における自分の状態フラグがSILENTフラグであるか否かを判断する。このステップS81では、現在の検出時刻で音声入力があり、かつ、前回の検出時刻で音声入力がなかったか否かを判断している。つまり、この対話装置12側のユーザが話し始めたタイミングであるか否かを判断している。

【0074】

ステップS81で“YES”であれば、今回話し始めるまでの2種類の間を発話フラグテーブルに基づいて計測する。具体的には、ステップS83で、自分が前に言葉を話し終えてから話し始めるまでの空白時間を計測する。また、ステップS85で、相手が言葉を話し終えてから自分が話し始めるまでの空白時間を計測する。そして、ステップS87で、計測データをサーバ14に送信する。これに応じて、サーバ14は間計測データを記憶する。ステップS87を終了し、または、ステップS81で“NO”である場合には、処理は図6のステップS11に戻る。

【0075】

このようにして、間の計測データの履歴をサーバ14で記憶していくことによって、対話者がどのような間を取りながら対話を行っているかをサーバ14で記録することができる。この間の履歴データと音声解析履歴データから、間のパターンを抽出することができる。

【0076】

続いて、図6のステップS11では、発話フラグテーブルに基づいて、 T 前のときの対話相手のフラグがSPEAKINGフラグであるか否かを判断する。このステップでは、相手が話しているのに、この対話装置12側の対話者も発話をしているのか否かを判断している。ステップS11で“YES”であれば、つまり、両対話者の発話が重複した場合には、ステップS13で、RECORDフラグはオンであるか否かを判断する。

【0077】

ステップS13で“NO”であれば、つまり、発話の重複が始まったばかりである場合には、ステップS15で、音声の録音を開始し、取得した音声データを音声ファイル化してメモリに記憶する。たとえば、音声データはPCM方式データであり、音声ファイルはWAVE形式であってよい。なお、音声データを送信前に適宜な方式で圧縮し、再生前に復号するようにしてよい。また、ステップS17で、メモリのRECORDフラグをオンにして、録音中であることを記憶する。

【0078】

なお、相手が先に話し始めている場合には、対話装置12は相手側から音声を受信してスピーカから出力しているので、この対話装置12側の対話者は、通常は無理に発話を続けずに、自分の発話を止めて相手の音声を聞くと考えられる。このため、録音される音声は非常に短時間のものになると考えられるので、この実施例では、音声は録音完了後に一括してサーバ14へ送信するようにしている。しかし、他の実施例では、その都度サーバ14に音声を送信するようにしてもよい。

【0079】

一方、ステップS13で“YES”であれば、つまり、既に録音を開始している場合には、ステップS19で、ステップS15で開始された録音を継続し、取得した音声データを音声ファイルに記憶する。

【0080】

ステップS17またはS19を終了すると、ステップS21で、サーバ14の発話フラグテーブルにRECORDINGフラグを記録する。具体的には、対話装置12は、時刻 t 、発話者(この対話装置12の識別情報)、対象(相手側対話装置12の識別情報)等

10

20

30

40

50

の情報とともに、録音中であることを示す情報（RECORDINGフラグ）をサーバ14に送信する。これに応じて、サーバ14は、受信した情報に基づいて、発話フラグテーブルに、時刻、発話者、対象およびRECORDINGフラグを記憶する。

【0081】

なお、システム10が3つ以上の対話装置12を含む場合、発話の対象（相手側対話装置12の識別情報）を入力装置の操作等によって選択できるようにしてもよい。

【0082】

さらに、ステップS23で、サーバ14の発話フラグテーブルのPLAYフラグに、T前の値に1を加算した値を記録する。具体的には、対話装置12は、時刻t、発話者および対象等の情報とともに、PLAYフラグの増加指示をサーバ14に送信する。これに応じて、サーバ14は、発話フラグテーブルのPLAYフラグの時刻tの1つ前の値を読み出して、この値に1を加算し、当該算出値を時刻tのPLAYフラグの値として記憶する。未再生の音声ファイルが残っていない状態で録音が始まったときは、PLAYフラグに1が記憶され、録音が継続中である限りPLAYフラグの値は時刻tの進行に合わせて1つずつ増加される。ステップS23を終了すると、処理は図9のステップS59に進む。

10

【0083】

一方、ステップS11で“NO”である場合には、処理は図7のステップS25に進む。つまり、この対話装置12側の対話者が発話している場合において、1つ前の時刻で相手側が発話の無い状態、または録音中であるときは、この対話装置12側の対話者の音声を相手に聞かせる。

20

【0084】

また、ステップS5で“NO”である場合には、つまり、この対話装置12側で対話者が発話していない場合には、処理は図8のステップS35に進む。

【0085】

図7のステップS25では、サーバ14の発話フラグテーブルにSPEAKINGフラグを記録する。具体的には、対話装置12は、時刻t、発話者、対象等の情報とともに、発話中であることを示す情報（SPEAKINGフラグ）をサーバ14に送信する。これに応じて、サーバ14は、受信した情報に基づいて、発話フラグテーブルに、時刻、発話者、対象およびSPEAKINGフラグを記憶する。

30

【0086】

続くステップS27で、取得した音声データを音声ファイル化して、メモリにローカルファイルとして記憶する。また、ステップS27で、取得した音声データとその再生指示を相手側の対話装置12に直接（すなわちサーバ14を介さずに）送信する。相手側の対話装置12は、音声データと再生指示を受信すると、当該音声データの再生処理を実行して、当該音声をスピーカ18または22から出力する。このようにして、この対話装置12側のみで発話が行われている場合、あるいは相手側が録音中である場合には、この対話装置12側の音声がローカルファイルに記録されつつ相手側に直接送信され、相手側の対話装置12で直ちに当該音声再生されて出力される。

【0087】

40

また、ステップS31では、メモリのSENDフラグをオンにして、送信中であることを記憶する。さらに、ステップS33では、サーバ14の発話フラグテーブルのPLAYフラグに、T前の値をそのまま記録する。具体的には、対話装置12は、時刻t、発話者および対象等の情報とともに、PLAYフラグの維持指示をサーバ14に送信する。これに応じて、サーバ14は、発話フラグテーブルのPLAYフラグの時刻tの1つ前の値を読み出して、この値を時刻tのPLAYフラグの値として記憶する。このように、録音中でない場合には、PLAYフラグの値として前回の値が維持される。ステップS33を終了すると、処理は図9のステップS59に進む。

【0088】

この対話装置12で音声入力が行われていない場合には、図8のステップS35で、R

50

RECORDフラグがオンであるか否かを判断する。ステップS35で“YES”であれば、つまり、1つ前の時刻まで録音が行われていた場合には、ステップS37で、音声ファイルへの音声の録音を終了する。また、ステップS39で、音声の録音中に入力装置を用いて入力された動作コマンドのコマンドファイルへの記録を終了する。さらに、ステップS41で、メモリのRECORDフラグをオフにする。

【0089】

そして、ステップS43で、録音した音声ファイルに対する音声解析処理を実行する。この音声解析処理の動作の一例が図11に詳細に示される。なお、図8のステップS53で実行される音声解析処理の動作も同じである。

【0090】

音声解析処理を開始すると、図11のステップS91で、メモリに録音された音声ファイルを読み込む。なお、図8のステップS53で実行される場合には、このステップS91では、ローカルファイルの音声データを読み込む。

【0091】

次に、ステップS93で、読み込んだ音源の基本周波数(ピッチ)および振幅を算出する。また、ステップS95では、音声データを音節に分割する処理を試みる。そして、ステップS97で、分割した音節が存在するか否かを判断する。ステップS97で“YES”であれば、続くステップS99で、当該音節の持続時間を算出する。さらに、当該音節の持続時間の平均を算出する。ステップS99を終了すると、ステップS97に戻って、分割した音節が残っている場合には、当該音節についてステップS99の処理を繰返す。ステップS97で“NO”であれば、ステップS101で、音声解析データをサーバ14に送信する。したがって、音声解析データは、基本周波数、振幅、および音節の平均持続時間等の情報を含む。この音声解析データは、たとえば、時刻、発話者、対象等の情報に対応付けられてサーバ14に送信される。これに応じて、サーバ14は、受信した音声解析データを音声解析履歴DB92に記憶する。このようにして、発話音声の特徴が抽出されて、その履歴が記録される。ステップS101を終了すると、この音声解析処理を終了して、図8のステップS45(ステップS43の場合)、またはステップS55(ステップS53の場合)へ戻る。

【0092】

ステップS43を終了すると、ステップS45で、録音した音声ファイルとコマンドファイルとをサーバ14に送信する。音声ファイルとコマンドファイルとは、たとえば時刻、発話者、対象等の情報に対応付けられてサーバ14に送信される。これに応じて、サーバ14は、受信した音声ファイルとコマンドファイルとをメモリの所定領域に保存する。発話フラグテーブルでは、録音された時刻の保存音声ファイル情報として、音声ファイルの記憶位置が登録されるとともに、同時刻の保存コマンドファイル情報として、コマンドファイルの記憶位置が登録される(図5参照)。ステップS45を終了すると処理はステップS55に進む。

【0093】

一方、ステップS35で“NO”であれば、ステップS47でSENDフラグがオンであるか否かを判断する。ステップS47で“YES”であれば、つまり、1つ前の時刻まで音声を相手側の対話装置12に送信していた場合には、ステップS49で、音声のローカルファイルへの記録を終了し、ステップS51で、メモリのSENDフラグをオフにする。そして、ステップS53で、ローカルファイルの音声データに対して、上述のような図11の音声解析処理を実行する。ステップS53を終了すると、または、ステップS47で“NO”である場合には、処理はステップS55へ進む。

【0094】

ステップS55では、サーバ14の発話フラグテーブルにSILENTフラグを記録する。具体的には、対話装置12は、時刻t、発話者、対象等の情報とともに、音声入力が無いことを示す情報(SILENTフラグ)をサーバ14に送信する。これに応じて、サーバ14は、受信した情報に基づいて、発話フラグテーブルに、時刻、発話者、対象およ

10

20

30

40

50

び SILENT フラグを記憶する。

【0095】

また、ステップ S 5 7 で、図 7 のステップ S 3 3 と同様にして、サーバの発話フラグテーブルの PLAY フラグに、T 前の値をそのまま記録する。ステップ S 5 7 を終了すると、処理は図 9 のステップ S 5 9 へ進む。

【0096】

図 9 のステップ S 5 9 から S 6 7 では、相手側の対話装置 1 2 がロボット 1 2 b である場合の処理である。したがって、相手側がロボット 1 2 b でない場合には、これらの処理は行われなくてよい。

【0097】

図 9 のステップ S 5 9 では、動作コマンドの入力をチェックする。具体的には、入力装置からの入力データを取得して、ロボット 1 2 b の身振りのための動作コマンドが選択されたか否かを判定する。たとえば、動作コマンドはディスプレイに選択可能なりストとして表示されてよい。なお、この対話装置 1 2 がロボット 1 2 b である場合には、入力装置とディスプレイを設ける必要がある。

【0098】

そして、ステップ S 6 1 で、選択された動作コマンドがあるかどうかを判断し、“YES”であれば、ステップ S 6 3 で、メモリの RECORD フラグがオンであるか否かを判断する。ステップ S 6 3 で“YES”であれば、つまり、音声録音中の場合には、ステップ S 6 5 で、動作コマンドをメモリのコマンドファイルに記憶する。このように、音声を録音している場合には、動作コマンドの入力も同時に記録して、録音終了後に上述のステップ S 4 5 でサーバ 1 4 に送信するようにしているので、両対話者の発話が重複した場合には、発話と身振りに対して同時に遅延を与えてから相手側で再生することができる。

【0099】

一方、ステップ S 6 3 で“NO”であれば、ステップ S 6 7 で、動作コマンドと再生指示とを相手側の対話装置 1 2 b に直接送信する。相手側の対話装置 1 2 b は、動作コマンドと再生指示を受信すると、当該動作コマンドに対応するプログラムおよびデータに従って動作し、その身振りを実行する。

【0100】

ステップ S 6 5 または S 6 7 を終了したとき、またはステップ S 6 1 で“NO”の場合には、ステップ S 6 9 で、所定時間 T (たとえば 1 フレーム) を加算することで時刻 (あるいはフレーム番号) t を更新する。そして、図 6 のステップ S 3 に戻って、次の時刻 t における処理を繰返す。このようにして、対話装置 1 2 では、この対話装置 1 2 側の対話者の発話の状態および相手側の発話の状態に応じた処理が実行される。

【0101】

図 1 2 にはサーバ 1 4 の継続促進処理における動作の一例が示される。また、図 1 3 には、サーバ 1 4 の遅延再生処理における動作の一例が示される。

【0102】

なお、サーバ 1 4 の他の処理、たとえば受信処理、発話フラグテーブルの作成処理および送信処理などのフロー図は省略する。サーバ 1 4 は上述のような各処理を並列的に実行している。サーバ 1 4 は、上述のように、対話装置 1 2 からデータを受信したときは、当該データをメモリに記憶し、必要に応じて当該データに対応する所定の処理を実行する。たとえば、サーバ 1 4 は、対話装置 1 2 から発話や処理の状態に関するデータを受信したときは発話フラグテーブルを作成する。音声ファイルおよび動作コマンドファイル等を受信したときは、これらのファイルを記憶するとともに、発話フラグテーブルに記憶位置を書き込む。音声解析データを受信したときは、当該データを音声解析履歴 DB 9 2 に記憶する。また、対話装置 1 2 から発話フラグテーブルの要求があったときは、当該対話装置 1 2 に発話フラグテーブルを送信する。

【0103】

図 1 2 に示す継続促進処理では、サーバ 1 4 の CPU は、ステップ S 1 1 1 で初期化を

10

20

30

40

50

実行し、たとえば変数 t に初期値 T を設定する。この初期値 T は発話フラグテーブルの時刻 t の最初の値 T であり、つまり、対話装置 1 2 における時刻 t の初期値 T である。したがって、継続促進処理は発話フラグテーブルの作成後に実行される。続くステップ S 1 1 3 から S 1 3 5 の処理をサーバ 1 4 の CPU は一定時間 T ごとに、たとえば 1 フレームごとに繰り返し実行する。

【 0 1 0 4 】

ステップ S 1 1 3 では、メモリの発話フラグテーブルを参照する。たとえば現時刻 t のデータを読み出す。そして、ステップ S 1 1 5 で、対話者同士で S I L E N T フラグであるか否かを判断する。たとえば、現時刻 t においてユーザと発話対象が互いに対になっている両対話者が存在しており、かつ、当該両対話者の状態フラグが S I L E N T フラグであることを判定する。たとえば図 5 では、時刻 $T + T$ のときがこの状態に相当する。

10

【 0 1 0 5 】

ステップ S 1 1 5 で “ Y E S ” であれば、つまり、対話において無音状態になっている場合には、ステップ S 1 1 7 で、空白時間を算出する。たとえば、現時刻 t 以前の発話フラグテーブルのデータを読み出して、現時刻 t から遡って両対話者のどちらかの状態フラグが S I L E N T フラグでなくなるまでに掛かった時間（またはフレーム数）を算出する。

【 0 1 0 6 】

続いて、ステップ S 1 1 9 で、音声解析履歴 DB 9 2 から対話者らの最新のデータを抽出する。具体的には、現時刻 t に最も近い時刻の発話者の音声解析データから、基本周波数、振幅および音節の平均持続時間等を読み出す。このように、ステップ S 1 1 7 と S 1 1 9 で、少なくとも空白時間と当該空白前の発話音声の特徴を含む間の状況が検出される。

20

【 0 1 0 7 】

そして、ステップ S 1 2 1 で、現在の間の状況と間パターンとの照合を実行して、ステップ S 1 2 3 で、現在の対話の間の状況にマッチする間パターンがあるか否かを判断する。上述の図 4 のように、間パターンデータ内には、空白時間 (t) および条件 (I) 設定されているので、このような間パターンに合う空白時間および発話音声の特徴（基本周波数、振幅、音節の平均持続時間など）を有する間の状況（すなわち、最終発話者の発話後の無音状態）が生じているか否かを判定する。マッチする間パターンがある場合には、当該間パターンに対応する間機能言葉を選択する。また、間パターンデータに設定されている最終発話者と発話者との関係（相手か自分か）に基づいて、間機能言葉を発話させる対話装置 1 2 を特定する。

30

【 0 1 0 8 】

ステップ S 1 2 3 で “ Y E S ” であれば、つまり、現在の対話における間の状況が、間パターンに基づく間を挿入すべき状況になっていると判定される場合には、ステップ S 1 2 5 で、選択した間機能言葉の音声ファイルをメモリの作業領域に読み出して、ピッチ、抑揚パターンを調整して、当該調整した間機能言葉の音声ファイルを生成する。これによって、発話者の発話の特徴（たとえば、高揚した口調、淡々とした発話など）に合わせた間機能言葉を出力することが可能になる。したがって、会話に合成音声が入力されても対話者に違和感をさほど覚えさせないようにすることができるし、また、それまでの会話の調子や流れを継続させることができる。

40

【 0 1 0 9 】

また、ステップ S 1 2 7 で、選択された間機能言葉に適した動作コマンドを選択する。この実施例では、間パターンデータにおいて、間機能言葉に対応する動作コマンドが登録されているので、当該動作コマンドを選択する。

【 0 1 1 0 】

そして、ステップ S 1 2 9 で、音声ファイルと動作コマンドファイルを、発話させる対話装置 1 2 に送信する。ファイル送信後、ステップ S 1 3 1 で、音声と動作の再生指示を同じ対話装置 1 2 に送信する。これによって、対話における無音領域に言葉や身振りを挿

50

入することができる。なお、その対話装置 1 2 は、音声ファイルの再生を実行し、当該音声出力する。また、対話装置 1 2 がロボット 1 2 b である場合には、さらに動作コマンドの再生を実行し、当該動作コマンドに対応する身振りを行う。

【 0 1 1 1 】

さらに、ステップ S 1 3 3 で、発話フラグテーブルにおいて、現時刻 t の状態フラグに INTERPOLATING フラグを上書きする (図 5 の時刻 $T + 2 T$ を参照)。これによって、以降のステップ S 1 1 7 では、当該時刻 t が無音であるとは見なされないようにすることができる。

【 0 1 1 2 】

一方、ステップ S 1 2 3 で “ NO ” である場合には、つまり、未だ、間パターンに従った間を与える必要がない場合には、処理はそのままステップ S 1 3 5 に進む。ステップ S 1 3 5 では、所定時間 T (たとえば 1 フレーム) を加算することで時刻 (あるいはフレーム番号) t を更新する。なお、このサーバ 1 4 における T は対話装置 1 2 における T と同一である。そして、ステップ S 1 1 3 に戻って、次の時刻 t における処理を繰り返す。このようにして、対話において無音が検出された場合には、必要に応じて言葉や身振りを挿入することによって、無音時間を適切な間に変えることができる。

【 0 1 1 3 】

図 1 3 に示す遅延再生処理では、サーバ 1 4 の CPU は、ステップ S 1 5 1 で初期化を実行する。たとえば、PLAYING フラグをオフにする。PLAYING フラグは録音された音声ファイルおよび動作コマンドファイルを再生中であるか否かを示す。また、図 1 2 の継続促進処理と同様に、変数 t に初期値 T を設定する。この初期値 T は発話フラグテーブルの時刻 t の最初の値 T であり、つまり、対話装置 1 2 における時刻 t の初期値 T である。したがって、この遅延再生処理も発話フラグテーブルの作成後に実行される。続くステップ S 1 5 3 から S 1 7 9 の処理をサーバ 1 4 の CPU は一定時間 T ごとに、たとえば 1 フレームごとに繰り返し実行する。

【 0 1 1 4 】

ステップ S 1 5 3 では、メモリの発話フラグテーブルを参照する。ステップ S 1 5 5 で、メモリの PLAYING フラグがオンであるか否かを判断する。ステップ S 1 1 5 で “ NO ” であれば、つまり、再生中ではない場合には、ステップ S 1 5 7 で、現時刻 t における両対話者のどちらかの状態フラグとして SPEAKING フラグがあるか否かを判断する。ステップ S 1 5 7 で “ YES ” の場合、一方が発話をしており、その音声が他方の対話装置 1 2 から出力されているはずである。したがって、遅延再生は行わず処理はステップ S 1 7 9 に進む。

【 0 1 1 5 】

一方、ステップ S 1 5 7 で “ NO ” であれば、つまり、両対話装置 1 2 で音声が出力されていない場合には、ステップ S 1 5 9 で、両対話者のどちらかの PLAY フラグが 1 以上であるか否かを判断する。ステップ S 1 5 9 で “ NO ” であれば、録音されたが未再生である音声ファイルが存在しないので、処理はそのままステップ S 1 7 9 に進む。

【 0 1 1 6 】

しかし、ステップ S 1 5 9 で “ YES ” であれば、つまり、録音されたが未再生の音声ファイルが残っている場合には、ステップ S 1 6 1 で、PLAY フラグが 1 である時刻 t が早いユーザを発話フラグテーブルから参照する。つまり、録音を開始した時刻が早いユーザを特定する。なお、録音の開始が両対話者で同時刻である場合には、予め設定しておいた優先順位 (たとえば $B > A$) に基づいて、ユーザを特定する。

【 0 1 1 7 】

続いて、ステップ S 1 6 4 で、再生のための設定を実行し、変数 F に 1 を設定し、変数 U に特定したユーザを設定する。変数 F は音声再生のためのフレームカウンタである。また、ステップ S 1 6 5 で、メモリの PLAYING フラグをオンにして、再生中であることを記憶する。そして、ステップ S 1 6 7 で、変数 U の PLAY フラグが変数 F の値である音声および動作を再生する。具体的には、当該音声ファイルを読み出して、当該ユーザ

10

20

30

40

50

の相手側の対話装置 1 2 に音声ファイルと再生指示とを送信する。なお、当該動作コマンドファイルも保存されている場合には、当該動作コマンドファイルも読み出して、音声ファイルと一緒に相手側の対話装置 1 2 に送信する。これに応じて、当該対話装置 1 2 は、音声ファイルおよび動作コマンドファイルを記憶するとともに、その再生を実行する。これによって、音声スピーカ 1 8 または 2 2 から出力され、動作コマンドもあった場合には、当該身振りも実行される。このようにして、録音されていた音声および記憶されていた動作の再生が開始される。

【 0 1 1 8 】

ステップ S 1 6 7 を終了すると、処理はステップ S 1 7 9 へ進む。ステップ S 1 7 9 では、時刻 t に所定時間 T が加算されて時刻 t が更新される。ステップ S 1 7 9 を終了すると、処理はステップ S 1 5 3 へ戻って、次の時刻 t における処理を繰り返す。

10

【 0 1 1 9 】

再生が開始されると、ステップ S 1 5 5 で “ Y E S ” と判断され、続くステップ S 1 6 9 で、時刻 t における変数 U の P L A Y フラグが変数 F の値に等しいか否かを判断する。上述のように、録音が終了した場合には、P L A Y フラグの値は前時刻の値を維持するので、このステップ S 1 6 9 では、再生中の音声ファイルの再生を完了したか否かを判定している。

【 0 1 2 0 】

ステップ S 1 6 9 で “ N O ” であれば、つまり、音声ファイルの再生が未だ完了していない場合には、ステップ S 1 7 1 で、変数 F をインクリメントする。その後、ステップ S 1 7 3 で、変数 U の P L A Y フラグが変数 F の値である音声および動作を再生する。これによって、上述のステップ S 1 6 7 と同様にデータが送信され、次のフレームの音声および動作が対話装置 1 2 で再生される。ステップ S 1 7 3 を終了すると、処理はステップ S 1 7 9 へ進む。

20

【 0 1 2 1 】

一方、ステップ S 1 6 9 で “ Y E S ” であれば、つまり、音声ファイルの再生を完了した場合には、ステップ S 1 7 5 で、メモリの P L A Y I N G フラグをオフにする。また、ステップ S 1 7 7 で、変数 U の P L A Y フラグの値を全て変数 F の値だけ減算する。なお、減算の結果、値が負になったとき、当該 P L A Y フラグの値は 0 に設定される。これによって、再生された変数 U および時刻 t の P L A Y フラグの値がすべて 0 になる。また、当該変数 U のユーザの未再生の音声ファイルが存在する場合には、当該ユーザの最も古く録音された音声ファイルのうち最も早い時刻の P L A Y フラグの値が 1 になる。したがって、次回は、当該未再生の音声を再生することが可能になる。ステップ S 1 7 7 を終了すると、処理はステップ S 1 7 9 へ進む。

30

【 0 1 2 2 】

このようにして、両対話者の発話の重複によって録音された音声および記録された動作コマンドを、後から再生することができる。

【 0 1 2 3 】

図 1 4 には、対話装置 1 2 の C P U の出力処理の動作の一例が示される。この出力処理は上述の図 6 から図 9 の入力処理と並列的に実行される。また、この出力処理は一定時間ごと、たとえば 1 フレームごとに繰り返し実行される。

40

【 0 1 2 4 】

ステップ S 1 9 1 では、音声を受信したか否かが判断され、“ Y E S ” であれば、ステップ S 1 9 3 で、受信した音声ファイルないし音声データをメモリに記憶する。

【 0 1 2 5 】

続いて、ステップ S 1 9 5 では、動作コマンドを受信したか否かが判断され、“ Y E S ” であれば、ステップ S 1 9 7 で、受信した動作コマンドファイルをメモリに記憶する。

【 0 1 2 6 】

続いて、ステップ S 1 9 9 では、再生指示を受信したか否かが判断され、“ Y E S ” であれば、ステップ S 2 0 1 で、音声を再生する。具体的には、対話装置 1 2 の C P U は、

50

受信した音声ファイルを再生を開始し、当該音声データを音声入出力ボードに与えてスピーカから当該音声を出力する。また、当該対話装置 12 が身体動作機能を有する対話装置 12 b である場合には、ステップ S 203 で、動作を再生する。具体的には、当該動作コマンドに従って対応する身振りを実行する。動作コマンドに対応する身振りを実行するためのプログラムおよび制御データは、対話装置 12 b のメモリ 70 に予め記憶されている。CPU 66 は動作コマンドに対応するプログラムに従って制御データをモータ制御ボード 72 に与えて、対応するモータを制御する。これによって対応する身体部位が動かされて所定の身振りが表現される。

【0127】

なお、上述の実施例では、両対話者の発話が重複したとき、後から発話された方の音声を録音して、その後どちらも発話しなくなってから、当該録音音声を相手側で出力するようにしていた。しかし、他の実施例では、両対話者の発話が重複したときには、後から発話された方の音声をキャンセルするようにしてもよい。

【0128】

また、上述の各実施例では、間機能言葉の音声データをサーバ 14 が記憶しておいて、サーバ 14 から対話装置 12 に送信するようにしていた。しかし、他の実施例では、間機能言葉の音声データを各対話装置 12 に予め記憶させておいて、サーバ 14 から再生すべき間機能言葉を指定する情報を送信するようにしてもよい。

【0129】

また、上述の各実施例では、システム 10 は、身体動作機能を有しない対話装置 12 a と身体動作機能を有する対話装置 12 b とを含んでいた。しかし、他の実施例では、身体動作機能を有しない対話装置 12 a のみが使用されてよく、この場合には、動作コマンド関連の処理が不要である。逆に、身体動作機能を有する対話装置 12 b のみが使用されてもよい。

【0130】

また、上述の各実施例では、システム 10 は対話装置 12 とは別に各対話装置 12 の音声取得状態および音声出力状態を示す情報（すなわち発話フラグテーブル）を管理するサーバ 14 を備えた。しかし、他の実施例では、サーバ 14 を別途に設けずに、サーバ 14 の機能（発話フラグテーブルの管理、継続促進処理、遅延再生処理など）を一方の対話装置 12 に備えさせるようにしてよいし、あるいは 2 つの対話装置 12 に分散して備えさせるようにしてもよい。

【図面の簡単な説明】

【0131】

【図 1】この発明の一実施例の遠隔地間対話システムの構成を示す図解図である。

【図 2】身体動作機能を有する対話装置の外観の一例を示す図解図である。

【図 3】図 2 の対話装置の電気的な構成の一例を示すブロック図である。

【図 4】間パターン DB に記憶される間パターンデータの一例を示す図解図である。

【図 5】サーバに記憶される発話フラグテーブルの一例を示す図解図である。

【図 6】対話装置の入力処理の動作の一例の一部を示すフロー図である。

【図 7】図 6 の続きの一部を示すフロー図である。

【図 8】図 6 の続きの一部を示すフロー図である。

【図 9】図 6、図 7 および図 8 の続きを示すフロー図である。

【図 10】図 6 の間計測処理の動作の一例を示すフロー図である。

【図 11】図 8 の音声解析処理の動作の一例を示すフロー図である。

【図 12】サーバの継続促進処理の動作の一例を示すフロー図である。

【図 13】サーバの遅延再生処理の動作の一例を示すフロー図である。

【図 14】対話装置の出力処理の動作の一例を示すフロー図である。

【符号の説明】

【0132】

10 ... 遠隔地間対話システム

10

20

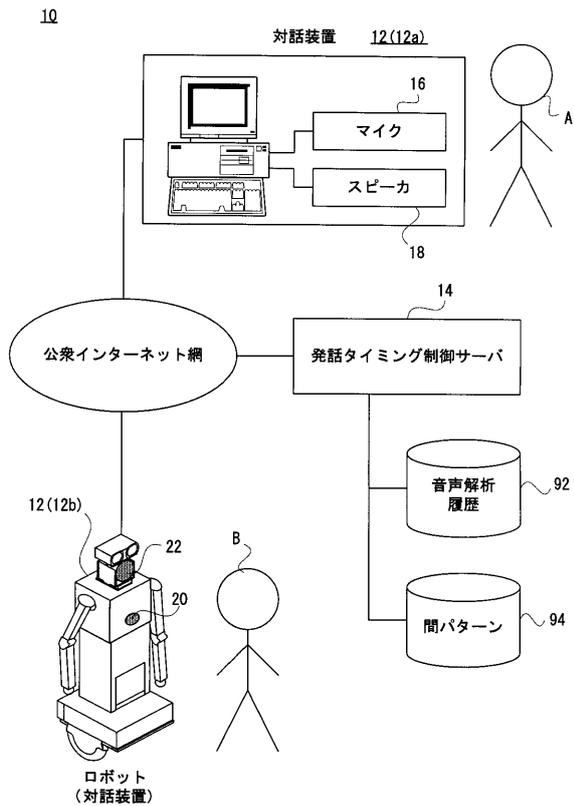
30

40

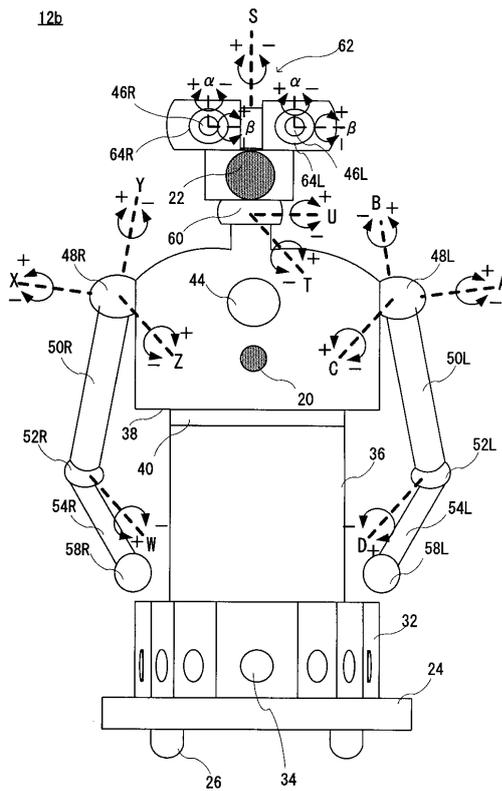
50

- 1 2 , 1 2 a , 1 2 b ... 対話装置
- 1 4 ... 発話タイミング制御サーバ
- 1 6 , 2 0 ... マイク
- 1 8 , 2 2 ... スピーカ
- 9 2 ... 音声解析履歴データベース
- 9 4 ... 間パターンデータベース

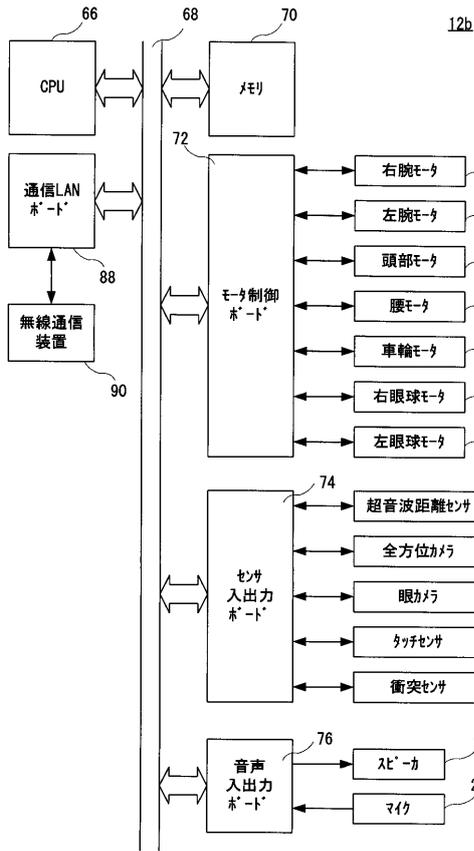
【 図 1 】



【 図 2 】



【図3】



【図4】

間パターンデータ

空白時間 (t) [sec]	最後の発話者	条件 (1)	間機能言葉	発話者	動作コマンド
$1 < t \leq 1.5$	A	$\omega > 150,$ $1000 < A \leq 1500,$ $30 < m \leq 50$	うんうん	B	うなづく
$0.5 < t \leq 2$	B	$\omega > 10,$ $50 < m \leq 100$	うーん	B	首傾げる
$t \leq 2$	A	$\omega > 300,$ $m \leq 30$	はいはい	B	うなづく
$t > 2$	B	$\omega > 80,$ $A \leq 1500$	えーと	A	視線を上方向ける
⋮	⋮	⋮	⋮	⋮	⋮

ω : 基本周波数 [Hz],
A : 振幅,
m : 音節の平均持続時間 [msec]

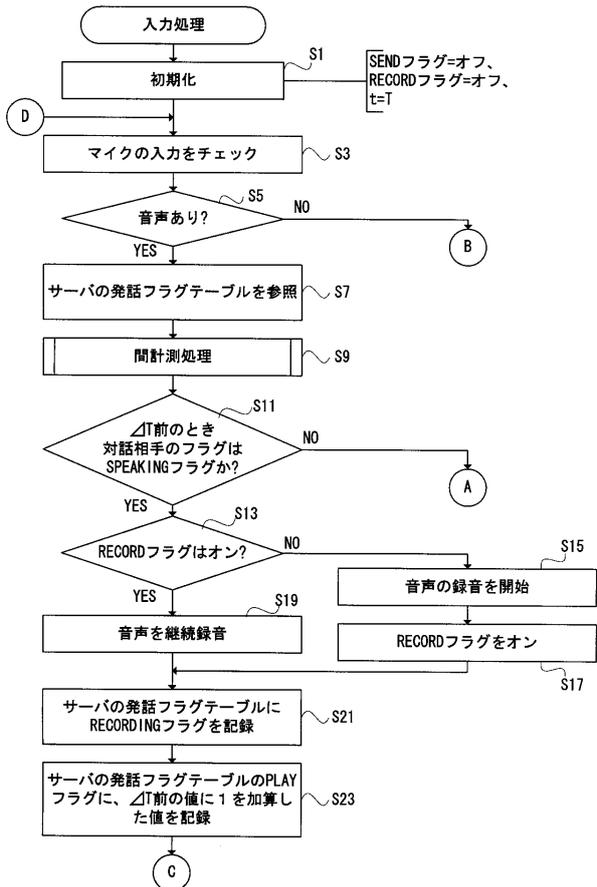
【図5】

発話フラグテーブル

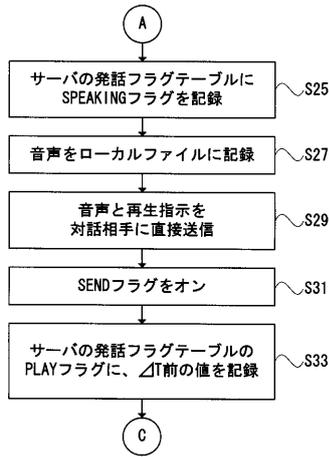
時刻 t	ユーザ	フラグ	対象	保存音声ファイル	保存コマンドファイル	PLAYフラグ
T	A	SPEAKING	B			0
T	B	SILENT	A			0
T+ΔT	A	SILENT	B			0
T+ΔT	B	SILENT	A			0
T+2ΔT	A	INTERPOLATING	B			0
T+2ΔT	B	INTERPOLATING	A			0
T+3ΔT	A	SILENT	B			0
T+3ΔT	B	SPEAKING	A			0
T+4ΔT	A	RECORDING	B	D:\voice%\v.wav	D:\command%\g.dat	1
T+4ΔT	B	SPEAKING	A			0
T+5ΔT	A	RECORDING	B	D:\voice%\v.wav	D:\command%\g.dat	2
T+5ΔT	B	SPEAKING	A			0
T+6ΔT	A	SILENT	B			2
T+6ΔT	B	SPEAKING	A			0
⋮	⋮	⋮	⋮	⋮	⋮	⋮

【図6】

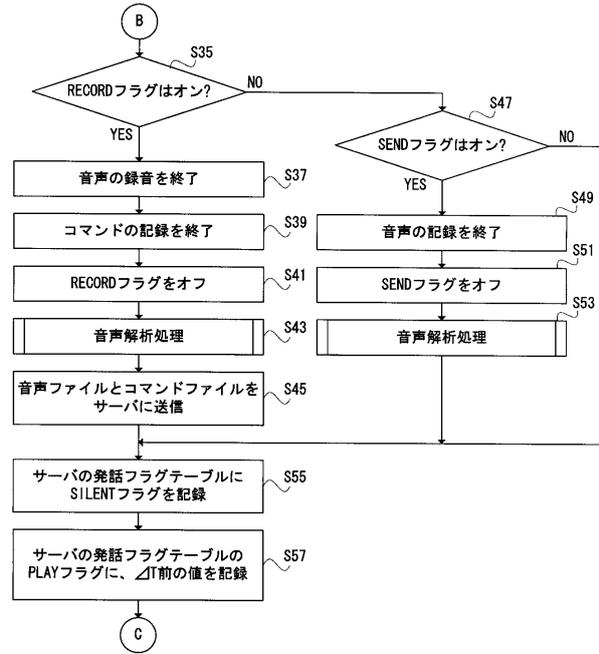
対話装置



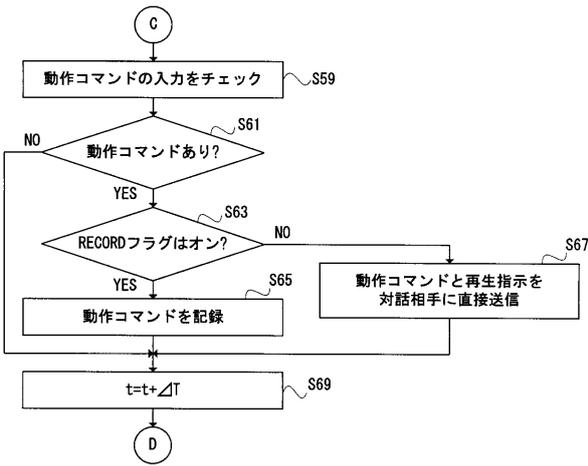
【図7】



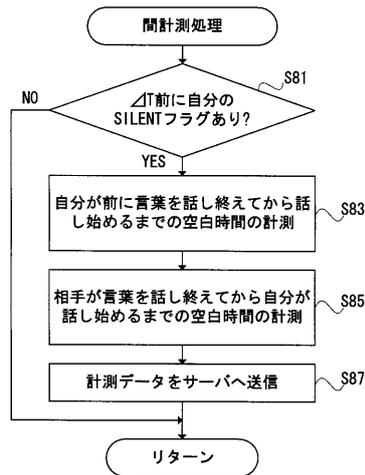
【図8】



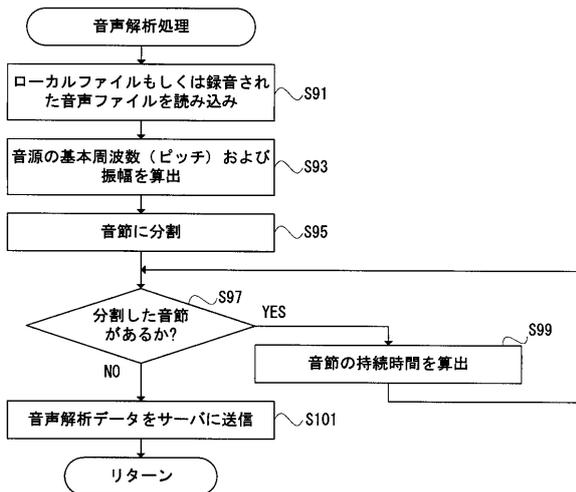
【図9】



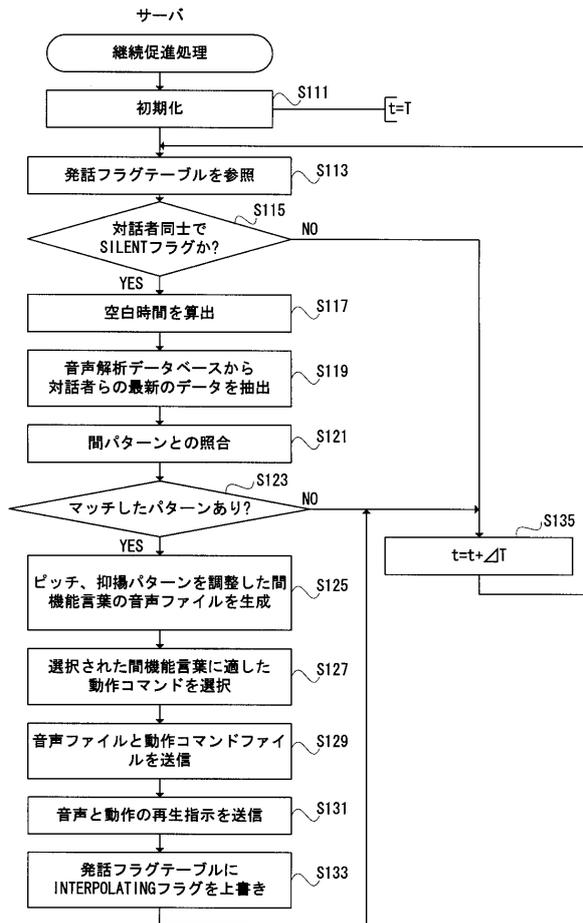
【図10】



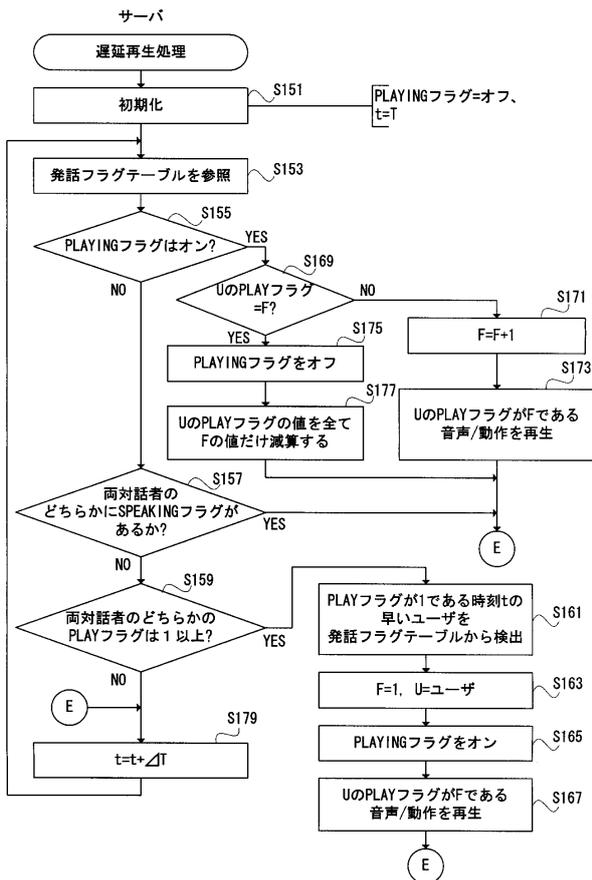
【図11】



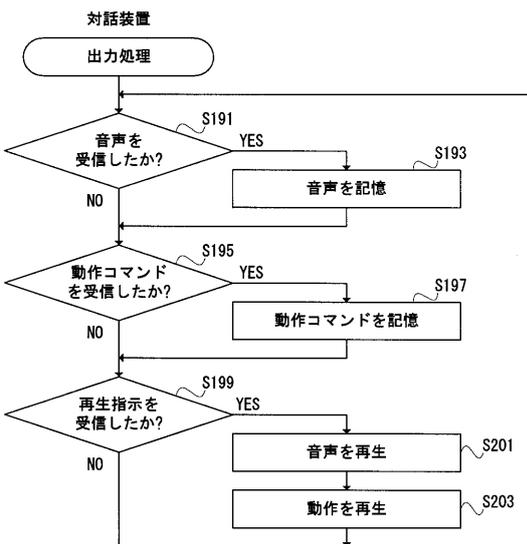
【図12】



【図13】



【図14】



フロントページの続き

(72)発明者 宮下 敬宏

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 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

審査官 松元 伸次

(56)参考文献 特開2004-064281(JP,A)

特開2005-012819(JP,A)

特開2003-304307(JP,A)

特開2001-134289(JP,A)

特開2000-349920(JP,A)

特開昭63-287126(JP,A)

(58)調査した分野(Int.Cl., DB名)

H04M 1/00、 1/24 - 1/253、
1/58 - 1/62、 1/66 - 3/00、
3/16 - 3/20、 3/38 - 3/58、
7/00 - 7/16、 11/00 - 11/10、 99/00