

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6686977号
(P6686977)

(45) 発行日 令和2年4月22日(2020.4.22)

(24) 登録日 令和2年4月6日(2020.4.6)

(51) Int. Cl.	F I
G 1 O L 21/0272 (2013.01)	G 1 O L 21/0272 1 0 0 Z
G 1 O L 15/00 (2013.01)	G 1 O L 15/00 2 0 0 H
G 1 O L 15/25 (2013.01)	G 1 O L 15/25

請求項の数 12 (全 25 頁)

(21) 出願番号	特願2017-123643 (P2017-123643)	(73) 特許権者	000001443
(22) 出願日	平成29年6月23日 (2017.6.23)		カシオ計算機株式会社
(65) 公開番号	特開2019-8134 (P2019-8134A)		東京都渋谷区本町1丁目6番2号
(43) 公開日	平成31年1月17日 (2019.1.17)	(74) 代理人	100074099
審査請求日	平成30年6月1日 (2018.6.1)		弁理士 大菅 義之
		(74) 代理人	100121083
			弁理士 青木 宏義
		(74) 代理人	100138391
			弁理士 天田 昌行
		(72) 発明者	山谷 崇史
			東京都羽村市栄町3丁目2番1号 カシオ
			計算機株式会社 羽村技術センター内
		(72) 発明者	中込 浩一
			東京都羽村市栄町3丁目2番1号 カシオ
			計算機株式会社 羽村技術センター内
			最終頁に続く

(54) 【発明の名称】 音源分離情報検出装置、ロボット、音源分離情報検出方法及びプログラム

(57) 【特許請求の範囲】

【請求項1】

音声を取得するために所定の指向性を有する音声取得手段と、
 前記音声取得手段により取得された前記音声から、所定の対象の信号音声の到来方向である第1方向を検出する第1方向検出手段と、
 前記音声取得手段により取得された前記音声から、雑音音声の到来方向である第2方向を検出する第2方向検出手段と、
 前記第1方向と前記第2方向とに基づいて、音源分離方向又は音源分離位置を検出する検出手段と、
前記音声取得手段が前記音声を取得するタイミングで前記所定の対象の口唇画像を取得する画像取得手段と、
前記口唇画像に基づいて、前記所定の対象の口唇の開口又は前記口唇の閉口を判定する判定手段と、
 を備え、
前記第1方向検出手段は、前記判定手段による前記口唇の開口の判定時に、前記音声取得手段により取得された前記音声を前記信号音声とし、
前記第2方向検出手段は、前記判定手段による前記口唇の閉口の判定時に、前記音声取得手段により取得された前記音声を前記雑音音声とする、
 ことを特徴とする音源分離情報検出装置。

【請求項2】

前記検出手段は、前記信号音声と前記雑音音声とから算出された信号対雑音比が閾値以下である場合、前記第1方向と前記第2方向とに基づいて、前記信号/雑音比が前記閾値を超える前記音源分離方向又は前記音源分離位置を検出する、ことを特徴とする請求項1に記載の音源分離情報検出装置。

【請求項3】

前記画像取得手段は、顔部画像を更に取得し、
前記口唇画像から前記所定の対象の前記口唇の移動量を取得する口唇移動量取得手段と

、
前記顔部画像から前記所定の対象の顔部の回転量を取得する顔部回転量取得手段と、
を更に備え、

前記判定手段は、前記口唇の移動量と前記顔部の回転量とに基づいて、前記所定の対象の前記口唇の開口又は前記口唇の閉口を判定する、
ことを特徴とする請求項1又は2に記載の音源分離情報検出装置。

【請求項4】

前記判定手段は、前記口唇の移動量のうちの前記口唇の開閉方向の移動量が第1の閾値を超え、且つ、前記口唇の移動量のうちの前記口唇の延伸方向の移動量が第2の閾値未満であり、且つ、前記顔部の回転量が第3の閾値未満であるときに、前記口唇の開口又は前記口唇の閉口を判定する、

ことを特徴とする請求項3に記載の音源分離情報検出装置。

【請求項5】

前記第1方向検出手段は、前記判定手段による前記口唇の開口の判定時に、前記信号音声の信号音声パワーに基づいて、前記第1方向を検出し、

前記第2方向検出手段は、前記判定手段による前記口唇の閉口の判定時に、前記雑音音声の雑音音声パワーに基づいて、前記第2方向を検出する、

ことを特徴とする請求項1乃至4の何れか1項に記載の音源分離情報検出装置。

【請求項6】

前記所定の対象にメッセージを報知する報知手段を更に備え、

前記報知手段は、前記所定の対象に現在位置から前記音源分離位置まで移動させるために、前記音源分離位置までの移動方向及び移動距離を含む前記メッセージを報知する、

ことを特徴とする請求項1乃至5の何れか1項に記載の音源分離情報検出装置。

【請求項7】

請求項1乃至6の何れか1項に記載の音源分離情報検出装置と、

自装置を移動する移動手段と、

前記自装置を動作する動作手段と、

前記音源分離情報検出装置、前記移動手段及び前記動作手段を制御する制御手段と、
を備える、

ことを特徴とするロボット。

【請求項8】

前記制御手段は、前記移動手段を制御して、前記音源分離位置に前記自装置を移動させる、

ことを特徴とする請求項7に記載のロボット。

【請求項9】

前記制御手段は、前記自装置が前記所定の対象とアイコンタクトを取りながら、又は、前記自装置が前記所定の対象の方を向きながら、前記音源分離位置に移動するように、前記動作手段を制御する、

ことを特徴とする請求項8に記載のロボット。

【請求項10】

前記制御手段は、前記自装置が前記音源分離位置に一気に移動するのではなく、少しだけ動いたり、回転のみをしたりして、前記音源分離位置に移動するように、前記移動手段及び前記動作手段を制御する、

10

20

30

40

50

ことを特徴とする請求項 8 又は 9 に記載のロボット。

【請求項 1 1】

音声を取得するために所定の指向性を有する音声取得手段により取得された前記音声から、所定の対象の信号音声の到来方向である第 1 方向を検出し、

前記音声取得手段により取得された前記音声から、雑音音声の到来方向である第 2 方向を検出し、

前記第 1 方向と前記第 2 方向とに基づいて、音源分離方向又は音源分離位置を検出し、

前記音声を取得するタイミングで前記所定の対象の口唇画像を取得し、

前記口唇画像に基づいて、前記所定の対象の口唇の開口又は前記口唇の閉口を判定し、

前記口唇の開口の判定時に、前記取得された前記音声を前記信号音声とし、

前記口唇の閉口の判定時に、前記取得された前記音声を前記雑音音声とする、

10

ことを含む、

ことを特徴とする音源分離情報検出方法。

【請求項 1 2】

音源分離情報検出装置のコンピュータを、

音声を取得するために所定の指向性を有する音声取得手段により取得された前記音声から、所定の対象の信号音声の到来方向である第 1 方向を検出し、

前記音声取得手段により取得された前記音声から、雑音音声の到来方向である第 2 方向を検出し、

前記第 1 方向と前記第 2 方向とに基づいて、音源分離方向又は音源分離位置を検出し、

前記音声を取得するタイミングで前記所定の対象の口唇画像を取得し、

前記口唇画像に基づいて、前記所定の対象の口唇の開口又は前記口唇の閉口を判定し、

前記口唇の開口の判定時に、前記取得された前記音声を前記信号音声とし、

前記口唇の閉口の判定時に、前記取得された前記音声を前記雑音音声とする、

20

ように機能させる、

ことを特徴とするプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音源分離情報検出装置、ロボット、音源分離情報検出方法及びプログラムに関する。

30

【背景技術】

【0002】

人間、動物等に模した形態を有し、人間と会話等のコミュニケーションをすることができるロボットが知られている。このようなロボットには、自装置に搭載されたマイクの出力に基づいてロボットの周囲に発生した音を検出し、その音が対象者の発声した音声であると判別すると、その対象者がいる方向にロボットの顔の向きあるいは体の向きを変え、その対象者に話しかけたり手を振ったりする等の動作をするものもある。

【0003】

かかるロボットの動作を実現するためには、対象者が発声する音声である信号音声（信号源）の方向又は位置を検出するために、ロボットの周囲に発生した音から信号音声以外の音声である不要な雑音音声（雑音源）を取り除いて対象者が発声する信号音声のみを分離させる、音源分離技術が必要となる。

40

【0004】

従来、信号音声対雑音音声比（S/N比）を高めるために音源分離技術の一種であるビームフォーミングをする技術（例えば特許文献 1）が知られている。

【先行技術文献】

【特許文献】

【0005】

【特許文献 1】特開 2005 - 253071 号公報

50

【発明の概要】

【発明が解決しようとする課題】

【0006】

しかしながら、従来の音源分離技術では、信号音声と雑音音声とが同方向から到来する場合には音源分離が困難であるという課題があり、そのような場合に、信号音声と雑音音声とを同時に検出してしまい、対象者の信号音声（信号源）の方向又は位置の検出において誤検出が発生するという問題点があった。

【0007】

本発明は、以上のような課題を解決するためのものであり、信号音声を雑音音声から分離できる音源分離情報を検出することが可能な音源分離情報検出装置、ロボット、音源分離情報検出方法及びプログラムを供給することを目的とする。

10

【課題を解決するための手段】

【0008】

前記目的を達成するため、本発明に係る音源分離情報検出装置の一様態は、
 音声を取得するために所定の指向性を有する音声取得手段と、
 前記音声取得手段により取得された前記音声から、所定の対象の信号音声の到来方向である第1方向を検出する第1方向検出手段と、

前記音声取得手段により取得された前記音声から、雑音音声の到来方向である第2方向を検出する第2方向検出手段と、

前記第1方向と前記第2方向とに基づいて、音源分離方向又は音源分離位置を検出する検出手段と、

20

前記音声取得手段が前記音声を取得するタイミングで前記所定の対象の口唇画像を取得する画像取得手段と、

前記口唇画像に基づいて、前記所定の対象の口唇の開口又は前記口唇の閉口を判定する判定手段と、

を備え、

前記第1方向検出手段は、前記判定手段による前記口唇の開口の判定時に、前記音声取得手段により取得された前記音声を前記信号音声とし、

前記第2方向検出手段は、前記判定手段による前記口唇の閉口の判定時に、前記音声取得手段により取得された前記音声を前記雑音音声とする、

30

ことを特徴とする。

【0009】

前記目的を達成するため、本発明に係るロボットの一様態は、
 前記音源分離情報検出装置と、
 自装置を移動する移動手段と、
前記自装置を動作する動作手段と、
前記音源分離情報検出装置、前記移動手段及び前記動作手段を制御する制御手段と、
 を備える、

ことを特徴とする。

【0010】

40

前記目的を達成するため、本発明に係る音源分離情報検出方法の一様態は、
 音声を取得するために所定の指向性を有する音声取得手段により取得された前記音声から、所定の対象の信号音声の到来方向である第1方向を検出し、

前記音声取得手段により取得された前記音声から、雑音音声の到来方向である第2方向を検出し、

前記第1方向と前記第2方向とに基づいて、音源分離方向又は音源分離位置を検出し、

前記音声を取得するタイミングで前記所定の対象の口唇画像を取得し、

前記口唇画像に基づいて、前記所定の対象の口唇の開口又は前記口唇の閉口を判定し、

前記口唇の開口の判定時に、前記取得された前記音声を前記信号音声とし、

前記口唇の閉口の判定時に、前記取得された前記音声を前記雑音音声とする、

50

ことを含む、
ことを特徴とする。

【0011】

前記目的を達成するため、本発明に係るプログラムの一様態は、
音源分離情報検出装置のコンピュータを、
音声を取得するために所定の指向性を有する音声取得手段により取得された前記音声から、所定の対象の信号音声の到来方向である第1方向を検出し、
前記音声取得手段により取得された前記音声から、雑音音声の到来方向である第2方向を検出し、

前記第1方向と前記第2方向とに基づいて、音源分離方向又は音源分離位置を検出し、
前記音声を取得するタイミングで前記所定の対象の口唇画像を取得し、

前記口唇画像に基づいて、前記所定の対象の口唇の開口又は前記口唇の閉口を判定し、
前記口唇の開口の判定時に、前記取得された前記音声を前記信号音声とし、
前記口唇の閉口の判定時に、前記取得された前記音声を前記雑音音声とする、

ように機能させる、
ことを特徴とする。

【発明の効果】

【0012】

本発明によれば、信号音声を雑音音声から分離できる音源分離情報を検出することが可能な音源分離情報検出装置、ロボット、音源分離情報検出方法及びプログラムを供給することが可能になる。

【図面の簡単な説明】

【0013】

【図1】本発明の実施の形態にかかるロボットの外觀図である。

【図2】ロボットの構成を示すブロック図である。

【図3】ロボット制御機能の構成を示すブロック図である。

【図4】図3のブロック図の構成の処理例を示すフローチャートである。

【図5】ラベル付けされた顔パーツ検出結果のフォーマット例を示す図である。

【図6】頭部の回転の自由度を模式的に表した図である。

【図7】雑音音声の音源到来方向推定処理の例を示すフローチャートである。

【図8】仮の音源位置とマイクの配置との一例を示す図である。

【図9】マイクアレイの指向特性の例を示す図である。

【図10】音源分離情報検出の説明図(その1)である。

【図11】音源分離情報検出の説明図(その2)である。

【図12】音源分離情報検出の説明図(その3)である。

【発明を実施するための形態】

【0014】

以下、本発明を実施するための形態について図面を参照しながら詳細に説明する。図1は、実施の形態に係るロボット100を正面から見た場合の外觀を模式的に示した図である。ロボット100は、頭部101と胴体102とを備えた人型のコミュニケーションロボットである。ロボット100は、例えば住宅内に設置され、所定の対象である住人等(以下「対象者」と記載)に呼びかけられると、呼びかけた対象者と会話する。

【0015】

図1に示すように、ロボット100の頭部101には、カメラ104(画像取得手段)と、マイクアレイ103(音声取得手段)と、スピーカ105(報知手段)と、センサ群106と、首関節駆動部107と、足回り駆動部108と、が設けられている。

【0016】

カメラ104は、頭部101の前面の下側、人の顔でいうところの鼻の位置に設けられている。カメラ104は、後述する制御部127の制御の下、撮像を行う。

【0017】

10

20

30

40

50

マイクアレイ 103 は、例えば 13 個のマイクからなる。13 個のマイクのうちの 8 個のマイクが、人の顔でいうところの額の高さの位置であって、頭部 101 の周りに等間隔で配置されている。これら 8 個のマイクよりも上側に、4 個のマイクが頭部 101 の周りに等間隔で配置されている。更に、1 個のマイクが頭部 101 の頭頂部に配置されている。マイクアレイ 103 はロボット 100 の周囲で発生した音を検出する。

【0018】

スピーカ 105 は、カメラ 104 よりも下側、人の顔でいうところの口の位置に設けられている。スピーカ 105 は、後述する制御部 127 の制御の下、各種の音声を出力する。

【0019】

センサ群 106 は、人の顔でいうところの目の位置と耳の位置とに設けられている。センサ群 106 は、加速度センサ、障害物検知センサ等を含み、ロボット 100 の姿勢制御や、安全性の確保のために使用される。

【0020】

首関節駆動部 107 は、頭部 101 と胴体 102 とを連結する部材である。頭部 101 は、破線で示される首関節駆動部 107 によって、胴体 102 に連結されている。首関節駆動部 107 は、複数のモータを含む。後述する制御部 127 がこれら複数のモータを駆動すると、ロボット 100 の頭部 101 が回転する。首関節駆動部 107 は、ロボット 100 の頭部 101 を回転させると共にその回転量を取得する、顔部回転量取得手段としての役割を有する。

【0021】

足回り駆動部 108 は、ロボット 100 を移動させる移動手段としての役割を有する。特に図示しないが、足回り駆動部 108 は、胴体 102 の下側に設けられた 4 つの車輪（ホイール）を含む。4 つの車輪のうち、2 つが胴体 102 の前側に、残り 2 つが後ろ側に配置されている。車輪として、例えば、オムニホイール、メカナムホイールが使用される。後述の制御部 201 は、足回り駆動部 108 の車輪を回転させることにより、ロボット 100 を移動させる。

【0022】

図 2 は、図 1 の外観を有するロボット 100 の制御系であるロボット制御システム 200 を示すブロック図である。図 2 において、図 1 と同じ参照番号を付した部分は図 1 と同じものである。図 2 において、胴体 102 内に設置される制御部 201 は、CPU（Central Processing Unit：中央演算処理装置）、RAM（Random Access Memory：ランダムアクセスメモリ）等を含む。制御部 201 は、頭部 101 内のマイクアレイ 103、カメラ 104、スピーカ 105、センサ群 106、胴体 102 内の首関節駆動部 107 及び足回り駆動部 108 と、それぞれ電氣的に接続され、RAM を作業領域として、後述する記憶部 202 に記憶されている制御プログラム 205 を読み出して実行することにより、前記各部を制御する。

【0023】

記憶部 202 は、ソリッドステートディスクドライブ、ハードディスクドライブ、フラッシュメモリ等を含み、胴体 102 の内部に設けられている。記憶部 202 は、制御部 201 によって実行される制御プログラム 205、マイクアレイ 103 が集音した音声データ、カメラ 104 が撮像した画像データ等を含む各種データを記憶する。記憶部 202 が記憶する制御プログラム 205 には、後述する音源分離情報検出プログラム、移動プログラム、及び対話プログラム等が含まれる。

【0024】

操作ボタン 203 は、胴体 102 の背中に設けられている（図 1 において不図示）。操作ボタン 203 は、ロボット 100 を操作するための各種のボタンであり、電源ボタン、スピーカ 105 の音量調節ボタン等を含む。

【0025】

電源部 204 は、胴体 102 に内蔵された充電電池であり、ロボット制御システム 200

10

20

30

40

50

の各部に電力を供給する。

【0026】

図3は、図2の制御部201が記憶部202内の制御プログラム205を実行する機能の一部として実現される対話機能の構成を示すブロック図である。なお、図3に示される各機能部は、制御部201内のFPGA(Field Programmable Array)等のハードウェアによって実現されてもよい。

【0027】

図3において、音声取得手段として機能する音声入力部301は、図1のマイクアレイ103を構成する各マイクから、音声を入力する。

【0028】

画像取得手段として機能する画像入力部304、顔検出部305、及び口パーツ検出部306は、音声入力部301が音声を取得するタイミングで、所定の対象である対象者の口唇画像を取得する。具体的には、画像入力部304が、図1のカメラ104から画像を入力する。次に、顔検出部305が、音声入力部301が例えば所定の閾値以上のパワーを有する音声を入力するタイミングで、入力された画像から顔領域を検出する。そして、口パーツ検出部306が、検出された顔領域から口パーツを検出し、口唇画像とする。

【0029】

判定手段として機能する口開閉判定部307は、口パーツ検出部306が出力する口唇画像に基づいて、対象者の口唇の開口又は口唇の閉口を判定する。

【0030】

音源到来方向推定部302は、口開閉判定部307が口唇の開口(口唇が開いている状態)を判定しているときには、第1方向検出手段として機能し、音声入力部301が入力する音声を信号音声として、口パーツ検出部306が出力する口唇画像及びその信号音声の信号音声パワーに基づいて、信号音声の到来方向である第1方向を推定する。

【0031】

一方、音源到来方向推定部302は、口開閉判定部307が口唇の閉口(口唇が閉じている状態)を判定しているときには、第2方向検出手段として機能し、音声入力部301が入力する音声を雑音音声として、その雑音音声の雑音音声パワーに基づいて、雑音音声の到来方向である第2方向を推定する。

【0032】

音源到来方向推定部302は、第2方向検出手段として機能するときの処理例として、音源定位手法の一手法であるMUSIC(Multiple Signal Classification)法に基づく処理を実行することにより、対象者以外の音源からの雑音音声の音源定位(雑音源の位置)を推定する。この処理の詳細については後述する。

【0033】

音源分離部303は、例えば下記文献1で示されているビームフォーミング技術に基づく演算処理を実行することにより、音源到来方向推定部302により現在得られている信号音声の到来方向である第1方向又は雑音音声の到来方向である第2方向を入力として、対象者が発声する信号音声を強調し又は信号音声以外の雑音音声を抑圧する音源分離の処理を実行する。

【0034】

<文献1>

浅野 太、“音源分離”、[online]、2011年11月受領、電子情報通信学会『知識の森』、[2017年6月15日検索]、インターネット

<URL:http://www.ieice-hbkb.org/files/02/02gun_06hen_02.pdf>

【0035】

具体的には、音源分離部303は、口開閉判定部307が口唇の開口を判定しているときには、上記ビームフォーミングの演算処理により、信号音声を音源到来方向推定部302により現在得られている第1方向にビームステアリング(強調)するビームステアリング演算処理を実行することにより、強調された信号音声を得て、それを音量算出部308

10

20

30

40

50

に出力する。

【0036】

一方、音源分離部303は、口開閉判定部307が口唇の閉口を判定しているときには、上記ビームフォーミングの演算処理により、雑音音声を音源到来方向推定部302により現在得られている第2方向にヌルステアリング(抑圧)するヌルステアリング演算処理を実行することにより、抑圧された雑音音声を得て、それを音量算出部308に出力する。

【0037】

なお、音源分離部303が実行する前記処理は、マイクアレイ103として所定の指向性を有する物理的な指向性マイクを用いて実現されてもよい。

10

【0038】

音量算出部308は、音源分離部303が出力するビームステアリング(強調)された信号音声又はヌルステアリング(抑圧)された雑音音声のそれぞれの音量を算出する。

【0039】

S/N比算出部309は、音量算出部308が算出する信号音声の音量と雑音音声の音量とに基づいて、信号対雑音比(以下「S/N比」と記載)を算出し、そのS/N比が閾値よりも大きいか否かを判定する。音源分離部303、音量算出部308、及びS/N比算出部309は、第1方向と第2方向とに基づいて、音源分離方向又は音源分離位置を検出する検出手段として機能する。

【0040】

20

S/N比算出部309での判定の結果、S/N比が閾値以下である場合には、図2の制御部201は、音声認識のための十分なS/N比が得られていないと判定する。この場合、制御部201は例えば、図1又は図2の足回り駆動部108を制御することにより、例えば対象者に対して一定の関係(例えば一定の距離又は一定の角度等)を維持しながら、ロボット100を移動させる。

【0041】

ロボット100の移動の後、制御部201は再び、図3のロボット制御機能を動作させ、上述と同様のS/N比の判定動作を実行させる。この結果、S/N比算出部309が算出するS/N比が閾値よりも大きくなると、図2の制御部201は、音声認識のための十分なS/N比が得られ、対象者に対するロボット100の位置関係が、信号音声を雑音音声から最も良く分離できる最適化された位置である音源分離位置になったと判定する(又は、対象者に対するロボット100の方向関係が、信号音声を雑音音声から最も良く分離できる最適化された方向である音源分離方向になったと判定する)。この場合、制御部201は、図3の音声認識部310に、音源分離部303が出力するビームステアリング(強調)された信号音声に対する音声認識処理を実行させることにより、対象者の発話内容を理解する。更に、制御部201は、この音声認識結果に応じて、対話アルゴリズムに従って、図3の発声部311から図1又は図2のスピーカ105を介して、対象者に対して発声を行って対話をする。

30

【0042】

図3において、音声認識部310は、既知の音声認識技術を使って音声認識処理を実行する。また、発声部311は、既知の音声合成技術を使って音声合成による発声処理を実行する。

40

【0043】

図4は、図3のブロック図の構成の処理例を示すフローチャートである。このフローチャートの処理例は、図3のブロック図の構成を実現する制御部201のハードウェアが実行する処理として、又は図2の制御部201が実行する制御プログラム205の処理として実現される。

【0044】

まず、図3の顔検出部305が、顔検出処理を実行する(ステップS401)。この顔検出処理では、音声入力部301が例えば所定の閾値以上のパワーを有する音声を入力す

50

るタイミングで、カメラ104から画像入力部304を介して入力された画像から、顔領域が検出される。顔検出処理としては、既知の顔検出技術を使用することができる。例えば、下記文献2に記載されている何れかの顔検出技術が適用されてよい。

【0045】

<文献2>

堀田 一弘、“小特集 顔認識技術 1.顔認識の研究動向”、[online]、2012年3月28日公開、映像情報メディア学会誌、Vol.64,No.4(2010),p.459-462、[2017年6月15日検索]、インターネット

<URL: https://www.jstage.jst.go.jp/article/itej/64/4/64_4_455/_pdf>

【0046】

次に、図3の口パーツ検出部306が、口パーツ検出処理を実行する(ステップS402)。口パーツ検出処理としては、既知の顔パーツ検出技術を使用することができる。例えば、下記文献3に記載されている何れかの顔パーツ検出技術が適用されてよい。

【0047】

<文献3>

littlewing、“WEBカメラで利用できる顔認識技術まとめ-その2”、[online]、2015年4月7日公開、[2017年6月15日検索]、インターネット

<URL: <http://littlewing.hatenablog.com/entry/2015/04/07/221856>>

【0048】

ステップS402の口パーツ検出処理により、まず例えばラベル付けされた座標値である顔パーツ検出結果が得られる。ラベル付けされた顔パーツ検出結果のフォーマット例としては、例えば図5に示されるように、下記文献4のFigure2として記載されている例を採用することができる。

【0049】

<文献4>

C.sagonas,“Facial point annotations”、[online]、[2017年6月15日検索]、インターネット

<URL: <https://ibug.doc.ic.ac.uk/resources/facial-point-annotations/>>

【0050】

ステップS402の口パーツ検出処理では、図5に例示される顔パーツ検出結果のうち例えば、ラベル49から68までが口パーツとして検出され、またラベル28から36までが鼻パーツとして検出される。

【0051】

次に、図3の口開閉判定部307は、ステップS402で算出された口パーツと鼻パーツとのラベル付けされた座標値(例えば図5のラベル49~68、ラベル28~36)を用いて、口唇の開口(口唇が開いているか)又は口唇の閉口(口唇が閉じているか)を検出する口開閉検出処理を実行する(ステップS403)。

【0052】

ステップS403で口開閉判定部307はまず、口唇の縦座標(顔の上下方向)の変化 y を算出する。今、ある時刻のフレーム $F(t)$ にて、下記(1)式の演算により、 y 座標量差分総計 $y(t)$ が算出される。

【0053】

$$y(t) = y_{y1} + y_{y2} \cdots (1)$$

【0054】

(1)式において、 y_{y1} は、上口唇(下側)と下口唇(上側)との y 座標量差分総計であり、図5の関係より、下記(2)式から(7)式までの累算演算により算出される。これらの式において、演算「+=」は、左辺の値に右辺の値を累算する演算を示す。また、関数「fabs()」は、括弧内の数値に対する絶対値を浮動小数で算出する関数である。また例えば、「data.y[61](t)」は、時刻 t のフレーム画像 $F(t)$ 内における図5のラベル61番の y 座標データ値を示す。他も同様である。

10

20

30

40

50

【 0 0 5 5 】

$$\begin{aligned}
 y y 1 + &= f a b s (d a t a . y [6 1] (t) \\
 &\quad - d a t a . y [6 7] (t)) \cdots (2) \\
 y y 1 + &= f a b s (d a t a . y [6 1] (t) \\
 &\quad - d a t a . y [5 8] (t)) \cdots (3) \\
 y y 1 + &= f a b s (d a t a . y [6 2] (t) \\
 &\quad - d a t a . y [6 6] (t)) \cdots (4) \\
 y y 1 + &= f a b s (d a t a . y [6 2] (t) \\
 &\quad - d a t a . y [5 7] (t)) \cdots (5) \\
 y y 1 + &= f a b s (d a t a . y [6 3] (t) \\
 &\quad - d a t a . y [6 5] (t)) \cdots (6) \\
 y y 1 + &= f a b s (d a t a . y [6 3] (t) \\
 &\quad - d a t a . y [5 6] (t)) \cdots (7)
 \end{aligned}$$

10

【 0 0 5 6 】

(1) 式において、 $y y 2$ は、鼻下と下口唇(上側)との y 座標量差分総計であり、図5の関係より、下記(8)式から(1 2)式までの演算により算出される。

【 0 0 5 7 】

$$\begin{aligned}
 y y 2 + &= f a b s (d a t a . y [3 1] (t) \\
 &\quad - d a t a . y [6 0] (t)) \cdots (8) \\
 y y 2 + &= f a b s (d a t a . y [3 2] (t) \\
 &\quad - d a t a . y [6 1] (t)) \cdots (9) \\
 y y 2 + &= f a b s (d a t a . y [3 3] (t) \\
 &\quad - d a t a . y [6 2] (t)) \cdots (1 0) \\
 y y 2 + &= f a b s (d a t a . y [3 4] (t) \\
 &\quad - d a t a . y [6 3] (t)) \cdots (1 1) \\
 y y 2 + &= f a b s (d a t a . y [3 4] (t) \\
 &\quad - d a t a . y [6 4] (t)) \cdots (1 2)
 \end{aligned}$$

20

【 0 0 5 8 】

図4のステップS 4 0 3で口開閉判定部3 0 7は次に、下記(1 3)式により、時刻 t のフレーム画像 $F (t)$ に対して(1)式の演算で算出した y 座標量差分総計 $y (t)$ と、1フレーム時刻前の時刻 $(t - 1)$ のフレーム画像 $F (t - 1)$ に対して(1)式と同様の演算で算出した y 座標量差分総計 $y (t - 1)$ との差分絶対値 y を求める。ここで、関数「 $a b s ()$ 」は、括弧内の数値に対する絶対値を整数で算出する関数である。

30

【 0 0 5 9 】

$$y = a b s (y (t) - y (t - 1)) \cdots (1 3)$$

【 0 0 6 0 】

(1 3)式で算出される y は、口唇の移動量を示しており、上口唇と下口唇とが離れる方向もしくは近づく方向に移動している時に大きくなる。即ち、口開閉判定部3 0 7は、口唇移動量取得手段として動作する。

【 0 0 6 1 】

図4のステップS 4 0 3で口開閉判定部3 0 7は、口唇の横座標(顔の左右方向)の変化 x についても、前記 y の場合と同様の演算で算出する。

40

【 0 0 6 2 】

即ち今、ある時刻のフレーム $F (t)$ にて、下記(1 4)式の演算によって、 x 座標量差分総計 $x (t)$ が算出される。(1 4)式で例えば、「 $d a t a . x [6 1] (t)$ 」は、時刻 t のフレーム画像 $F (t)$ 内における図5のラベル6 1番の x 座標データ値を示す。他も同様である。

【 0 0 6 3 】

$$\begin{aligned}
 x (t) &= d a t a . x [6 1] (t) + d a t a . x [6 2] (t) \\
 &\quad + d a t a . x [6 3] (t) + d a t a . x [6 7] (t)
 \end{aligned}$$

50

$$+ \text{data} . x [6 6] (t) + \text{data} . x [6 5] (t)$$

$$\dots (1 4)$$

【 0 0 6 4 】

次に、下記 (1 5) 式により、時刻 t のフレーム画像 $F (t)$ に対して (1 4) 式の演算で算出した x 座標量差分総計 $x (t)$ と、1フレーム時刻前の時刻 $(t - 1)$ のフレーム画像 $F (t - 1)$ に対し (1 4) 式と同様の演算で算出した x 座標量差分総計 $x (t - 1)$ との差分絶対値 x が算出される。

【 0 0 6 5 】

$$x = \text{abs} (x (t) - x (t - 1)) \dots (1 5)$$

【 0 0 6 6 】

(1 5) 式で算出される x の値は、 y の場合と同様に口唇の移動量を示しており、口唇が左右どちらかに移動している時に大きくなる。この場合も口開閉判定部 3 0 7 は、口唇移動量取得手段として動作する。

【 0 0 6 7 】

図 4 のステップ S 4 0 3 で口開閉判定部 3 0 7 は続いて、図 1 の頭部 1 0 1 の回転判定を行う。口開閉判定部 3 0 7 は、図 1 又は図 2 の首関節駆動部 1 0 7 から制御部 2 0 1 に入力する信号に基づいて、フレーム時刻 t のフレーム画像 $F (t)$ と、その 1 時刻前のフレーム時刻 $(t - 1)$ のフレーム画像 $F (t - 1)$ における、頭部姿勢の差分 $r o l l$ 、 $y a w$ 及び $p i t c h$ を、下記 (1 6) 式、(1 7) 式及び (1 8) 式により算出する。

【 0 0 6 8 】

$$r o l l = \text{abs} (F (t) r o l l - F (t - 1) r o l l) \dots (1 6)$$

$$y a w = \text{abs} (F (t) y a w - F (t - 1) y a w) \dots (1 7)$$

$$p i t c h = \text{abs} (F (t) p i t c h - F (t - 1) p i t c h)$$

$$\dots (1 8)$$

【 0 0 6 9 】

ここで例えば、 $F (t) r o l l$ は、時刻 t のフレーム画像 $F (t)$ に対応して図 1 又は図 2 の首関節駆動部 1 0 7 から制御部 2 0 1 に入力するロール角度値 $F (t - 1) r o l l$ は、時刻 $(t - 1)$ のフレーム画像 $F (t - 1)$ に対応して図 1 又は図 2 の首関節駆動部 1 0 7 から制御部 2 0 1 に入力するロール角度値である。ヨー角度値 $F (t) y a w$ 及び $F (t - 1) y a w$ 、ピッチ角度値 $F (t) p i t c h$ 及び $F (t - 1) p i t c h$ についても、それぞれ同様である。図 6 は、図 1 のロボット 1 0 0 の頭部 1 0 1 の回転の自由度を模式的に表した図である。図 1 又は図 2 の首関節駆動部 1 0 7 により、ロボット 1 0 0 の頭部 1 0 1 は、胴体 1 0 2 に対して、ピッチ軸 $X m$ の軸回り、ロール軸 $Z m$ の軸回り、ヨー軸 $Y m$ の軸回りにそれぞれ回転可能である。首関節駆動部 1 0 7 は、ピッチ軸 $X m$ の軸回りのピッチ角度値、ロール軸 $Z m$ の軸回りのロール角度値及びヨー軸 $Y m$ の軸回りのヨー角度値をそれぞれ、上記のようにして制御部 2 0 1 に出力する。

【 0 0 7 0 】

図 4 のステップ S 4 0 3 で口開閉判定部 3 0 7 は、上記 (1 6) 式、(1 7) 式及び (1 8) 式の演算の結果、ロール角度差分値 $r o l l$ 、ヨー角度差分値 $y a w$ 及びピッチ角度差分値 $p i t c h$ を、頭部 1 0 1 の回転角度として算出する。この場合、口開閉判定部 3 0 7 は、頭部 1 0 1 = 口唇画像の回転量を取得する口唇回転量取得手段として動作する。

【 0 0 7 1 】

なお、頭部 1 0 1 の回転角度の推定方式としては様々な手法が知られており、上記以外の技術が採用されてもよい。

【 0 0 7 2 】

図 4 のステップ S 4 0 3 で口開閉判定部 3 0 7 は、以上のようにして、口唇の縦座標の変化 y と、横座標の変化 x と、ロボット 1 0 0 の頭部 1 0 1 の回転角度としてロール角度差分値 $r o l l$ 、ヨー角度差分値 $y a w$ 及びピッチ角度差分値 $p i t c h$ に基

10

20

30

40

50

づいて、以下のルールにより口唇の開閉判定を行う。即ち、口開閉判定部307は、下記(19)式の論理式で示される条件が満たされたときに、口唇の開口(口唇が開いている状態)を判定し、その条件が満たされないときに、口唇の閉口(口唇が閉じている状態)を判定する。なお、(19)式において、第1の閾値である y_th 、第2の閾値である x_th 、並びに、第3の閾値群である $roll_th$ 、 yaw_th 及び $pitch_th$ はそれぞれ、 y 、 x 、 $roll$ 、 yaw 及び $pitch$ の判定閾値である。

【0073】

```

y > y_th  &&
x < x_th  &&
roll < roll_th  &&
yaw < yaw_th  &&
pitch < pitch_th

```

10

・・・(19)

【0074】

即ち、口開閉判定部307は、上口唇と下口唇とが離れる方向もしくは近づく方向に移動しており、口唇の横方向移動量は少なく、かつロボット100の頭部101があまり回転していない場合に、口唇の開口を判定する。 y だけでなく、 x 、 $roll$ 、 yaw 、及び $pitch$ も口唇の開閉判定に用いることにより、イヤイヤ(左右に首を振る)、考えるために首を傾げるといった動作でも、誤判定を起りにくくすることができる。

20

【0075】

図4の説明に戻り、上記ステップS403での一連の処理により口開閉判定部307により口唇の開口が判定されると、以下のステップS404からステップS406までの一連の処理が実行される。

【0076】

まず、図3の音源到来方向推定部302が、信号音声の到来方向の推定処理として、図3の顔検出部305により検出されている顔画像(=口唇画像)の口唇方向に基づいて、ロボット100(のカメラ104)に対する口唇方向角度 S_ang を算出する処理を実行する(ステップS404)。

【0077】

続いて、図3の音源分離部303が、例えば前述した文献1に記載されているビームフォーミングの演算処理により、ステップS404で算出された口唇方向角度 S_ang の方向(第1方向)にビームステアリング(強調)するビームステアリング演算処理を実行することにより、強調された信号音声を得る(ステップS405)。

30

【0078】

そして、図3の音量算出部308が、ステップS405で得られたビームステアリング(強調)された信号音声の音量 S_{pow} を算出する(ステップS406)。

【0079】

一方、ステップS403での一連の処理により口開閉判定部307により口唇の閉口が判定されると、以下のステップS407からステップS409までの一連の処理が実行される。

40

【0080】

まず、図3の音源到来方向推定部302が、音源定位手法の一手法であるMUSIC法に基づく処理を実行することにより、対象者以外の音源からの雑音音声の音源定位(雑音源の位置)を推定してノイズ方向角度 N_ang を決定する処理を実行する(ステップS407)。この処理の詳細については、後述する。

【0081】

続いて、図3の音源分離部303が、例えば前述した文献1に記載されているビームフォーミングの演算処理により、ステップS407で算出されたノイズ方向角度 N_ang の方向(第2方向)にヌルステアリング(抑圧)するヌルステアリング演算処理を実行す

50

ることにより、抑圧された雑音音声を得る（ステップS408）。

【0082】

そして、図3の音量算出部308が、ステップS408で得られたヌルステアリング（抑圧）された雑音音声の音量 N_{pow} を算出する（ステップS409）。

【0083】

その後、図3のS/N比算出部309が、ステップS406で算出された信号音声の音量 S_{pow} とステップS409で算出された雑音音声の音量 N_{pow} とに基づいて、下記（20）式の演算に基づいて、S/N比を算出する。

【0084】

$$S/N比 = S_{pow} / N_{pow} \cdots (20)$$

10

【0085】

更に、S/N比算出部309が、下記（21）式の判定演算に基づいて、算出したS/N比が閾値 s_{n_th} よりも大きいかが判定する（ステップS410）。

【0086】

$$S/N比 > s_{n_th} \cdots (21)$$

【0087】

ステップS410の判定がNOの場合には、図2の制御部201は、音声認識のための十分なS/N比が得られていないと判定する。この場合、制御部201は例えば、図1又は図2の足回り駆動部108を制御することにより、例えば対象者に対して一定の関係（例えば一定の距離又は一定の角度等）を維持しながら、ロボット100を移動させる（ステップS411）。移動処理の詳細については、後述する。

20

【0088】

ロボット100の移動の後再び、図4のステップS401からステップS409までの一連の制御処理が実行され、ステップS410のS/N比の判定が行われる。

【0089】

やがて、ステップS410の判定がYESになると、図2の制御部201は、音声認識のための十分なS/N比が得られ、対象者に対するロボット100の位置関係が、信号音声を雑音音声から最も良く分離できる最適化された位置である音源分離位置になったと判定する。この場合、制御部201は、図3の音声認識部310に、音源分離部303が出力するビームステアリング（強調）された信号音声に対する音声認識処理を実行させることにより、対象者の発話内容を理解する。更に、制御部201は、この音声認識結果に応じて、対話アルゴリズムに従って、図3の発声部311から図1又は図2のスピーカ105を介して、対象者に対して発声を行って対話をする（以上、ステップS412）。対話終了後、図2の制御部201は、図4のフローチャートで示される制御処理を終了する。

30

【0090】

図7は、図4のステップS403での一連の処理により口開閉判定部307により口唇の閉口が判定された場合に、ステップS407で図3の音源到来方向推定部302によりMUSIC法に基づいて実行される、対象者以外の音源からの雑音音声の音源定位（雑音源の位置）を推定してノイズ方向角度 N_{ang} を決定する処理の詳細例を示すフローチャートである。

40

【0091】

まず、図1又は図2のマイクアレイ103に入力された音声、時間周波数変換される（ステップS701）。ここでは例えば、時間周波数変換演算処理として、STFT（Short-Time Fourier Transform：短時間フーリエ変換）が実行される。

【0092】

音源数をNとすると、第n番目の音源の信号 S_n は、下記（22）式で表せる。なお、 ω は角周波数、 f はフレーム番号である（以下の説明でも同様）。

【0093】

$$S_n(\omega, f) \quad (n = 1, 2, \dots, N) \cdots (22)$$

50

【0094】

図1又は図2のマイクアレイ103の各マイクで観測される信号は、マイクアレイ103におけるマイクの数 M とすると、下記(23)式で表せる。

【0095】

$$X_m(\omega, f) \quad (m = 1, 2, \dots, M) \dots (23)$$

【0096】

音源から出た音は、空気を伝わってマイクアレイ103のマイクで観測されるが、そのときの伝達関数を $H_{nm}(\omega)$ とすると、音源の信号を表す数式に、伝達関数を乗じることによって、マイクアレイ103の各マイクで観測される信号を求めることができる。 m 番目のマイクで観測される信号 $X_m(\omega, f)$ は下記(24)式のように表される。

10

【0097】

【数1】

$$X_m(\omega, f) = \sum_{n=1}^N S_n(\omega, f) H_{nm}(\omega) \dots (24)$$

【0098】

ロボット100は、マイクアレイ103としてマイクを複数有しているので、マイクアレイ103全体で観測される信号 $x(\omega, f)$ は下記(25)式で表すことができる。

【0099】

20

【数2】

$$x(\omega, f) = \begin{bmatrix} X_1(\omega, f) \\ X_2(\omega, f) \\ \vdots \\ X_M(\omega, f) \end{bmatrix} \dots (25)$$

30

【0100】

同様に、全音源の信号 $s(\omega, f)$ も下記(26)式で表すことができる。

【0101】

【数3】

$$s(\omega, f) = \begin{bmatrix} S_1(\omega, f) \\ S_2(\omega, f) \\ \vdots \\ S_N(\omega, f) \end{bmatrix} \dots (26)$$

40

【0102】

同様に、第 n 番目の音源の伝達関数 $h_n(\omega)$ は下記(27)式で表すことができる。

【0103】

【数4】

$$h_n(\omega) = \begin{bmatrix} H_{n1}(\omega) \\ H_{n2}(\omega) \\ \vdots \\ H_{nM}(\omega) \end{bmatrix} \dots (27)$$

10

【0104】

全ての伝達関数を下記(28)式のように表記する。

【0105】

$$h(\) = [h_1(\), h_2(\), \dots, h_N(\)] \dots (28)$$

【0106】

(28)式で表される伝達関数を、前述した(24)式に適用すると、下記(29)式のように表される。

【0107】

$$x(\ , f) = h(\) s(\ , f) \dots (29)$$

20

【0108】

$h_n(\)$ は音源位置毎に独立であり、ある程度のフレーム数(例えば、フレーム数をLとする)で見れば $S_n(\ , f)$ は無相関とみなせるので、 $x(\ , f)$ は音源数NをRANKとする超平面を構成する。このとき、距離で正規化した音量が大きな音源の伝達関数方向に分布が広がりやすい。そこで、部分空間とゼロ空間とに分解することを考える。

【0109】

再び図7を参照する。次の(30)式に示されるように、相関行列が計算される(ステップS702)。ここで、「*」は複素共役転置を表す。

【0110】

【数5】

$$R(\omega, f) = \sum_{l=0}^{L-1} x(\omega, f+l)x^*(\omega, f+l) \dots (30)$$

30

【0111】

続いて、固有値分解が実行される(ステップS703)。ここで、固有値 $\lambda_m(\ , f)$ と固有ベクトル $e_m(\ , f)$ とは固有値が降順になるように並べ替えられているものとする

40

【0112】

原理的には、 $h_n(\)$ は部分空間の固有ベクトル $e_m(\ , f)$ ($m=1 \sim N$)の重み付け加算から復元できるが、実際には復元が困難であるためゼロ空間を構成する固有ベクトル $e_m(\ , f)$ ($m=N+1 \sim M$)が $h_n(\)$ と直交することを使って音源定位を実現する。

【0113】

しかし、雑音音声の音源は例えば建物室内を移動する可能性があるため、音源位置を予め知ることはできず、音源位置の伝達関数を予め取得しておくことは難しい。このため、仮の音源位置が決められ、仮の音源位置の伝達関数が予め用意されて、音源定位が行われ

50

る。

【 0 1 1 4 】

図 8 は、仮の音源位置とマイクの配置との一例を示す図である。図 8 では、太線の円がロボット 1 0 0 の頭 1 1 0 を表し、太線上の黒丸がマイクアレイ 1 0 3 のマイクを表す。なお、ここでは、便宜上図 1 のマイクアレイ 1 0 3 の 1 3 個のマイクの全てを表示していない。ロボット 1 0 0 の回りには 4 個の仮の音源位置があるものとする。

【 0 1 1 5 】

マイクアレイ 1 0 3 の複数のマイクは、ロボット 1 0 0 の頭 1 1 0 に配置されていることから、円周に沿って配置されているとみなすことができる。X 軸の正の向きと、各マイクが成す円の中心（ロボット 1 0 0 の頭 1 1 0 の中心位置に相当）と仮の音源 1 ~ 4 とをそれぞれ結んだ線と、がなす角度を 1、 2、 3、 4 として、それぞれの伝達関数 $h_{\theta}(\omega)$ を予め計算しておく。

10

【 0 1 1 6 】

図 8 では、音源が 4 個の例を示したが、音源数が N 個の場合、 1、 2、 …、 N のそれぞれの伝達関数 $h_{\theta}(\omega)$ を予め計算しておけばよい。或いは、仮の音源位置の伝達関数を用意するのではなく、幾何的な情報をもとに予め伝達関数を計算しておいてもよい。

【 0 1 1 7 】

再び図 7 を参照する。下記 (3 1) 式を使用して、周波数帯毎の M U S I C スペクトルが計算される (ステップ S 7 0 4) 。

20

【 0 1 1 8 】

【 数 6 】

$$M_{\theta}(\omega, f) = \frac{h_{\theta}^*(\omega)h_{\theta}(\omega)}{\sum_{m=N+1}^M |h_{\theta}^*(\omega)e_m(\omega, f)|^2} \dots (31)$$

【 0 1 1 9 】

ここで、(3 1) 式の分母は、ノイズや誤差、S T F T の周波数帯間の信号漏洩の影響等からゼロにはならない。また、音源の方向と予め決めた角度 (1、 2、 …、 N) の何れかが近い場合、つまり $h_n(\omega)$ と $h_{\theta}(\omega)$ とが近い場合、(3 1) 式の値は極端に大きなものになる。図 8 に示す例では、雑音音声の音源と仮の音源の位置とが近いたため、 2 の伝達関数を使用した場合、(3 1) 式の値が極端に大きくなることが想定される。

30

【 0 1 2 0 】

次に、統合した M U S I C のパワーを求めるため、下記 (3 2) 式の演算により、周波数帯毎の M U S I C スペクトルが重み付け加算される (ステップ S 7 0 5) 。

【 0 1 2 1 】

【 数 7 】

$$M(f) = \sum_{\omega} w(\omega)M(\omega, f) \dots (32)$$

40

【 0 1 2 2 】

重み付け係数は、固有値 $\lambda_m(\omega, f)$ が大きいほど大きくすれば、 $S_n(\omega, f)$ に含まれるパワーに応じた計算をすることもできる。この場合は $S_n(\omega, f)$ に殆どパワーがない場合の悪影響を軽減できる。

【 0 1 2 3 】

最後に、パワースペクトルから適切なピーク (極大値) が選択される (ステップ S 7 0

50

6)。具体的には、まず、複数のピークが算出され、その中から適切なピークが選択されて、選択されたピークにおける θ が図4のステップS407で説明した雑音音声の音源方向のノイズ方向角度 N_ang とされる。ここで、ピークを求めるのは以下のような理由による。本来の音源方向の θ のパワーが必ずしも一番大きいとは限らず、本来の音源方向に近い θ のパワーは総じて大きくなるので、音源方向は複数のピークの何れかに正解があるからである。その後、図7のフローチャートの処理が終了して、図4のステップS407の雑音音声の音源到来方向推定処理が終了する。

【0124】

以上の説明では、雑音音声の音源到来方向として平面を仮定して説明したが、3次元空間を仮定しても上記説明は成り立つ。

10

【0125】

図9は、図1又は図2のマイクアレイ103の指向特性の例を示す図、図10から図12は、音源分離方向検出の説明図である。図9において、マイクアレイ103は、120度付近で、各周波数においてまんべんなくマイナスゲインが得られている。従って、下記(33)式のように、図4のステップS404で算出される口唇方向角度である対象者の信号音声方向 S_ang と、ステップS407で算出されるノイズ方向角度 N_ang との差分の絶対値が120度付近になる音源分離方向が、最も良い音源分離が期待できる方向となる。

【0126】

$$abs(S_ang - N_ang) \cdots (33)$$

20

【0127】

図4のステップS410の判定がNO S411 S401として実行される処理により実現されるアルゴリズムとしては、ロボット100の位置毎に、前述した(20)式により算出される信号音声対雑音音声のS/N比が前述した(21)式の判定演算により閾値 s_nth を超えたか否かが判定されながら、前述した図4のステップS411でのロボット100の移動処理が繰り返され、S/N比が閾値 s_nth を超えたと判定された地点が、信号音声と雑音音声との最適な分離位置、即ち音源分離位置とされる。

【0128】

なお、S/N比が閾値 s_nth を超えた時点ではなく、閾値を超えた後にS/N比が最高となる地点が音源分離位置とされてもよい。

30

【0129】

例えば、図10は、 $S_ang = 0$ 度、 $N_ang = 10$ 度と算出された状態の例を示している。ロボット100のカメラ104(図1)から見ると、対象者の右10度方向にノイズの音源が存在していることになる。この状態から、図2の制御部201は、図1又は図2の足回り駆動部108を制御することにより、ロボット100を、対象者を中心にして例えば右方向(図10の方向A)に移動させる。左方向への移動が行われてももちろん良い。ただし、図10の例の場合は、右移動の方が音源分離位置に最短距離で近づくことができる。図11は、上記移動後のロボット100と対象者とノイズ音源の位置関係を示す図である。このような移動が繰り返されることにより、ロボット100は最終的に、図12に示される音源分離位置まで移動をして、移動を完了する。この音源分離位置は、マイクアレイ103が図9に示される指向特性を有する場合に、(33)式で算出される信号音声方向 S_ang とノイズ方向角度 N_ang との差分の絶対値が120度付近になる位置である。

40

【0130】

上述の動作において、図2の制御部201は、移動開始時に「聞き取りやすい位置に移動するね」などの音声を、図3の発声部311から発声させることにより、移動中は対話を中止してもらえらるような文言を喋らせることが望ましい。また、移動中にも対話できるようにしてもよい。

【0131】

上述した図4のフローチャートで例示される制御処理において、S/N比算出部309

50

でのステップS410の判定の結果、S/N比が閾値 s_{n_th} 以下である場合に、制御部201は例えば、図3の発声部311を介して図1又は図2のスピーカ105から、対象者に対して、「僕を中心にして・・・度ほど回転するように移動してください。」というような意味の発声を行って、対象者に移動を促すような制御が行われてもよい。

【0132】

また、上述のような発声を行いながら、継続的に取得したノイズ方向角度 N_{ang} が都合の良い角度になるまで、「もう少し」や「ストップ」などの発声を行って対象者に指示をするような制御が行われてもよい。

【0133】

例えば建物室内のマップ情報を利用できる場合には、対象者やノイズの2次元又は3次元の音源位置をマップ上で推定し、その推定結果に基づいて音源分離位置に移動するような制御が実施されてもよい。音源位置のマップは、ノイズ音源になるべく近づいてその位置を特定して登録するようにしてもよい。

【0134】

一方、音源位置のマップが無い場合には、ロボット100の移動中に獲得したノイズ方向とその時の位置とロボット100本体の向きとから、ノイズの音源位置を推定するようにしてもよい。この場合、観測点が2点以上あれば音源位置が決められる。推定方向にある程度の誤差を持たせて、より多くの観測点から推定が行われるようにしてもよい。

【0135】

更に、上記のようなマップ情報を使ったノイズの音源位置の推定結果に基づいて、「あと・・・度回転して」というような発声を行って対象者に指示をするような制御が行われてもよい。

【0136】

上述した実施形態において、ロボット100が移動するときに、ロボット100がそばを向きながら移動したり、ロボット100が勝手に動いていると、対象者が違和感を感じてしまうため、対象者が違和感を感じないように、移動することが望ましい。例えば、対象者とアイコンタクトを取ったり、対象者の方を向きながら移動することが望ましい。また、音源分離位置まで一気に移動するのではなく、少しだけ動いたり、回転のみをしてもよい。

【0137】

以上の実施形態によれば、信号音声 S が雑音音声 N から最も良く分離した状態で音源分離が行える最適化された音源分離情報(音源分離方向又は音源分離位置)を検出することが可能となる。これにより、対象者の音声以外の他の音声を排除して、音声認識の誤認識を減らすことが可能となる。

【0138】

以上説明した実施形態において、図2の制御部201が記憶部202に記憶され図4や図7のフローチャートの処理例で示される制御プログラム205を実行することにより図3で示される機能を実現する場合、制御プログラム205は、例えば外部記憶装置や可搬記録媒体に記録して配布してもよく、あるいは特には図示しない無線や有線の通信インタフェースを介してネットワークから取得できるようにしてもよい。

【0139】

以上の実施形態に関して、更に以下の付記を開示する。

(付記1)

音声を取得するために所定の指向性を有する音声取得手段と、

前記音声取得手段により取得された所定の対象の信号音声から、前記信号音声の到来方向である第1方向を検出する第1方向検出手段と、

前記音声取得手段により取得された雑音音声から、前記雑音音声の到来方向である第2方向を検出する第2方向検出手段と、

前記第1方向と前記第2方向とに基づいて、音源分離方向又は音源分離位置を検出する検出手段と、

10

20

30

40

50

を備える、
ことを特徴とする音源分離情報検出装置。

(付記 2)

前記検出手段は、前記信号音声と前記雑音音声とから算出された信号対雑音比が閾値以下である場合、前記第 1 方向と前記第 2 方向とに基づいて、前記信号対雑音比が前記閾値を超える前記音源分離方向又は前記音源分離位置を検出する、
ことを特徴とする付記 1 に記載の音源分離情報検出装置。

(付記 3)

前記音声取得手段が前記音声を取得するタイミングで前記所定の対象の口唇画像を取得する画像取得手段と、

前記口唇画像に基づいて、前記所定の対象の口唇の開口又は前記口唇の閉口を判定する判定手段と、

を更に備え、

前記第 1 方向検出手段は、前記判定手段による前記口唇の開口の判定時に、前記音声取得手段により取得された前記音声を前記信号音声とし、

前記第 2 方向検出手段は、前記判定手段による前記口唇の閉口の判定時に、前記音声取得手段により取得された前記音声を前記雑音音声とする、

ことを特徴とする付記 1 又は 2 に記載の音源分離情報検出装置。

(付記 4)

前記画像取得手段は、顔部画像を更に取得し、

前記口唇画像から前記所定の対象の前記口唇の移動量を取得する口唇移動量取得手段と

、

前記顔部画像から前記所定の対象の顔部の回転量を取得する顔部回転量取得手段と、
を更に備え、

前記判定手段は、前記口唇の移動量と前記顔部の回転量とに基づいて、前記所定の対象の前記口唇の開口又は前記口唇の閉口を判定する、

ことを特徴とする付記 3 に記載の音源分離情報検出装置。

(付記 5)

前記判定手段は、前記口唇の移動量のうちの前記口唇の開閉方向の移動量が第 1 の閾値を超え、且つ、前記口唇の移動量のうちの前記口唇の延伸方向の移動量が第 2 の閾値未満であり、且つ、前記顔部の回転量が第 3 の閾値未満であるときに、前記口唇の開口又は前記口唇の閉口を判定する、

ことを特徴とする付記 4 に記載の音源分離情報検出装置。

(付記 6)

前記第 1 方向検出手段は、前記判定手段による前記口唇の開口の判定時に、前記信号音声の信号音声パワーに基づいて、前記第 1 方向を検出し、

前記第 2 方向検出手段は、前記判定手段による前記口唇の閉口の判定時に、前記雑音音声の雑音音声パワーに基づいて、前記第 2 方向を検出する、

ことを特徴とする付記 3 乃至 5 の何れか 1 つに記載の音源分離情報検出装置。

(付記 7)

前記検出手段は、前記信号対雑音比が前記閾値を超えて最大となる方向を前記音源分離方向とする、又は、前記信号対雑音比が前記閾値を超えて最大となる位置を前記音源分離位置とする、

ことを特徴とする付記 2 乃至 6 の何れか 1 つに記載の音源分離情報検出装置。

(付記 8)

前記検出手段は、前記信号対雑音比が前記閾値を超える場合、現在方向を前記音源分離方向とする、又は、現在位置を前記音源分離位置とする、

ことを特徴とする付記 2 乃至 6 の何れか 1 つに記載の音源分離情報検出装置。

(付記 9)

前記所定の対象にメッセージを報知する報知手段を更に備え、

10

20

30

40

50

前記報知手段は、前記所定の対象に現在位置から前記音源分離位置まで移動させるために、前記音源分離位置までの移動方向及び移動距離を含む前記メッセージを報知する、ことを特徴とする付記 1 乃至 8 の何れか 1 つに記載の音源分離情報検出装置。

(付記 10)

前記所定の対象は人又は動物である、ことを特徴とする付記 1 乃至 9 の何れか 1 つに記載の音源分離情報検出装置。

(付記 11)

付記 1 乃至 10 の何れかに記載の音源分離情報検出装置と、
自装置を移動させる移動手段と、
前記音源分離情報検出装置及び前記移動手段を制御する制御手段と、
を備える、
ことを特徴とするロボット。

10

(付記 12)

前記制御手段は、前記移動手段を制御して、前記音源分離位置に前記自装置を移動させる、
ことを特徴とする付記 11 に記載のロボット。

(付記 13)

音声を取得するために所定の指向性を有する音声取得手段により取得された所定の対象の信号音声から、前記信号音声の到来方向である第 1 方向を検出し、

前記音声取得手段により取得された雑音音声から、前記雑音音声の到来方向である第 2 方向を検出し、

20

前記第 1 方向と前記第 2 方向とに基づいて、音源分離方向又は音源分離位置を検出する、

ことを含む、

ことを特徴とする音源分離情報検出方法。

(付記 14)

音源分離情報検出装置のコンピュータを、
音声を取得するために所定の指向性を有する音声取得手段により取得された所定の対象の信号音声から、前記信号音声の到来方向である第 1 方向を検出し、

前記音声取得手段により取得された雑音音声から、前記雑音音声の到来方向である第 2 方向を検出し、

30

前記第 1 方向と前記第 2 方向とに基づいて、音源分離方向又は音源分離位置を検出する、

ように機能させる、

ことを特徴とするプログラム。

【符号の説明】

【0140】

100 ロボット

101 頭部

102 胴体

40

103 マイクアレイ

104 カメラ

105 スピーカ

106 センサ群

107 首関節駆動部

108 足回り駆動部

200 ロボット制御システム

201 制御部

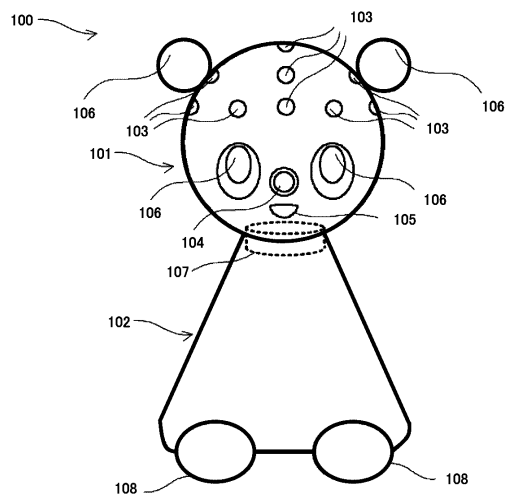
202 記憶部

203 操作ボタン

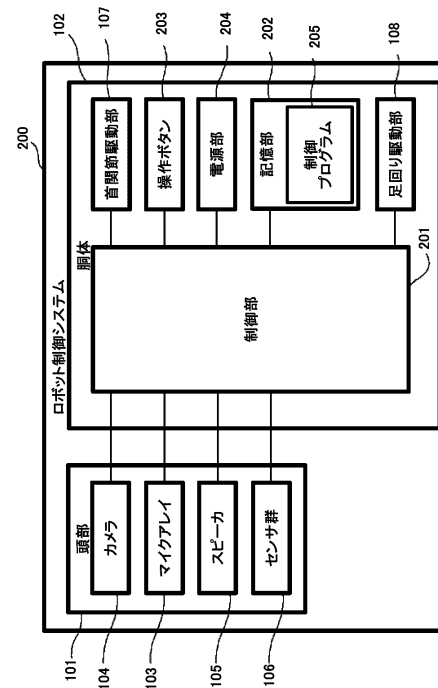
50

- 204 電源部
- 205 制御プログラム
- 301 音声入力部
- 302 音源到来方向推定部
- 303 音源分離部
- 304 画像入力部
- 305 顔検出部
- 306 口パーツ検出部
- 307 口開閉判定部
- 308 音量算出部
- 309 S/N比算出部
- 310 音声認識部
- 311 発声部

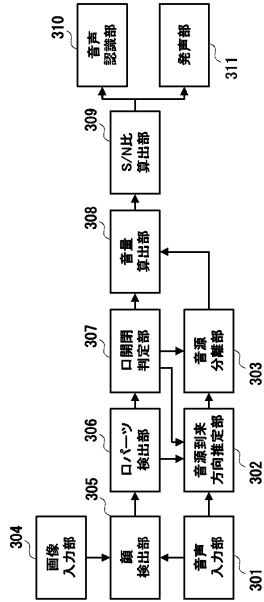
【図1】



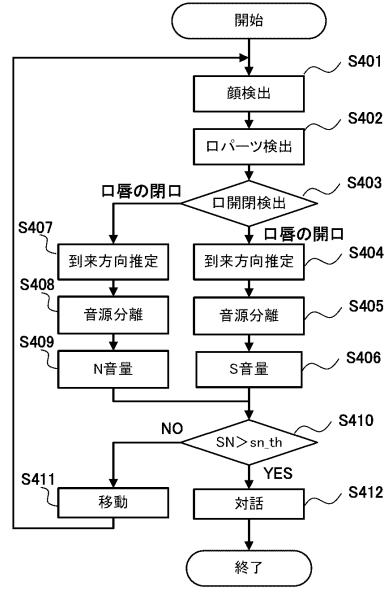
【図2】



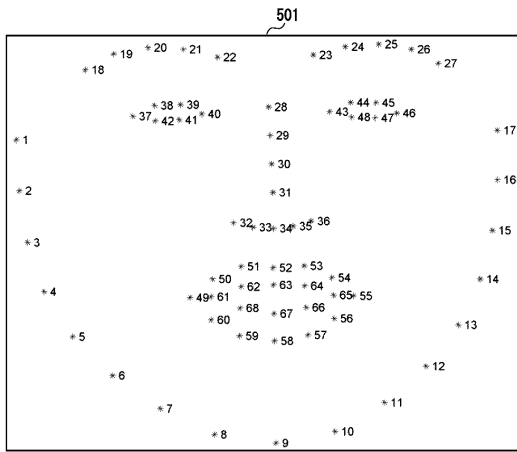
【図3】



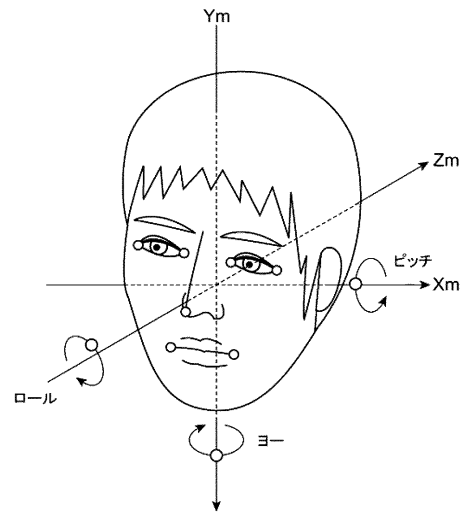
【図4】



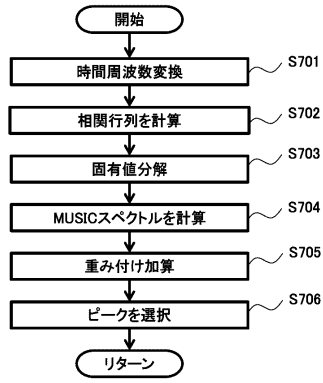
【図5】



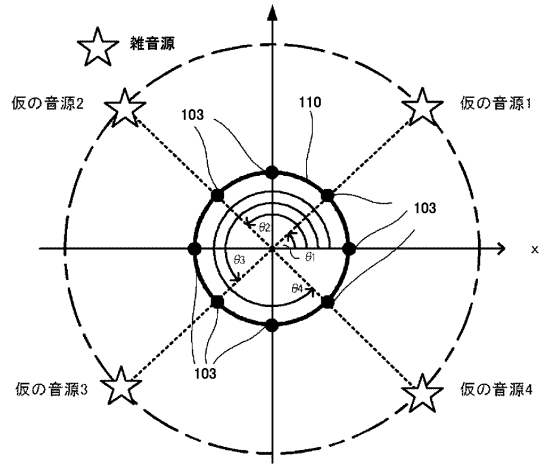
【図6】



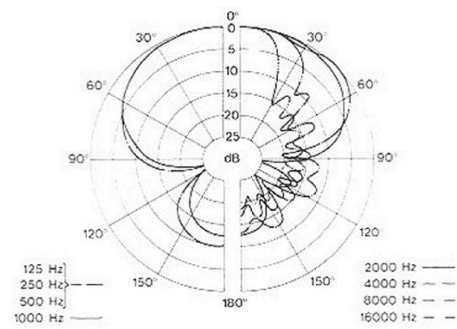
【図7】



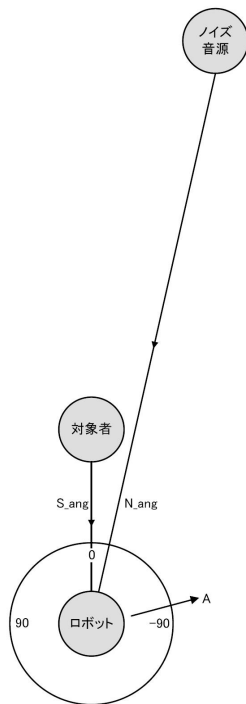
【図8】



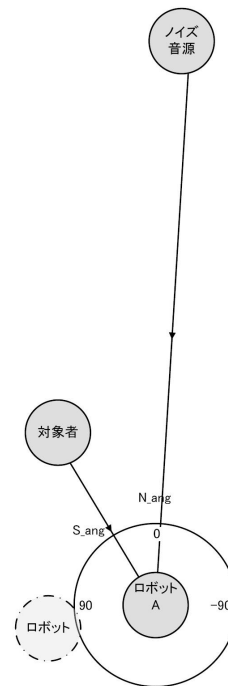
【図9】



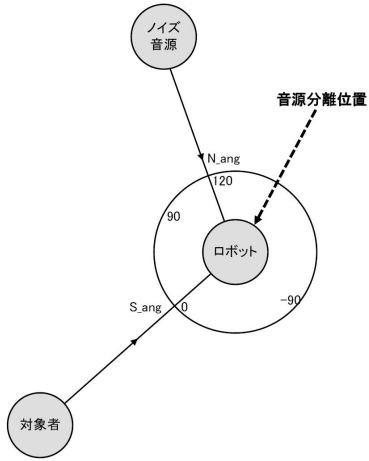
【図10】



【図11】



【図 12】



フロントページの続き

審査官 山下 剛史

- (56)参考文献 特表2005-529421(JP,A)
特開2006-181651(JP,A)
特開2014-207589(JP,A)
特開2005-253071(JP,A)
藤田善弘, パーソナルロボットR100, 日本ロボット学会誌, 2000年 3月, Vol.18, No.
.2, p.198-199

- (58)調査した分野(Int.Cl., DB名)
G10L 13/00 - 99/00