

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6673276号
(P6673276)

(45) 発行日 令和2年3月25日 (2020.3.25)

(24) 登録日 令和2年3月9日 (2020.3.9)

(51) Int.Cl.	F I
G 1 O L 25/84 (2013.01)	G 1 O L 25/84
G 1 O L 25/51 (2013.01)	G 1 O L 25/51 4 O O
G 1 O L 15/00 (2013.01)	G 1 O L 15/00 2 O O H

請求項の数 9 (全 26 頁)

(21) 出願番号	特願2017-62756 (P2017-62756)	(73) 特許権者	000001443
(22) 出願日	平成29年3月28日 (2017.3.28)		カシオ計算機株式会社
(65) 公開番号	特開2018-165759 (P2018-165759A)		東京都渋谷区本町1丁目6番2号
(43) 公開日	平成30年10月25日 (2018.10.25)	(74) 代理人	100095407
審査請求日	平成30年12月7日 (2018.12.7)		弁理士 木村 満
		(72) 発明者	島田 敬輔
			東京都羽村市栄町3-2-1 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3-2-1 カシオ計算機株式会社 羽村技術センター内
		(72) 発明者	山谷 崇史
			東京都羽村市栄町3-2-1 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】 音声検出装置、音声検出方法、及びプログラム

(57) 【特許請求の範囲】

【請求項1】

音声を検出する音声検出手段と、
 前記音声検出手段により検出された音声である検出音声の音声発生源が特定の音声発生源であるか否かを判別する第1判別手段と、
 前記第1判別手段の判別結果に基づいて自機を制御する制御手段と、
 前記検出音声が発生した方向を判別する第2判別手段と、
 前記特定の音声発生源以外の他の音声発生源の位置を示す情報を含む音声発生源位置情報を記憶した記憶部と、
 前記第2判別手段による判別結果と前記記憶された音声発生源位置情報とに基づいて、
 前記自機に対する前記検出音声が発生した方向に前記他の音声発生源が存在するか否かを判別する第3判別手段と、を備え、
 前記制御手段は、前記第3判別手段により前記検出音声が発生した方向に前記他の音声発生源が存在しないと判別されている場合に、前記自機の動作を制御する、
 音声検出装置。

【請求項2】

前記制御手段は、前記第1判別手段により、前記検出された音声の音声発生源が前記特定の音声発生源であると判別された場合、前記自機の位置及び姿勢の少なくとも一方に関する制御を実行する、

請求項1に記載の音声検出装置。

【請求項 3】

撮像部をさらに備え、

前記第 3 判別手段により前記検出音声が発生した方向に前記他の音声発生源が存在しないと判別されている場合に、前記制御手段は、前記自機の動作の制御として、前記撮像部の撮像方向を前記第 2 判別手段が判別した方向に向けるように、前記自機の位置及び姿勢の少なくとも一方に関する制御を実行する、

請求項 1 または 2 に記載の音声検出装置。

【請求項 4】

前記音声発生源位置情報は、前記自機の周囲の複数の位置の各々に、前記特定の音声発生源以外の音声発生源が存在する確率を示す情報を含む、

請求項 1 から 3 のいずれか 1 項に記載の音声検出装置。

【請求項 5】

前記制御手段により前記自機を移動させている間に、前記撮像部が撮像した画像から認識された音声発生源の位置を示す情報を、前記音声発生源位置情報に追加する、

請求項 3 または請求項 3 に従属する請求項 4 に記載の音声検出装置。

【請求項 6】

前記第 1 判別手段は、前記音声検出手段により検出された前記音声の前記自機宛てに発せられた音声か否かを判別し、前記自機宛ての音声であると判別した場合、前記音声検出手段が検出した音声の音声発生源が前記特定の音声発生源であるか否かを判別する、

請求項 1 から 5 のいずれか 1 項に記載の音声検出装置。

【請求項 7】

音声を検出する音声検出手段と、

前記音声検出手段により検出された音声の音声発生源が特定の音声発生源であるか否かを判別する判別手段と、

前記判別手段の判別結果に基づいて自機を制御する制御手段と、

撮像部と、

前記特定の音声発生源以外の音声発生源であって、登録された音声発生源の位置を示す情報を含む音声発生源位置情報をあらかじめ記憶した記憶部と、を備え、

前記判別手段は、前記音声検出手段により検出された前記音声の音声発生源の位置を判別し、前記判別された位置が、前記音声発生源位置情報に含まれる前記登録された音声発生源の位置であるか否かを判別し、

前記判別手段により判別された位置が、前記音声発生源位置情報に含まれる前記登録された音声発生源の位置でないと判別した場合に、前記制御手段は、前記撮像部の撮像方向を前記判別手段が判別した位置に向けるように、自機の位置、姿勢の少なくとも一方を変える、

音声検出装置。

【請求項 8】

ロボットに搭載されたコンピュータが音声を検出する音声検出方法であって、

音声を検出する音声検出ステップと、

前記音声検出ステップで検出された音声である検出音声の音声発生源が特定の音声発生源であるか否かを判別する第 1 判別ステップと、

前記第 1 判別ステップの判別結果に基づいて、前記ロボットの動作を制御する制御ステップと、

前記検出音声が発生した方向を判別する第 2 判別ステップと、

前記第 2 判別ステップによる判別結果と、記憶部に記憶された、前記特定の音声発生源以外の他の音声発生源の位置を示す情報を含む音声発生源位置情報とに基づいて、前記ロボットに対する前記検出音声が発生した方向に前記他の音声発生源が存在するか否かを判別する第 3 判別ステップと、を備え、

前記制御ステップでは、前記第 3 判別ステップにより前記検出音声が発生した方向に前記他の音声発生源が存在しないと判別されている場合に、前記ロボットの動作を制御する

10

20

30

40

50

、
音声検出方法。
 【請求項 9】
 ロボットに搭載されたコンピュータに、
 音声を検出する音声検出機能と、
 前記音声検出機能により検出された音声である検出音声の音声発生源が特定の音声発生源であるか否かを判別する第 1 判別機能と、
 前記第 1 判別機能により判別された判別結果に基づいて、前記ロボットの動作を制御する制御機能と、
 前記検出音声が発生した方向を判別する第 2 判別機能と、
 前記第 2 判別機能による判別結果と、記憶部に記憶された、前記特定の音声発生源以外の他の音声発生源の位置を示す情報を含む音声発生源位置情報とに基づいて、前記ロボットに対する前記検出音声が発生した方向に前記他の音声発生源が存在するか否かを判別する第 3 判別機能と、を実現させ、
 前記制御機能は、前記第 3 判別機能により前記検出音声が発生した方向に前記他の音声発生源が存在しないと判別されている場合に、前記ロボットの動作を制御する、

10

プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

20

本発明は、音声検出装置、音声検出方法、及びプログラムに関する。

【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、自機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、人の声であると判別すると、人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、その人に話しかける、手を振る等の動作をするものもある。

【0003】

特許文献 1 には、ロボットが、マイクロホンに閾値以上の振幅の音が入力されることにより、音イベントが発生したことを検出し、音源方向を推定して、推定した音源方向に振り向くことが記載されている。

30

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2003 - 266351 号公報

【非特許文献】

【0005】

【非特許文献 1】Andrew J. Davison, "Real-Time Simultaneous Localization and Mapping with a Single Camera", Proceedings of the 9th IEEE International Conference on Computer Vision Volume 2, 2003, pp. 1403 - 1410

40

【非特許文献 2】Richard Hartley, Andrew Zisserman, "Multiple View Geometry in Computer Vision", Second Edition, Cambridge University Press, March 2004, chapter 9

【非特許文献 3】Csurka, G., Dance, C. R., Fan, L., Willamowski, J. and Bray, C.: Visual categorization with bags of keypoints, ECCV Internat

50

ional Workshop on Statistical Learning in Computer Vision (2004)

【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、特許文献1に記載されているロボットは、音イベントを検出すると振り向くので、実際にロボットに対して人から発せられた音だけではなく、例えば、テレビ、ラジオ等の電子機器のスピーカから出力される音声にも反応してしまうことが予想される。

【0007】

本発明は、上記実情を鑑みてなされたものであり、ロボットに実際の人から直接発せられた音声か電子機器のスピーカから出力された音声かを判別させることで、ロボットの無駄な動作を減らすことを目的とする。

【課題を解決するための手段】

【0008】

上記目的を達成するため、本発明に係る音声検出装置は、
 音声を検出する音声検出手段と、
 前記音声検出手段により検出された音声である検出音声の音声発生源が特定の音声発生源であるか否かを判別する第1判別手段と、
 前記第1判別手段の判別結果に基づいて自機を制御する制御手段と、
 前記検出音声が発生した方向を判別する第2判別手段と、
 前記特定の音声発生源以外の他の音声発生源の位置を示す情報を含む音声発生源位置情報を記憶した記憶部と、

前記第2判別手段による判別結果と前記記憶された音声発生源位置情報とに基づいて、前記自機に対する前記検出音声が発生した方向に前記他の音声発生源が存在するか否かを判別する第3判別手段と、を備え、

前記制御手段は、前記第3判別手段により前記検出音声が発生した方向に前記他の音声発生源が存在しないと判別されている場合に、前記自機の動作を制御する。

【発明の効果】

【0009】

本発明によれば、ロボットに実際の人から直接発せられた音声か電子機器のスピーカから出力された音声かを判別させることで、ロボットの無駄な動作を減らすことができる。

【図面の簡単な説明】

【0010】

【図1】本発明の実施の形態1にかかるロボットの外觀図である。

【図2】ロボットの頭の自由度を説明するための図である。

【図3】ロボットの構成を示すブロック図である。

【図4】部屋内のロボットとユーザの位置の一例を示す図である。

【図5】地図作成処理のフローチャートである。

【図6】呼びかけ応答処理のフローチャートである。

【図7】音源定位の処理のフローチャートである。

【図8】仮の音源の位置を説明するための図である。

【図9】自機位置推定の処理のフローチャートである。

【図10】実施の形態2にかかるロボットの記憶部の構成を示すブロック図である。

【図11】呼びかけ移動処理のフローチャートである。

【図12】顔位置推定の処理のフローチャートである。

【発明を実施するための形態】

【0011】

(実施の形態1)

以下、図面を参照しながら本発明の実施の形態1について説明する。図1は、実施の形

10

20

30

40

50

態１に係るロボット１００を正面から見た場合の外観を模式的に示した図である。ロボット１００は、頭１１０と胴体１２０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、住宅内に設置されており、住人に呼びかけられると、呼びかけた住人と会話する。

【００１２】

図１に示すように、ロボット１００の頭１１０には、カメラ１１１と、マイク１１２と、スピーカ１１３と、が設けられている。

【００１３】

カメラ１１１（撮像手段）は、頭１１０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１１１は、後述する制御部１２７の制御の下、撮像を行う。

10

【００１４】

マイク１１２は、１３個のマイクを含む。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１１０の周りに等間隔で配置されている。これら８個のマイクより上側に、４個のマイクが頭１１０の周りに等間隔で配置されている。さらに、１個のマイクが頭１１０の頭頂部に配置されている。マイク１１２はロボット１００の周囲で発生した音を検出する。マイク１１２は、後述の制御部１２７と協働して、音声検出手段としての役割を果たす。

【００１５】

スピーカ１１３は、カメラ１１１より下側、人の顔でいうところの口に相当する位置に設けられている。スピーカ１１３は、後述する制御部１２７の制御の下、各種の音声を出

20

【００１６】

首関節１２１は、頭１１０と胴体１２０とを連結する部材である。頭１１０は、破線で示される首関節１２１によって、胴体１２０に連結されている。首関節１２１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１１０が回転する。図２にロボット１００の頭１１０の回転の自由度を模式的に表した図を示す。首関節１２１により、ロボット１００の頭１１０は、胴体１２０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 の軸回り回転可能である。首関節１２１は、後述の足回り部１２６とともに、後述の制御部１２７と協働して、ロボット１００の各部位の動作を制御することで、自機の位置、姿勢の少なくとも一方を変える制御手段としての役割を果たす。

30

【００１７】

図３を参照する。上述の構成に加え、ロボット１００は、操作ボタン１２２と、センサ群１２３と、電源部１２４と、記憶部１２５と、足回り部１２６と、制御部１２７と、を備える。

【００１８】

操作ボタン１２２は、胴体１２０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１２２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１１３の音量調節ボタン等を含む。

【００１９】

図１に示すように、センサ群１２３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１２３は、距離センサ、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。

40

【００２０】

図３を参照する。電源部１２４は、胴体１２０に内蔵された充電電池であり、ロボット１００の各部に電力を供給する。

【００２１】

記憶部１２５は、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１２０の内部に設けられている。記憶部１２５は、後述の制御部１２７によって実行されるプログラム、カメラ１１１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１２５

50

が記憶するプログラムには、後述の呼びかけ応答処理に係る呼びかけ応答プログラム 1251、地図作成処理に係る地図作成プログラム 1252が含まれる。さらに、記憶部 125には、後述のSLAM (Simultaneous Localization And Mapping) 法で作成される部屋の地図であるSLAM地図 1253、撮像画像の特徴点等を格納するフレームデータベース 1254、後述のラベリングの音声発生確率が定義された音声発生確率データベース 1255が含まれる。

【0022】

足回り部 126は、胴体 120の下側に設けられた4つの車輪 (ホイール) を含む。図 1に示すように、4つの車輪のうち、2つが胴体 120の前側に、残り2つが後側 (不図示) が配置されている。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 127が足回り部 126の車輪を回転させると、ロボット 100は移動する。足回り部 126は、前述の首関節 121とともに、後述の制御部 127と協働して、ロボット 100の各部位の動作を制御することで、自機の位置、姿勢の少なくとも一方を変える制御手段としての役割を果たす。

【0023】

さらに、足回り部 126の車輪にはロータリエンコーダが設けられている。ロータリエンコーダで車輪の回転数を計測し、車輪の直径や車輪間の距離等の幾何学的関係を利用することで並進移動量及び回転量を計算できる。

【0024】

図3を参照する。制御部 127は、CPU (Central Processing Unit)、RAM (Random Access Memory) 等で構成される。制御部 127は、上述のロボット 100の各部に接続されており、RAMをワークスペースとして、記憶部 125に記憶されたプログラムを実行することにより、ロボット 100の各部を制御する。

【0025】

本実施の形態においては、制御部 127は、ロボット 100の各部位の動作を制御するため、前述の首関節 121、足回り部 126を制御することで、自機の位置、姿勢の少なくとも一方を変える制御手段の役割を果たす。

【0026】

さらに、制御部 127は、足回り部 126の車輪に設けられたロータリエンコーダの回転数から、自機の位置 (移動開始時の位置を基準とした自機の位置) を計測することができる。例えば、車輪の直径を D 、回転数を R (足回り部 126のロータリエンコーダにより測定) とすると、その車輪の接地部分での並進移動量は $D \cdot R$ となる。また、車輪の直径を D 、車輪間の距離を I 、右車輪の回転数を R_R 、左車輪の回転数を R_L とすると、向き変更の回転量は (右回転を正とすると) $360^\circ \times D \times (R_L - R_R) / (2 \times I)$ となる。この並進移動量や回転量を逐次足し合わせていくことで、自機位置 (移動開始時の位置及び向きを基準とした位置及び向き) を計測することができる。このように、制御部 127は、オドメトリとしても機能する。

【0027】

上述のように、ロボット 100は、住人 (ユーザ) に呼びかけられると会話するので、呼びかけられたことを判別すると、呼びかけた住人 (ユーザ) の顔検出を行う必要がある。以下、ロボット 100が行う顔検出の処理を説明する。ここでは、ユーザの呼びかけに応答する一連の処理 (呼びかけ応答処理) の中で、ロボット 100がユーザの顔検出を行う例を説明する。図4に示すように、部屋RM内にロボット 100とユーザPがあり、ロボット 100とユーザPとが正対していない場合に、ユーザPがロボット 100に呼びかける場面を想定する。

【0028】

本実施の形態においては、部屋RM内に存在する音源の位置が登録された地図 (音声発生源位置情報) があらかじめ作成されている。ロボット 100の制御部 127は、人の声がしたことを検出したときに、まず、その音の音源の方向を判別する。そして、制御部 127は、音源の方向と自機 (ロボット 100) の位置とあらかじめ作成されている部屋R

10

20

30

40

50

Mの中の地図とに基づき、判別した音源の方向に、人以外の音源が存在するか否かを判別し、存在する否かに応じて、振り向くか振り向かないかを判別する。

【0029】

呼びかけ応答処理に先立ってあらかじめ作成される実空間（ここでは部屋RM）内の地図の作成方法を説明する。制御部127の制御の下、ロボット100は、毎日決められた時刻に、部屋の中を動き回りながら、撮像し、撮像画像に基づいて部屋の地図を作成し、作成した地図を記憶部125に格納する。

【0030】

地図の作成には、SLAM法を採用する。SLAM法は、実空間の地図を作成するための手法のひとつである。この手法では、カメラの撮影する動画像の複数フレームから、同一の特徴点を追跡することで、自機の3次元位置（カメラ位置）と特徴点の3次元位置（これが集まって地図の情報を構成する）とを交互または同時に推定する処理を行う。SLAM法の詳細は、非特許文献1に記載されている。

【0031】

以下、図5のフローチャートを参照しながら、制御部127が行うSLAM法を採用した地図作成処理を説明する。制御部127は、記憶部125に記憶されている地図作成プログラム1252を実行することによって、以下の処理を実現する。

【0032】

まず、制御部127は、撮像画像を取得し、撮像画像の二次元特徴点（2D特徴点）を抽出する（ステップS11）。2D特徴点とは、画像中のエッジ部分など、画像内の特徴的な部分であり、SIFT（Scale-Invariant Feature Transform）やSURF（Speed-Up Robust Features）等のアルゴリズムを用いて取得することができる。

【0033】

具体的には、ステップS11において、制御部127は、カメラ111を制御して、撮像を行う。そして、撮像した画像から2D特徴点を抽出する。さらに、前述のようにオドメトリとしても機能する制御部127は、足回り部126のロータリエンコーダを使用して、自機（ロボット100）の現在位置を計測する。制御部127は、2D特徴点と、自機の現在位置と、を撮像画像と対応づけて記憶部125に記憶する。

【0034】

制御部127は、地図作成処理の開始後に撮像した画像が2枚以上であるか否かを判別する（ステップS12）。2枚未満であると判別すると、（ステップS12；N）、制御部127は、足回り部126を制御して、自機を所定の距離だけ移動し（ステップS19）、再びステップS11へ戻る。

【0035】

一方、撮像した画像が2枚以上であると判別した場合（ステップS12；Yes）、制御部127は、2つの画像の2D特徴点の対応を取得する（ステップS13）。2つの画像は、例えば、今回撮像した画像と、直前に撮像した画像である。

【0036】

ステップS13で取得した2つの画像の対応する特徴点（対応特徴点）の個数が、閾値未満であるか否かを判別する（ステップS14）。これは、取得した特徴点の個数が少ないと、後述のTwo-view Structure from Motion法での計算ができないためである。

【0037】

2つの画像の対応する特徴点の個数が、閾値未満であると判別した場合（ステップS14；No）、制御部127は、足回り部126を制御して、自機を所定の距離だけ移動し（ステップS19）、再びステップS11へ戻る。

【0038】

一方、2つの画像の対応する特徴点の個数が、閾値以上であると判別した場合（ステップS14；Yes）、制御部127は、2つの画像間の姿勢を推定する（ステップS15）。

10

20

30

40

50

【0039】

具体的には、ステップS15において、Two-view Structure from Motion法を用いて、2つの画像の間で対応する2D特徴点の2次元座標(2D座標)と、2つの画像のそれぞれの撮影位置(撮影時の自機の位置)の間の距離とから、2つの画像間の姿勢(それぞれの画像を取得した位置の差分(並進ベクトル t)及び向きの差分(回転行列 R))を推定する。この推定は、非特許文献2に記載されているように、エプipoラ拘束式により、対応する特徴点から基礎行列 E を求め、基礎行列 E を並進ベクトル t と回転行列 R とに分解することによって得られる。

【0040】

続いて、制御部127は、2つの画像の間で対応する2D特徴点(2Dの対応特徴点)の3次元座標(3D座標)を推定する(ステップS16)。具体的には、これは、ステップS15で算出した2つの画像間の姿勢を表す値と、2つの画像の間で対応する2D特徴点の2D座標と、を用いて推定する。

10

【0041】

制御部127は、ステップS16で推定した推定値をデータベースに登録する(ステップS17)。具体的には、制御部127は、ステップS16で求めた「2Dの対応特徴点の3D座標(X、Y、Z)」と、「2D特徴点の特徴量」(例えばSIFT等で得た特徴量)と、を記憶部125のSLAM地図1253に登録する。

【0042】

また、制御部127は、記憶部125のフレームデータベース1254に、画像の情報として、「SLAM地図内での画像の姿勢」(その画像を撮像したときの自機のSLAM座標内での位置(並進ベクトル t)及び向き(回転行列 R))と、「抽出した全ての2D特徴点」と、「すべての2D特徴点の中で3D位置(3D座標)が既知の点」と、「キーフレーム自体の特徴」と、を記憶部125のフレームデータベース1254に登録する。

20

【0043】

ここで、キーフレームとは、処理の対象となる撮像画像のことである。キーフレーム自体の特徴とは、キーフレーム間の画像類似度を求める処理を効率化するためのデータであり、画像中の2D特徴点のヒストグラム等を用いてもよいし、画像自体を「キーフレーム自体の特徴」としてもよい。

【0044】

制御部127は、処理が終了であると判別すると(ステップS18; Yes)、地図作成処理を終了する。一方、処理が終了でないと判別すると(ステップS18; No)、足回り部126を制御して、自機を所定の距離だけ移動し(ステップS19)、再びステップS11へ戻る。以上が地図作成処理である。

30

【0045】

さらに、上述のように作成したSLAM地図1253に、部屋RM内のそれぞれの位置における障害物が存在する確率を示す障害物情報として確率変数を付加してもよい。障害物情報の確率変数の値は、その値が高いほど、その位置に障害物がある可能性が高いことを表している。障害物情報の確率変数は、例えば、SLAM地図1253の作成処理におけるデータベース登録(図5のステップS17)のタイミングで、SLAM地図1253に付加することができる。

40

【0046】

さらに、本実施の形態においては、上述のように作成したSLAM地図1253に、人以外の音声発生源情報を付加したものを使用して、ロボット100が、検出した音がか人であるか否かを判別する。

【0047】

人以外の音声発生源情報は、例えば、SLAM地図1253の作成処理におけるデータベース登録(図5のステップS17)のタイミングで、SLAM地図1253に付加することができる。

【0048】

50

音声発生源の特定は、例えば以下のような方法で行う。S L A M地図1 2 5 3の作成時にロボット1 0 0が部屋R M内を動き回り撮像した画像に対して、一般画像認識（画像に含まれる物体を一般的な名称で認識する処理）を行い、音声発生源か否かのラベリングする方法を用いてもよい。画像内で音声発生源としてラベリングされた領域に存在する2 D特徴点に対応する地図内の地点に対して第1の値（第2の値より大きい値）を登録する。また、それ以外の2 D特徴点に対応する地点には、第2の値（第1の値より小さい値）を登録する。具体的には、ロボット1 0 0が通過した地点には、第2の値を登録し、ロボット1 0 0が通過した際に、接触センサ、距離センサ等により障害物に接触したと判別した地点には、第1の値を登録する。

【0 0 4 9】

10

上述の例では、確率変数を2値とすることを説明した。あるいは、一般画像認識結果の尤度に、ラベリングの音声発生確率を乗じた値を確率変数としてもよい。

【0 0 5 0】

ラベリングの音声発生確率、ここでは、部屋R M内のそれぞれの位置における音声発生源の確率を示す情報（確率変数）はあらかじめ記憶部1 2 5の音声発生確率データベース1 2 5 5に登録されているものとする。確率変数の値は、値が高いほど、当該位置に人以外の音声発生源が存在する可能性が高いことを示す。

【0 0 5 1】

音声発生確率データベース1 2 5 5に登録されているラベリングの音声発生確率として、例えば、換気扇：0 . 8、ドア：0 . 5、観葉植物：0といった値が登録されている。換気扇は、動作中にそれなりの音を出し、ドアは、開け閉めする人により出る音の大きさに差があり、置かれているだけの観葉植物については、音は発生しないといった観点で、このような値が規定される。

20

【0 0 5 2】

また、ラベリングの音声発生確率は、時刻、季節、気温等に応じて、複数の値を規定しておいてもよい。季節に応じたラベリングの場合、例えば、夏：0 . 8、冬：0とする。夏場であれば、窓を開けることが多いため、室内でも室外で発生した音が聞こえることがあり、冬場の窓を閉め切った状態であれば、室外の音はほぼ聴こえないからである。

【0 0 5 3】

また、一般画像認識結果の尤度を使用するのは次のような理由による。一般画像認識を使用した場合、どのような画像に対しても、認識の精度が高いというわけではない。一般画像認識結果の尤度を用いることで、一般画像認識が誤認識した場合の影響を減らすことができる。

30

【0 0 5 4】

また、一般画像認識ではなく他の手法を用いてもよい。非特許文献3に記載されているBag-of-featuresという手法がある。この手法は、画像中の物体がどのカテゴリに属するかを求める画像分類問題の手法である。

【0 0 5 5】

あるいは、一般画像認識ではなく、ユーザが指定した音声発生源の領域、音声発生源となる物体を示す情報を、作成したS L A M地図1 2 5 3に追加してもよい。この場合、例えば、ロボット1 0 0は、タッチパネル、ディスプレイ等の表示装置と、タッチパネル、キーボード等の入力装置を備え、ユーザに対して作成したS L A M地図1 2 5 3を提示して、ユーザに音声発生源を入力されるようにしてもよい。

40

【0 0 5 6】

あるいは、ロボット1 0 0は、S L A M地図1 2 5 3後に、部屋R内を動き回り、部屋R M内にある物体を指さしして、ユーザに、当該物体が音声発生源であるか等を尋ねてもよい。ユーザの回答に基づく音声発生源の情報を、S L A M地図1 2 5 3に追加することができる。

【0 0 5 7】

あるいは、S L A M地図1 2 5 3後に、部屋R M内の物体を撮像し、撮像画像を表示装

50

置に表示し、ユーザに、当該物体が音声発生源であるか等を尋ねてもよい。この場合も、ユーザの回答に基づく音声発生源の情報を、SLAM地図1253に追加することができる。

【0058】

次に、音を検出した場合に、地図を使用して、検出した音の音源が人であるか否かを判別し、判別結果に応じて応答する呼びかけ応答処理を説明する。呼びかけ応答処理の開始に先立って、上述の地図作成処理はすでに実行されているものとし、SLAM地図1253、フレームデータベース1254、音声発生確率データベース1255には、適宜の情報が登録済みであるとする。

【0059】

制御部127は、記憶部125の呼びかけ応答プログラム1251を実行することで、以下の呼びかけ応答処理を行い、検出した音声発生源が特定の音声発生源（ここでは人間）であるか否かを判別する判別手段として機能する。

【0060】

図6のフローチャートを参照しながら、呼びかけ応答処理を説明する。制御部127は、ロボット100の周辺である程度の大きさの音を検出したか否かを判別する（ステップS101）。具体的には、制御部127は、1つ以上のマイク112に所定の閾値以上の振幅の音が入力されたか否かを判別する。なお、所定の大きさとは、マイク112の感度によるものとする。

【0061】

マイク112により所定の大きさの音を検出できない場合（ステップS101；No）、制御部127は、音を検出するまで待ち受ける。

【0062】

一方、ある程度の大きさの音を検出したと判別した場合（ステップS101；Yes）、制御部127は、マイク112により検出した音が人間の声か否かを判別する（ステップS102）。具体的には、制御部127は、ステップS101で検出した音が特定の周波数帯域の音であるか否かを判別する。ステップS101で検出した音が人間の声でない場合（ステップS102；No）、制御部127はステップS101へ戻り、音を検出するまで待ち受ける。

【0063】

一方、人間の声であると判別すると（ステップS102；Yes）、制御部127は、音源の位置（ここではユーザPの音が発せられた位置）を求めるため、音声定位を行う（ステップS103）。ここでは、音源の位置を推定するため、音源定位のひとつの手法であるMUSIC（Multiple Signal Classification）を採用することとする。なお、音源定位の最中に音源であるユーザPは移動せず、静止しているものとする。

【0064】

図7を参照して音源定位を説明する。まず、マイク112に入力された音声を時間周波数変換する（ステップS10301）。ここでは、時間周波数変換として、STFT（Short-Time Fourier Transform）（短時間フーリエ変換）を行う。

【0065】

音源数をNとすると、第n番目の音源の信号 S_n は、下記式（1）で表せる。

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

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）。

【0066】

マイク112で観測される信号は、マイク112の数をMとすると、下記式（2）で表せる。

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

【0067】

10

20

30

40

50

音源から出た音は、空気を伝わってマイク 1 1 2 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で、マイク 1 1 2 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 1 1 2 で観測される信号 $X_m(\omega, f)$ は下記式 (3) のように表される。

【 0 0 6 8 】

【 数 1 】

$$X_m(\omega, f) = \sum_{n=1}^N S_n(\omega, f) H_{nm}(\omega) \quad \dots \quad (3)$$

10

【 0 0 6 9 】

ロボット 1 0 0 は、マイク 1 1 2 を複数有しているので、マイク 1 1 2 全体で観測される信号 $x(\omega, f)$ は下記式 (4) で表すことができる。

【 0 0 7 0 】

【 数 2 】

$$x(\omega, f) = \begin{bmatrix} X_1(\omega, f) \\ X_2(\omega, f) \\ \vdots \\ X_M(\omega, f) \end{bmatrix} \quad \dots \quad (4)$$

20

【 0 0 7 1 】

同様に、全音源の信号 $s(\omega, f)$ も下記式 (5) で表すことができる。

【 0 0 7 2 】

【 数 3 】

$$s(\omega, f) = \begin{bmatrix} S_1(\omega, f) \\ S_2(\omega, f) \\ \vdots \\ S_N(\omega, f) \end{bmatrix} \quad \dots \quad (5)$$

40

【 0 0 7 3 】

同様に、第 n 番目の音源の伝達関数 $h_n(\omega)$ は下記式 (6) で表すことができる。

【 0 0 7 4 】

50

【数 4】

$$h_n(\omega) = \begin{bmatrix} H_{n1}(\omega) \\ H_{n2}(\omega) \\ \vdots \\ H_{nM}(\omega) \end{bmatrix} \cdots \quad (6)$$

10

【0075】

全ての伝達関数を下記式(7)のように表記する。

$$h(\quad) = [h_1(\quad), h_2(\quad), \dots, h_N(\quad)] \cdots (7)$$

【0076】

上記の式(7)で表される伝達関数を、上述の式(3)に適用すると、下記式(8)の
ように表される。

20

$$x(\quad, f) = h(\quad) s(\quad, f) \cdots (8)$$

【0077】

$h_n(\quad)$ 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 $s_n(\quad, f)$ は無相関とみなせるので、 $x(\quad, f)$ 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に分解することを考える
。

【0078】

再び図7を参照する。次の式(9)に示すように相関行列を計算する(ステップS10
302)。ここで、*は複素共役転置を意味する。

30

【0079】

【数5】

$$R(\omega, f) = \sum_{l=0}^{L-1} x(\omega, f+l) x^*(\omega, f+l) \cdots \quad (9)$$

【0080】

続いて、固有値分解する(ステップS10303)。ここで、固有値 $\lambda_m(\quad, f)$ と
固有ベクトル $e_m(\quad, f)$ 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

【0081】

原理的には、 $h_n(\quad)$ は部分空間の固有ベクトル $e_m(\quad, f)$ ($m = 1 \sim N$) 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 $e_m(\quad, f)$ ($m = N+1 \sim M$) が $h_n(\quad)$ と直交することを使って音源定位
を実現する。

【0082】

しかし、音源であるユーザPが部屋RM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、

50

仮の音源位置を決め、仮の音源位置の伝達関数をあらかじめ用意しておき、音源定位を行う。

【 0 0 8 3 】

図 8 に、仮の音源位置とマイクの配置の一例を示す。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 1 1 2 を表す。なお、ここでは、便宜上 1 3 個のマイク 1 1 2 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 0 8 4 】

複数のマイク 1 1 2 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、マイク 1 1 2 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、それぞれの伝達関数 h () を予め計算しておく。

【 0 0 8 5 】

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 () を予め計算しておけばよい。また、あるいは、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 0 8 6 】

再び図 7 を参照する。次の式 (1 0) を使用して、周波数帯毎の M U S I C スペクトルを計算する (ステップ S 1 0 3 0 4) 。

【 0 0 8 7 】

【 数 6 】

$$M_{\theta}(\omega, f) = \frac{h^*_{\theta}(\omega)h_{\theta}(\omega)}{\sum_{m=N+1}^M |h^*_{\theta}(\omega)e_m(\omega, f)|^2} \cdots (10)$$

【 0 0 8 8 】

ここで、式 (1 0) 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、音源の方向とあらかじめ決めた角度 (1、 2、 ...

N) のいずれかが近い場合、つまり h_n () と h () が近い場合、式 (1 0) の値は極端に大きなものになる。図 8 に示す例では、音源である人と仮の音源 2 の位置が近いため、 2 の伝達関数を使用した場合、式 (1 0) の値が極端に大きくなることが想定される。

【 0 0 8 9 】

そして、統合した M U S I C のパワーを求めるため、式 (1 1) に示すように周波数帯毎の M U S I C スペクトルを重み付け加算する (ステップ S 1 0 3 0 5) 。

【 0 0 9 0 】

【 数 7 】

$$M(f) = \sum_{\omega} w(\omega)M(\omega, f) \cdots (11)$$

【 0 0 9 1 】

重み付け係数は、固有値 m (, f) が大きいほど大きくすれば、 S_n (, f) に

10

20

30

40

50

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 $S_n(\quad, f)$ に殆どパワーがない場合の悪影響を軽減できる。

【0092】

続いて、パワースペクトルから適切なピーク（極大値）を選択する（ステップS10306）。具体的には、まず、複数のピークを求め、その中から適切なピークを選択し、選択したピークにおける θ を音源方向とする。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。

【0093】

また、テレビが点いている、ドアホンが鳴る等の部屋RM内に他の音源がある場合でも、多くの場合、人は、テレビ、ドアホン等の周囲の音より大きな声でロボット100に呼びかけると考えられる。よって、人の声のパワーの方が、人以外のテレビ、ドアホン等の音源から発せられる音のパワーより大きくなることが想定される。よって、単純にパワーが最大となる仮の音源位置を示す θ を音源方向として選択しても問題はない。ただし、周囲の環境などによっては、パワーが最大となる仮の音源位置ではなく、パワーが2番目あるいはそれ以降となる仮の音源位置を、音源方向と選択することが適切な場合もある。このようにして、制御部127は、音源方向、ここでは、ロボット100の位置から見たユーザPがいる方向、を判別することができる。

【0094】

音源定位の処理は以上である。ここでは、平面を仮定して説明したが、3次元を仮定しても上記説明は成り立つ。

【0095】

再び図6を参照する。ステップS103の音源定位を実行して音源方向を判別すると、制御部127は、音源方向を示す情報として、ロボット100の向いている方向に対する音源の方向を示す角度 θ を記憶部125に記憶する。続いて、制御部127は、ステップS104へ進み、撮影画像と、地図（SLAM地図1253、フレームデータベース1254）を用いて自機位置推定の処理を実行する。

【0096】

図9を参照して、自機位置の推定の処理を説明する。制御部127は、カメラ111により撮像された画像の二次元特徴点（2D特徴点）を抽出する（ステップS10401）。具体的には、制御部127は、カメラ111を制御して撮像し、撮像した画像から2D特徴点を抽出する。

【0097】

続いて、制御部127は、記憶部125のフレームデータベース1254を参照して、フレームデータベース1254に登録されている以前のフレームの情報から、その画像の情報に含まれている2D特徴点のうち、3D位置が既知である2D特徴点を取得し、取得した2D特徴点から、ステップS10401で抽出した2D特徴点と、対応が取れる特徴点を抽出する（ステップS10402）。ここで、3D位置が既知であるとは、即ち、2D特徴点がSLAM地図に登録されていることを意味する。

【0098】

制御部127は、ステップS10402で抽出した対応が取れる特徴点の個数が、閾値以上であるか否かを判別する（ステップS10403）。閾値未満であると判別した場合（ステップS10403；No）、制御部127は、足回り部126を制御して、自機を所定の距離だけ移動し（ステップS10406）、再びステップS10401へ戻る。

【0099】

一方、ステップS10402で抽出した対応特徴点の個数が、閾値以上であると判別した場合（ステップS10403；Yes）、制御部127は、記憶部125のSLAM地図1253から、ステップS10402で抽出した対応特徴点それぞれの3D座標（ X_i, Y_i, Z_i ）を取得する（ステップS10404）。

10

20

30

40

50

【 0 1 0 0 】

続いて、制御部 1 2 7 は、自機の姿勢を推定する（ステップ S 1 0 4 0 5）。ここでは、制御部 1 2 7 は、対応特徴点の S L A M 地図上の 3 D 位置と、対応特徴点のフレーム座標（2 D 座標）の関係から自機の姿勢（並進ベクトル t 及び回転行列 R で表される自機の位置及び向き）を推定する。

【 0 1 0 1 】

具体的には、今撮像した画像に含まれている対応特徴点のフレーム座標を (u_i, v_i) とし、その対応特徴点の 3 D 座標を (X_i, Y_i, Z_i) とする（ i は 1 から対応特徴点の数までの値を取る）。ここで、各対応特徴点の 3 D 位置 (X_i, Y_i, Z_i) を下記式（1 2）によってフレーム座標系に投影した値 (u_{xi}, v_{xi}) とフレーム座標 (u_i, v_i) とは理想的には一致する。

$$(u_{xi} \ v_{xi} \ 1)' \sim A(R | t)(X_i \ Y_i \ Z_i \ 1)' \dots (12)$$

【 0 1 0 2 】

しかし、実際には (X_i, Y_i, Z_i) にも (u_i, v_i) にも誤差が含まれているため、 (u_{xi}, v_{xi}) と (u_i, v_i) とが一致することはめったにない。そして、未知数は R と t （3 次元空間ではそれぞれ 3 次元となり、 $3 + 3 = 6$ が未知数の個数である）だけなのに、数式は対応特徴点の個数の 2 倍存在する（対応特徴点一つに対して、フレーム座標の u, v それぞれに対する式が存在するため）ことになるため、過剰条件の連立一次方程式になり、上述したように最小二乗法で求めることになる。

【 0 1 0 3 】

具体的には、制御部 1 2 7 は、以下の式（1 3）のコスト関数 $E1$ を最小化する姿勢（並進ベクトル t 及び回転行列 R ）を求める。

【 0 1 0 4 】

【 数 8 】

$$E1 = \sum_{i=1}^{\text{対応特徴点数}} ((u_i - u_{xi})^2 + (v_i - v_{xi})^2) \dots (13)$$

【 0 1 0 5 】

このように求めた値が、S L A M 法で求めた S L A M 座標での自機の姿勢（並進ベクトル t 及び回転行列 R で表される自機の位置及び向き）を示す値である。このようにして算出した値により自機の姿勢が推定される。以上が自機位置推定の処理である。

【 0 1 0 6 】

再び、図 6 を参照する。制御部 1 2 7 は、ステップ S 1 0 4 の自機位置の推定の処理が終わると、ステップ S 1 0 5 へ進み、S L A M 地図 1 2 5 3 と音声発生確率データベース 1 2 5 5 とを参照して、ステップ S 1 0 4 で推定した自機位置から、ステップ S 1 0 3 で求めた音源方向に、人以外の音声発生源が存在する確率を取得する（ステップ S 1 0 5）。ここでは、音源方向の各点の確率の平均を求め、求めた平均を人以外の音声発生源が存在する確率としてもよい。あるいは、音源方向の各点の確率について最大値を人以外の音声発生源が存在する確率としてもよい。

【 0 1 0 7 】

次に、制御部 1 2 7 は、ステップ S 1 0 5 で求めた人以外の音声発生源が存在する確率が閾値以上であるか否かを判別する（ステップ S 1 0 6）。人以外の音声発生源が存在する確率が閾値以上であると判別した場合（ステップ S 1 0 6 ; Y e s）、制御部 1 2 7 は、音源方向の音源は人以外であると判別して、首関節 1 2 1 を回転駆動させず、再びステップ S 1 0 1 へ戻り、音の入力を待ち受ける。

【 0 1 0 8 】

一方、人以外の音声発生源が存在する確率が閾値未満であると判別した場合（ステップ S 1 0 6 ; N o）、制御部 1 2 7 は、ステップ S 1 0 7 へ進む。

【 0 1 0 9 】

続いて制御部 1 2 7 は、頭 1 1 0 の回転をさせるため、首関節 1 2 1 を回転駆動させる（ステップ S 1 0 7）。ここで、制御部 1 2 7 は、ロボット 1 0 0 の頭 1 1 0 の正面（カメラ 1 1 1 のレンズ面）が音源（ユーザ P）の方向に向くまで、頭 1 1 0 を回転する。具体的には、制御部 1 2 7 は、記憶部 1 2 5 に記憶されている音源定位により求められた角度に基づいて、求めた角度だけ頭 1 1 0 を回転し、その後、回転駆動を停止する。このようにして、カメラ 1 1 1 のレンズ面を音源（ユーザ P）がいる方向に向ける。

【 0 1 1 0 】

回転駆動を停止した後、制御部 1 2 7 は、顔検出の処理を実行する（ステップ S 1 0 8）。まず、制御部 1 2 7 は、カメラ 1 1 1 を制御して撮像し、撮像した画像に対して以下の処理を施すことで、顔検出処理を実行する。

10

【 0 1 1 1 】

制御部 1 2 7 は、まず、ピラミッド画像を作成する。ピラミッド画像とは、元画像を一定の比率で縮小を繰り返して作成した一連の画像群であり、ピラミッド画像の各階層に対して、固定サイズの顔検出器を適用することで様々なサイズ（つまり距離に相当）の顔を検出することができる。ここでは、回転によるカメラの見え方は対象までの距離によって変わるので、ピラミッド画像を使用して顔検出を行う。

【 0 1 1 2 】

まず、顔探索対象を最初の階層に設定する。ここでは縮小前の元の画像とする。最初の検出窓を設定する。初期位置は例えば左上の隅とする。設定した検出窓に対して、固定サイズの顔検出器を適用する。この階層でのスライドによる探索が完了したかを判定する。スライドによる探索が完了でないなら、検索窓をスライドさせ、再度顔検出を行う。スライドによる探索が完了ならば、ピラミッド画像のすべての階層での処理が完了したかの判定を行う。すべての階層での処理が完了でないなら、階層を移動し、移動先の階層でもスライドによる顔検出を行う。すべての階層での処理が完了したならば、顔検出の処理を終了する。

20

【 0 1 1 3 】

なお、ロボット 1 0 0 から近い顔画像は、画角に入りきらない場合があることと、全体の計算負荷の割合が小さいことを考慮して、縮小率の大きい階層の顔探索はしないほうがより望ましい。

30

【 0 1 1 4 】

顔検出処理により、撮像画像から顔を検出することができなかった場合（ステップ S 1 0 8；N）、制御部 1 2 7 は、再びステップ S 1 0 1 に戻る。

【 0 1 1 5 】

一方、顔検出が成功すると（ステップ S 1 0 8；Y e s）、続いて、制御部 1 2 7 は、ユーザ P がロボット 1 0 0 に注目しているかどうかを判別する（ステップ S 1 0 9）。具体的には、制御部 1 2 7 は、カメラ 1 1 1 を制御して、ユーザ P を撮像し、撮像した画像からユーザ P の顔が、ロボット 1 0 0 の方を向いているか否かを判別する。ユーザ P がロボット 1 0 0 に注目していないと判別すると（ステップ S 1 0 9；N o）、再びステップ S 1 0 1 へ戻り、音の入力を待ち受ける。

40

【 0 1 1 6 】

一方、ユーザ P の顔が、ロボット 1 0 0 の方を向いていると判別すると（ステップ S 1 0 9；Y e s）、制御部 1 2 7 は、ユーザ P に近づくように所定の距離だけ移動し（ステップ S 1 1 0）、ユーザ P との距離が決められた距離以下となったかを判別する（ステップ S 1 1 1）。このユーザ P とロボット 1 0 0 との間の決められた距離は、ロボット 1 0 0 が、ユーザ P が発声する内容を音声認識することができる程度の距離である。制御部 1 2 7 は、ロボット 1 0 0 とユーザ P との距離が決められた距離以下ではないと判別した場合に（ステップ S 1 1 1；N o）、再びステップ S 1 1 0 に戻る。

【 0 1 1 7 】

一方、制御部 1 2 7 は、ユーザ P との距離が所定の距離となったと判別した場合に（ス

50

テップ S 1 1 1 ; Y e s)、ユーザ P と対話する (ステップ S 1 1 2)。例えば、制御部 1 2 7 は、スピーカ 1 1 3 を制御して、ユーザ P に対して、例えば、「何かご用ですか？」と話しかけ、また、マイク 1 1 2 から入力したユーザの発言を音声解析し、解析した内容に基づいて、なんらかの音声をスピーカ 1 1 3 から出力する。

【 0 1 1 8 】

以上、説明したように、本実施の形態においては、ロボット 1 0 0 は、あらかじめ作成した S L A M 地図 1 2 5 3 に基づき、判別した音源方向に人以外の音源がある場合、人に呼ばれたのではないと判別する。よって、人以外の音源であるテレビ、ラジオ等から人の声が聞こえた場合であっても、振り向かないので、無駄な動作を減らすことができる。

【 0 1 1 9 】

なお、上述の説明においては、回転駆動は y a w を前提で説明したが、他の方向の回転があっても成立する。

【 0 1 2 0 】

実施の形態 1 においては、ロボット 1 0 0 は、ユーザ P の方向へ近づくよう、単に移動したが、ロボット 1 0 0 は、S L A M 地図 1 2 5 3 を使用して、部屋 R M 内を移動し、ユーザに近づいてもよい。

【 0 1 2 1 】

(実施の形態 2)

実施の形態 2 においては、ユーザ P から呼びかけられたロボット 1 0 0 が、S L A M 地図 1 2 5 3 を使用して、移動経路を作成し、移動経路に沿って移動する。ロボット 1 0 0 が備える構成は、実施の形態 1 と同様である。以下の説明においては、実施の形態 2 に特有の構成を中心に説明する。

【 0 1 2 2 】

実施の形態 1 と同様に、あらかじめ S L A M 地図 1 2 5 3 が作成されているものとする。

【 0 1 2 3 】

実施の形態 2 においては、図 1 0 に示すように、記憶部 1 2 5 には後述の呼びかけ移動処理のための呼びかけ移動プログラム 1 2 5 6 が記憶されているものとする。制御部 1 2 7 は、呼びかけ移動プログラム 1 2 5 6 を実行することによって、以下の処理を行う。

【 0 1 2 4 】

図 1 1 にユーザ P に呼びかけられたときに、ロボット 1 0 0 がユーザ P のいる場所まで移動する処理 (呼びかけ移動処理) のフローチャートを示す。なお、上述の呼びかけ応答処理と同様であるので、ここでは、所定の大きさの音を検出し、検出した音が人間の声であると判別したと仮定して、説明を行う。

【 0 1 2 5 】

制御部 1 2 7 は、撮像した画像と S L A M 地図 1 2 5 3 を用いて自機位置推定の処理を実行する (ステップ S 2 0 1)。自機位置推定の処理については図 9 を参照して説明したため、ここでは、説明を省略する。

【 0 1 2 6 】

続いて、制御部 1 2 7 は、S L A M 法によるユーザ P の顔の位置を推定する処理 (顔位置推定の処理) を実行する (ステップ S 2 0 2)。図 1 2 を参照して、顔位置推定の処理を説明する。制御部 1 2 7 は、カメラ 1 1 1 を制御して撮像し、撮像した画像から二次元特徴点 (2 D 特徴点) を抽出する (ステップ S 2 0 2 0 1)。特徴抽出には S I F T や S U R F 等のアルゴリズムを用いる。

【 0 1 2 7 】

制御部 1 2 7 は、ステップ S 2 0 2 0 1 で抽出した 2 D 特徴点のうち、撮像した画像の顔の領域内の特徴点 (2 D 顔特徴点) を抽出する (ステップ S 2 0 2 0 2)。顔領域内に特徴点がない場合は、顔パーツ検出の結果を特徴点として用いる。

【 0 1 2 8 】

制御部 1 2 7 は、顔位置推定の処理開始後に撮像した画像が 2 枚以上であるか否かを判

10

20

30

40

50

別する（ステップS20203）。2枚未満であると判別すると、（ステップS20203；N）、制御部127は、足回り部126を制御して、自機を所定の距離だけ移動し（ステップS20208）、再びステップS20201へ戻る。

【0129】

一方、撮像した画像が2枚以上であると判別した場合（ステップS20203；Yes）、制御部127は、2つの画像の2D顔特徴点の対応を取得する（ステップS20204）。

【0130】

制御部127は、ステップ20202で抽出した対応する2D顔特徴点の個数が、閾値以上であるか否かを判別する（ステップS20205）。閾値未満であると判別した場合（ステップS20205；No）、制御部127は、足回り部126を制御して、自機を所定の距離だけ移動し（ステップS20208）、再びステップS20201へ戻る。

【0131】

一方、2つの画像の対応する2D顔特徴点の個数が、閾値以上であると判別した場合（ステップS20205；Yes）、制御部127は、2つの画像間の姿勢を推定する（ステップS20206）。

【0132】

具体的には、2つの画像の間で対応する2D顔特徴点の2次元座標（2D座標）と、2つの画像のそれぞれの撮影位置（撮影時の自機の位置）の間の距離と、に対して、Two-view Structure from Motion法を用いて、2つの画像間の姿勢（それぞれの画像を取得した位置の差分（並進ベクトル t ）及び向きの差分（回転行列 R ））を推定する。

【0133】

続いて、制御部127は、2つの画像の間で対応する2D顔特徴点の3次元座標（3D座標）を推定する（ステップS20207）。具体的には、これは、ステップS20206で算出した2つの画像間の姿勢を表す値と、2つの画像の間で対応する2D顔特徴点の2D座標と、を用いて推定する。以上が、顔位置推定の処理である。

【0134】

図11を再び参照する。ステップS203に進み、制御部127は、自機位置からユーザPの顔位置までの経路の作成を行う（ステップS203）。

【0135】

実施の形態2においては、SLAM地図1253に、部屋RM内のそれぞれの位置における障害物が存在する確率を示す障害物情報の確率変数を付加したものを使用する。障害物情報の確率変数の値は、その値が高いほど、その位置に障害物がある可能性が高いことを表している。

【0136】

経路の作成は、まず、記憶部125からSLAM地図1253を読み出し、SLAM地図1253上にランダムにノードを配置する（ノード情報を追加する）。このとき、自機（ロボット100）と同じ高さにノードを配置する。また、障害物情報の確率変数が閾値以上である位置（点）を中心とした一定の範囲内には、ノードを配置しない。

【0137】

なお、高さは、重力方向のオフセット値を用いて推定する。具体的には、ロボット100が過去に移動した位置から面推定を行い、法線ベクトル（重力方向）を求めて、自機位置と法線ベクトルの内積を求め自機位置の高さとする。経路中のノードも同様にして値を求める。自機位置の高さを表す値と、ノードの高さを表す値の差が決められた閾値以内であれば、自機位置の高さとノードの高さが同じであるとみなす。

【0138】

配置したノードについて、当該ノードを中心とした一定の範囲内に存在する他のノードとをつなぐ。これを、ランダムに配置した全てのノードについて行う。このようにして、グラフ構造を作る。

10

20

30

40

50

【0139】

ステップS201の自機位置推定で推定した自機の位置の一番近くに存在するノード、ステップS202顔位置推定で推定したユーザPの顔の一番近くに存在するノード、をそれぞれ選択する。そして、ダイクストラ法により、選択した2つのノード間の最短経路を求める。

【0140】

その後、求めた最短経路に従って、移動する(ステップS204)。以上が、実施の形態2にかかる呼びかけ移動処理である。

【0141】

(変形例)

実施の形態2の呼びかけ移動処理では、2つの画像間の姿勢の推定(図12のステップS20206)に、Two-view Structure from Motion法を用いた。姿勢の推定はこれに限られない。

【0142】

上述の自機位置推定の処理における姿勢の推定(図9のステップS10405)で行ったように姿勢を推定してもよい。この方法の方が、精度が高く、計算に要する時間も少ない。また、あるいは、被写体の顔のサイズが標準的な顔サイズであると仮定し、顔検出結果(顔のサイズ、位置)とカメラパラメータ(画角、焦点距離)を用いて、自機と顔間の相対的な姿勢を推定し、SLAM法により求めた自機の姿勢のSLAM地図上における推定結果を用いて、SLAM地図上における顔の姿勢を算出してもよい。また、あるいは、ロボット100に距離センサを設け、距離センサを使用して、ロボット100とユーザPの顔の間の距離を測定してもよい。

【0143】

上記の実施の形態2のSLAM地図1253を使用した移動処理は、実施の形態1における呼びかけ応答処理時の、図6のステップS110の移動の際にも応用可能である。

【0144】

SLAM地図1253の精度を上げる方法として、次のようなものがある。フレームデータベースにある程度の撮像した画像のデータが蓄積されたところで、3D位置が既知で無い特徴点の対応を再探索し、3D位置を計算してもよい。

【0145】

また、バンドルアジャストメント処理を行い、キーフレーム姿勢とMap点の3D位置の精度を向上させることができる。バンドルアジャストメント処理とは、カメラ姿勢(キーフレーム姿勢)とMap点の3D位置とを同時に推定する非線形最適化法である。この方法を使用することで、SLAM地図上の点を、画像上に投影させたときに発生する誤差が最小になるような最適化を行うことができる。

【0146】

また、ループクロージング処理を行ってもよい。ループクロージング処理とは、以前にきたことのある同じ場所に戻ってきたことを認識した場合に、以前その場所にいた時の姿勢の値と現在の姿勢の値とのずれを用いて、以前にきた時から今までの軌跡中の画像や、関連するMap点の3D位置を修正することである。

【0147】

制御部127は、音源から発せられる音が、ロボット100に向けられたものか否かを判別し、ロボット100に向けられたものであると判別した場合だけ、音検出を行い、その音が人間の声であるか否かを判別し、人間の声であると判別した場合に、上述の処理により振り返り判定を行ってもよい。この場合、例えば、マイク112に含まれる13個のマイクとして単一指向性マイクを使用することで、音源から発せられる音の方向を精度良く判別することができる。

【0148】

本発明は、上記実施形態に限定されず、本発明の要旨を逸脱しない部分での種々の修正は勿論可能である。

10

20

30

40

50

【 0 1 4 9 】

上述の実施の形態では、ロボット 1 0 0、ユーザ P とともに屋内（部屋 R M 内）にいる例を説明したが、屋外であっても同様に、本発明を採用して、ロボットは振り向き判定を行うことができる。

【 0 1 5 0 】

上述の顔検出では、ピラミッド画像の階層を順次移動して、顔検出を行ったが、制御部 1 2 7 は、マイク 1 1 2 の入力音声の大きさ（振幅の大きさ）に基づいて、ロボット 1 0 0 から音源までの距離を推定し、推定した距離に基づいて、ピラミッド画像の全ての階層について顔検出を行わないようにしてもよい。例えば、ユーザ P が近くにいると判別した場合、ある程度小さく縮小したピラミッド画像を使用する必要はない。

10

【 0 1 5 1 】

また、制御部 1 2 7 は、マイク 1 1 2 の入力音声を、そのときの、人か否かの判別結果とともに、記憶部 1 2 5 に記憶しておいてもよい。再度、同じ音を検出したとき、人か否かの判別が容易となるからである。

【 0 1 5 2 】

また、ユーザが、あらかじめ、ロボット 1 0 0 のマイク 1 1 2 の入力とないうる人以外の音のデータを記憶させてもよい。例えば、インターホンの音、電話の呼び出し音である。よって、ロボット 1 0 0 は、当該音声が入った場合には、人以外であると判別することができる。

20

【 0 1 5 3 】

また、ロボット 1 0 0 が屋外にいる場合、あらかじめ、周囲を撮像して、撮像した画像を画像認識しておくことが好ましい。屋外の場合、音源になりうるものの数が室内に比べて多くなることが想定されるからである。例えば、公園内であれば、大型スピーカが設置されていることがあり、あらかじめ、撮像画像から大型スピーカを画像認識しておき、音源として記憶することで、ロボット 1 0 0 が、振り向き判定をしやすくなる。

【 0 1 5 4 】

上述の実施の形態では、音声発生源が人間であるか否かを判別する構成を説明した。しかし、判別する特定の音声発生源は、人間だけに限られない。音声発生源の判別の対象に、人間のように自らの意思で話す人工知能を搭載したロボットを含めることができる。本発明を採用することで、人間に加え、人間のように自らの意思で話す人工知能を搭載したロボットの音声についても、同様に判別することができる。

30

【 0 1 5 5 】

また、本発明に係る顔認識装置は、専用のシステムによらず、通常のコンピュータシステムを用いて実現可能である。例えば、ネットワークに接続されているコンピュータに、上記動作を実行するためのプログラムを、コンピュータシステムが読み取り可能な記録媒体（CD-ROM（Compact Disc Read Only Memory）、MO（Magnetooptical）等）に格納して配布し、当該プログラムをコンピュータシステムにインストールすることにより、上述の処理を実行する顔認識装置を構成してもよい。

【 0 1 5 6 】

40

また、コンピュータにプログラムを提供する方法は任意である。例えば、プログラムは、通信回線の掲示板（BBS（Bulletin Board System））にアップロードされ、通信回線を介してコンピュータに配信されてもよい。また、プログラムは、プログラムを表す信号により搬送波を変調した変調波により伝送され、この変調波を受信した装置が変調波を復調してプログラムを復元するようにしてもよい。そして、コンピュータは、このプログラムを起動して、OS（Operating System）の制御のもと、他のアプリケーションと同様に実行する。これにより、コンピュータは、上述の処理を実行する顔認識装置として機能する。

【 0 1 5 7 】

この発明は、この発明の広義の精神と範囲を逸脱することなく、様々な実施の形態及び

50

変形が可能とされるものである。また、上述した実施の形態は、この発明を説明するためのものであり、この発明の範囲を限定するものではない。すなわち、この発明の範囲は、実施の形態ではなく、請求の範囲によって示される。そして、請求の範囲内及びそれと同等の発明の意義の範囲内で施される様々な変形が、この発明の範囲内とみなされる。この発明の範囲内とみなされる。以下に、本願出願の当初の特許請求の範囲に記載された発明を付記する。

【 0 1 5 8 】

(付 記)

(付 記 1)

音声を検出する音声検出手段と、

前記音声検出手段が検出した音声の音声発生源が特定の音声発生源であるか否かを判別する判別手段と、

前記判別手段の判別結果に基づいて自機を制御する制御手段と、

を備える音声検出装置。

10

【 0 1 5 9 】

(付 記 2)

前記制御手段は、前記判別手段が、前記音声検出手段が検出した音声の音声発生源が前記特定の音声発生源であると判別した場合、自機の位置、姿勢の少なくとも一方を変えるよう自機を制御する、

付記 1 に記載の音声検出装置。

20

【 0 1 6 0 】

(付 記 3)

撮像部と、

前記特定の音声発生源以外の音声発生源であって、登録された音声発生源の位置を示す情報を含む音声発生源位置情報があらかじめ記憶された記憶部と、

をさらに備え、

前記判別手段は、前記音声検出手段が検出した音声の音声発生源の位置を判別し、判別した位置が、前記音声発生源位置情報に含まれる前記登録された音声発生源の位置であるか否かを判別し、

前記判別手段が判別した位置が、前記音声発生源位置情報に含まれる前記登録された音声発生源の位置でないと判別した場合に、前記制御手段は、前記撮像部の撮像方向を前記判別手段が判別した位置に向けるように、自機の位置、姿勢の少なくとも一方を変える、

付記 1 または 2 に記載の音声検出装置。

30

【 0 1 6 1 】

(付 記 4)

前記音声発生源位置情報は、さらに、前記登録された音声発生源の位置に、前記特定の音声発生源以外の音声発生源が存在する確率を示す情報を含む、

付記 3 に記載の音声検出装置。

【 0 1 6 2 】

(付 記 5)

前記制御手段により自機を移動させている間に、前記撮像部が撮像した画像から認識された音声発生源の位置を示す情報を、前記音声発生源位置情報に追加する、

付記 3 または 4 に記載の音声検出装置。

40

【 0 1 6 3 】

(付 記 6)

前記判別手段は、前記音声検出手段が検出した音声が自機宛てに発せられた音声か否かを判別し、自機宛ての音声であると判別した場合、前記音声検出手段が検出した音声の音声発生源が前記特定の音声発生源であるか否かを判別する、

付記 1 から 5 のいずれか 1 つに記載の音声検出装置。

【 0 1 6 4 】

50

(付記 7)

ロボットに搭載されたコンピュータが音声を検出する音声検出方法であって、
 音声を検出する音声検出ステップと、
 前記音声検出ステップで検出された音声の音声発生源が特定の音声発生源であるか否かを判別する判別ステップと、
 前記判別ステップの判別結果に基づいて、前記ロボットの動作を制御する制御ステップと、
 を備える音声検出方法。

【0165】

(付記 8)

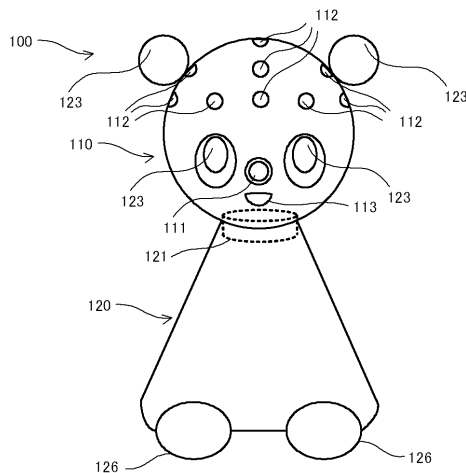
ロボットに搭載されたコンピュータに、
 音声を検出する音声検出機能と、
 前記音声検出機能により検出された音声の音声発生源が特定の音声発生源であるか否かを判別する判別機能と、
 前記判別機能により判別された判別結果に基づいて、前記ロボットの動作を制御する制御機能と、
 を実現させるプログラム。

【符号の説明】

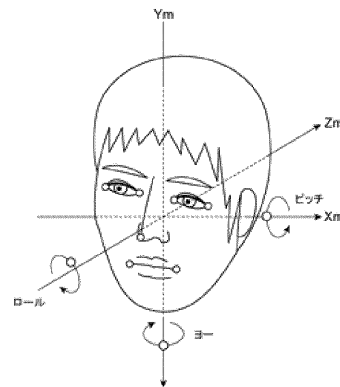
【0166】

100...ロボット、110...頭、111...カメラ、112...マイク、113...スピーカ、
 120...胴体、121...首関節、122...操作ボタン、123...センサ群、124...電源部、
 125...記憶部、126...足回り部、127...制御部、128...呼びかけ応答プログラム、
 129...地図作成プログラム、130...SLAM地図、131...フレームデータベース、
 132...音声発生確率データベース、133...呼びかけ移動プログラム、RM...部屋

【図 1】



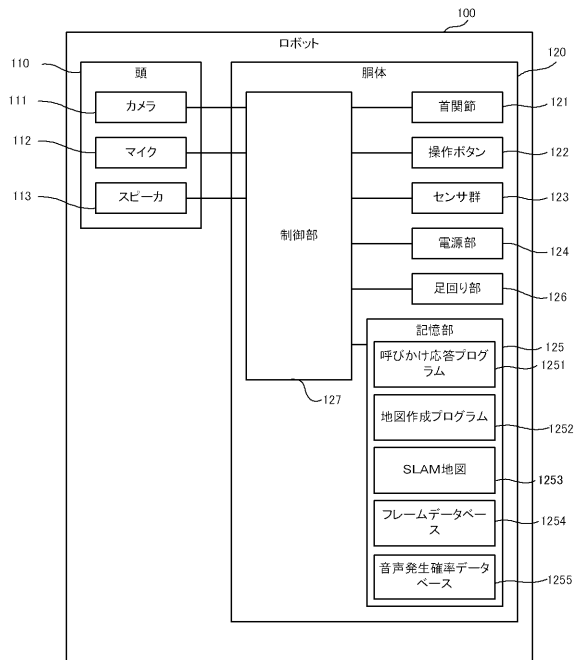
【図 2】



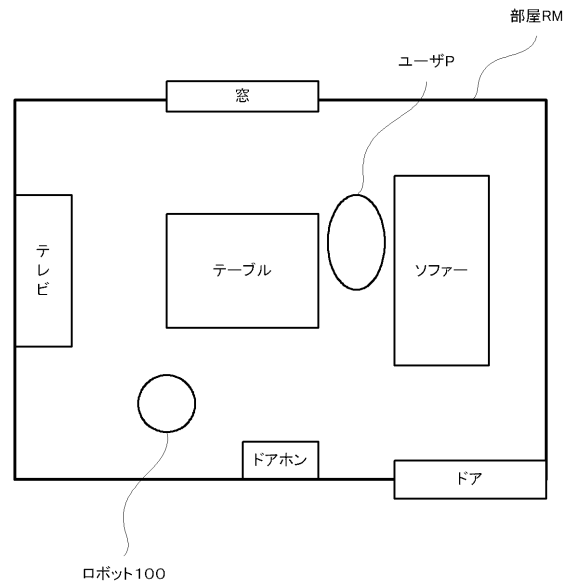
10

20

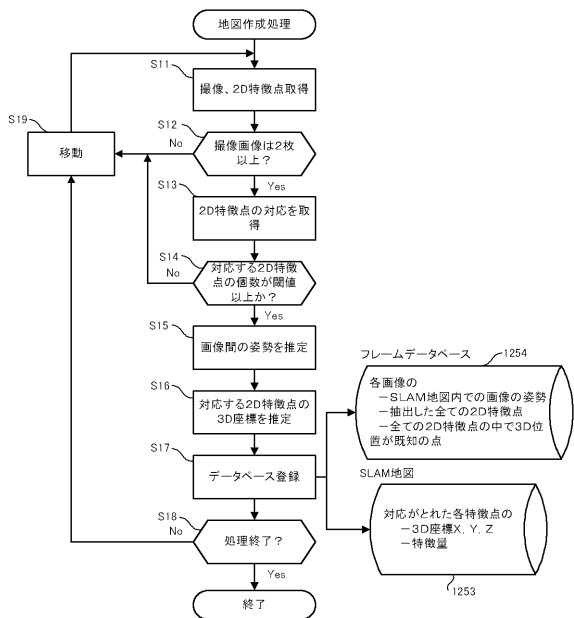
【図3】



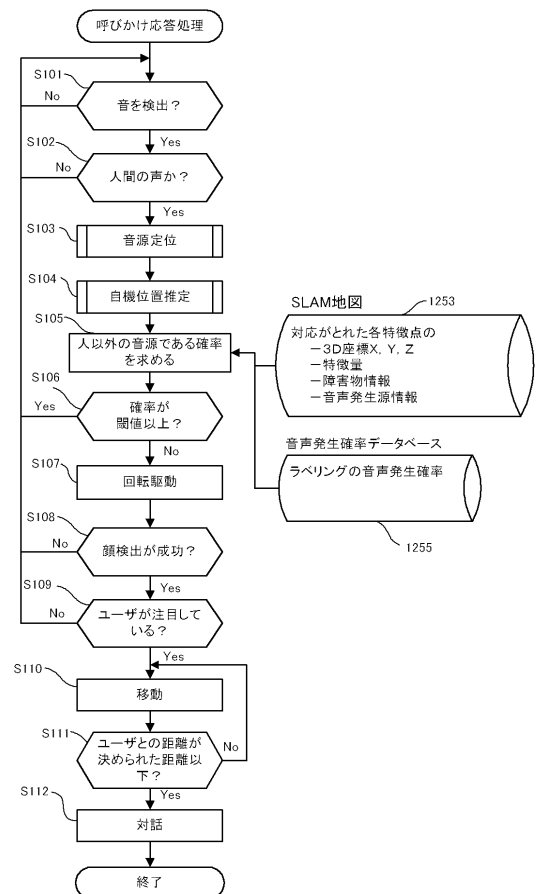
【図4】



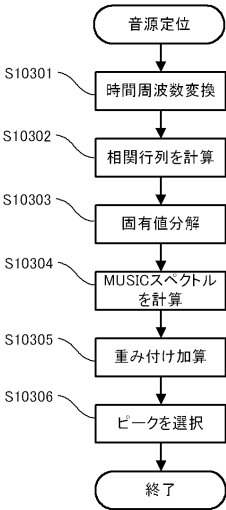
【図5】



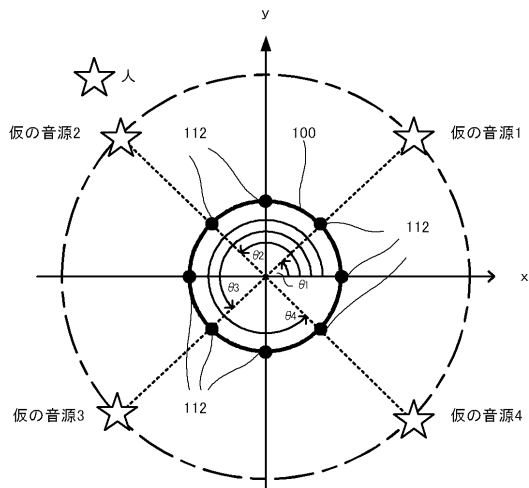
【図6】



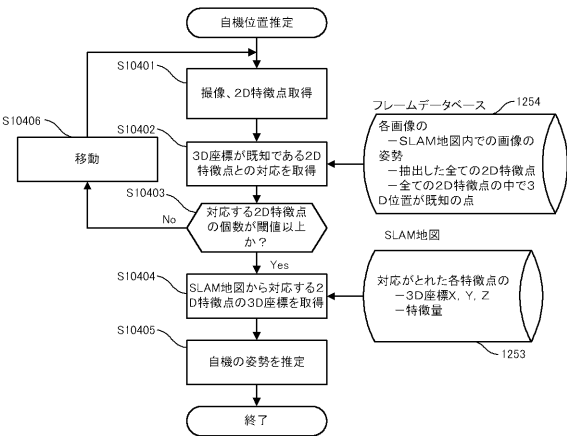
【図 7】



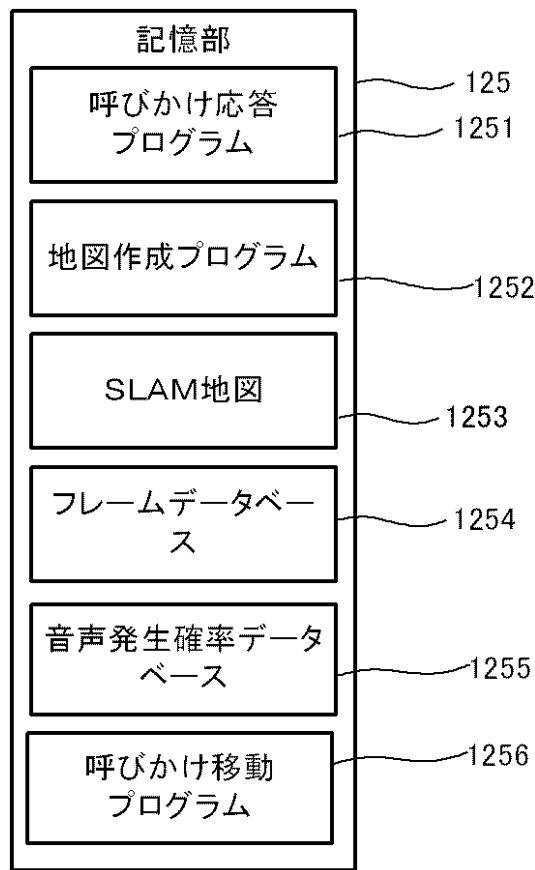
【図 8】



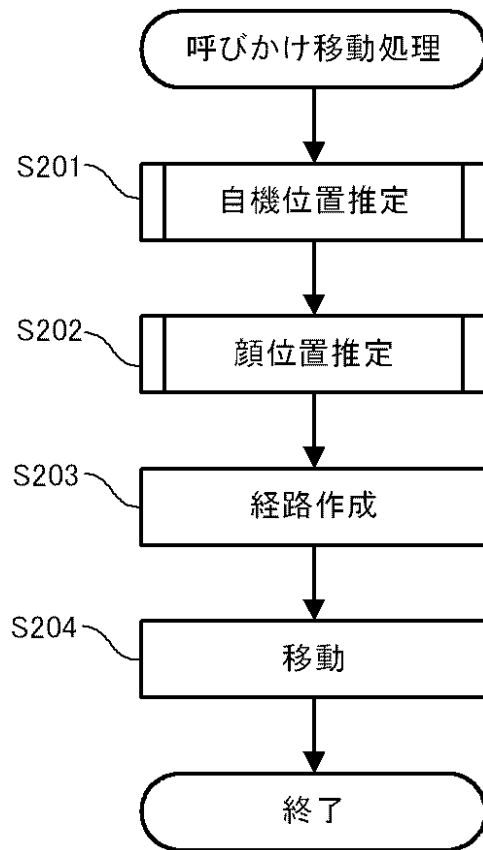
【図 9】



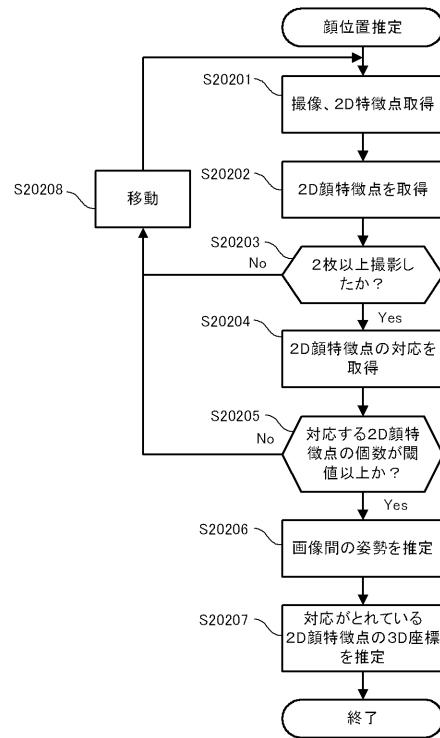
【図 10】



【図 1 1】



【図 1 2】



フロントページの続き

審査官 山下 剛史

(56)参考文献 特開 2 0 0 3 - 2 6 6 3 5 1 (J P , A)
特開 2 0 1 4 - 1 3 7 2 2 6 (J P , A)
特開 2 0 0 7 - 2 2 1 3 0 0 (J P , A)
特開 2 0 0 3 - 6 2 7 7 7 (J P , A)
国際公開第 2 0 1 4 / 1 6 7 7 0 0 (W O , A 1)
特開 2 0 0 5 - 2 5 0 2 3 3 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)
G 1 0 L 1 3 / 0 0 - 9 9 / 0 0
A 6 3 H 1 / 0 0 - 3 7 / 0 0