

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2020-38315
(P2020-38315A)

(43) 公開日 令和2年3月12日(2020.3.12)

(51) Int. Cl.		F I		テーマコード (参考)
G10L 21/0272 (2013.01)		G10L 21/0272	100B	5D220
G10L 21/028 (2013.01)		G10L 21/028	B	
H04R 3/00 (2006.01)		H04R 3/00	320	
G10L 25/51 (2013.01)		G10L 25/51	400	

審査請求 未請求 請求項の数 15 O L (全 16 頁)

(21) 出願番号	特願2018-165875 (P2018-165875)	(71) 出願人	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成30年9月5日(2018.9.5)	(74) 代理人	110001689 青稜特許業務法人
		(72) 発明者	池下 林太郎 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		(72) 発明者	藤岡 拓也 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		(72) 発明者	堀口 翔太 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		Fターム(参考)	5D220 BA06 BC05

(54) 【発明の名称】 音声情報処理装置および方法

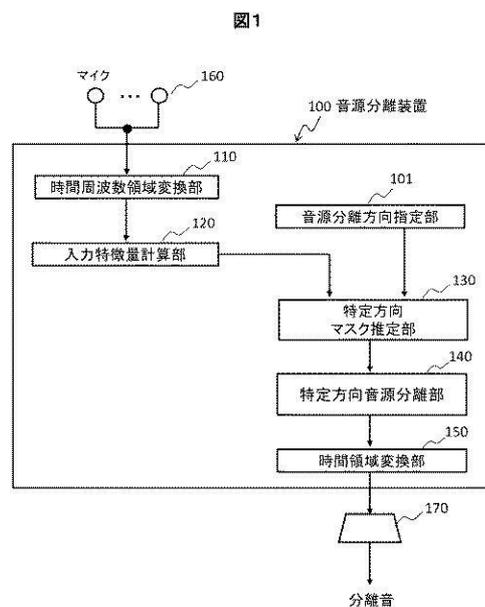
(57) 【要約】

【課題】時間周波数領域における音源のスパース性が成立しない場合においても、高性能に音源分離を行うことを目的とする。

【解決手段】

入力装置と、出力装置と、処理部を備える音声情報処理装置である。処理部は、入力装置から入力された音声信号を時間周波数領域に変換する、時間周波数領域変換部と、時間周波数領域から特徴量を計算する、特徴量計算部と、特徴量から特定方向の音を分離するための第1の時間周波数マスクを推定する、第1の時間周波数マスク推定部と、を備える。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

入力装置と、出力装置と、処理部を備え、
前記処理部は、
前記入力装置から入力された音声信号を時間周波数領域に変換する、時間周波数領域変換部と、
前記時間周波数領域から特徴量を計算する、特徴量計算部と、
前記特徴量から特定方向の音を分離するための第 1 の時間周波数マスクを推定する、第 1 の時間周波数マスク推定部と、
を備える音声情報処理装置。

10

【請求項 2】

前記第 1 の時間周波数マスク推定部に、1 または複数方向を指定する音源分離方向指定部を備える、
請求項 1 記載の音声情報処理装置。

【請求項 3】

時間周波数マスクの類似度、あるいは、時間周波数マスクで抽出された音声信号の類似度を判定する、時間周波数マスクマッチング部を備える、
請求項 1 記載の音声情報処理装置。

【請求項 4】

前記特徴量から特定の音源を分離するための第 2 の時間周波数マスクを推定する、第 2 の時間周波数マスク推定部を備える、
請求項 3 記載の音声情報処理装置。

20

【請求項 5】

前記時間周波数マスクマッチング部は、前記第 1 の時間周波数マスクと前記第 2 の時間周波数マスクの類似度を判定する、
請求項 4 記載の音声情報処理装置。

【請求項 6】

前記第 1 の時間周波数マスク推定部は、異なる方向に対する複数の第 1 の時間周波数マスクを推定し、
前記第 2 の時間周波数マスク推定部は、異なる音源に対する複数の第 2 の時間周波数マスクを推定し、
前記時間周波数マスクマッチング部は、前記複数の前記第 1 の時間周波数マスクと複数の前記第 2 の時間周波数マスクの類似度をそれぞれ判定し、指定した方向において、方向に対する第 1 の時間周波数マスクと、それと最も類似度の高い第 2 の時間周波数マスクとの類似度を音源存在スコアとして出力する、
請求項 5 記載の音声情報処理装置。

30

【請求項 7】

前記音源存在スコアが、予め定めた閾値以上であるとき、指定した方向に音源が存在すると判定し、音源方向推定結果を出力する音源方向推定部を備える、
請求項 6 記載の音声情報処理装置。

40

【請求項 8】

前記第 1 の時間周波数マスクを用いて音源分離フィルタを形成し、音源分離を実行する特定方向音源分離部を備え、
前記特定方向音源分離部は、前記音源方向推定部が、指定した方向に音源が存在すると判定した場合、方向に対する第 1 の時間周波数マスクを用いて音源分離フィルタを形成し、音源分離を実行する、
請求項 7 記載の音声情報処理装置。

【請求項 9】

前記特定方向音源分離部が出力する音源分離結果を、時間領域の分離信号に変換して出力する時間領域変換部を備える、

50

請求項 8 記載の音声情報処理装置。

【請求項 10】

前記第 1 の時間周波数マスク推定部は、異なる方向に対する複数の第 1 の時間周波数マスクを推定し、

前記第 2 の時間周波数マスク推定部は、異なる音源に対する複数の第 2 の時間周波数マスクを推定し、

前記第 2 の時間周波数マスク推定部が出力する複数の第 2 の時間周波数マスクを加算し、合成された第 2 の時間周波数マスクを出力する時間周波数マスク加算部を備え、

前記時間周波数マスクマッチング部は、前記複数の前記第 1 の時間周波数マスクと前記合成された第 2 の時間周波数マスクの類似度をそれぞれ判定し、前記類似度に基づいて所定の方向 に複数の音源が存在する可能性を示す複数音源存在スコアを出力する、

10

請求項 5 記載の音声情報処理装置。

【請求項 11】

移動機構と、

前記複数音源存在スコアが、予め定めた閾値以上であるとき、所定の方向 に音源が複数存在すると判定し、複数音源同一方向判定結果を出力する、複数音源同一方向判定部を備え、

前記移動機構は、前記複数音源同一方向判定結果が出力された場合に、各音源の方向が別れる位置に移動する、

請求項 10 記載の音声情報処理装置。

20

【請求項 12】

前記第 1 の時間周波数マスク推定部から出力される、異なる複数の時間における所定の方向 に対する第 1 の時間周波数マスクを記憶する特定方向マスク貯蔵部を備え、

前記時間周波数マスクマッチング部は、前記特定方向マスク貯蔵部に記憶された前記第 1 の時間周波数マスク同士の類似度を判定する、

請求項 3 記載の音声情報処理装置。

【請求項 13】

音声合成機能と、

時間周波数マスクマッチング部が出力する前記類似度が、予め定めた閾値以上であるとき、所定の方向 から同一の発話が 2 度あったと判定し、同一発話判定結果を出力する同一発話判定部を備え、

30

前記音声合成機能は、前記同一発話判定結果が出力された場合に、予め定められた内容の発話を行なう、

請求項 12 記載の音声情報処理装置。

【請求項 14】

入力装置と、出力装置と、処理部を備える情報処理装置を用いた音声情報処理方法であって、

前記処理部が、前記入力装置から入力された音声信号を時間周波数領域に変換する、時間周波数領域変換処理、

前記処理部が、前記時間周波数領域から特徴量を計算する、特徴量計算処理、

40

前記処理部が、前記特徴量から特定方向の音を分離するための第 1 の時間周波数マスクを推定する、第 1 の時間周波数マスク推定処理、

前記処理部が、前記特徴量から特定の音源を分離するための第 2 の時間周波数マスクを推定する、第 2 の時間周波数マスク推定処理、

前記処理部が、前記第 1 の時間周波数マスクと前記第 2 の時間周波数マスクのマッチングを行なう周波数マスクマッチング処理、

を行なう音声情報処理方法。

【請求項 15】

前記処理部が、複数の前記第 2 の時間周波数マスクの合成を行なって、合成時間周波数マスクを出力するマスク加算処理を行ない、

50

前記周波数マスクマッチング処理では、前記第 1 の時間周波数マスクと前記合成時間周波数マスクのマッチングを行なう、

請求項 1 4 記載の音声情報処理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声情報処理技術に関し、音源分離あるいは音源方向推定にかかる技術に関するものである。

【背景技術】

【0002】

音源分離技術とは、複数音源が混合した観測信号から、混合前の個々の源信号を推定する技術のことである。近年、音声信号の時間周波数スペクトルの特徴を事前に学習しておくことで音源分離を達成する、教師あり機械学習に基づく音源分離手法の研究が盛んに行われている。

【0003】

非特許文献 1 には、「マイクロホンで観測した観測信号を時間周波数領域の信号に変換し、観測信号の各時間周波数ピンをニューラルネットワークを用いて高次元空間のベクトルに写像し、写像されたベクトルをクラスタリングすることによって各時間周波数ピンを音源毎のクラスターに分類し、音源分離を達成する」技術が開示されている。

【先行技術文献】

【非特許文献】

【0004】

【非特許文献 1】 J. R. Hershey, Z. Chen, J. L. Roux, and S. Watanabe, "Deep Clustering: Discriminative embeddings for segmentation and separation," IEEE International Conference on Acoustics, Speech and Signal Processing, 2016, pp. 31-35.

【発明の概要】

【発明が解決しようとする課題】

【0005】

非特許文献 1 の音源分離方法は、時間周波数ピン (bin) をクラスタリングすることで、音源分離を達成できる。しかしながら、時間周波数ピンのクラスタリングに基づく音源分離方法では、空間に存在する音源の数が多いなどといった、時間周波数領域における音源のスパース性 (各時間周波数ピンにおいて原信号のうち 1 つだけが支配的であること) が成立しない場合において、音源分離性能が低いという問題があった。

【課題を解決するための手段】

【0006】

本発明の好ましい一側面は、入力装置と、出力装置と、処理部を備える音声情報処理装置である。処理部は、入力装置から入力された音声信号を時間周波数領域に変換する、時間周波数領域変換部と、時間周波数領域から特徴量を計算する、特徴量計算部と、特徴量から特定方向の音を分離するための第 1 の時間周波数マスクを推定する、第 1 の時間周波数マスク推定部と、を備える。

【0007】

好ましい具体例では、時間周波数マスクの類似度、あるいは、時間周波数マスクで抽出された音声信号の類似度を判定する、時間周波数マスクマッチング部を備える。

【0008】

別の好ましい具体例では、特徴量から特定の音源を分離するための第 2 の時間周波数マスクを推定する、第 2 の時間周波数マスク推定部を備える。

【0009】

さらに好ましい具体例では、時間周波数マスクマッチング部は、第 1 の時間周波数マスクと第 2 の時間周波数マスクの類似度を判定する。

【0010】

10

20

30

40

50

他のさらに好ましい具体例では、時間周波数マスクマッチング部は、異なる時間区間における第1の時間周波数マスクの類似度を判定する。

【0011】

本発明の他の好ましい一側面は、入力装置と、出力装置と、処理部を備える情報処理装置を用いた音声情報処理方法である。処理部が、入力装置から入力された音声信号を時間周波数領域に変換する、時間周波数領域変換処理、時間周波数領域から特徴量を計算する、特徴量計算処理、特徴量から特定方向の音を分離するための第1の時間周波数マスクを推定する、第1の時間周波数マスク推定処理、特徴量から特定の音源を分離するための第2の時間周波数マスクを推定する、第2の時間周波数マスク推定処理、第1の時間周波数マスクと第2の時間周波数マスクのマッチングを行なう周波数マスクマッチング処理、の各処理を行なう。

10

【0012】

他の具体的な例では、第1、第2の時間周波数マスク推定部は、教師有り学習で学習されたニューラルネットワークで構成される。

【発明の効果】

【0013】

本発明によれば、時間周波数領域におけるスパース性が成立しない音源に対しても、高い分離性能を有する音源分離方法を提供できる。

【図面の簡単な説明】

【0014】

20

【図1】本発明の第1実施形態に関わるブロック図。

【図2】本発明の第1実施形態に関わるフローチャート。

【図3】本発明の第2実施形態に関わるブロック図。

【図4】本発明の第2実施形態に関わるフローチャート。

【図5】本発明の実施形態に関わる時間周波数マスクの概念図。

【図6】本発明の第3実施形態に関わるブロック図。

【図7】本発明の第3実施形態に関わるフローチャート。

【図8】本発明の第3実施形態に関わる時間周波数マスクの概念図。

【図9】本発明の第4実施形態に関わるブロック図。

【図10】本発明の第4実施形態に関わるフローチャート。

30

【発明を実施するための形態】

【0015】

実施の形態について、図面を用いて詳細に説明する。ただし、本発明は以下に示す実施の形態の記載内容に限定して解釈されるものではない。本発明の思想ないし趣旨から逸脱しない範囲で、その具体的構成を変更し得ることは当業者であれば容易に理解される。

【0016】

以下に説明する発明の構成において、同一部分又は同様な機能を有する部分には同一の符号を異なる図面間で共通して用い、重複する説明は省略することがある。

【0017】

同一あるいは同様な機能を有する要素が複数ある場合には、同一の符号に異なる添字を付して説明する場合がある。ただし、複数の要素を区別する必要がない場合には、添字を省略して説明する場合がある。

40

【0018】

本明細書等における「第1」、「第2」、「第3」などの表記は、構成要素を識別するために付するものであり、必ずしも、数、順序、もしくはその内容を限定するものではない。また、構成要素の識別のための番号は文脈毎に用いられ、一つの文脈で用いた番号が、他の文脈で必ずしも同一の構成を示すとは限らない。また、ある番号で識別された構成要素が、他の番号で識別された構成要素の機能を兼ねることを妨げるものではない。

【0019】

図面等において示す各構成の位置、大きさ、形状、範囲などは、発明の理解を容易にす

50

るため、実際の位置、大きさ、形状、範囲などを表していない場合がある。このため、本発明は、必ずしも、図面等に開示された位置、大きさ、形状、範囲などに限定されない。

【0020】

本明細書で引用した刊行物、特許および特許出願は、そのまま本明細書の説明の一部を構成する。

【0021】

本明細書において単数形で表される構成要素は、特段文脈で明らかに示されない限り、複数形を含むものとする。

【0022】

以下で詳細に説明される実施例の概要を説明する。実施例に関わる代表的な音源分離方法ないし装置は、マイクロホンアレイで観測された複数音源の混合信号から、指定した方向に対する音源分離結果を与える時間周波数マスクを計算するニューラルネットワークを用いる。ここで、時間周波数マスクとは、分離対象とする信号が支配的である時間周波数ピンを推定し、その時間周波数成分をマスク処理により分離するマスクである（明細書・図面等において、「時間周波数マスク」を単に「マスク」ということがある）。

10

【0023】

より具体的な例では、指定した方向に対する音源分離結果を与える第1の時間周波数マスクと、空間の複数の音源（便宜上「全音源」と称することがあるが、必ずしも全ての音源を含む必要はない）に対する音源分離結果を与える第2の時間周波数マスクとを計算し、2つの時間周波数マスクを用いることにより、精度のよい音源分離を行なう。より具体的な例では、計算された2つの時間周波数マスクをマッチングすることによって、所定の方向に音源が存在するか否かを判定する。

20

【0024】

例えば第1の時間周波数マスクは、（マイクから見た）方向X、方向Y、方向Zを識別するマスクであり、これは既知の方向X、方向Y、方向Zからの音声を教師データとして、ニューラルネットワークを学習することで生成することができる。また、第2の時間周波数マスクは、音源A、音源B、音源Cを識別するマスクであり、これは既知の音源A、音源B、音源Cからの音声を教師データとして、ニューラルネットワークを学習することで生成することができる。

【実施例1】

30

【0025】

図1と図2を用いて、本発明の第一実施形態に関わる音源分離装置100を説明する。本実施例の音源分離装置100は、入力装置、出力装置、記憶装置、および処理装置を含む、通常のコンピュータ（例えばサーバ）で構成した。本実施例では計算や制御等の機能は、記憶装置に格納されたプログラムが処理装置によって実行されることで、定められた処理を他のハードウェアと協働して実現される。計算機などが実行するプログラム、その機能、あるいはその機能を実現する手段を、「～部」等と呼ぶ場合がある。また記憶装置と処理装置を総称して「処理部」と呼ぶ場合がある。実施例2以降で説明される音源分離装置も同様である。

【0026】

40

また、音源分離装置100の構成は、単体のコンピュータで構成してもよいし、あるいは、入力装置、出力装置、処理装置、記憶装置の任意の部分が、ネットワークで接続された他のコンピュータで構成されてもよい。また、本実施例中、プログラムで構成した機能と同等の機能は、FPGA（Field Programmable Gate Array）、ASIC（Application Specific Integrated Circuit）などのハードウェアでも実現できる。

【0027】

図1は第1実施形態に関わる音源分離装置100のブロック図である。音源分離装置100は、音源分離方向指定部101と、時間周波数領域変換部110と、入力特徴量計算部120と、特定方向時間周波数マスク推定部130と、特定方向音源分離部140と、時間領域変換部150から構成される。これらは、既述のように、処理装置によって実行

50

されるプログラムとして実装されている。また、入力装置として例えばマイクロホンアレイ160を、分離音の出力装置として例えばスピーカ170を備える。

【0028】

図2も参照しつつ、音源分離装置100の動作例を説明する。時間周波数領域変換部110は、公知の短時間フーリエ変換などにより、マイクロホンアレイ160を用いて観測した混合信号（例えば横軸を時間、縦軸を振幅（音圧）で表現される）の時間周波数表現を計算して出力する（ステップS201）。公知のように、時間周波数表現としては、例えば横軸を時間、縦軸を周波数とし、周波数の強度をコントラストで示すものがある。本実施例では、時間周波数表現がマイクロホンアレイ160のマイクの数得られることになる。

10

【0029】

入力特徴量計算部120は、時間周波数領域変換部110が出力した観測信号の時間周波数表現から音源分離装置が用いる特徴量を計算し、特定方向時間周波数マスク推定部130に出力する（ステップS202）。特徴量としては、例えば、観測信号の振幅スペクトルとマイクロホン間の位相差を束ねたものを用いることができる。所定の音声について、入力特徴量計算部120からはマイクロホンアレイ160のマイクの数だけ、観測信号から得た振幅スペクトルと位相差の情報の組からなる情報が出力され、これらを纏めて当該音声の特徴量とする。一般には、マイクロホンアレイ160からのアナログ信号をデジタル信号に変換して後続の処理を行なう。この場合、周波数や時間等の量子化の粒度は任意である。

20

【0030】

位相差の情報については、マイクロホンアレイの各マイクの位置が異なることにより、各マイクに音波が到達する時刻が異なるので、得られる観測信号に位相差が生じる。位相差の情報から方向情報を得ることができる。位相差の情報は音源の方向を特定するために重要である。

【0031】

音源分離方向指定部101は、音声端末デバイスから見た、音源分離したい方向を指定する。具体的にはマイクロホンアレイ160から見た方向である（ステップS203）。特定方向時間周波数マスク推定部130は、入力特徴量計算部120が出力する観測信号の特徴量を用いて、音源分離方向指定部101が指定するデバイスから見た音源分離したい方向 に対する時間周波数マスクを推定する（ステップS204）。

30

【0032】

ステップS204における指定した方向 に対する時間周波数マスクの推定には、例えばニューラルネットワークを用いることができる。ニューラルネットワークは、観測信号の特徴量を入力したとき、指定した方向 に対する音源を抽出する時間周波数マスクを出力する。ここで、指定した方向 に対する音源を抽出するとは、物理的な音源そのものを特定するのではなく、指定した方向 から到来する音を特定することを意味する。具体例としては、特定方向時間周波数マスク推定部130は学習済みのディープニューラルネットワーク（DNN）で構成することができる。DNNは、方向 が既知の音源からの観測信号を教師データとして用いて、方向からの音を選択的に抽出する時間周波数マスクを出力とするように教師有り学習をしておく。

40

【0033】

特定方向音源分離部140は、特定方向時間周波数マスク推定部130が出力する、指定した方向 に対する時間周波数マスクを用いて、指定した方向 に対する音源分離フィルタを形成し、音源分離を実行する（ステップS205）。音源分離フィルタは、例えば各時間周波数ビンに対して1または0の重みを付けるものでもよいし、1~0の間で離散的な値を取るものでも良い。

【0034】

なお、特定方向時間周波数マスク推定部130が出力する時間周波数マスクをそのまま用いるのではなく、目的に合ったフィルタとなるように、信号処理技術を用いて再形成

50

(修正, 調整)し, 再形成したフィルタを用いて, 特定方向からの音を抜き出す処理を行ってもよい。目的にかなったフィルタとは, 例えば音声を歪ませないようにマスキングを調整したり, 不要な周波数成分を除去するなど, ユーザが任意に設定することができる。

【0035】

時間領域変換部150は, 特定方向音源分離部140が出力する, 指定した方向 に対する音源分離結果を, 時間領域の分離信号に変換して出力する。出力は, 例えばスピーカ170から音声信号として出力される。

【0036】

上記の特定方向に対する音源分離処理は, 音源分離方向指定部101が指定する方向として, 様々な方向を指定することで, 任意方向に対して実行できる。例えば, 音源分離方向指定部101が指定する方向 として, 予め数式1

【0037】

【数1】

$$\theta \in \left\{ \frac{2k\pi}{36} \mid k = 0, 1, \dots, 35 \right\}$$

【0038】

のように全方位を10度刻みに指定しておくことができる。このために, 本実施例では, DNNに各方向 に対して専用の層を準備しておき, ステップS203で選択するように構成される。数式1の例に従うと, kごとに36個の層を準備する。

【0039】

他の実装方法としては, 全方向の に対して共通の層を準備しておき, 補助ネットワークで層の重みパラメータを変化させることにより, 指定した方向 に対応した層に変換することも可能である。

【0040】

上記の実施例では, 特定方向時間周波数マスク推定部130に入力する特徴量として, 観測信号の振幅スペクトルとマイクロホン間の位相差の両方を用いた。他の手法としては, 観測信号の振幅スペクトルを用いず, マイクロホン間の位相差のみを特徴量として用いても良い。このようにすると, 音の到来方向の情報だけを抽出できるので, 音源の種類による影響を小さくすることができる。

【0041】

時間周波数スペクトルを特徴量に含める場合には, 音源の種類が支配的にならないように, 様々な音源の種類と, 様々な方向の組み合わせを, 十分多くの数用意した学習データを用いて, 教師有り学習を行なうことが望ましい。この方法により, 音源の種類により方向の推定結果が影響されることを緩和することができる。

【0042】

一方, 時間周波数スペクトルを用いる場合の利点として以下が考えられる。例えば, マイクロホンアレイの性質によっては, 「右方向の音源からの音がマイク1に小さく入り, マイク2に大きく入る」といったことが起こり得る。位相差の特徴量のみでは, このような音源方向によって各マイクに届く音の大きさが異なるという情報を利用することができない。

【0043】

実施例1の構成によれば, スパース性が不成立の場合でも, 指定した方向 に対しては高い分離性能が得られる。

【実施例2】

【0044】

図3~図5を用いて, 本発明の第二実施形態に関わる音源分離装置300を説明する。

10

20

30

40

50

第二実施形態は、例えば会議室等において、複数の発言者を識別するとともにその方向（あるいは位置）を特定するのに好適である。また、本実施例は２種類の時間周波数マスクをマッチングすることにより、高精度の音源分離、音源方向推定が可能となる。

【 0 0 4 5 】

第二実施例の音源分離装置 3 0 0 は、図 3 における全音源時間周波数マスク推定部 3 1 0 と、時間周波数マスクマッチング部 3 2 0 と、音源方向推定部 3 3 0 が加わることを除けば、図 1 に示した第一実施形態の音源分離装置 1 0 0 と同じ構成であるので、以下では、全音源時間周波数マスク推定部 3 1 0 と、時間周波数マスクマッチング部 3 2 0 と、音源方向推定部 3 3 0 についてのみ説明し、他の説明を省略する。

【 0 0 4 6 】

図 4 に示した第二実施形態の処理フローも、ステップ S 4 0 1 と、ステップ S 4 0 2 と、ステップ S 4 0 3 が加わることを除けば、図 2 に示した第一実施形態の処理フローと同一であるため、以下ではステップ S 4 0 1 と、ステップ S 4 0 2 と、ステップ S 4 0 3 についてのみ説明し、他の説明を省略する。

【 0 0 4 7 】

全音源時間周波数マスク推定部 3 1 0 は、入力特徴量計算部 1 2 0 が出力する観測信号の特徴量を用いて、空間に存在する全音源に対して、各音源の分離結果に対応する、時間周波数マスクを推定する（ステップ S 4 0 1）。ステップ S 4 0 1 における全音源に対する時間周波数マスクの推定には、例えば、非特許文献 1 で提案された深層クラスタリングに基づく時間周波数マスクの推定方法を用いることができる。全音源時間周波数マスク推定部 3 1 0 は、たとえば DNN で構成することができる。DNN は、既知の音源からの観測信号を教師データとして、公知の教師有り学習で学習させることができる。

【 0 0 4 8 】

本実施例では、全音源時間周波数マスク推定部 3 1 0 に入力する特徴量は、特定方向時間周波数マスク推定部 1 3 0 に入力する特徴量と同じ特徴量とし、観測信号の振幅スペクトルとマイクロホン間の位相差の両方を用いた特徴量とした。ただし、別々の特徴量を用いても良い。また、全音源時間周波数マスク推定部 3 1 0 へ入力する特徴量は、マイクロホン間の位相差を省略してもよい。

【 0 0 4 9 】

時間周波数マスクマッチング部 3 2 0 は、全音源時間周波数マスク推定部 3 1 0 が出力する全音源に対する時間周波数マスクの推定値と、特定方向時間周波数マスク推定部 1 3 0 が出力する指定した方向 に対する時間周波数マスクの推定値を入力として受け取り、各音源 n の時間周波数マスクの推定値と指定した方向 に対する時間周波数マスクの推定値の類似度を計算する（ステップ S 4 0 2）。時間周波数マスクの類似度 d_1 としては、例えば数式 2

【 0 0 5 0 】

【 数 2 】

$$d_1(y_n, y_\theta) = -\|y_n - y_\theta\|_F^2$$

【 0 0 5 1 】

を用いればよい。ここで、 y_n は全音源時間周波数マスク推定部 3 1 0 が出力する音源 n の時間周波数マスクの推定値、 y_θ は特定方向時間周波数マスク推定部 1 3 0 が出力する方向 に関する時間周波数マスクの推定値、右辺の「 $\|\cdot\|_F$ 」は行列のフロベニウスノルムである。

【 0 0 5 2 】

ここで類似度 d_1 は、特定の音源 n に対する時間周波数マスクと特定の方向 に対する

10

20

30

50

時間周波数マスクの類似度を示しており、これが類似（あるいは一致）した場合には、特定の音源が特定の方向に存在する可能性が高い。前述のように特定の音源に対する時間周波数マスクは、スパース性が不成立の場合には精度が低くなる可能性はあるが、おおよその音源分離は可能である。一方、特定の方向に対する時間周波数マスクは、定めた方向に対する分離性能は高いが、その方向に音源があるかは不確かである。すなわち、その方向にはノイズ源や壁などの反射物がある可能性もある。よって、両方の時間周波数マスクを利用することで、音源を特定する精度を向上させることが期待できる。

【 0 0 5 3 】

図 5 に、特定の音源 n に対する時間周波数マスクと特定の方向 θ に対する時間周波数マスクの類似度を概念的に示す。単純化するため、マスクは各ピンに対して 0（白）と 1（黒）の 2 値を取るものとしている。図 5 には、音源 1 ~ 音源 4 に対する時間周波数マスクと、方向 1 ~ 方向 4 に対する時間周波数マスクが概念的に示されている。これらの類似度を計算することにより、例えば音源 2 のマスクと方向 2 のマスクがマッチすることが分かる。この場合、音源 2 が方向 2 にあり、方向 2 からの音は音源 2 の音である可能性が高い。

10

【 0 0 5 4 】

上記のような理論に基づいて、さらに、時間周波数マスクマッチング部 3 2 0 は、指定した方向 θ に音源が存在するか否かを表す、音源存在スコアを計算し出力する（ステップ S 4 0 2）。音源存在スコアとしては、例えば、数式 3。

【 0 0 5 5 】

【 数 3 】

20

$$\ell_1(\theta) = \max_n d_1(y_n, y_\theta)$$

を用いれば良い

【 0 0 5 6 】

数式 3 の音源存在スコア ℓ_1 は、ある方向 θ において、 θ 方向に対する時間周波数マスクと、それと最も類似度の高い音源に対する時間周波数マスクとの類似度を示している。音源方向推定部 3 3 0 は、時間周波数マスクマッチング部 3 2 0 が出力する、指定した方向 θ に対する音源存在スコア ℓ_1 が、予め定めた閾値以上であるとき、指定した方向 θ に音源が存在すると判定し、音源方向推定結果を出力する（ステップ S 4 0 3）。

30

【 0 0 5 7 】

ステップ S 4 0 3 で方向 θ に音源が存在すると判定した場合には、時間領域変換部 1 5 0 は、特定方向音源分離部 1 4 0 が出力する方向 θ に対する音源分離結果を、時間領域の分離信号に変換して出力する。出力は、例えばスピーカ 1 7 0 から音声信号（分離音）として出力される。このとき、図 3 のようにあわせて音源方向推定結果を表示することにより、音源の特定がさらに容易となる。例えば、会議室において、音源方向に該当する座席を表示することにより、音源である発言者を特定することができる。

40

【 0 0 5 8 】

会議室のように、音源が複数想定される場合には、ステップ S 2 0 3 で θ を 0 ~ 3 6 0 度順次指定し、全方位に渡って同様の処理を行えばよい。

【 0 0 5 9 】

本実施例では、分離音は、特定方向時間周波数マスク推定部 1 3 0 で生成した θ 方向に対する時間周波数マスクを用いて得ているが、全音源時間周波数マスク推定部 3 1 0 で生成した、音源に対する時間周波数マスクで得たものを用いても良い。あるいは、両者を選択可能としても良い。音源に対する時間周波数マスクとしては、例えば θ 方向に対する時間周波数マスクと最も類似度の高い音源に対する時間周波数マスクを用いる。ただし、音

50

源が多い場合等、スパース性が成立しない場合は、一般に、方向に対する時間周波数マスクを用いたほうが精度が良い。

【0060】

本実施例では、マスクマッチング部では2つのマスク同士を比較している。他の手法として、マスクを用いて抽出した音源信号同士を比較して類似性を評価しても同様の効果が得られる。ただし、この場合には両方のマスクにより音源分離を行なう必要がある。

【実施例3】

【0061】

図6～図8を用いて、本発明の第三実施形態に関わる音源分離装置500を説明する。第三実施形態は、複数の音源が同一方向に存在することを判定可能な例である。本実施例は、例えばロボット等に応用することもできる。

10

【0062】

第三実施例の音源分離装置500は、図6における時間周波数マスク加算部510と、時間周波数マスクマッチング部520と、複数音源同一方向判定部530が加わることを除けば、図3に示した第二実施形態の音源分離装置300と同じ構成であるので、以下では、時間周波数マスク加算部510と、時間周波数マスクマッチング部520と、複数音源同一方向判定部530についてのみ説明し、他の説明を省略する。

【0063】

図7に示した第三実施形態の処理フローも、ステップS601からステップS604が加わることを除けば、図4に示した第二実施形態の処理フローと同一であるため、以下ではステップS601からステップS604についてのみ説明し、他の説明を省略する。

20

【0064】

時間周波数マスク加算部510は、全音源時間周波数マスク推定部310が出力する全音源に対する時間周波数マスクを入力として受け取り、複数の音源の時間周波数マスクを加算し、時間周波数マスクマッチング部520に出力する(ステップS601)。例えば、音源数がNで各音源の時間周波数マスクを y_1, \dots, y_N とするとき、時間周波数マスク加算部510は、数式4

【0065】

【数4】

$$E := \{I \subseteq \{1, \dots, N\} \mid |I| \geq 2\}$$

【0066】

と数式5

【0067】

【数5】

$$\left\{ \sum_{n \in I} y_n \mid I \in E \right\}$$

で表される合計 $2^N - 1$ (N + 1)個の加算時間周波数マスクを計算する。

【0068】

ここで時間周波数マスクの加算とは、例えばマスクAとマスクBを加算したマスクA + Bがある場合、マスクA + Bを使って抽出した音源が、マスクAを使って抽出した音源と、マスクBを使って抽出した音源との加算になるようなマスクを合成する処理をいう。

50

【 0 0 6 9 】

時間周波数マスクマッチング部 5 2 0 は、時間周波数マスク加算部 5 1 0 が出力する加算時間周波数マスクと、特定方向時間周波数マスク推定部 1 3 0 が出力する指定した方向に対する時間周波数マスクを入力として受け取り、上記 2 つのマスクの類似度を計算する（ステップ S 6 0 2）。時間周波数マスクの類似度 d_2 としては、例えば数式 6

【 0 0 7 0 】

【数 6】

$$d_2(\sum_{n \in I} y_n, y_\theta) = -\|\sum_{n \in I} y_n - y_\theta\|_F^2$$

【 0 0 7 1 】

を用いれば良い。ここで、 y_n は全音源時間周波数マスク推定部 3 1 0 が出力する音源 n の時間周波数マスクの推定値、 I は数式 4 の要素、 y_θ は特定方向時間周波数マスク推定部 1 3 0 が出力する方向に対する時間周波数マスクの推定値である。類似度 d_2 が大きいということは、所定の方向に複数の音源が存在する可能性を示している。さらに、時間周波数マスクマッチング部 5 2 0 は、指定した方向に複数の音源が存在するか否かを表す、複数音源存在スコアを計算し出力する（ステップ S 6 0 2）。複数音源存在スコア l_2 としては、例えば、数式 7

20

【 0 0 7 2 】

【数 7】

$$l_2(\theta) = \max_{I \subseteq E} d_2(\sum_{n \in I} y_n, y_\theta)$$

を用いれば良い。ここで、 E は数式 4 で定義される。

30

【 0 0 7 3 】

図 8 で、所定の方向に複数の音源が存在する場合の判定の原理を概念的に説明する。図 8 には、音源 1 と音源 2 に対する時間周波数マスクから抽出された音源と、各方向からの音源が概念的に示されている。音源 1 と音源 2 に対する 2 つのマスクを加算することにより、音源 1 + 音源 2 に対する時間周波数マスクが生成される。加算した時間周波数マスクによる音源と所定方向の時間周波数マスクによる音源の類似度を計算することにより、例えば音源 1 + 音源 2 と方向 2 の音源がマッチすることが分かる。この場合、音源 1 と音源 2 が方向 2 にあり、方向 2 からの音は音源 1 および音源 2 の音である可能性が高い。以上の説明では、理解を容易にするため抽出した音源のマッチングを行なっているが、既に述べたように、音源を抽出するためのマスクのマッチングを行なっても同様の効果がある。

40

【 0 0 7 4 】

上記の理論を基にして、複数音源同一方向判定部 5 3 0 は、時間周波数マスクマッチング部 5 2 0 が出力する、指定した方向に対する複数音源存在スコア l_2 が、予め定めた閾値以上であるとき、指定した方向に音源が複数存在すると判定し、判定結果を出力する（ステップ S 6 0 3）。

【 0 0 7 5 】

マイクロホンアレイ 1 6 0 を備えた音声アシスタント端末、例えばモーターと操舵機構を備える移動機構を備えるロボットが、音声認識を行なってユーザに応答する用途を想定すると、一方向に複数の音源（複数のユーザ）が存在し、同時に発話する場合、音声認識が困難になる場合が考えられる。そこで、指定した方向に複数の音源が存在すると判定

50

されたとき，マイクロホンアレイを備えた音声アシスタント端末は，同一方向にいる複数の音源が別々の方向になるように移動機構 540 を制御し，移動する（ステップ S604）。

【0076】

そのためには，ステップ S603、S203、S204、S602 のループ処理を繰り返し，音源が複数存在すると判定された方向の時間周波数マスクに対する複数音源存在スコア l_2 が，閾値を下回るまで移動を行なえばよい。

【実施例 4】

【0077】

図 9 と図 10 を用いて，本発明の第四実施形態に関わる音源分離装置 700 を説明する。第四実施形態は，同一方向に存在する音源が複数回発話することを判定可能な例である。本実施例は，例えばロボット等に応用することもできる。

10

【0078】

第四実施例の音源分離装置 700 は，図 9 における特定方向時間周波数マスク貯蔵部 710 と，時間周波数マスクマッチング部 720 と，同一発話判定部 730 が加わることを除けば，図 1 に示した第一実施形態の音源分離装置 100 と同じ構成であるので，以下では，特定方向時間周波数マスク貯蔵部 710 と，時間周波数マスクマッチング部 720 と，同一発話判定部 730 についてのみ説明し，他の説明を省略する。

【0079】

図 10 に示した第四実施形態の処理フローも，ステップ S801 からステップ S803 が加わることを除けば，図 2 に示した第一実施形態の処理フローと同一であるため，以下ではステップ S801 からステップ S803 についてのみ説明し，他の説明を省略する。

20

【0080】

特定方向時間周波数マスク貯蔵部 710 は，特定方向時間周波数マスク推定部 130 が毎時刻出力する指定した方向に関する時間周波数マスクを保存する（ステップ S801）。毎時刻出力する時間周波数マスクは，それぞれ異なる所定時間区間の観測信号を反映している。

【0081】

時間周波数マスクマッチング部 720 は，特定方向時間周波数マスク貯蔵部 710 に保存された，指定した方向に関する 2 つの時刻の時間周波数マスクを入力として受け取り，類似度を計算する（ステップ S802）。時間周波数マスクの類似度 d_3 としては，例えば数式 8

30

【0082】

【数 8】

$$d_3(y(t_1, \theta), y(t_2, \theta)) = - \|y(t_1, \theta) - y(t_2, \theta)\|_F^2$$

【0083】

を用いれば良い。ここで， $y(t_1, \theta)$ と $y(t_2, \theta)$ は，異なる 2 つの時刻 t_1 と t_2 における指定した方向に関する時間周波数マスクの推定値である。他にも，時間周波数マスクの類似度として，時間周波数マスク $y(t_2, \theta)$ を時間方向に伸縮した結果である $y'(t_2, \theta)$ と， $y(t_1, \theta)$ の間の類似度を計算することにも良い。

【0084】

同一発話判定部 730 は，時間周波数マスクマッチング部 720 が出力する時間周波数マスクの類似度 d_3 が，予め定めた閾値以上であるとき，指定した方向から同一の発話が 2 度あったと判定し，判定結果を出力する（ステップ S802）。

50

【0085】

マイクロホンアレイ160を備えた音声アシスタント端末、例えばロボットが、音声認識を行なってユーザに応答する用途を想定すると、同一のユーザが同じ発話を複数回繰り返す場合、ロボットが音声認識できず応答ができていない場合が考えられる。そこで、指定した方向から同一の発話が2度あったと判定されたとき、マイクロホンアレイを備えた音声アシスタント端末は、方向に向きを変え、音声応答を行う(ステップS803)。このとき、音声アシスタント端末は必ずしもユーザの発話の音声認識ができていなくても、同一方向から同一の発話が2度あったことをトリガーとして、予め定められた応答を行なうことができる。例えば、応答のための音声合成部740を備え、同一発話判定部730の判定結果により、所定の発話を行なうように音声合成部740を制御する。この構成により、ユーザの音声認識ができていなくても、音声合成部740が「よく聞こえませんでした。もう一度話してください。」のように発話することにより、円滑な対応が可能となる。

10

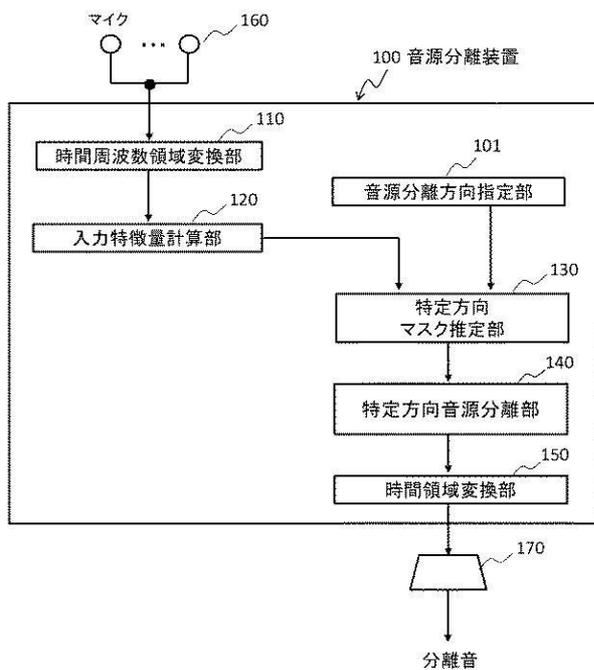
【符号の説明】

【0086】

音源分離装置100，音源分離方向指定部101，時間周波数領域変換部110，入力特徴量計算部120，特定方向時間周波数マスク推定部130，特定方向音源分離部140，時間領域変換部150

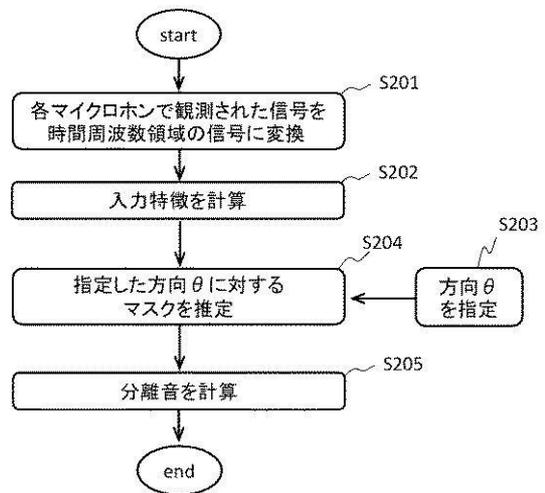
【図1】

図1

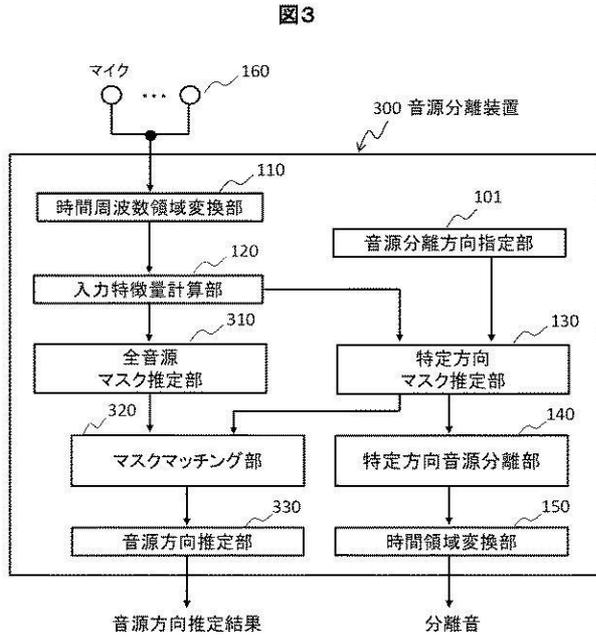


【図2】

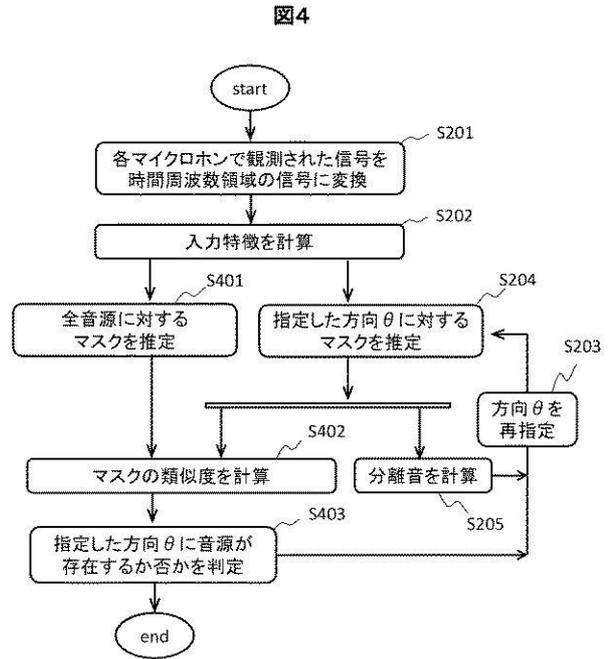
図2



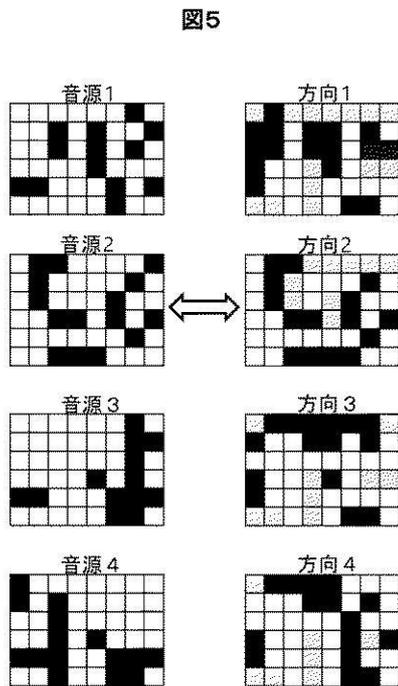
【 図 3 】



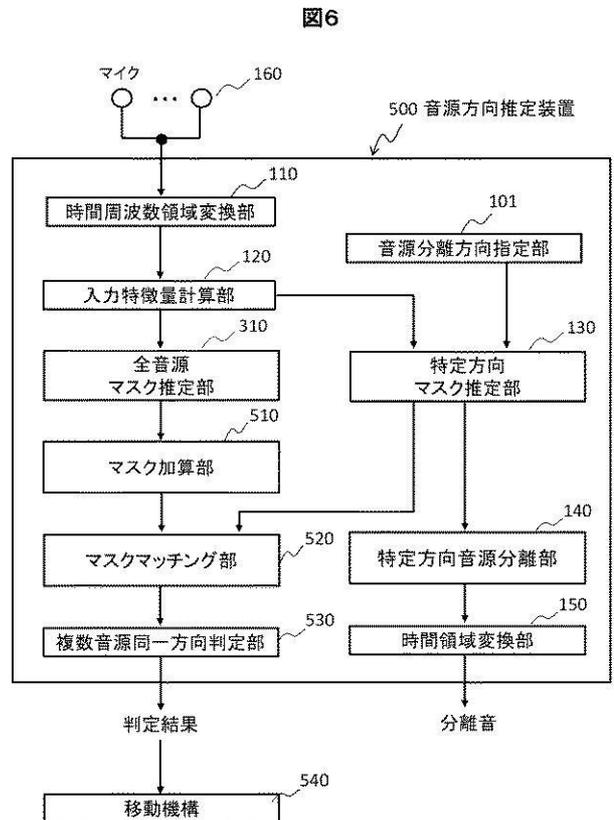
【 図 4 】



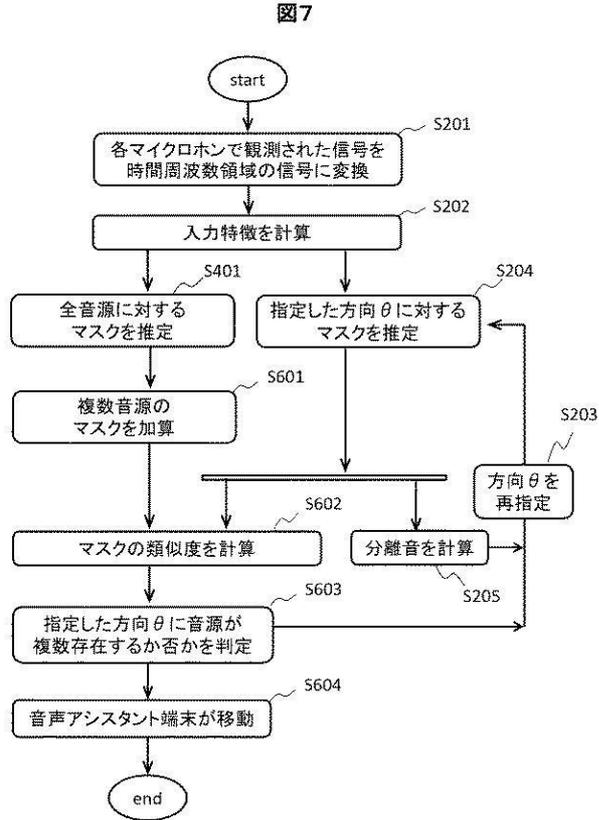
【 図 5 】



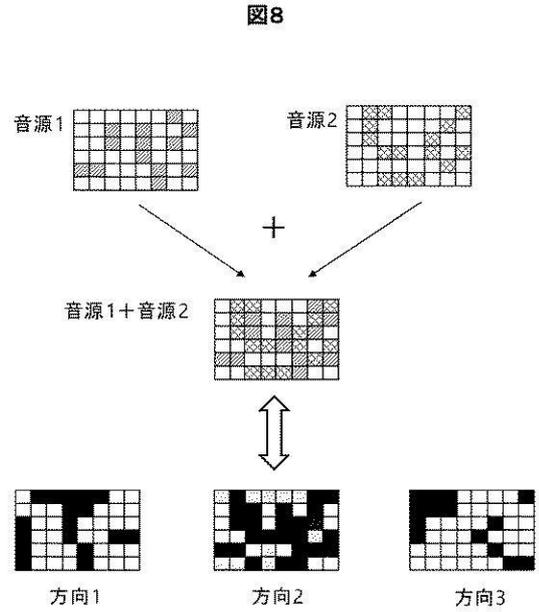
【 図 6 】



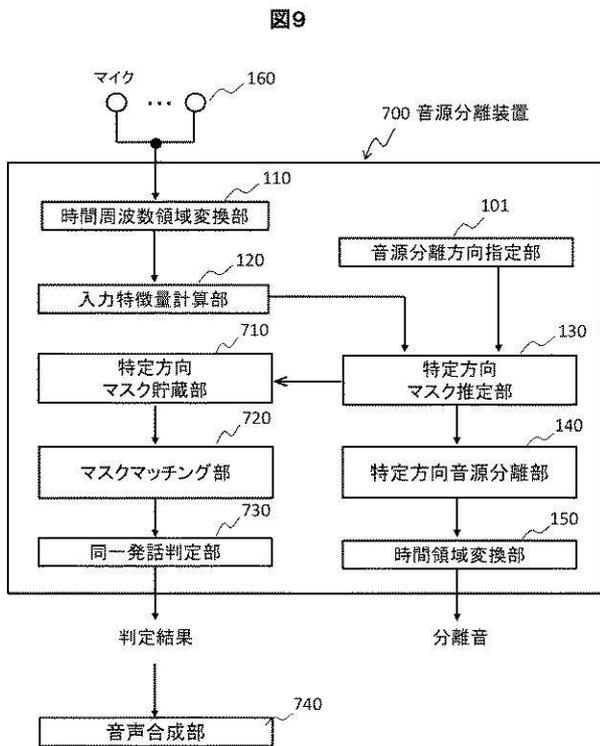
【 図 7 】



【 図 8 】



【 図 9 】



【 図 10 】

