

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6059112号
(P6059112)

(45) 発行日 平成29年1月11日(2017.1.11)

(24) 登録日 平成28年12月16日(2016.12.16)

(51) Int.Cl. F I
G 1 O L 21/0308 (2013.01) G 1 O L 21/0308 Z
G 1 O L 21/028 (2013.01) G 1 O L 21/028 B

請求項の数 7 (全 17 頁)

(21) 出願番号	特願2013-171079 (P2013-171079)	(73) 特許権者	000004226
(22) 出願日	平成25年8月21日 (2013.8.21)		日本電信電話株式会社
(65) 公開番号	特開2015-40934 (P2015-40934A)		東京都千代田区大手町一丁目5番1号
(43) 公開日	平成27年3月2日 (2015.3.2)	(74) 代理人	100121706
審査請求日	平成27年6月29日 (2015.6.29)		弁理士 中尾 直樹
		(74) 代理人	100128705
			弁理士 中村 幸雄
		(74) 代理人	100147773
			弁理士 義村 宗洋
		(72) 発明者	木下 慶介
			東京都千代田区大手町二丁目3番1号 日
			本電信電話株式会社内
		(72) 発明者	中谷 智広
			東京都千代田区大手町二丁目3番1号 日
			本電信電話株式会社内

最終頁に続く

(54) 【発明の名称】 音源分離装置とその方法とプログラム

(57) 【特許請求の範囲】

【請求項1】

複数の音源から発せられる音源信号を複数のマイクロホンで收音した複数チャンネルの観測信号と、上記複数のマイクロホンの各々で観測される上記複数の音源の各々からの信号の音圧が異なると仮定した観測信号のモデルを用いて、各マイクロホンごとに各音源に関する音源存在事後確率を推定するマイク別音源存在事後確率推定部と、

上記複数チャンネルの観測信号と、上記音源存在事後確率を入力として、観測信号のモデルパラメータを推定するモデルパラメータ推定部と、

上記複数チャンネルの観測信号と、上記音源存在事後確率と、上記モデルパラメータと、を入力として上記各マイクロホンごとに上記各音源からの到来信号を推定して出力する出力音推定部と、

を具備する音源分離装置。

【請求項2】

請求項1に記載した音源分離装置において、

上記観測信号のモデルは、

m番目のマイクロホンで観測される信号 $o_{t,f}^{(m)}$ (但し、tは時間のインデックス、fは周波数のインデックスとする)が、上記複数の音源の各々から到来し当該m番目のマイクロホンで観測される到来信号のうち、最大の音圧を持つ到来信号と同値となるよう定義されたモデルであり、

上記到来信号のモデルは、

m番目のマイクロホンで観測されるi番目の音源からの到来信号 $x_{t,f}(i,m)$ を

i番目の音源のクリーン音声信号 $s_{t,f}(i)$ と、

i番目の音源からm番目のマイクロホンに到来する信号の音圧に対応する伝達関数 $f_{t,f}(i,m)$ と、

i番目の音源からm番目のマイクロホンに到来する信号とm番目のマイクロホンで観測されるi番目の音源からの信号との差に対応するエラー項 $e_{t,f}(i,m)$ と、

により定義した確率モデルであり、

上記モデルパラメータは、上記音源のクリーン音声信号 $s_{t,f}(i)$ と上記伝達関数 $f_{t,f}(i,m)$ と上記エラー項 $e_{t,f}(i,m)$ の分散 $\sigma_{t,f}^2(i,m)$ とである、
ことを特徴とする音源分離装置。

10

【請求項3】

請求項2に記載した音源分離装置において、

更に、記憶部と反復処理部とを備え、

上記記憶部は上記観測信号のモデルパラメータ $\hat{\Lambda}(i)$ を記憶するものであり、

上記マイク別音源存在事後確率推定部は、上記マイクロホンmごとの観測信号 $o_{t,f}(m)$ と上記記憶部に記憶されたモデルパラメータ $\hat{\Lambda}(i)$ とを入力として、当該マイクロホンmごとの観測信号 $o_{t,f}(m)$ とモデルパラメータ $\hat{\Lambda}(i)$ とを上記観測信号のモデルに当てはめたときの上記観測信号 $o_{t,f}(m)$ と上記観測信号のモデルパラメータ $\hat{\Lambda}(i)$ との同時確率に基づいて、上記マイクロホンmと音源iごとに音源存在事後確率 $\hat{M}_{t,f}(i,m)$ を推定するものであり、

20

上記モデルパラメータ推定部は、上記マイクロホンmごとの観測信号 $o_{t,f}(m)$ と上記記憶部に記憶されたモデルパラメータ $\hat{\Lambda}(i)$ と上記音源存在事後確率 $\hat{M}_{t,f}(i,m)$ とを入力として、当該マイクロホンmごとの観測信号 $o_{t,f}(m)$ とモデルパラメータ $\hat{\Lambda}(i)$ とを上記観測信号のモデルに当てはめたときの上記観測信号 $o_{t,f}(m)$ と上記観測信号のモデルパラメータ $\hat{\Lambda}(i)$ との同時確率の対数に、上記音源存在事後確率 $\hat{M}_{t,f}(i,m)$ に対応する重みを乗じた値を、全ての観測信号について足し合わせた重み付き和が大きくなるように、上記記憶部に記憶された伝達関数 $f_{t,f}(i,m)$ とエラー項 $e_{t,f}(i,m)$ の分散 $\sigma_{t,f}^2(i,m)$ とクリーン音声信号 $s_{t,f}(i)$ とを更新するものであり、

30

上記反復処理部は、所定の基準を満たすまで、上記マイク別音源存在事後確率推定部と上記モデルパラメータ推定部の処理を繰り返すものであり、

上記出力音推定部は、上記複数チャネルの観測信号と上記音源存在事後確率と上記記憶部に記憶されたパラメータ $\hat{\Lambda}(i)$ とを入力として上記音源iごとの到来信号 $x_{t,f}(i,m)$ を計算するもの、

であることを特徴とする音源分離装置。

【請求項4】

複数の音源から発せられる音源信号を複数のマイクロホンで収録した複数チャネルの観測信号と、上記複数のマイクロホンの各々で観測される上記複数の音源の各々からの信号の音圧が異なると仮定した観測信号のモデルを用いて、各マイクロホンごとに各音源に関する音源存在事後確率を推定するマイク別音源存在事後確率推定過程と、

40

上記複数チャネルの観測信号と、上記音源存在事後確率を入力として、観測信号のモデルパラメータを推定するモデルパラメータ推定過程と、

上記複数チャネルの観測信号と、上記音源存在事後確率と、上記モデルパラメータと、を入力として上記各マイクロホンごとに上記各音源からの到来信号を推定して出力する出力音推定過程と、

を備える音源分離方法。

【請求項5】

請求項4に記載した音源分離方法において、

上記観測信号のモデルは、

50

m 番目のマイクロホンで観測される信号 $o_{t, f}^{(m)}$ (但し、 t は時間のインデックス、 f は周波数のインデックスとする) が、上記複数の音源の各々から到来し当該 m 番目のマイクロホンで観測される到来信号のうち、最大の音圧を持つ到来信号と同値となるよう定義されたモデルであり、

上記到来信号のモデルは、

m 番目のマイクロホンで観測される i 番目の音源からの到来信号 $x_{t, f}^{(i, m)}$ を

i 番目の音源のクリーン音声信号 $s_{t, f}^{(i)}$ と、

i 番目の音源から m 番目のマイクロホンに到来する信号の音圧に対応する伝達関数 $f_{(i, m)}$ と、

i 番目の音源から m 番目のマイクロホンに到来する信号と m 番目のマイクロホンで観測される i 番目の音源からの信号との差に対応するエラー項 $e_{t, f}^{(i, m)}$ と、

により定義した確率モデルであり、

上記モデルパラメータは、上記音源のクリーン音声信号 $s_{t, f}^{(i)}$ と上記伝達関数 $f_{(i, m)}$ と上記エラー項 $e_{t, f}^{(i, m)}$ の分散 $\sigma_{t, f}^{(i, m)}$ とである、ことを特徴とする音源分離方法。

【請求項 6】

請求項 5 に記載した音源分離方法において、

更に、反復処理過程を備え、

上記マイク別音源存在事後確率推定過程は、上記マイクロホン m ごとの観測信号 $o_{t, f}^{(m)}$ と記憶部に記憶されたモデルパラメータ $\hat{\Lambda}^{(i)}$ とを入力として、当該マイクロホン m ごとの観測信号 $o_{t, f}^{(m)}$ とモデルパラメータ $\hat{\Lambda}^{(i)}$ とを上記観測信号のモデルに当てはめたときの上記観測信号 $o_{t, f}^{(m)}$ と上記観測信号のモデルパラメータ $\hat{\Lambda}^{(i)}$ との同時確率に基づいて、上記マイクロホン m と音源 i ごとに音源存在事後確率 $\hat{M}_{t, f}^{(i, m)}$ を推定するものであり、

上記モデルパラメータ推定過程は、上記マイクロホン m ごとの観測信号 $o_{t, f}^{(m)}$ と上記記憶部に記憶されたモデルパラメータ $\hat{\Lambda}^{(i)}$ と上記音源存在事後確率 $\hat{M}_{t, f}^{(i, m)}$ とを入力として、当該マイクロホン m ごとの観測信号 $o_{t, f}^{(m)}$ とモデルパラメータ $\hat{\Lambda}^{(i)}$ とを上記観測信号のモデルに当てはめたときの上記観測信号 $o_{t, f}^{(m)}$ と上記観測信号のモデルパラメータ $\hat{\Lambda}^{(i)}$ との同時確率の対数に、上記音源存在事後確率 $\hat{M}_{t, f}^{(i, m)}$ に対応する重みを乗じた値を、全ての観測信号について足し合わせた重み付き和が大きくなるように、上記記憶部に記憶された伝達関数 $f_{(i, m)}$ とエラー項 $e_{t, f}^{(i, m)}$ の分散 $\sigma_{t, f}^{(i, m)}$ とクリーン音声信号 $s_{t, f}^{(i)}$ とを更新するものであり、

上記反復処理過程は、所定の基準を満たすまで、上記マイク別音源存在事後確率推定過程と上記モデルパラメータ推定過程の処理を繰り返すものであり、

上記出力音推定過程は、上記複数チャネルの観測信号と上記音源存在事後確率と上記記憶部に記憶されたパラメータ $\hat{\Lambda}^{(i)}$ とを入力として上記音源 i ごとの到来信号 $x_{t, f}^{(i, m)}$ を計算する過程、

であることを特徴とする音源分離方法。

【請求項 7】

請求項 4 乃至 6 の何れかに記載した音源分離方法を、コンピュータで処理するためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、入力信号に複数の目的信号が含まれている場合において、各目的信号を精度良く抽出する音源分離装置と、その方法とプログラムに関する。

【背景技術】

【0002】

10

20

30

40

50

複数の目的音源が存在する環境で音響信号を收音すると、しばしば目的信号同士が互いに重なり合った混合信号が観測される。この時、注目している目的音源が音声信号である場合、その他の音源信号がその目的信号に重畳した影響により、目的音声の明瞭度は大きく低下してしまう。その結果、本来の目的音声信号（以下、目的信号）の性質を抽出することが困難となり、自動音声認識（以下、音声認識）システムの認識率も著しく低下する。よって認識率の低下を防ぐためには、複数の目的信号をそれぞれ分離することで、目的信号の明瞭度を回復する工夫（方法）が必要である。

【0003】

この複数の目的信号をそれぞれ分離する要素技術は、さまざまな音響信号処理システムに用いることが可能である。例えば、実環境下で收音された音から目的信号を抽出して聞き取り易さを向上させる補聴器、目的信号を抽出することで音声の明瞭度を向上させるTV会議システム、実環境で用いられる音声認識システム、機械制御インターフェースにおける機械と人間との対話装置、楽曲を検索したり採譜したりする音楽情報処理システムなどに利用することが出来る。

10

【0004】

図7に、例えば非特許文献1に開示されている従来の音源分離装置900の機能構成を示してその動作を簡単に説明する。音源分離装置900は、全マイク共通音源存在事後確率推定部90、フィルタリング部91、を備える。

【0005】

全マイク共通音源存在事後確率推定部90は、複数の音源から発せられる音源信号を複数のマイクロホンで收音した複数チャンネルの観測信号を入力として、当該各観測信号の各時間周波数ピンを特徴付ける特徴ベクトルを算出し、その特徴ベクトルを分類することで各音源に関する存在確率を計算する。フィルタリング部91は、複数のマイクロホンで收音した複数チャンネルの観測信号に、上記存在確率を乗算することで音源信号を回復する。

20

【先行技術文献】

【非特許文献】

【0006】

【非特許文献1】H. Sawada, S. Araki, and S. Makino, "Underdetermined convolutive blind source separation via frequency bin-wise clustering and permutation alignment," IEEE Trans. Audio, Speech and Lang. Process., vol. 19, pp.516-527, March 2011.

30

【発明の概要】

【発明が解決しようとする課題】

【0007】

しかし、複数のマイクロホンが空間的に大きく分散された形で配置されていると、各マイクロホンで観測されるある音源の音圧は同程度にならない。極端な場合は、ある音源はあるマイクロホンにおいて実質的に観測不可能な状況も起こり得る。このような状況では、各マイクロホンで異なる音源存在確率（アクティビティパターン）を仮定することが妥当である。しかし、従来の方法では、マイクロホン別に音源存在確率を計算することができないため、分散マイクロホンアレイ環境において、効率的な音源分離を行うことができない課題があった。

40

【0008】

この発明は、このような課題に鑑みてなされたものであり、分散マイクロホンアレイ環境においても効率的に音源分離を行うことができる音源分離装置とその方法とプログラムを提供することを目的とする。

【課題を解決するための手段】

【0009】

この発明の音源分離装置は、マイク別音源存在事後確率推定部と、モデルパラメータ推定部と、出力音推定部と、を具備する。マイク別音源存在事後確率推定部は、複数の音源から発せられる音源信号を複数のマイクロホンで收音した複数チャンネルの観測信号と、上

50

記複数のマイクロホンの各々で観測される上記複数の音源の各々からの信号の音圧が異なると仮定した観測信号のモデルを用いて、各マイクロホンごとに各音源に関する音源存在事後確率を推定する。モデルパラメータ推定部は、複数チャンネルの観測信号と、音源存在事後確率を入力として、観測信号のモデルパラメータを推定する。出力音推定部は、複数チャンネルの観測信号と、音源存在事後確率と、モデルパラメータと、を入力として各マイクロホンごとに各音源からの到来信号を推定して出力する。

【発明の効果】

【0010】

この発明の音源分離装置によれば、複数のマイクロホンごとに各音源に関して推定した音源存在事後確率を用いて、音源ごとに音源からの到来信号（音源イメージ）を推定するので分散マイクロホンアレイ環境においても効率的に音源分離を行うことができる。評価実験で確認した具体的な効果については後述する。

10

【図面の簡単な説明】

【0011】

【図1】この発明の音源分離装置100の機能構成例を示す図。

【図2】音源分離装置100の動作フローを示す図。

【図3】この発明のEMアルゴリズムとNewton-Raphson法を用いる音源分離装置100の機能構成例を示す図。

【図4】モデルパラメータ最適化の動作フローを示す図。

【図5】評価実験に使用した音響環境を示す図。

20

【図6】評価実験結果を示す図

【図7】従来の音声分離装置900の機能構成例を示す図。

【発明を実施するための形態】

【0012】

以下、この発明の実施の形態を図面を参照して説明する。複数の図面中同一のものには同じ参照符号を付し、説明は繰り返さない。実施例の説明の前に、観測信号をモデル化する。

【0013】

〔観測信号のモデル化〕

複数の点音源（1, 2, ... N_i）から発音する音声を、複数のマイクロホン（1, 2, ... N_m）のm番目のマイクロホンで観測した場合、i番目の音源から到来する信号 $x_{t,f}^{(i,m)}$ は、時間周波数領域において以下のように表される。t（t = 1, ... N_t）, f（f = 1, ..., N_f）は、時間と周波数のインデックスである。

30

【0014】

【数1】

$$\begin{aligned} x_{t,f}^{(i,m)} &= \log \left| S_{t,f}^{(i)} H_f^{(i,m)} \right|^2 + e_{t,f}^{(i,m)}, \\ &= \log \left| S_{t,f}^{(i)} \right|^2 + \log \left| H_f^{(i,m)} \right|^2 + e_{t,f}^{(i,m)}, \\ &= s_{t,f}^{(i)} + \beta_f^{(i,m)} + e_{t,f}^{(i,m)}, \end{aligned} \quad (1)$$

40

【0015】

ここで $S_{t,f}^{(i)}$ と $s_{t,f}^{(i)}$ は、それぞれi番目の音源からのクリーン音声信号の短時間フーリエ変換領域での信号と対数パワー領域での信号に相当し、それぞれマイク位置非依存のパラメータである。また、 $H_f^{(i,m)}$ と $h_f^{(i,m)}$ は、同様に短時間フーリエ変換領域と対数パワースペクトル領域での伝達関数に相当する。

【0016】

以降の説明では、変数 $h_f^{(i,m)}$ はマイク位置依存・音源時不変ゲインと称する。i番目の音源から到来する信号 $x_{t,f}^{(i,m)}$ を音源イメージと称する。 $e_{t,f}^{(i,m)}$

50

i, m) はエラー項であり、 $x_{t,f}^{(i,m)}$ と $\log |S_{t,f}^{(i)}| H_f^{(i,m)}$ の差であり、例えば伝達関数の揺らぎを表す。このエラー項 $e_{t,f}^{(i,m)}$ は、平均 0、分散 $\sigma_{t,f}^{(i,m)}$ の白色信号であると仮定する。

【0017】

以上の定義に従うと、 i 番目の音源からのクリーン音声信号 $s_{t,f}^{(i)}$ とその音源イメージ $x_{t,f}^{(i,m)}$ との関係は、ガウス分布の確率密度関数として次のようにモデル化することができる。

【0018】

【数 2】

$$p(x_{t,f}^{(i,m)}; \theta^{(i)}) = N(s_{t,f}^{(i)} + \beta_f^{(i,m)}, \sigma_f^{(i,m)}) \quad (2)$$

10

【0019】

ここで、 $\theta^{(i)}$ はモデルパラメータ式を表す。N は正規分布 (Normal distribution) を意味する。

【0020】

次に、LogMax 近似を用いて、複数の点音源が存在する環境における m 番目のマイクロホンで収録した観測信号 $o_{t,f}^{(m)}$ をモデル化する。その近似を用いれば、次式に示すように観測信号 $o_{t,f}^{(m)}$ は、全点音源の中で最大の音圧を持つ支配的な音源信号の値と同値となる。

【0021】

【数 3】

$$o_{t,f}^{(m)} = \max \{x_{t,f}^{(1,m)}, \dots, x_{t,f}^{(N_i,m)}\} \quad (3)$$

20

【0022】

このモデル化では支配的ではない音源は、観測信号の対数パワースペクトル以下の値であれば、任意の値を取ることができる。上記した LogMax 近似モデルは、次式に示すように確率的に定式化される。

【0023】

【数 4】

$$p(o_{t,f}^{(m)} | I_{t,f}^{(m)}, x_{t,f}^{(1,m)}, \dots, x_{t,f}^{(N_i,m)}) = \delta(o_{t,f}^{(m)} - x_{t,f}^{(I_{t,f}^{(m)},m)}) \quad (4)$$

30

$$p(I_{t,f}^{(m)} | x_{t,f}^{(1,m)}, \dots, x_{t,f}^{(N_i,m)}) = \begin{cases} 1 & \text{if } I_{t,f}^{(m)} = \arg \max_i x_{t,f}^{(i,m)} \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

【0024】

ここで、 $I_{t,f}^{(m)}$ は、 m 番目のマイクロホンの観測信号の各時間周波数ビンにおける支配的な音源の音源インデックスを表し、 $\delta(\cdot)$ はディラックのデルタ関数を表す。以降の説明では、変数 $I_{t,f}^{(m)}$ は支配的音源インデックス (DSI: Dominant Source Index) と称し、簡単のために添え字は省略する。

40

【0025】

式 (3) は、 m 番目のマイクロホンにおける観測信号 $o_{t,f}^{(m)}$ が、そのマイクロホンにおける支配的な音源イメージと同値であることを表している。ここで、マイクロホンごとに異なる音声のアクティビティパターン、つまり支配的音源インデックス DSI が割り当てられていることに注意されたい。

【0026】

上記した確率モデルを用いると観測信号 $o_{t,f}^{(m)}$ と I (支配的音源インデックス DSI) の同時確率は次式のように導出される。

【0027】

50

【数5】

$$p(o_{t,f}^{(m)}, I; \theta) = p(x_{t,f}^{(I,m)} = o_{t,f}^{(m)}; \theta^{(I)}) \\ \times \prod_{i \neq I} \int_{-\infty}^{o_{t,f}^{(m)}} p(x_{t,f}^{(i,m)}; \theta^{(i)}) dx_{t,f}^{(i,m)}. \quad (6)$$

【0028】

なお、 $\theta^{(i)}$ は各音源 i に関するパラメータを表し、 θ はすべての音源に関するパラメータを表す。すなわち、式(6)は、観測信号 $o_{t,f}^{(m)}$ と I (支配的音源インデックス DSI) を含むモデルパラメータの同時確率である。各音源の音源イメージ $x_{t,f}^{(i,m)}$ と観測信号の確率モデルを、上記したようにモデル化した前提で、以下の実施例を説明する。なお、以降の説明では、上述のLogMax近似モデル(式(4))を、「LogMax観測モデル」あるいは「観測信号の確率モデル」として参照する。

10

【0029】

〔この発明の考え〕

この発明の音源分離方法は、上記した音源イメージ $x_{t,f}^{(i,m)}$ に含まれる重要なパラメータに着目することで、複数のマイクロホンごとに異なるアクティビティパタンの推定を可能にする。

【0030】

この発明の音源分離方法を特徴付ける重要なパラメータは、支配的音源インデックス DSI である。支配的音源インデックス DSI は、各音源の各マイクロホンにおけるアクティビティパターンを示しているため、このパラメータを推定できれば、各マイクロホンごとに異なるアクティビティパターンを推定することが直接的に可能となる。

20

【0031】

この支配的音源インデックス DSI に加えて、当該パラメータを暗に支える形となっている時不変のマイク位置依存・音源時不変ゲイン $g_{t,f}^{(i,m)}$ と、時変のマイク非依存・音源対数パワースペクトル $s_{t,f}^{(i)}$ を用いる(式(1)参照)。

【0032】

これらのパラメータを用いることで、アクティビティパターンが推定できる原理を簡単に説明する。例えば、仮にある音源が m 番目のマイクロホンに高い SNR で到来すると、 SNR に対応するパラメータであるマイク位置依存・音源時不変ゲイン $g_{t,f}^{(i,m)}$ は相対的に高い値を取る傾向にあり、その音源はLogMax観測モデルの元で支配的な音源として観測される。

30

【0033】

ある時間周波数ビンにおいて支配的な音源として陽に観測された信号は、その音源の対数パワースペクトルを推定することを可能にする。一方で、ある音源が m 番目のマイクロホンに低い SNR で到来すると、マイク位置依存・音源時不変ゲイン $g_{t,f}^{(i,m)}$ は相対的に低い値を取る傾向にあり、その音源はLogMax観測モデルの元で非支配的な音源となる。LogMax観測モデルの元では、非支配的な音源のスペクトルは陽には観測されないため、その音源の対数パワースペクトルの推定は行われぬ。

40

【0034】

このようにこの発明では、各音源の対数パワースペクトルの推定を行うのに SNR の高い、一般的には音源に近いマイクロホンの観測信号を主に用いるようになる。その結果、複数のマイクロホンからの情報を効果的に加味しながら、各マイクロホンごとに異なるアクティビティパタンの推定が可能となる。

【0035】

具体的な実施例では、支配的音源インデックス DSI を潜在変数とした期待値最大化法(EMアルゴリズム)を用いてアクティビティパタンの推定を行う。Eステップ(期待値)では、支配的音源インデックス DSI に関する事後確率を更新し、どの音源がどのマイ

50

クロホンのどの時間周波数ビンで支配的かという情報を推定する。Mステップ(更新)では、その事後確率に基づいて、各音源のマイク位置依存・音源時不変ゲイン $f(i, m)$ とマイク非依存・音源対数パワースペクトル $s_{t, f}(i)$ とエラー項 $e_{t, f}(i, m)$ の分散 $\sigma_{t, f}(i, m)$ を更新する。

【実施例1】

【0036】

図1に、この発明の音源分離装置100の機能構成例を示す。その動作フローを図2に示す。音源分離装置100は、マイク別音源存在事後確率推定部10と、モデルパラメータ推定部20と、出力音推定部30と、を具備する。音源分離装置100の各部の機能は、例えばROM、RAM、CPU等で構成されるコンピュータに所定のプログラムが読み込まれて、CPUがそのプログラムを実行することで実現されるものである。

10

【0037】

マイク別音源存在事後確率推定部10は、複数の音源から発せられる音源信号を複数のマイクロホンで収録した複数チャンネルの観測信号 $o_{t, f}(m)$ と、マイクロホンの各々で観測される上記複数の音源 i の各々からの信号の音圧が異なると仮定した観測信号のモデルを用いて、各マイクロホン m ごとに各音源 i に関する音源存在事後確率 $\hat{M}_{t, f}(i, m)$ を推定する(ステップS10)。ここで、観測信号のモデルは、 m 番目のマイクロホンで観測される信号 $o_{t, f}(i, m)$ が、複数の音源の各々から到来し当該 m 番目のマイクロホンで観測される到来信号のうち、最大の音圧を持つ到来信号と同値となるように定義されたモデル(LogMax観測モデル、式(4))である。また、到来信号のモデルは、 m 番目のマイクロホンで観測される i 番目の音源の音源イメージ $x_{t, f}(i, m)$ が、 i 番目の音源のマイク非依存・音源対数パワースペクトル $s_{t, f}(i)$ と、 i 番目の音源から m 番目のマイクロホンに到来する信号の音圧に対応するマイク位置依存・音源時不変ゲイン $f(i, m)$ と、 i 番目の音源から m 番目のマイクロホンに到来する信号と m 番目のマイクロホンで観測される i 番目の音源からの信号との差に対応するエラー項 $e_{t, f}(i, m)$ と、により定義した確率モデルである(式(1))。

20

【0038】

なお、マイク非依存・音源対数パワースペクトル $s_{t, f}(i)$ は、マイクロホンに依存しない音源からのクリーン音声信号と称しても良いものである。また、マイク位置依存・音源時不変ゲイン $f(i, m)$ は、音源とマイクロホン位置によって変化する値であり、伝達関数と称しても良いものである。なお、 $\hat{\cdot}$ 等の表記は、図及び式中に表記されているように変数の直上に位置するのが正しい表記である。

30

【0039】

モデルパラメータ推定部20は、複数チャンネルの観測信号 $o_{t, f}(m)$ と、マイク別音源存在事後確率推定部10で推定した音源存在事後確率 $\hat{M}_{t, f}(i, m)$ を入力として、観測信号のモデルパラメータ $\hat{\cdot}(i)$ を推定する(ステップS20)。モデルパラメータ $\hat{\cdot}(i)$ は、マイク非依存・音源対数パワースペクトル $s_{t, f}(i)$ と、マイク位置依存・音源時不変ゲイン $f(i, m)$ と、エラー項 $e_{t, f}(i, m)$ の分散 $\sigma_{t, f}(i, m)$ と、である。

【0040】

出力音推定部30は、複数チャンネルの観測信号 $o_{t, f}(m)$ と、マイク別音源存在事後確率推定部10で推定した音源存在事後確率 $\hat{M}_{t, f}(i, m)$ と、モデルパラメータ推定部20で推定したモデルパラメータ $\hat{\cdot}(i)$ と、を入力として各マイクロホン m ごとに各音源 i に関する音源イメージ $x_{t, f}(i, m)$ を推定して出力する(ステップS30)。

40

【0041】

以上説明したように動作する音源分離装置100は、複数の各マイクロホン m において各音源 i ごとに推定した音源存在事後確率 $\hat{M}_{t, f}(i, m)$ を用いて、音源 i ごとの音源イメージ $x_{t, f}(i, m)$ を推定するので分散マイクロホンアレイ環境においても効率的に音源分離を行うことができる。以降において、音源分離装置100の動作を更に

50

詳しく説明する。

【0042】

音源分離装置100は、最大事後確率(MAP)基準で効果的にモデルパラメータ $\hat{\theta}^{(i)}$ の推定を行う。この実施例では、支配的音源インデックスDSIを潜在変数とみなして、モデルパラメータ $\hat{\theta}^{(i)} = (s_{t,f}^{(i)}, \{f^{(i,m)}, t_{t,f}^{(i,m)}\})$ を推定する。効率的な最大事後確率パラメータ推定を行うために、この実施例ではEMアルゴリズムを用い以下の補助関数を繰り返し最大化する。

【0043】

【数6】

$$Q(\theta|\hat{\theta}) = E\{\log p(\{o_{t,f}^{(m)}\}, \{I_{t,f}^{(m)}\}, \theta|\hat{\theta})\} \quad 10$$

$$= \sum_t \sum_m \sum_f \sum_i \hat{M}_{t,f}^{(i,m)} \log p(o_{t,f}^{(m)}, I_{t,f}^{(m)} = i; \theta),$$

$$= \sum_m \sum_i Q^{(i,m)}(\theta^{(i)}|\hat{\theta}),$$

$$Q^{(i,m)}(\theta^{(i)}|\hat{\theta}) = \sum_t \sum_f [\hat{M}_{t,f}^{(i,m)} p(x_{t,f}^{(i,m)} = o_{t,f}^{(m)}; \theta^{(i)}) \quad 20$$

$$+ (1 - \hat{M}_{t,f}^{(i,m)}) \int_{-\infty}^{o_{t,f}^{(m)}} p(x_{t,f}^{(i,m)}; \theta^{(i)}) dx_{t,f}^{(i,m)}] \quad (7)$$

【0044】

ここで、 $\hat{\theta}^{(i)}$ はモデルパラメータの事前推定値、 $\theta^{(i)}$ はモデルパラメータの推定値を表す。また、式(7)における $p(x_{t,f}^{(i,m)}; \theta^{(i)})$ は、式(2)で定義されている通り、モデルパラメータの事前推定値 $\hat{\theta}^{(i)}$ から算出することができる。なお、事前推定値 $\hat{\theta}^{(i)}$ は予め与えられているものとする。すなわち、上述の補助関数 $Q(\theta|\hat{\theta})$ は、観測信号 $o_{t,f}^{(m)}$ と支配的音源インデックスDSIを含むモデルパラメータの事前推定値との同時確率 $p(o_{t,f}^{(m)}, I_{t,f}^{(m)} = i; \theta^{(i)})$ に、音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ に対応する重みを乗じた値を、全ての観測信号について足し合わせた重み付き和である。EMアルゴリズムでは、この補助関数の値が大きくなるように、モデルパラメータを更新する。

30

【0045】

各マイクロホンmにおける音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ は次式で表せる。

【0046】

【数7】

$$\hat{M}_{t,f}^{(i,m)} = \frac{p(o_{t,f}^{(m)}, I_{t,f}^{(m)} = i; \hat{\theta})}{\sum_{i'} p(o_{t,f}^{(m)}, I_{t,f}^{(m)} = i'; \hat{\theta})} \quad (8) \quad 40$$

【0047】

式(7)は、第二項の複雑性により、解析的に最大化することができない。そこで、この実施例では、Newton-Raphson法を用いて効率的に補助関数を最大化する。

【0048】

図3に、EMアルゴリズムとNewton-Raphson法を用いる音源分離装置100の機能構成例を示す。音源分離装置100は、音源分離装置100の構成に加えて、更に記憶部40と、反復処理部50と、を備える。モデルパラメータ推定部20は、マイク位置依存・音源時不変ゲイン推定手段201と、マイク非依存・音源対数パワースペクトル推定手

50

段 202 と、を含む。

【0049】

パラメータの最適化手順は、マイク別音源存在事後確率推定部 10 とモデルパラメータ推定部 20 と記憶部 40 と反復処理部 50 と、で行う。図 4 に、パラメータの最適化手順の動作フローを示す。

【0050】

記憶部 40 には、モデルパラメータ $\hat{\theta}^{(i)} = (\hat{s}_{t,f}^{(i)}, \hat{\beta}_f^{(i,m)}, \hat{\sigma}_f^{(i,m)})$ の初期値 と、更新された値とが記憶される。記憶部 40 は、更新されたモデルパラメータ $\hat{\theta}^{(i)}$ のみを記憶し、初期値 はその値を必要とする各部に予め定数として持たせるようにしても良い。

10

【0051】

マイク別音源存在事後確率推定部 10 は、複数のマイクロホンごとの観測信号 $o_{t,f}^{(m)}$ と、記憶部 40 に記憶されたモデルパラメータ $\hat{\theta}^{(i)} = (\hat{s}_{t,f}^{(i)}, \hat{\beta}_f^{(i,m)}, \hat{\sigma}_f^{(i,m)})$ とを入力として、各マイクロホンごとに、式 (8) により、各音源 i に関する音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ を計算する (ステップ S10)。すなわち、マイク別音源存在事後確率推定部 10 は、観測信号 $o_{t,f}^{(m)}$ とモデルパラメータ $\hat{\theta}^{(i)}$ とを観測信号のモデルに当てはめたときの、観測信号 $o_{t,f}^{(m)}$ とモデルパラメータ $\hat{\theta}^{(i)}$ との同時確率に基づいて、音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ を計算する。この処理は、EM アルゴリズムの E ステップに当たる。

20

【0052】

マイク位置依存・音源時不変ゲイン推定手段 201 は、複数のマイクロホンごとの観測信号 $o_{t,f}^{(m)}$ と、マイク別音源存在事後確率推定部 10 で計算した音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ と、記憶部 40 に記憶されたモデルパラメータ $\hat{\theta}^{(i)}$ のマイク非依存・音源対数パワースペクトル $\hat{s}_{t,f}^{(i)}$ を入力として、次式でマイク位置依存・音源時不変ゲイン $\hat{\beta}_f^{(i,m)}$ と分散 $\hat{\sigma}_f^{(i,m)}$ を計算して、記憶部 40 に記憶されている当該パラメータの値を更新する (ステップ S201)。なお、以下の式では、条件 $o_{t,f}^{(m)} > (\hat{s}_{t,f}^{(i)} + \hat{\beta}_f^{(i,m)})$ が満たされる場合は、 $\hat{\beta}_f^{(i,m)} = \hat{M}_{t,f}^{(i,m)}$ とし、満たされない場合は $\hat{\beta}_f^{(i,m)} = 1$ とする。

30

【0053】

【数 8】

$$\hat{\beta}_f^{(i,m)} \leftarrow \hat{\beta}_f^{(i,m)} - \left(\frac{\partial^2 Q^{(i,m)}(\theta^{(i)} | \hat{\theta})}{\partial \beta_f^{(i,m)^2}} \right)^{-1} \left(\frac{\partial Q^{(i,m)}(\theta^{(i)} | \hat{\theta})}{\partial \beta_f^{(i,m)}} \right), \quad (9)$$

$$\hat{\sigma}_f^{(i,m)} \leftarrow \frac{\sum_t \hat{K}_{t,f}^{(i,m)} (o_{t,f}^{(i,m)} - (\hat{s}_{t,f}^{(i)} + \hat{\beta}_f^{(i,m)}))^2}{\sum_t \hat{M}_{t,f}^{(i,m)}} \quad (10)$$

40

【0054】

マイク非依存・音源対数パワースペクトル推定手段 202 は、マイクロホン m ごとの観測信号 $o_{t,f}^{(m)}$ と、記憶部 40 に記憶されたモデルパラメータ $\hat{\theta}^{(i)}$ と、マイク別音源存在事後確率推定部 10 で計算した音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ を入力として、複数のマイクロホン m との間で共通となる i 番目の音源からのクリーン音声信号 $s_{t,f}^{(i)}$ を次式で計算して、記憶部 40 に記憶されている当該パラメータの値を更新する (ステップ S202)。ステップ S201 と S202 の処理 (ステップ S20) は、EM アルゴリズムの M ステップに当たる。

【0055】

50

【数9】

$$\hat{s}_{t,f}^{(i)} \leftarrow \hat{s}_{t,f}^{(i)} - \left(\frac{\partial^2 Q^{(i,m)}(\theta^i | \hat{\theta})}{\partial s_{t,f}^{(i)2}} \right)^{-1} \left(\frac{\partial Q^{(i,m)}(\theta^i | \hat{\theta})}{\partial s_{t,f}^{(i)}} \right), \quad (11)$$

【0056】

また、 $\hat{s}_{t,f}^{(i)}$ と $\hat{\mu}_f^{(i,m)}$ の更新式は類似していることが分かる。これらの更新式の違いは平均化処理にあり、 $\hat{s}_{t,f}^{(i)}$ はマイクロホン番号に関する平均として計算され、一方で $\hat{\mu}_f^{(i,m)}$ は、時間インデックスに関する平均として計算される。

10

【0057】

なお、式(9)における補助関数は、式(7)で定義される補助関数と式(12)で計算される値に重み w を乗じたものを加算した値とする。これは、あるマイクロホンにおいて全く支配的にならない音源 (LogMax観測モデルの元では陽には全く観測されない音源) があると、マイク位置依存・音源時不変ゲイン $\hat{\mu}_f^{(i,m)}$ の最適解は無限小となってしまう推定処理全体が不安定になる。前述のように、マイク非依存・音源対数パワースペクトル $\hat{s}_{t,f}^{(i)}$ に関して以下のような正規化項 (事前分布) 203 を定義し、補助関数に重み w で加算すれば、このような問題を回避することができる。

【0058】

20

【数10】

$$\log p(\{s_{t,f}^{(i)}\}) = \sum_f \log N(s_{t,f}^{(i)}; \bar{\mu}_f^{(i)}, \bar{\sigma}_f^{(i)}) \quad (12)$$

【0059】

正規化項203は、記憶部40に予め記憶させておいても良いし、図3に示すようにモデルパラメータ推定部20の内部に定数として持たせるようにしても良い。

【0060】

以上のように、モデルパラメータ推定部20では、式(7)の補助関数、つまり、観測信号 $o_{t,f}^{(m)}$ と現在のモデルパラメータ推定値 $\hat{\mu}_f^{(i)}$ を観測モデルに当てはめたときの、観測信号 $o_{t,f}^{(m)}$ と支配的音源インデックス $D S I$ を含むモデルパラメータ推定値 $\hat{\mu}_f^{(i)}$ との同時確率 $p(o_{t,f}^{(m)}, I_{t,f}^{(m)} = i; \hat{\mu}_f^{(i)})$ に、音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ に対応する重みを乗じた値を、全ての観測信号について足し合わせた重み付き和が大きくなるように、モデルパラメータ (マイク位置依存・音源時不変ゲイン $\hat{\mu}_f^{(i,m)}$ と分散 $\hat{\sigma}_{t,f}^{(i,m)}$ とマイク非依存・音源対数パワースペクトル $\hat{s}_{t,f}^{(i)}$) を更新する (式(9)~(11))。

30

【0061】

反復処理部50は、所定の基準を満たすまでEステップとMステップを繰り返す (ステップS51)。所定の基準としては、例えば更新前のモデルパラメータ $\hat{\mu}_f^{(i)}$ 及び各音源に関する音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ から計算される式(7)に示したQ関数 (補助関数) の値と、更新後のモデルパラメータ及び各音源に関する音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ から計算されるQ関数の値との差が所定の閾値未満となった時を、所定の基準を満たしたと判定する方法や、予め定めた繰り返す回数に達した場合に所定の基準を満たしたと判定する方法が考えられる。繰り返し処理を行うことで補助関数を最大化することができる。

40

【0062】

所定の基準を満たすと、出力音推定部30は、複数のマイクロホンごとの観測信号 $o_{t,f}^{(m)}$ と、マイク別音源存在事後確率推定部10で計算した音源存在事後確率 $\hat{M}_{t,f}^{(i,m)}$ と、記憶部40に記憶されたモデルパラメータ $\hat{\mu}_f^{(i)}$ と、を入力として、m番目のマイクロホンにおけるi番目の音源イメージ $\hat{x}_{t,f}^{(i,m)}$ を計算し

50

て出力する。EMアルゴリズムを用いてパラメータ推定を行うと最小二乗誤差推定で音源イメージ $\hat{x}_{t,f}^{(i,m)}$ を求めることが可能となる。推定される音源イメージ $\hat{x}_{t,f}^{(i,m)}$ は、次式で表される。

【0063】

【数11】

$$\hat{x}_{t,f}^{(i,m)} = M_{t,f}^{(i,m)} o_{t,f}^{(m)} + \left(1 - M_{t,f}^{(i,m)}\right) \frac{\int_{-\infty}^{o_{t,f}^{(m)}} \tilde{x}_{t,f}^{(i,m)} p\left(x_{t,f}^{(i,m)}; \hat{\theta}^{(i)}\right) dx_{t,f}^{(i,m)}}{\int_{-\infty}^{o_{t,f}^{(m)}} p\left(x_{t,f}^{(i,m)}; \hat{\theta}^{(i)}\right) dx_{t,f}^{(i,m)}}. \quad (13)$$

ここで、 $\tilde{x}_{t,f}^{(i,m)}$ は $\tilde{x}_{t,f}^{(i,m)} = s_{t,f}^{(i)} + \beta_f^{(i,m)}$ とした。

【0064】

〔評価実験〕

この発明の音源分離装置100の性能を評価する目的で評価実験を行った。実験条件は次の通りとした。

【0065】

図5に、シミュレーションに用いた音響環境を示す。部屋のサイズは10m(W)×5m(D)×5m(H)であり、残響時間は100msである。この音響環境を鏡像法(参考文献1: J. B. Allen and D. A. Berkeley, "Image method for efficiently simulating small-room acoustics," J. Acoust. Soc. Am., vol. 65(4), pp. 943-950, 1979.)を用いてシミュレーションした。

【0066】

音響環境としては4つの環境を模擬した。第1音響環境と第2音響環境は、3人の話者が半径80cmの円状に等間隔を開けて座り、同時会話する状況を想定した。第1音響環境は、3つのマイクロホンが半径10cmの同心円状に配置されている状況とし、第2音響環境は、同じマイクロホンが半径50cmの同心円状に配置されている状況とした。図3において、第1音響環境と第2音響環境は一方の2人の話者とマイクロホンのグループが存在しない状態である。

【0067】

第3音響環境と第4音響環境は、3人の話者と2人の話者の2つのグループが同じ部屋で会話している状況を想定した。第3音響環境は、5つのマイクロホンが半径10cmの同心円状に配置されている状況とし、第4音響環境は、同じマイクロホンが半径50cmの同心円状に配置されている状況とした。

【0068】

第1番目と第2番目の音響環境においては3音源の分離を行った。第3番目と第4番目の音響環境においては5音源の分離を行った。この発明と比較する従来法は、すべてのマイクロホンにおいて共通の音源アクティビティパターンを仮定して、ソフトマスクを用いた音源分離を行う非特許文献1に示された方法とした。従来法では、各音源に最も近いマイク観測信号にソフトマスク処理を行い、分離信号を算出した。

【0069】

この発明の方法では、EMアルゴリズムの初期値として従来方法の処理結果を使用した。式(12)に示した正規化項の計算にも従来法の処理結果を用いた。正規化項の重みは $w = 0.00001$ とした。

【0070】

評価指標としてはケプストラム距離を用いた。ケプストラム距離は、比較対象信号と各音源に最も近いマイクロホンにおける各音源イメージの距離とした。評価音声としては、

10

20

30

40

50

TIMIT (参考文献2: W. Fisher, G.R. Doddington, and K. M. Goudie-Marshall, "The DARPA speech recognition research database: specifications and status," in Proc. DARPA workshop on Speech Recognition, 7986, pp. 96-99.) から無作為に抽出した音声を用い、各音響環境において計20個の異なる混合音声を用意し、結果はそれらの平均値として算出した。

【0071】

図6に、評価実験の結果を示す。横軸は音響環境、縦軸はケプストラム距離(dB)である。音響環境ごとに観測信号と従来法と本発明のケプストラム距離を示す。ここで、観測信号のケプストラム距離の算出のためには、各話者に最も近いマイクロホンの観測信号を用いており、最近傍マイクロホンを既知とした際のマイクロホン選択処理の結果に相当する。

10

【0072】

第1音響環境における結果では、従来法でもケプストラム距離を減らしているが、本発明は更にケプストラム距離を減らすことができている。これは、この発明の方法がケプストラム領域と類似する対数パワースペクトル領域にてパラメータ最適推定を行っているためと考えられる。

【0073】

第2～第4音響環境では、従来法による性能改善を確認することができない。従来法はケプストラム距離尺度で性能が劣化しており、過抑圧などにより歪が増大していることが予想される。本発明の方法では、全ての音響環境において、効果的にケプストラム距離を減少させることができた。このように本発明の音源分離装置100によれば、分散マイクロホンアレイ環境においても効率的に音源分離を行うことが確認できた。

20

【0074】

上記した音声分離装置100における処理手段をコンピュータによって実現する場合、各装置が有すべき機能の処理内容はプログラムによって記述される。そして、このプログラムをコンピュータで実行することにより、各装置における処理手段がコンピュータ上で実現される。

【0075】

なお、効率的に最大事後確率パラメータ推定を行う目的で、EMアルゴリズムNewton-Raphson法を用いた音源分離装置100について説明を行ったが、この発明はこの実施例に限定されない。例えば最大事後確率パラメータ推定を行うのに、EMアルゴリズムを用いる必要はない。全ての組み合わせを探索する全組み合わせ探索法を用いても、この発明の技術思想の範囲に含まれる。

30

【0076】

この処理内容を記述したプログラムは、コンピュータで読み取り可能な記録媒体に記録しておくことができる。コンピュータで読み取り可能な記録媒体としては、例えば、磁気記録装置、光ディスク、光磁気記録媒体、半導体メモリ等のようなものでもよい。具体的には、例えば、磁気記録装置として、ハードディスク装置、フレキシブルディスク、磁気テープ等を、光ディスクとして、DVD(Digital Versatile Disc)、DVD-RAM(Random Access Memory)、CD-ROM(Compact Disc Read Only Memory)、CD-R(Recordable)/RW(ReWritable)等を、光磁気記録媒体として、MO(Magneto Optical disc)等を、半導体メモリとしてEEPROM(Electronically Erasable and Programmable-Read Only Memory)等を用いることができる。

40

【0077】

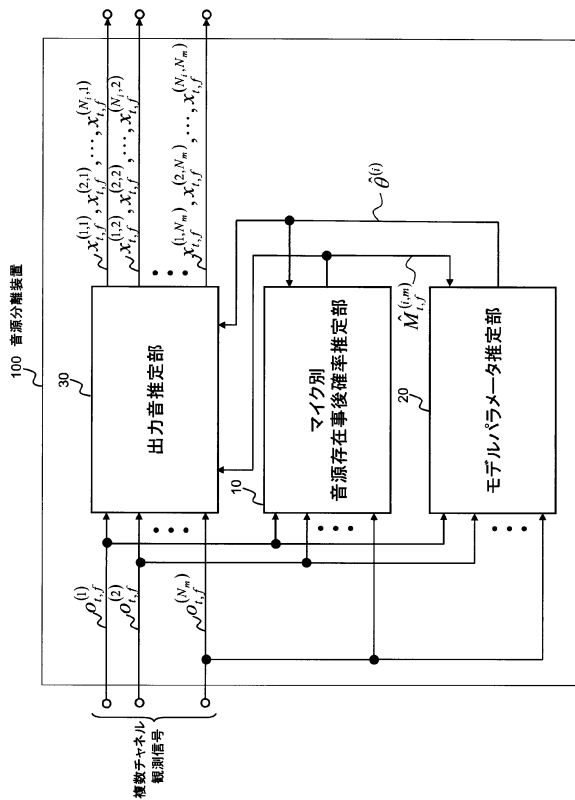
また、このプログラムの流通は、例えば、そのプログラムを記録したDVD、CD-ROM等の可搬型記録媒体を販売、譲渡、貸与等することによって行う。さらに、このプログラムをサーバコンピュータの記録装置に格納しておき、ネットワークを介して、サーバコンピュータから他のコンピュータにそのプログラムを転送することにより、このプログラムを流通させる構成としてもよい。

【0078】

50

また、各手段は、コンピュータ上で所定のプログラムを実行させることにより構成することにもよいし、これらの処理内容の少なくとも一部をハードウェア的に実現することとしてもよい。

【図1】



【図2】

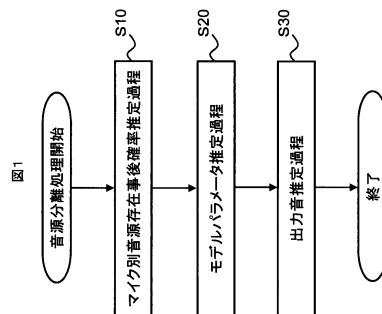
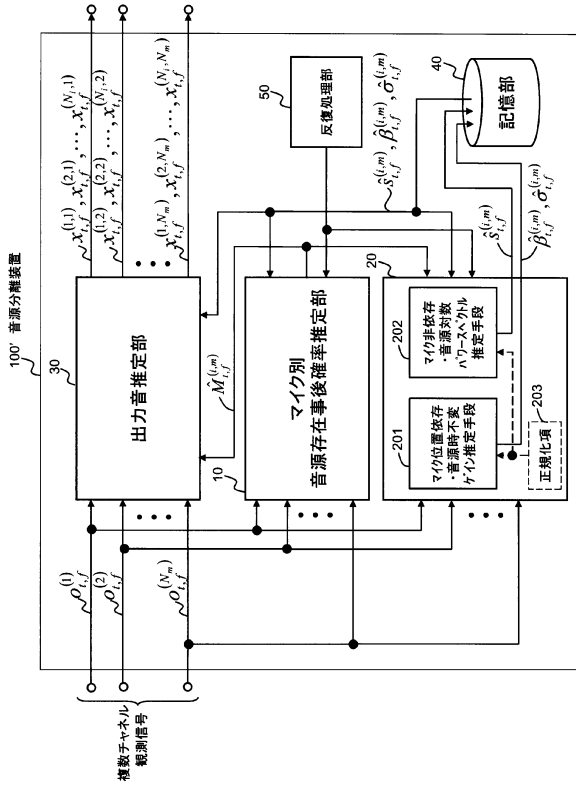
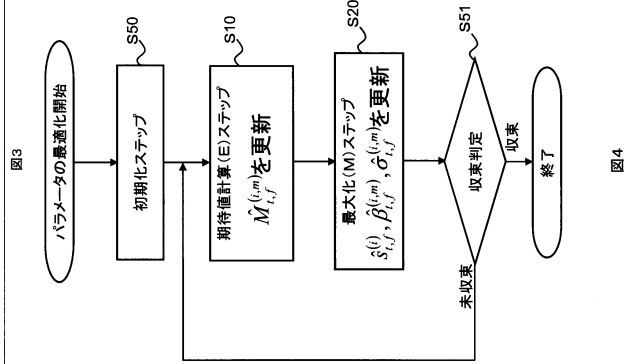


図2

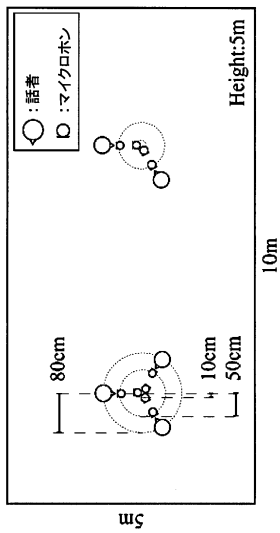
【 図 3 】



【 図 4 】



【 図 5 】



【 図 6 】

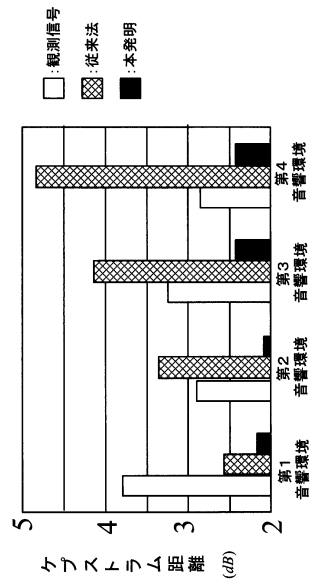


図5

図6

【図7】

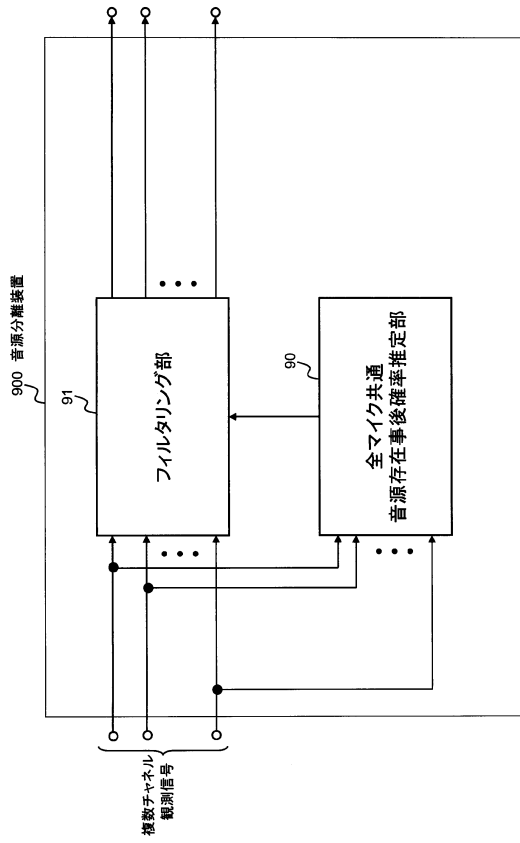


図7

フロントページの続き

審査官 安田 勇太

- (56)参考文献 特開2013-054258(JP,A)
特開2008-079256(JP,A)
ソウデン メRez, ノード内・ノード間情報の統合に基づく分散マイクアレイ音源分離, 日本音響学会 2013年 春季研究発表会講演論文集, 日本, 2013年 3月

- (58)調査した分野(Int.Cl., DB名)
G10L 21/0208 - 21/0308