

(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2014-89249

(P2014-89249A)

(43) 公開日 平成26年5月15日(2014.5.15)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 O L 21/028 (2013.01)	G 1 O L 21/02 2 O 1 D	5 D O 1 8
G 1 O L 21/0308 (2013.01)	G 1 O L 21/02 2 O 3 Z	5 D 2 2 0
H O 4 R 1/40 (2006.01)	H O 4 R 1/40 3 2 O A	
H O 4 R 3/00 (2006.01)	H O 4 R 3/00 3 2 O	

審査請求 未請求 請求項の数 5 O L (全 13 頁)

(21) 出願番号	特願2012-237835 (P2012-237835)	(71) 出願人	000006013 三菱電機株式会社 東京都千代田区丸の内二丁目7番3号
(22) 出願日	平成24年10月29日(2012.10.29)	(74) 代理人	100123434 弁理士 田澤 英昭
		(74) 代理人	100101133 弁理士 濱田 初音
		(74) 代理人	100173934 弁理士 久米 輝代
		(74) 代理人	100156351 弁理士 河村 秀央
		(72) 発明者	太刀岡 勇氣 東京都千代田区丸の内二丁目7番3号 三菱電機株式会社内
		Fターム(参考)	5D018 BB21 5D220 BA01 BC05

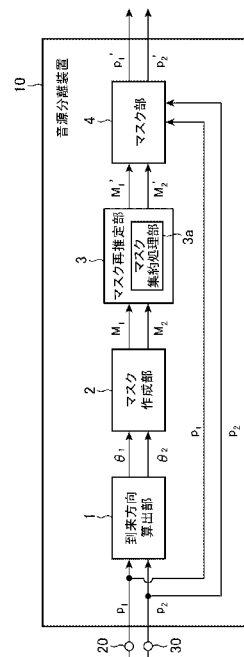
(54) 【発明の名称】音源分離装置

(57) 【要約】

【課題】TDOAから観測音の到来方向を算出して作成したマスクを音声らしさを用いて修正し、マスクのパラミュテーションを解決する。

【解決手段】各音源から出力された音声の到来方向を算出する到来方向算出部1と、到来方向算出部1が算出した各到来方向の時間周波数平面において、観測信号スペクトルから複数の音源のうち対応する音源から出力された目的音声以外の観測音の観測信号スペクトルをマスキングするマスクを作成するマスク作成部2と、音声の特徴に基づいて、各マスクについて目的音声と当該目的音声以外の観測音との分離性能を検証し、検証結果に基づいてマスクを再推定するマスク再推定部3と、マスク再推定部3が再推定した各マスクを用いて、観測信号スペクトルから目的音声以外の観測音の観測信号スペクトルをマスキングし、目的音声の観測信号スペクトルを取得するマスク部4とを備える。

【選択図】図2



【特許請求の範囲】**【請求項 1】**

複数の音源からの音声が入力された観測音を時間周波数領域に変換した観測信号スペクトルから、前記各音源から出力された音声の到来方向を算出する到来方向算出部と、

前記到来方向算出部が算出した各到来方向の時間周波数平面において、前記観測信号スペクトルから前記複数の音源のうち対応する音源から出力された目的音声以外の観測音の観測信号スペクトルをマスキングするマスクを作成するマスク作成部と、

前記音声の特徴に基づいて、前記マスク作成部が作成した各マスクについて、前記目的音声と当該目的音声以外の観測音との分離性能を検証し、検証結果に基づいて前記マスクを再推定するマスク再推定部と、

前記マスク再推定部が再推定した各マスクを用いて、前記観測信号スペクトルから前記目的音声以外の観測音の観測信号スペクトルをマスキングし、前記目的音声の観測信号スペクトルを取得するマスク部とを備えた音源分離装置。

10

【請求項 2】

前記マスク部は、前記マスク作成部が作成した各マスクを用いて、前記観測信号スペクトルから前記目的音声以外の観測音の観測信号スペクトルをマスキングし、前記目的音声の観測信号スペクトルを取得し、

前記マスク部において前記マスク再推定部が再推定したマスクを用いて取得した前記目的音声の観測信号スペクトル、および前記マスク部において前記マスク作成部が作成したマスクを用いて取得した前記目的音声の観測信号スペクトルについて、音声モデルに対するそれぞれの尤度を算出する尤度算出部と、

前記尤度算出部が算出した尤度に基づいて、前記マスク再推定部が再推定したマスク、または前記マスク作成部が作成したマスクのいずれか一方を選択し、選択したマスクを用いてマスキングした前記目的音声の観測信号スペクトルを取得するマスク選択部とを備えたことを特徴とする請求項 1 記載の音源分離装置。

20

【請求項 3】

前記マスク再推定部は、前記音声の時間的連続性、前記音声の倍音構造、または前記音声の話者特性に基づいて、前記マスク作成部が作成した各マスクについて、前記目的音声と当該目的音声以外の観測音との分離性能を検証し、検証結果に基づいて同一の音源から出力された前記目的音声の観測信号スペクトルは同一のマスクのマスキングによって取得されるよう前記マスク作成部が作成した各マスクを再推定するマスク集約処理部を備えることを特徴とする請求項 2 記載の音源分離装置。

30

【請求項 4】

前記マスク再推定部は、前記音声の時間的連続性、前記音声の倍音構造、または前記音声の話者特性に基づいて、前記マスク作成部が作成した各マスクについて、前記目的音声と当該目的音声以外の観測音との分離性能を検証し、検証結果に基づいて異なる音源から出力された前記目的音声の観測信号スペクトルはそれぞれ異なるマスクのマスキングによって取得されるように、前記マスク作成部が作成したマスクを再推定するマスク分離処理部を備えることを特徴とする請求項 2 または請求項 3 記載の音源分離装置。

【請求項 5】

前記マスク再推定部は、前記マスク作成部が作成した各マスクの信頼度が低い音声帯域において、前記マスク作成部が作成した複数のマスクを組み合わせてなるマスクの組み合わせを生成するマスク交叉部を備え、

前記マスク部は、前記マスク交叉部が生成したマスクの組み合わせで指定されたマスクを用いて前記観測信号スペクトルから前記目的音声以外の観測音の観測信号スペクトルをマスキングし、前記目的音声の観測信号スペクトルを取得し、

前記尤度算出部は、前記マスク部が取得した各マスクの組み合わせによって取得された前記目的音声の観測信号スペクトルについて、前記音声モデルに対するそれぞれの尤度を算出し、

前記マスク選択部は、前記尤度算出部が算出した尤度に基づいて、前記マスク交叉部が

40

50

生成したマスクの組み合わせのうち最も尤度の高いマスクの組み合わせを選択し、選択した組み合わせのマスクを用いてマスクングした前記目的音声の観測信号スペクトルを取得することを特徴とする請求項 2 記載の音源分離装置。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、複数の音源からの音声信号が混在した観測信号から、それぞれの音源に対応する分離信号を得る音源分離装置に関するものである。

【背景技術】

【0002】

複数人の音声混ざった音声信号を分離して、各人の音声信号を取り出す技術は音声認識技術の適用範囲拡大に寄与する。音源毎の音声信号の分離方法としては、マイクの死角を対象外の話者に向けたビームフォーミング（以下、BFと称する）による方法や、独立成分分析（ICA：independent Component analysis）により混合行列を推定する方法が用いられている。また近年は、時間周波数平面上のスペクトルで音声がスパースなことを利用して、対象話者以外の成分をマスクするバイナリマスクによる分離方法が用いられている。

【0003】

一方で、BFはノイズの抑圧には優れているが、混成音声の分離にはあまり有効でない。また、ICAは残響や騒音の影響で性能が低下する。さらに、BFやICAによる分離方法では、マイクの数音源数以上でなければならないという制約がある。これに対して、バイナリマスクにはこのような制限がないため、適用先が広く、有望であると言える。

【0004】

バイナリマスクにもいくつかの手法があるが、ここでは時間・周波数binにおける音声の到来時間差（TDOA）に着目して分類を行う方法について述べる。

2つのマイクで観測された音声信号の短時間フーリエ変換後の時間周波数平面(t, f)におけるスペクトルを p_1 , p_2 とすると、各スペクトルの位相差は以下の式(1)で表される。

$$\alpha = \arg [p_1(f, t)/p_2(f, t)] \quad (-\pi \leq \theta \leq \pi) \quad \dots (1)$$

【0005】

さらに、式(1)から各スペクトルの成分の時間差と音波の到来方向が、以下の式(2)により求められる。

$$= 1/2 f \sin^{-1}(c/l_m) \quad \dots (2)$$

cは音速、 l_m はマイク間隔である。音波の到来方向を別手法で推定する、もしくはクラスタリングすることにより、音源の方向別に(t, f)領域でのマスクを作成する。

【0006】

例えば到来角 θ_1 の第1の信号に対するマスクが $M_1(t, f)$ であった場合、以下の式(3)のように推定される。

$$\begin{aligned} M_1(t, f) &= 1 && \text{for } |\theta - \theta_1| \leq \theta_t \\ M_1(t, f) &= \varepsilon && \text{for } |\theta - \theta_1| > \theta_t \quad \dots (3) \end{aligned}$$

θ_t は許容誤差、 ε は十分小さい数である。

推定されたマスクを用いてマスクされた以下の式(4)で示すスペクトルを、逆フーリエ変換してマスク後の信号を得る。

$$p'_1(f, t) = M_1(t, f) p_1(f, t) \quad \dots (4)$$

【0007】

10

20

30

40

50

TDOAによるバイナリマスクを用いた従来の音源分離装置は、例えば上述した式(2)に基づいて到来方向を算出する手段、上述した式(3)に基づいてマスクを作成する手段、および上述した式(4)に基づいて音声スペクトルをマスクすることにより音声分離スペクトルを得る手段によって構成される。

【0008】

しかし、バイナリマスクによる分離方法では、バイナリマスクの推定において、マスクを時間および周波数binといった少ない情報から推定するため、推定の精度が誤差の影響を受けやすいという問題があった。特にTDOAから観測音の到来方向を算出してマスクを作成する方法では、マイクの間隔に比して、波長の長い低周波成分の場合には位相差が付きにくいことから、波長の短い高周波成分の場合には空間的エイリアシングの影響でマスクの推定精度が低下するという問題があった。

10

【0009】

そこで、バイナリマスクを用いた音源分離方法において、音声の特徴を生かしてマスクの誤判定を抑制する技術として、例えば特許文献1および特許文献2に開示されているものがある。特許文献1には、ある周波数binに隣接する複数の周波数binのスペクトル成分に対する時間変化を連結する手法が開示されている。特許文献2には、音源分離のためのバイナリマスクングにおいて、パワースペクトルからマスクパターンを生成する手法が開示されている。

【0010】

また、バイナリマスクの妥当性を、音声モデルを用いて検証する技術として、例えば特許文献3から特許文献5に開示されているものがある。

20

特許文献3には、ブラインド音声分離にEMアルゴリズムを適用し、最大尤度を与える音源方向と、各時間周波数成分への各音源の寄与率をEMアルゴリズムによって推定する手法が開示されている。特許文献4には、信号分離において、事後確率の類似度を指標として観測信号のクラスタリングを行う手法が開示されている。特許文献5には、音源分離装置において、確率モデルのモデルパラメタと各音源の存在確率を用いて有効音源を抽出する手法が開示されている。

【先行技術文献】

【特許文献】

【0011】

30

【特許文献1】特開2008-026625号公報

【特許文献2】特開2010-239424号公報

【特許文献3】特開2008-145610号公報

【特許文献4】特開2009-053349号公報

【特許文献5】特開2011-164467号公報

【発明の概要】

【発明が解決しようとする課題】

【0012】

しかしながら、上述した特許文献1および特許文献2に開示された技術では、マスクのスパース性を利用していないため、滑らかではあるが分離性能の低い非合理的なマスクを生じるという課題があった。また特許文献3に開示された技術では、マスクがスパースになるような基準が設けられておらず、分離性能の低い非合理的なマスクを生じるという課題があった。また、特許文献4および特許文献5に開示された技術では、音声らしさを基準として用いていないため、分離音に聴感上や音声認識にとって悪影響を及ぼすひずみが入りやすいという課題があった。

40

【0013】

この発明は、上記のような課題を解決するためになされたもので、TDOAから観測音の到来方向を算出して作成したマスクを音声らしさをを用いて修正し、マスクのパーミュテーションを解決する音源分離装置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】

50

【 0 0 1 4 】

この発明に係る音源分離装置は、複数の音源からの音声が入力された観測音を時間周波数領域に変換した観測信号スペクトルから、各音源から出力された音声の到来方向を算出する到来方向算出部と、到来方向算出部が算出した各到来方向の時間周波数平面において、観測信号スペクトルから複数の音源のうち対応する音源から出力された目的音声以外の観測音の観測信号スペクトルをマスクするマスクを作成するマスク作成部と、音声の特徴に基づいて、マスク作成部が作成した各マスクについて、目的音声と当該目的音声以外の観測音との分離性能を検証し、検証結果に基づいてマスクを再推定するマスク再推定部と、マスク再推定部が再推定した各マスクを用いて、観測信号スペクトルから目的音声以外の観測音の観測信号スペクトルをマスクし、目的音声の観測信号スペクトルを取得するマスク部とを備えるものである。

10

【 発明の効果 】

【 0 0 1 5 】

この発明によれば、分離性能の高いマスクを作成することができ、明瞭な目的音声を取得することができる。

【 図面の簡単な説明 】

【 0 0 1 6 】

【 図 1 】 実施の形態 1 による音源分離装置のマスク再推定処理を示す説明図である。

【 図 2 】 実施の形態 1 による音源分離装置の構成を示すブロック図である。

【 図 3 】 実施の形態 2 による音源分離装置の構成を示すブロック図である。

20

【 図 4 】 実施の形態 3 による音源分離装置の構成を示すブロック図である。

【 図 5 】 実施の形態 4 による音源分離装置の構成を示すブロック図である。

【 図 6 】 16 kHz サンプリングでの波形とスペクトログラムを示す図である。

【 図 7 】 実施の形態 5 による音源分離装置の構成を示すブロック図である。

【 発明を実施するための形態 】

【 0 0 1 7 】

実施の形態 1 .

混合前の音声を用いて、それぞれの音源に対応する分離信号を得るためのマスク（理想マスク）を作成して観察すると、 $M_1(t, f)=1$ となる (t, f) は、ある程度まとまっている傾向にある。すなわち時間・周波数方向にはスペクトルは局所的にはある程度の連続性がある。ところが TDOA により作成したマスクは孤立点が多い。これは TDOA の推定誤差の影響で、1つの音源からの音が異なるマスクに分類されてしまうためである。そこで、この実施の形態 1 では、 $M_1(t, f)=1$ となる (t, f) を近い範囲にまとめることで、より分離性能の高いマスクを作成する。

30

【 0 0 1 8 】

なお以下では、説明の簡単化のため上述した式（3）において ϵ を 0 とする。また、2つの音源から出力されて混合された音声信号を分離する場合を例に説明する。なお、本発明の構成は、3つ以上の音源から出力されて混合された音声信号を分離する場合にも適用可能である。

【 0 0 1 9 】

40

ここである (t, f) において $M_1(t, f)=1, M_2(t, f)=0$ であったとする。この時、時間周波数平面でのスペクトルの局所的な連続性を考慮すると、 $M_2(t, f)$ の周囲の点に1が多く、 $M_1(t, f)$ の周囲の点に0が多い場合は、推定誤りであって実は $M_1(t, f)=0, M_2(t, f)=1$ である可能性が高い。ここでは例えば、密集度の指標として以下の式（5）を用いることとする。

$$\sigma_i(t, f) = \sum_{(|t-t'| \leq \Delta t, |f-f'| \leq \Delta f)} M_i(t', f') \quad \dots (5)$$

t, f はそれぞれ時間・周波数領域での近接範囲を示す。 $M_1(t, f) < M_2(t, f)$ であった場合には $M_1(t, f)=0, M_2(t, f)=1$ とする。この操作を時間周波数平面に対して行う。さらにそれを繰り返すことで、2つのマスクのうち正しいマスクに集約させることができる

50

。言い換えると、マスクの密集度を高め凝縮させることができる。

【 0 0 2 0 】

時間周波数平面上でマスクが分散している場合、スパース性が低くなり音声らしさが失われる。そのため、近接範囲のマスクの状況を参考にして、対象マスクの { 0 , 1 } を切り替えることによりできるだけ小さい範囲にマスクをまとめることができる。この処理を繰り返すことにより、時間周波数平面上でマスクが局所的に分布するようになり、マスクの密集度を高める前と比較してマスクのスパース性を向上させることができる。

【 0 0 2 1 】

次に、マスクの密集度を高める処理を具体的に説明する。

図 1 は、この発明の実施の形態 1 による音源分離装置のマスク再推定処理を示す説明図である。図 1 (a) はマスクの初期状態を示し、図 1 (b) はマスク再推定処理を 1 回行った状態を示し、図 1 (c) はマスク再推定処理を 2 回行った状態を示している。

図 1 で示す表の列方向は時間を変化させた領域であり、行方向は周波数を変化させた領域である。

図 1 (a) の初期状態において、マスク M_1 の領域 A は $M_1(t, f)=1$ であり、領域 B は $M_1(t+1, f)=1$ である。一方、マスク M_2 の領域 A' は $M_2(t, f)=0$ であり、領域 B' は $M_2(t+1, f)=0$ である。領域 A , A' では密集度が $\rho_1 < \rho_2$ であることから、 $M_1(t, f)=0$, $M_2(t, f)=1$ と再推定される。一方、領域 B , B' では密集度が $\rho_1 > \rho_2$ であることから、 $M_1(t, f)=1$, $M_2(t, f)=0$ と再推定される。

【 0 0 2 2 】

図 1 (b) の再推定処理 1 回目の状態において、マスク M_1 の領域 A は $M_1(t, f)=0$ であり、領域 B は $M_1(t+1, f)=1$ である。一方、マスク M_2 の領域 A' は $M_2(t, f)=1$ であり、領域 B' は $M_2(t+1, f)=0$ である。領域 A , A' では密集度が $\rho_1 < \rho_2$ であることから、 $M_1(t, f)=0$, $M_2(t, f)=1$ と再推定される。同様に、領域 B , B' においても密集度が $\rho_1 < \rho_2$ であることから、 $M_1(t+1, f)=0$, $M_2(t+1, f)=1$ と再推定される。

【 0 0 2 3 】

図 1 (c) の 2 回目の状態において、マスク M_1 の領域 A は $M_1(t, f)=0$ であり、領域 B は $M_1(t+1, f)=0$ である。一方、マスク M_2 の領域 A' は $M_2(t, f)=1$ であり、領域 B' は $M_2(t+1, f)=1$ である。領域 A , A' および領域 B , B' では共に密集度が $\rho_1 < \rho_2$ であり、上述した再推定結果である $M_1(t, f)=0$, $M_2(t, f)=1$ と変化はない。

【 0 0 2 4 】

図 1 (a) で示した初期状態と、図 1 (c) で示した 2 回の再推定処理を行った状態との密度 ρ により、マスク M_1 , M_2 の密集度が高まったと判定される。図 1 の例では、再推定処理を 2 回行う構成を示したが、再推定処理を行う回数はあらかじめ設定しておいてもよいし、密集度 ρ の変化が閾値以下になった時に処理を終了するように構成してもよい。

また、図 1 の例では領域 A , A' , B , B' に対して再推定処理を行う構成を示したが、図 1 で示したその他全ての領域を構成する全ての要素に対して近接要素の影響を勘案して再推定処理を行う。

【 0 0 2 5 】

図 2 は、この発明の実施の形態 1 による音源分離装置の構成を示すブロック図である。

音源分離装置 10 は、到来方向算出部 1、マスク作成部 2、マスク再推定部 3 およびマスク部 4 で構成されている。

到来方向算出部 1 は、第 1 のマイク 20 および第 2 のマイク 30 でそれぞれ観測された 2 つの混合音声信号の時間周波数領域 (t , f) におけるスペクトルから、第 1 のマイク 20 および第 2 のマイク 30 からの音波の到来角 θ_1 , θ_2 を算出する。到来方向 θ の算出は、上述した式 (1) および式 (2) を用いて行われる。

【 0 0 2 6 】

マスク作成部 2 は、到来方向算出部 1 が算出した音波の到来方向 θ をクラスタリングする、または音波の到来方向を異なる手法で推定した結果を取得することにより、第 1 のマ

10

20

30

40

50

イク 20 の時間周波数領域 (t, f) のマスク M_1 および第 2 のマイク 30 の時間周波数領域 (t, f) のマスク M_2 を作成する。例えば、到来角 θ_1 の第 1 のマイク 20 からの第 1 の音声信号に対するマスク M_1 の時間周波数領域 $M_1(t, f)$ は、上述した式 (3) で示したように推定される。

【0027】

マスク再推定部 3 のマスク集約処理部 3a は、上述した式 (5) の基準に従って、マスク間で推定誤りを解消するため、例えば $M_1(t, f) = 1$ となる時間周波数 (t, f) 領域を所定の範囲内にまとめるようにマスク M_1, M_2 の再推定処理を行い、マスクの集約を行う。マスク M_1, M_2 の再推定処理として、図 1 で示した再推定処理を適用する。再推定されたマスク M_1', M_2' は、マスク部 4 に出力される。

10

【0028】

マスク部 4 は、マスク再推定部 3 から入力されたマスク M_1, M_2 を用いて上述した式 (4) に基づいて、第 1 のマイク 20 および第 2 のマイク 30 で観測された音声スペクトル p_1, p_2 をマスクし、分離音声スペクトル p_1', p_2' を得る。

【0029】

以上のように、この実施の形態 1 によれば、マスク間での推定誤りを解消するために時間周波数平面でのスペクトルの局所的な連続性を考慮してマスクの密集度を高めて適切なマスクへの集約を行うマスク集約処理部 3a を備えたマスク再推定部 3 を備えるように構成したので、1 つの音源に寄与すべきマスクが分散していた場合に、1 つのマスクの集約することができ、分離音声に他の音源の音声は漏れ出すことが抑制され、対象話者の音声を明瞭化することができる。

20

【0030】

実施の形態 2 .

この実施の形態 2 では、複数の音源に寄与するマスクが、1 つの分離音声スペクトルに寄与する 1 つのマスクにまとめられている場合に、適切なマスクに分離する構成を示す。

図 3 は、この発明の実施の形態 2 による音源分離装置の構成を示すブロック図である。なお、以下では、実施の形態 1 による音源分離装置 10 の構成要素と同一または相当する部分には実施の形態 1 で使用した符号と同一の符号を付して説明を省略または簡略化する。

【0031】

この実施の形態 2 では、マスク再推定部 3 をマスク分離処理部 3b で構成している。マスク分離処理部 3b には、例えば異なる周波数 bin でのマスクの共起確率を用いる方法を適用する。ある周波数 f_1 と f_2 [Hz] の共起する確率を、学習データによりあらかじめ求めておく。共起確率 p_s としては、 $p_s(f_1, f_2) = N(f_1, f_2) / N_t$ を用いることができる。ここで N_t は総フレーム数、 $N(f_1, f_2)$ は $M(t, f_1) = M(t, f_2) = 1$ であったフレーム数である。これにより、共起しやすい周波数の組み合わせを知ることができる。例えば、音声は倍音構造を持つので f と $n * f$ は共起しやすい。 f は任意の周波数 [Hz]、 n は自然数である。ここで、 f と f_1 が共起する確率は低いとする (f_1 は f でない任意の周波数 [Hz])。ここでマスク M_1 において $M_1(t, f) = M_1(t, f_1) = 1$ であり、マスク M_2 において $M_2(t, f_1) = 0$ 且つ $M_2(t, n * f_1) = 1$ であった場合 $M_1(t, f_1) = 0$ $M_2(t, f_1) = 1$ である確率が高いと考えられる。このように倍音構造の利用では、 F_0 に対してその倍音成分のマスクを確認して音声らしさを推定することによりマスクを再度推定する。

30

40

【0032】

以上のようにこの実施の形態 2 によれば、マスク分離処理部 3b が共起しやすい周波数の組み合わせに注目して音声らしさを推定するように構成したので、複数の音源に寄与するマスクを適切なマスクに分離することができる。これにより、本来複数の音源に寄与すべきマスクが 1 つにまとめられていた場合に、適切なマスクに分離されるため、分離音声の雑音が抑制され、対象話者の音声を明瞭化することができる。

【0033】

実施の形態 3 .

50

この実施の形態 3 では、実施の形態 2 で再推定されたマスクが、再推定前のマスクよりも妥当であるか否か音声モデルを用いて検証する構成を示す。

図 4 は、この発明の実施の形態 3 による音源分離装置の構成を示すブロック図である。なお、以下では、実施の形態 2 による音源分離装置 10 の構成要素と同一または相当する部分には実施の形態 1 で使用した符号と同一の符号を付して説明を省略または簡略化する。

【0034】

実施の形態 2 で示した音源分離装置 10 に対して、マスク部 4 の後段に尤度算出部 5、音声モデル記憶部 6 およびマスク選択部 7 を追加して設けている。

上述した実施の形態 2 の処理を行うことにより、マスク作成部 2 が作成した元のマスク M_1 、 M_2 と、マスク再推定部 3 のマスク分離処理部 3b により再推定されたマスク M_1' 、 M_2' の 2 通りのマスクが得られる。

マスク部 4 は、マスク再推定部 3 から入力されたマスク M_1' 、 M_2' を用いて上述した式 (4) に基づいて、第 1 のマイク 20 および第 2 のマイク 30 で観測された音声スペクトル p_1 、 p_2 をマスキングし、分離音声スペクトル p_1' 、 p_2' を得る。さらにマスク部 4 は、マスク作成部 2 が作成した元のマスク M_1 、 M_2 を用いて、第 1 のマイク 20 および第 2 のマイク 30 で観測された音声スペクトル p_1 、 p_2 をマスキングし、音声スペクトル p_1'' 、 p_2'' を得る。

【0035】

これら全ての分離音声スペクトル p_1' 、 p_2' 、 p_1'' 、 p_2'' について、尤度算出部 5 が音声モデル記憶部 6 に記憶された音声モデルに対する尤度をフレーム単位で計算する。通常、複数の話者のそれぞれの発話内容は異なるので、異なる話者の音声が入混在した場合には、異なる話者の音声が入混在したスペクトルよりも、単一の話者による音声のスペクトルの方が音声モデルに対する尤度が高くなり音声らしいと判断されることになる。例えば、以下の式 (6) の GMM (Gaussian Mixture Model) によりモデル化された音声のモデル中の最大尤度を求めることで、音声らしさを判断することができる。

$$p(x) = \sum_{k=1}^K \pi_k N(\mu_k, \Sigma_k) \quad \dots \quad (6)$$

式 (6) において、 N は平均 μ_k 、共分散 Σ_k 、混合率 π_k の正規分布である。

【0036】

マスク選択部 7 は、尤度算出部 5 が算出した尤度を参照し、分離音声スペクトル p_1' 、 p_2' 、 p_1'' 、 p_2'' のうち最も音声らしい組み合わせを選択し、選択した分離音声スペクトルに対応したマスクを選択する。これにより、音声らしさの高い分離音声を出力するマスクを選択することができる。マスク選択部 7 が選択したマスクを用いて、再度音声スペクトル p_1 、 p_2 をマスキングしてもよいが、マスク部 4 の処理により分離音声スペクトル p_1' 、 p_2' 、 p_1'' 、 p_2'' が既に得られているので、対応するマスクの分離音声スペクトルを選択して最終的な分離音声スペクトルを得ることができる。

尤度算出部 5 およびマスク選択部 7 を備えたことにより、元のマスク M_1 、 M_2 と再推定されたマスク M_1' 、 M_2' のうち音声らしさの高いマスクを選択することができる。

【0037】

以上のように、この実施の形態 3 によれば、マスク分離処理部 3b が再推定したマスクおよびマスク作成部 2 が作成したマスクを用いてマスキングされた分離音声スペクトル p_1' 、 p_2' 、 p_1'' 、 p_2'' の尤度を計算する尤度算出部 5 と、算出された尤度に基づいて音声らしい組み合わせとなる分離音声スペクトルに対応したマスクを選択するマスク選択部 7 とを備えるように構成したので、マスク分離処理部 3b によるマスクの分離の誤りを検出し、適切なマスクの選択を行うことができる。

【0038】

10

20

30

40

50

実施の形態 4 .

上述したように、TDOAによって作られたマスクは低周波数域と高周波数域において性能が低い。低周波数域では位相の変化が小さいため、誤差が生じるためである。また高周波数域でもマイクの間隔よりも短い波長の音波が到来した場合には、位相が2回転したものと区別がつかない空間的エイリアシングの影響で推定精度が低くなる。この実施の形態4では、信頼性の低い周波数域の分離結果を、音声モデルを用いて補正する構成を示す。

【0039】

図5は、この発明の実施の形態4による音源分離装置の構成を示すブロック図である。なお、以下では、実施の形態3による音源分離装置10の構成要素と同一または相当する部分には実施の形態3で使用した符号と同一の符号を付して説明を省略または簡略化する。

この実施の形態4では、マスク再推定部3をマスク交叉部3cで構成している。マスク交叉部3cは、マスクの性能が低い領域（以下、低信頼領域と称する）においてパーミュテーションが起こっているものとして、マスク作成部2が作成したマスクをそれぞれ交叉させて得られるマスクの組み合わせを生成する。マスク部4は、マスク交叉部3cが生成したマスクの組み合わせを用いて、上述した式(4)に基づいて、第1のマイク20および第2のマイク30で観測された音声スペクトル p_1 、 p_2 をマスクングし、分離音声スペクトルを得る。

【0040】

次に、具体例を挙げながら実施の形態4の音源分離装置10の処理内容を説明する。以下では、高周波数域での場合を例に説明を行うが、低周波数域であっても同様に適用することができる。

マイクアレイを用いたTDOAの精度は、音波の半波長がマイク間隔以下になると低下する。例えば6cmの間隔のアレイを用いた場合、3kHz以上で空間的エイリアシングが起こることになる。図6は、16kHzサンプリングでの波形とスペクトログラムを示す図である。図6において0Hz~4kHzを高信頼領域、4kHz~8kHzを低信頼領域とする。まず低信頼領域をいくつかの帯域に分割する。例えば、低信頼領域である4kHz~8kHzを、4kHz~6kHzと6kHz~8kHzの2つの領域に分割した場合を想定する。

【0041】

ここで、マスク作成部2はマスク M_1 、 M_2 を作成する。低信頼領域ではパーミュテーションが起こっているものとして、マスク交叉部3cは以下に示す(a)から(d)に示す4通りのマスクのかけ方の組み合わせを生成する。

(a) ($M_1 - M_1 - M_1$, $M_2 - M_2 - M_2$)

(b) ($M_1 - M_1 - M_2$, $M_2 - M_2 - M_1$)

(c) ($M_1 - M_2 - M_1$, $M_2 - M_1 - M_2$)

(d) ($M_1 - M_2 - M_2$, $M_2 - M_1 - M_1$)

【0042】

マスク部4は、上述の(a)~(d)の4通りのマスクのかけ方を用いて8つの分離音声スペクトルを作成する。

そこで、この4通りのマスクのかけ方を用いて8つの分離音声スペクトルを作成する。例えば(a)のマスクのかけ方により2つの分離音声スペクトルが生成されるので、それぞれ p_{a-1} 、 p_{a-2} と呼ぶ。上述した実施の形態3と比較して分離音声スペクトルの生成数が増加する。これは、低信頼領域を2つの領域に分割しているためである。

【0043】

尤度算出部5は、マスク部4が作成した8つの分離音声スペクトルに対して、音声モデル記憶部6に記憶された音声モデルに対する尤度をフレーム単位で計算する。実施の形態3と同様に、異なる話者の音声混在した分離音声スペクトルよりも、単一の話者による音声スペクトルの方が音声らしいと判断されることになる。

10

20

30

40

50

音声モデルとしては、例えばモノフォンやトライフォンといった単位での音声のGMMが考えられる。上述した式(6)で示したGMMによりモデル化された音声のモデル中の最大尤度を求めることで、音声らしさを判断することができる。

【0044】

その他にもLPC係数などスペクトル包絡の滑らかさなどを基準とすることもできる。分離音声スペクトル p_{a-1} に対するモデル中の最大尤度と、分離音声スペクトル p_{a-2} に対するモデル中の最大尤度を加算したものを $L_{(1)}$ とする。同様に(b)から(d)のマスクに対する分離音声スペクトルに対する $L_{(2)}$ 、 $L_{(3)}$ 、 $L_{(4)}$ を算出し、最も尤度の高いマスクの組み合わせを選択することで高信頼領域と接続のよい分離音声信号を選び出すことができる。

10

【0045】

マスク選択部7は、尤度算出部5が算出した尤度を参照し、最も音声らしい組み合わせの分離音声スペクトルを選択し、選択した分離音声スペクトルに対応したマスクを選択する。これにより、高信頼領域の情報を活用すると共に、低信頼領域のパーミュテーションを解決することができる。

【0046】

以上のように、この実施の形態4によれば、低信頼領域において、マスク作成部2が作成したマスクをそれぞれ交叉させて得られるマスクの組み合わせを生成するマスク交叉部3cと、マスク交叉部3cが生成したマスクの組み合わせに基づいて音声スペクトルをマスクングして分離音声スペクトルを生成するマスク部4と、生成された全ての分離音声スペクトルの尤度を計算する尤度算出部5と、算出された尤度に基づいて最も音声らしい組み合わせとなる分離音声スペクトルに対応したマスクを選択するマスク選択部7とを備えるように構成したので、音声らしさの高いマスクを選択することができる。また、高信頼領域と接続のよい分離音声信号を選び出すことができると共に、低信頼領域のパーミュテーションを解決することができる。さらにマスク作成部2が作成したマスクの分離の誤りを検証することができ、適切なマスクを選択することができる。

20

【0047】

実施の形態5

マスクの再推定処理には、様々な方法が適用可能であり、パラメータの調整の余地もある。また、音素や話者によって分離性能の高いマスク推定方法が異なる場合も存在する。そこで、この実施の形態5では、マスク再推定部3がマスクの再推定を行う複数の構成を備える例を示す。

30

【0048】

図7は、この発明の実施の形態6の音源分離装置の構成を示すブロック図である。図7の例では、実施の形態1で示したマスク集約処理部3aおよび実施の形態2で示したマスク分離処理部3bを用いてマスク再推定部3を構成している。マスク部4は、マスク集約処理部3aおよびマスク分離処理部3bにより再推定されたマスクを用いて、音声スペクトル p_1 、 p_2 をマスクングする。なお、尤度算出部5、音声モデル記憶部6およびマスク選択部7の動作は上述した実施の形態3および実施の形態4と同一であるため、説明を省略する。

40

【0049】

以上のように、この実施の形態5によれば、マスク作成部2が作成したマスクの再推定処理を行うマスク集約処理部3aおよびマスク分離処理部3bを備えるように構成したので、複数のマスク再推定方法により、密集度が高められたマスク、あるいは適切な分離が行われたマスクから、最も音声らしい分離スペクトルが得られるマスクを選択することができる。最適なマスクの再推定処理を選択することができる。

【0050】

なお、上述した実施の形態1から実施の形態5では、TDOAに着目した構成を示したが、TDOA以外のその他の手法に関しても、マスクを作成可能であれば、本願発明の構成を適用することができる。

50

【0051】

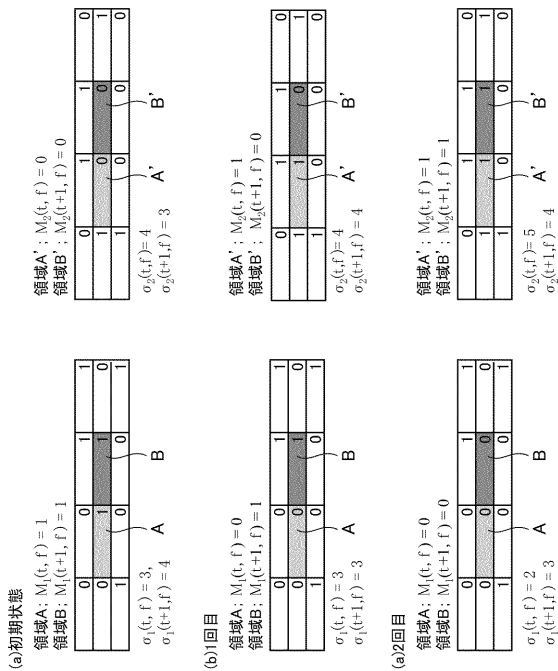
なお、本願発明はその発明の範囲内において、各実施の形態の自由な組み合わせ、あるいは各実施の形態の任意の構成要素の変形、もしくは各実施の形態において任意の構成要素の省略が可能である。

【符号の説明】

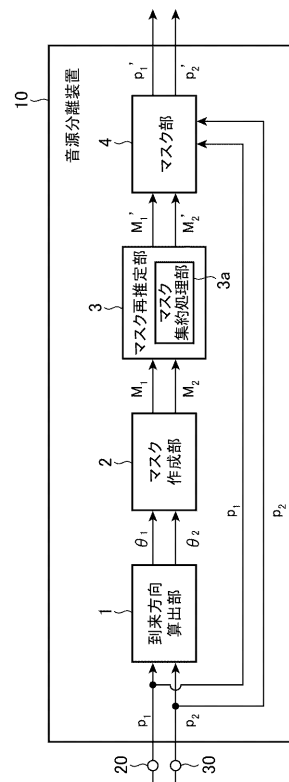
【0052】

1 到来方向算出部、2 マスク作成部、3 マスク再推定部、3 a マスク集約処理部、3 b マスク分離処理部、3 c マスク交叉部、4 マスク部、5 尤度算出部、6 音声モデル記憶部、7 マスク選択部、10 音源分離装置、20 第1のマイク、30 第2のマイク。

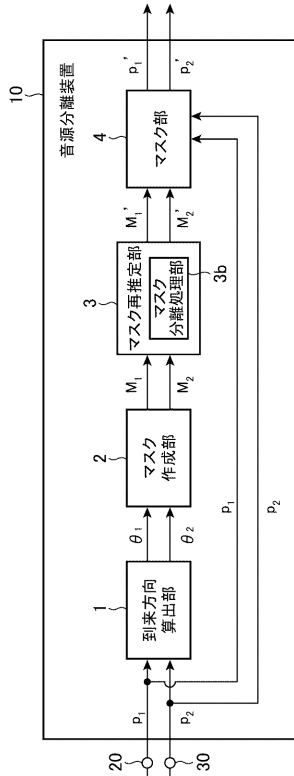
【図1】



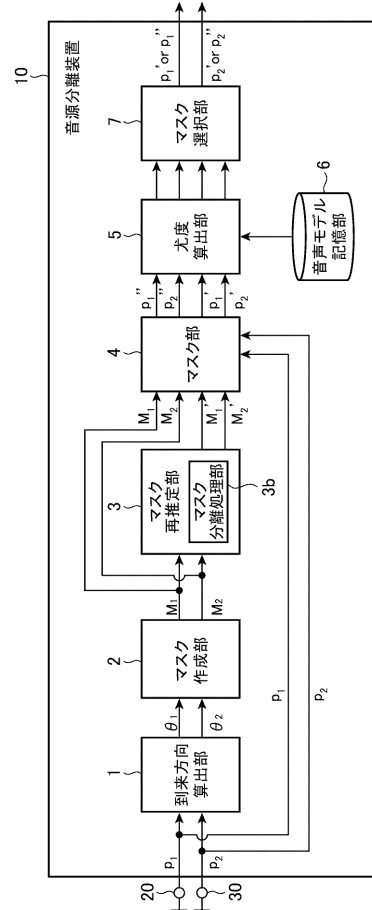
【図2】



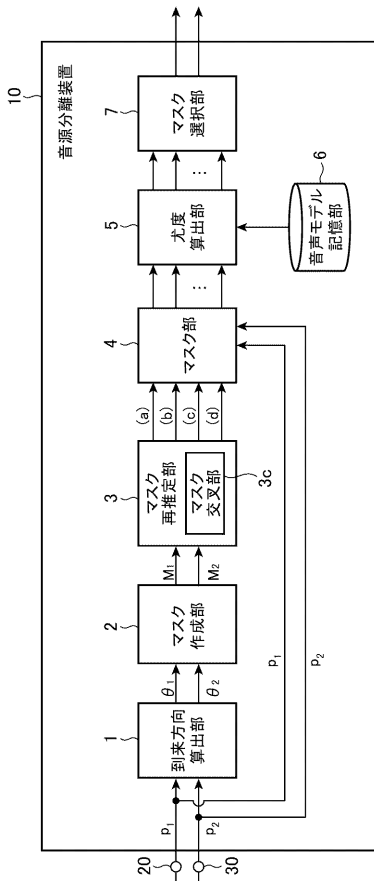
【 図 3 】



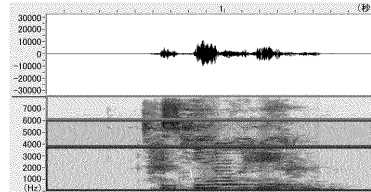
【 図 4 】



【 図 5 】



【 図 6 】



【 図 7 】

