

(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2016-156938

(P2016-156938A)

(43) 公開日 平成28年9月1日(2016.9.1)

(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 21/028 B	
G 1 O L 21/0272 (2013.01)	G 1 O L 21/0272 1 O O Z	
G 1 O L 21/0308 (2013.01)	G 1 O L 21/0308 Z	

審査請求 未請求 請求項の数 13 O L (全 17 頁)

(21) 出願番号 (22) 出願日 申請有り	特願2015-34339 (P2015-34339) 平成27年2月24日 (2015.2.24)	(71) 出願人 504132272 国立大学法人京都大学 京都府京都市左京区吉田本町36番地1 (74) 代理人 100091443 弁理士 西浦 ▲嗣▼晴 (72) 発明者 池宮 由楽 京都府京都市左京区吉田本町36番地1 国立大学法人京都大学内 (72) 発明者 吉井 和佳 京都府京都市左京区吉田本町36番地1 国立大学法人京都大学内 (72) 発明者 糸山 克寿 京都府京都市左京区吉田本町36番地1 国立大学法人京都大学内
-----------------------------------	--	--

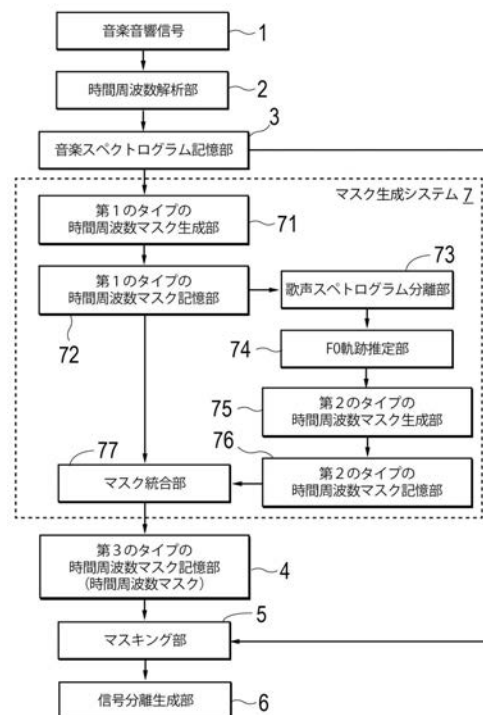
(54) 【発明の名称】 歌声信号分離方法及びシステム

(57) 【要約】

【課題】 歌声信号と伴奏音信号とを含む音楽音響信号から歌声信号を分離する精度を従来よりも改善することができる歌声信号分離方法及びシステムを提供することにある。

【解決手段】 調波構造から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、出現可能性から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン及び前記出現可能性から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、調波構造から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン音楽スペクトログラムからマスクングする機能を有する時間周波数マスクを準備する。次にこの時間周波数マスクを音楽スペクトログラムに適用して分離用歌声スペクトログラムを生成する。そして分離用歌声スペクトログラムに基づいて歌声信号を分離生成する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

歌声信号と伴奏音信号とを含む音楽音響信号から前記歌声信号を分離する歌声信号分離方法であって、

前記音楽音響信号を、時間周波数解析を行って音楽スペクトログラムに変換する変換ステップと、

調波構造から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、出現可能性から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン及び前記出現可能性から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、前記調波構造から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビンを前記音楽スペクトログラムからマスキングする機能を有する時間周波数マスクを準備するマスク準備ステップと、

前記時間周波数マスクを前記音楽スペクトログラムに適用して分離用歌声スペクトログラムを生成するマスキングステップと、

前記分離用歌声スペクトログラムに基づいて前記歌声信号を分離生成する分離生成ステップとからなることを特徴とする歌声信号分離方法。

【請求項 2】

前記マスク準備ステップでは、

前記音楽スペクトログラムを、ロバスト主成分分析を用いて低ランク行列とスパース行列とに分解し、

前記低ランク行列と前記スパース行列の比較に基づいて、前記音楽スペクトログラムを歌声が出現している可能性が高いスペクトルビンを含むスパース行列と歌声が出現している可能性が低いスペクトルビンを含む低ランク行列とに分離する機能を有する第 1 のタイプの時間周波数マスクを生成し、

第 1 のタイプの時間周波数マスクを前記音楽スペクトログラムに適用して歌声が出現している可能性が高いスペクトルビンを含む歌声スペクトログラムを分離し、

分離された前記歌声スペクトログラムに対して歌声基本周波数 F_0 を推定して歌声基本周波数 F_0 軌跡を推定し、

前記歌声基本周波数 F_0 軌跡に基づいて、前記歌声基本周波数 F_0 と倍音周辺以外のスペクトルビンをマスキングする機能を有する第 2 のタイプの時間周波数マスクを生成し、

前記第 1 のタイプの時間周波数マスクと前記第 2 のタイプの時間周波数マスクとを統合して、前記第 2 のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、前記第 1 のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビン及び前記第 1 のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、前記第 2 のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビンをマスキングする機能を有する第 3 のタイプの時間周波数マスクを前記時間周波数マスクとして準備することを特徴とする請求項 1 に記載の歌声信号分離方法。

【請求項 3】

前記第 1 のタイプの時間周波数マスクと前記第 2 のタイプの時間周波数マスクとの統合とは、前記第 1 のタイプの時間周波数マスクの選択領域と前記第 2 のタイプの第 2 の時間周波数マスクの選択領域との論理積をとることである請求項 2 に記載の歌声信号分離方法。

【請求項 4】

前記第 1 のタイプの時間周波数マスクと前記第 2 のタイプの時間周波数マスクとの統合とは、前記第 1 のタイプの時間周波数マスクの選択領域と前記第 2 のタイプの時間周波数マスクの選択領域との論理積をとって、仮統合時間周波数マスクを生成し、前記仮統合時間周波数マスクから、歌が無い区間を推定して、推定された時間フレームの全要素を 0 にすることにより前記第 3 のタイプの時間周波数マスクとすることである請求項 1 に記載の歌声信号分離方法。

10

20

30

40

50

【請求項 5】

前記第 1 のタイプの時間周波数マスクと前記第 2 のタイプの時間周波数マスクとの統合とは、前記第 1 のタイプの時間周波数マスクの選択領域と前記第 2 のタイプの時間周波数マスクの選択領域との論理積をとって、仮統合時間周波数マスクを生成し、前記仮統合時間周波数マスクから、歌が無い区間を推定して、推定された時間フレームの全要素を 0 にし、且つ前記第 1 のタイプの時間周波数マスクから子音を通過させる要素を得て該要素を前記仮統合時間周波数マスクに反映することである請求項 1 に記載の歌声信号分離方法。

【請求項 6】

前記第 1 のタイプの時間周波数マスク、前記第 2 のタイプの時間周波数マスク及び第 3 のタイプの時間周波数マスクは、それぞれバイナリマスクである請求項 3, 4 または 5 に記載の歌声信号分離方法。

10

【請求項 7】

前記分離生成ステップでは、分離用歌声スペクトログラムを時間領域に逆変換することにより歌声信号を分離生成することを特徴とする請求項 1 に記載の歌声信号分離方法。

【請求項 8】

前記各ステップを 1 以上のプロセッサで実施することを特徴とする請求項 1 乃至 7 に記載の歌声信号分離方法。

【請求項 9】

歌声信号と伴奏音信号とを含む音楽音響信号から前記歌声信号を分離する歌声信号分離システムであって、

20

前記音楽音響信号を、時間周波数解析を行って音楽スペクトログラムに変換する時間周波数解析部と、

調波構造から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、出現可能性から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン及び前記出現可能性から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、前記調波構造から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビンを前記音楽スペクトログラムからマスクする機能を有する時間周波数マスクを用いて、前記時間周波数マスクを前記音楽スペクトログラムに適用して分離用歌声スペクトログラムを生成するマスク部と、

前記分離用歌声スペクトログラムに基づいて前記歌声信号を分離生成する信号分離生成部とからなることを特徴とする歌声信号分離システム。

30

【請求項 10】

前記音楽スペクトログラムを、ロバスト主成分分析を用いて低ランク行列とスパース行列とに分解し、前記低ランク行列と前記スパース行列の比較に基づいて、前記音楽スペクトログラムを歌声が出現している可能性が高いスペクトルビンを含むスパース行列と歌声が出現している可能性が低いスペクトルビンを含む低ランク行列とに分離する機能を有する第 1 のタイプの時間周波数マスクを生成する第 1 のタイプの時間周波数マスク生成部と、

前記第 1 のタイプの時間周波数マスクを記憶する第 1 のタイプの時間周波数マスク記憶部と、

40

前記第 1 のタイプの時間周波数マスクを前記音楽スペクトログラムに適用して歌声が出現している可能性が高いスペクトルビンを含む歌声スペクトログラムを分離する歌声スペクトログラム分離部と、

分離された前記歌声スペクトログラムに対して歌声基本周波数 F_0 を推定して歌声基本周波数 F_0 軌跡を推定する F_0 軌跡推定部と、

前記歌声基本周波数 F_0 軌跡に基づいて作成されて、前記歌声基本周波数 F_0 と倍音周波数以外のスペクトルビンをマスクする機能を有する第 2 のタイプの時間周波数マスクを生成する第 2 のタイプの時間周波数マスク生成部と、

前記第 2 のタイプの時間周波数マスクを記憶する第 2 のタイプの時間周波数マスク記憶部と、

50

前記第 1 のタイプの時間周波数マスクと前記第 2 のタイプの時間周波数マスクとを統合して作成された、前記第 2 のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、前記第 1 のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビン及び前記第 1 のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、前記第 2 のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビンをもスキミングする機能を有する第 3 のタイプの時間周波数マスクを前記時間周波数マスクとして統合するマスク統合部とからなるマスク生成システムによって、前記時間周波数マスクが生成されたものである請求項 9 に記載の歌声信号分離システム。

10

【請求項 11】

前記第 1 のタイプの時間周波数マスク、前記第 2 のタイプの時間周波数マスク及び第 3 のタイプの時間周波数マスクは、それぞれバイナリマスクである請求項 10 に記載の歌声信号分離システム。

【請求項 12】

前記信号分離生成部は、分離用歌声スペクトログラムを時間領域に逆変換することにより歌声信号を分離生成することを特徴とする請求項 9 に記載の歌声信号分離システム。

【請求項 13】

上記構成要件は、1 以上のプロセッサとメモリによって実現されている請求項 9 乃至 12 のいずれか 1 項に記載の歌声信号分離システム。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、歌声信号と伴奏音信号とを含む音楽音響信号から歌声信号を分離する歌声信号分離方法及びシステムに関するものである。

【背景技術】

【0002】

非特許文献 1 [Combining Modeling of Singing Voice and Background Music for Automatic Separation of Musical Mixtures, ISMIR (2013)] には、歌声と伴奏を分離する従来の技術の一例が開示されている。

30

【0003】

例えば、非特許文献 2 [Mixture of Gaussian Process Experts for Predicting Sung Melodic Contour with Expressive Dynamic Fluctuations, ICASSP (2014)] は、歌声の F0 軌跡を不連続な楽譜成分と微細な変動成分の重ね合わせとして表現する確率モデルを用いて、任意の楽譜から歌声の F0 軌跡を生成する手法を提案している。同様のモデルは、非特許文献 3 [混合ガウス過程に基づく歌声音量軌跡の生成過程モデル, 情処研報 (2013)] において、歌声の音量軌跡に対しても適用されている。

【先行技術文献】

【非特許文献】

【0004】

【非特許文献 1】Rafii, Z., Germain, F. G., Sun, D. L., and Mysore, G. J.: Combining Modeling of Singing Voice and Background Music for Automatic Separation of Musical Mixtures, ISMIR (2013)

40

【非特許文献 2】Ohishi, Y., Mochihashi, D., Kameoka, H., and Kashino, K.: Mixture of Gaussian Process Experts for Predicting Sung Melodic Contour with Expressive Dynamic Fluctuations, ICASSP (2014)

【非特許文献 3】大石康智, 持橋大地, 亀岡弘和, 柏野邦夫: 混合ガウス過程に基づく歌声音量軌跡の生成過程モデル, 情処研報 (2013)

【発明の概要】

【発明が解決しようとする課題】

50

【0005】

混合音中の歌声に対する編集システムを実現するには、高精度な歌声・伴奏音分離と歌声のF0推定が必要である。しかしながら従来の技術では、両タスクの相互依存性を考慮して、精度を一挙に改善することができるものはなかった。

【0006】

本発明の目的は、歌声信号と伴奏音信号とを含む音楽音響信号から歌声信号を分離する精度を従来よりも改善することができる歌声信号分離方法及びシステムを提供することにある。

【0007】

本発明の他の目的は、高精度な歌声・伴奏音分離と歌声のF0推定の相互依存性を考慮して、歌声信号と伴奏音信号とを含む音楽音響信号から歌声信号を分離する精度を一挙に改善することができる歌声信号分離方法及びシステムを提供することにある。

【課題を解決するための手段】

【0008】

本発明は、歌声信号と伴奏音信号とを含む音楽音響信号から歌声信号を分離する歌声信号分離方法及びシステムを改良の対象とする。本発明の方法では、まず音楽音響信号を、時間周波数解析を行って音楽スペクトログラムに変換する(変換ステップ)。また調波構造から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、出現可能性から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン及び前記出現可能性から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、調波構造から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビンを音楽スペクトログラムからマスキングする機能を有する時間周波数マスクを準備する(マスク準備ステップ)。

【0009】

次にこの時間周波数マスクを音楽スペクトログラムに適用して分離用歌声スペクトログラムを生成する(マスキングステップ)。そして分離用歌声スペクトログラムに基づいて歌声信号を分離生成する(分離生成ステップ)。上記のような時間周波数マスクを用いると、歌声信号と伴奏音信号とを含む音楽音響信号から歌声信号を従来よりも精度よく分離することができる。

【0010】

マスク準備ステップでは、具体的には、次のようにして時間周波数マスクを準備する。まず音楽スペクトログラムを、ロバスト主成分分析を用いて低ランク行列とスパース行列とに分解する。次に低ランク行列とスパース行列の比較に基づいて、音楽スペクトログラムを歌声が出現している可能性が高いスペクトルビンを含むスパース行列と歌声が出現している可能性が低いスペクトルビンを含む低ランク行列とに分離する機能を有する第1のタイプの時間周波数マスクを生成する。次に第1のタイプの時間周波数マスクを音楽スペクトログラムに適用して歌声が出現している可能性が高いスペクトルビンを含む歌声スペクトログラムを分離する。そして分離された歌声スペクトログラムに対して歌声基本周波数F0を推定して歌声基本周波数F0軌跡を推定する。次に歌声基本周波数F0軌跡に基づいて、歌声基本周波数F0と倍音周辺以外のスペクトルビンをマスキングする機能を有する第2のタイプの時間周波数マスクを生成する。ここで「歌声基本周波数F0と倍音周辺」とは、歌声基本周波数F0のピークとその倍音のピークを中心として、予め定めた周波数幅に入る周波数である。この周波数幅は、歌声基本周波数F0とその倍音のスペクトルの形状から自動的に定めることもできる。

【0011】

最後に、第1のタイプの時間周波数マスクと第2のタイプの時間周波数マスクとを統合して、第2のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、第1のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビン及び第1のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、第2のタイプの

10

20

30

40

50

時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビン进行マスクする機能を有する第3のタイプの時間周波数マスクを時間周波数マスクとして準備する。この具体的方法では、基本的にはロバスト主成分分析 (Robust Principal Component Analysis; RPCA) を用いてスペクトログラム上で歌声・伴奏音分離を行う。歌声のF0情報を用いれば、不要な伴奏音を抑制することができる。一方、混合音に対して歌声のF0推定を行うよりも、分離した歌声に対してF0推定を行う方がずっと容易である。

【0012】

第1のタイプの時間周波数マスクと第2のタイプの時間周波数マスクとの統合とは、両マスクの機能を優れた機能を併用可能にすることを意味し、例えば、第1のタイプの時間周波数マスクと第2のタイプの時間周波数マスクとの統合とは、第1のタイプの時間周波数マスクの選択領域と前記第2のタイプの第2の時間周波数マスクの選択領域との論理積をとることにより両マスクを統合することができる。

10

【0013】

また統合の他の例では、第1のタイプの時間周波数マスクの選択領域と第2のタイプの時間周波数マスクの選択領域との論理積をとって、仮統合時間周波数マスクを生成し、仮統合時間周波数マスクから、歌が無い区間を推定して、推定された時間フレームの全要素を0にすることにより第3のタイプの時間周波数マスクとすることができる。この時間周波数マスクでは、歌が無い区間を推定して、推定された時間フレームの全要素を0にするため、さらに分離精度を高めることができる。

20

【0014】

さらに統合の他の例では、第1のタイプの時間周波数マスクの選択領域と第2のタイプの時間周波数マスクの選択領域との論理積をとって、仮統合時間周波数マスクを生成し、仮統合時間周波数マスクから、歌が無い区間を推定して、推定された時間フレームの全要素を0にし、且つ第1のタイプの時間周波数マスクから子音を通過させる要素を得て該要素を仮統合時間周波数マスクに反映する。このようにするとさらに分離精度を高めることができる。

【0015】

なお第1のタイプの時間周波数マスク、第2のタイプの時間周波数マスク及び第3のタイプの時間周波数マスクは、それぞれバイナリマスクであるのが好ましい。バイナリマスクを用いると、1と0の組み合わせによりマスクが構成されるため、歌声と伴奏がくっきり分かれ、伴奏音側に歌声が残る可能性はほとんどなくなる。

30

【0016】

分離生成ステップでは、分離用歌声スペクトログラムを時間領域に逆変換することにより歌声信号を分離生成することができる。そして各ステップは、1以上のプロセッサで実施することができる。

【0017】

本発明の方法を実施する本発明の歌声信号分離システムは、時間周波数解析部と、マスク部と、信号分離生成部とから構成される。

【0018】

時間周波数解析部は、音楽音響信号を、時間周波数解析を行って音楽スペクトログラムに変換する。マスク部は、調波構造から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、出現可能性から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン及び出現可能性から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、調波構造から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン进行音楽スペクトログラムからマスクする機能を有する時間周波数マスクを用いて、時間周波数マスクを音楽スペクトログラムに適用して分離用歌声スペクトログラムを生成する。そして信号分離生成部は、分離用歌声スペクトログラムに基づいて歌声信号を分離生成する。

40

【0019】

50

時間周波数マスクはマスク生成システムによって生成される。マスク生成システムは、第1のタイプの時間周波数マスク生成部と、第1のタイプの時間周波数マスク記憶部と、歌声スペクトログラム分離部と、F0軌跡推定部と、第2のタイプの時間周波数マスク生成部と、第2のタイプの時間周波数マスク記憶部と、マスク統合部とから構成される。

【0020】

第1のタイプの時間周波数マスク生成部は、音楽スペクトログラムを、ロバスト主成分分析を用いて低ランク行列とスパース行列とに分解し、低ランク行列とスパース行列の比較に基づいて、音楽スペクトログラムを歌声が出現している可能性が高いスペクトルビンを含むスパース行列と歌声が出現している可能性が低いスペクトルビンを含む低ランク行列とに分離する機能を有する第1のタイプの時間周波数マスクを生成する。第1のタイプの時間周波数マスク記憶部は、第1のタイプの時間周波数マスクを記憶する。歌声スペクトログラム分離部は、第1のタイプの時間周波数マスクを音楽スペクトログラムに適用して歌声が出現している可能性が高いスペクトルビンを含む歌声スペクトログラムを分離する。F0軌跡推定部は、分離された歌声スペクトログラムに対して歌声基本周波数F0を推定して歌声基本周波数F0軌跡を推定する。第2のタイプの時間周波数マスク生成部は、歌声基本周波数F0軌跡に基づいて、歌声基本周波数F0と倍音周辺以外のスペクトルビンをマスクする機能を有する第2のタイプの時間周波数マスクを生成する。第2のタイプの時間周波数マスク記憶部は、第2のタイプの時間周波数マスクを記憶する。マスク統合部は、第1のタイプの時間周波数マスクと第2のタイプの時間周波数マスクとに基づき、第2のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、第1のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビン及び第1のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、第2のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビンをマスクする機能を有する第3のタイプの時間周波数マスクを時間周波数マスクとして統合する。

10

20

【0021】

信号分離生成部は、分離用歌声スペクトログラムを時間領域に逆変換することにより歌声信号を分離生成する。なお上記構成要件は、1以上のプロセッサとメモリによって実現するのが好ましい。

30

【図面の簡単な説明】

【0022】

【図1】本発明の歌声信号分離方法を実施する歌声信号分離システムの一例の構成を示すブロック図である。

【図2】図1の実施の形態の歌声信号分離システムをコンピュータ(1以上のプロセッサと1以上のメモリを含む)で実施する際に使用されるソフトウェアのアルゴリズムを示すフローチャートである。

【図3】入力音楽音曲信号を時間周波数解析して得る音楽スペクトログラムの一例を示す図である。

【図4】音楽スペクトログラムからロバスト主成分分析によりスパース行列と低ランク行列とに分析した結果の一例と、両行列の各要素の値を比較して得た第1のタイプの時間周波数分析マスクとしてバイナリマスクの例を示す図である。

40

【図5】図4の表示内容の理解を高めるために、音楽スペクトログラムの一部を拡大し、またスパース行列と低ランク行列の一部を拡大し、さらに第1のタイプの時間周波数分析マスクとしてバイナリマスクの一部を拡大した図を示している。

【図6】F0軌跡推定部によって推定された歌声基本周波数F0軌跡から前記第2のタイプの時間周波数マスク(バイナリマスク)を生成する過程の一例を示す図である。

【図7】第1のタイプの時間周波数マスク(ロバスト主成分分析によるバイナリマスク)と第2のタイプの時間周波数マスク(歌声基本周波数F0による調波構造のバイナリマスク)を統合する場合の一例を画像で示す図である。

50

【図 8】図 7 の画像の理解を高めるために、図 7 に示した複数の画像の一部をそれぞれ拡大して示す図である。

【図 9】マスキング部における処理を画像で示すための図である。

【図 10】図 9 の画像の理解を高めるために、図 9 に示した複数の画像の一部をそれぞれ拡大して示す図である。

【図 11】(A)乃至(D)は、マスキング部によるマスキング処理の状況を示す波形図である。

【図 12】歌声信号の再合成を説明するために用いる図である。

【図 13】マスクの統合の他の例を示す概念図である。

【図 14】マスクの統合のさらに他の例を示す概念図である。

10

【発明を実施するための形態】

【0023】

以下図面を参照して、本発明の歌声信号分離方法及びシステムの実施の形態の一例を詳細に説明する。図 1 は、本発明の歌声信号分離方法を実施する歌声信号分離システムの一例の構成を示すブロック図である。図 2 は、図 1 の実施の形態の歌声信号分離システムをコンピュータ(1以上のプロセッサと1以上のメモリを含む)で実施する際に使用されるソフトウェアのアルゴリズムを示すフローチャートである。

【0024】

本発明の方法を実施する本発明の歌声信号分離システムは、時間周波数解析部 2 と、音楽スペクトログラム記憶部 3 と、第 3 のタイプの時間周波数マスク記憶部 4 と、マスキング部 5 と、信号分離生成部 6 とから構成される。図 1 には、第 3 のタイプの時間周波数マスク記憶部 4 に記憶する時間周波数マスクを生成するためのマスク生成システム 7 も併せて記載してある。説明の都合上、本実施の形態の説明の途中でマスク生成システム 7 についても説明する。

20

【0025】

時間周波数解析部 2 は、歌声信号と伴奏音信号とを含む音楽音響信号 1 を、時間周波数解析を行って音楽スペクトログラム(行列)に変換する(ステップ S T 1)。まず、短時間フーリエ変換(Short-Term Fourier Transform; STFT)あるいは定 Q 変換を用いて入力音楽音響信号の時間周波数解析を行う。定 Q 変換については、「Schorkhuber, C. and Klapuri, A.: Constant-Q Transform Toolbox for Music Processing, SMC Conference (2010)」に詳しく記載されている。

30

【0026】

実用上、全時間サンプル n における対数スペクトルピッチを求めるのではなく、例えば 10 [msec] などの時間幅で切り出す。以後分かりやすさのため、時間インデックス、周波数インデックスをそれぞれ t, f とし、音楽スペクトログラムを $X(t, f)$ と記述する。図 3 には、入力音楽音曲信号を時間周波数解析して得る音楽スペクトログラムの一例を示している。

【0027】

第 3 のタイプの時間周波数マスク記憶部 4 には、マスク生成システム 7 で作成した時間周波数マスク(統合マスク)として第 3 のタイプの時間周波数マスクが記憶される。マスク生成システム 7 は、第 1 のタイプの時間周波数マスク生成部 7 1 と、第 1 のタイプの時間周波数マスク記憶部 7 2 と、歌声スペクトログラム分離部 7 3 と、F0 軌跡推定部 7 4 と、第 2 のタイプの時間周波数マスク生成部 7 5 と、第 2 のタイプの時間周波数マスク記憶部 7 6 と、マスク統合部 7 7 とから構成される。第 1 のタイプの時間周波数マスク生成部 7 1 は、音楽スペクトログラム記憶部 3 に記憶した音楽スペクトログラム中の音楽スペクトログラムを、ロバスト主成分分析を用いて低ランク行列とスパース行列とに分解し(ステップ S T 2)、低ランク行列とスパース行列の比較に基づいて、音楽スペクトログラムを歌声が出現している可能性が高いスペクトルピッチを含むスパース行列と歌声が出現している可能性が低いスペクトルピッチを含む低ランク行列とに分離する機能を有する第 1 のタイプの時間周波数マスクを生成する(ステップ S T 3)。そして第 1 のタイプの時間周

40

50

波数マスクは、第1のタイプの時間周波数マスク記憶部72記憶される。

【0028】

ロバスト主成分分析は、与えられた行列(2次元配列)を低ランク行列とスパース行列とに分解する手法であり、次式で定式化される。

【0029】

【数1】

$$\text{minimize } \|L\|_* + \lambda \|S\|_1 \quad (\text{subject to } L + S = X)$$

10

ここで、 X 、 L 、 S はそれぞれ入力行列、低ランク行列およびスパース行列であり、 $\|\cdot\|_*$ と $\|\cdot\|_1$ はそれぞれ核ノルムとL1ノルム、 λ は低ランク性とスパース性のトレードオフパラメータを表す。一般に時間変化するデータ集合などを入力とし、頻出する成分(各フレームで繰り返し現れる成分)が低ランク行列に、それ以外の成分(各フレームに稀にしか現れない成分)がスパース行列に分解される。

【0030】

音楽スペクトログラムを入力行列 X と見なしてロバスト主成分分析を適用すると、繰り返し演奏されるため何度も出現する伴奏音(ドラムやギター)のスペクトルピッチは低ランク行列 L へ、それ以外の歌声などの時間的な変動が大きいスペクトルピッチはスパース行列 S へ分解される。本実施の形態では、分析結果から第1のタイプの時間周波数分析マスクとしてバイナリマスクを作成する。

20

【0031】

【数2】

$$M_r(t, f) = \begin{cases} 1 & |S(t, f)| > |L(t, f)| \\ 0 & \text{otherwise} \end{cases}$$

このバイナリマスクからなる第1のタイプの時間周波数分析マスクを音楽スペクトログラム $X(t, f)$ へ適用することで歌声スペクトログラムが分離できる。

30

【0032】

なお図4には、音楽スペクトログラムからロバスト主成分分析により分析した結果のスパース行列(歌声)と低ランク行列(伴奏)とに分析した結果の一例と、両行列の各要素の値を比較して得た第1のタイプの時間周波数分析マスクとしてバイナリマスクの例を示している。図5は、図4の表示内容の理解を高めるために、音楽スペクトログラムの一部を拡大し、またスパース行列(歌声)と低ランク行列(伴奏)の一部を拡大し、さらに第1のタイプの時間周波数分析マスクとしてバイナリマスクの一部を拡大した図を示している。

【0033】

40

歌声スペクトログラム分離部73は、第1のタイプの時間周波数マスク(バイナリマスク)を音楽スペクトログラムに適用して歌声が出現している可能性が高いスペクトルピッチを含む歌声スペクトログラムを分離する(ステップST4)。F0軌跡推定部74は、分離された歌声スペクトログラムに対して歌声基本周波数F0を推定して歌声基本周波数F0軌跡を推定する(ステップST5)。

【0034】

具体的には、ロバスト主成分分析により分離された歌声スペクトログラム $X_s^{\text{Pca}}(t, f)$ から、Subharmonic Summation(SHS)を用いて歌声のF0軌跡を推定する。SHSについては、「Hermes, D. J.: Measurement of pitch by subharmonic summation, J. Acoust. Soc. Am., Vol. 83, No. 1, pp.257-264 (online), DOI: 10.1121/1.396427 (198

50

8)」に詳しく説明されている。SHSは計算コストの低さとノイズへの頑健性を兼ね備えた音高推定法であり、スペクトルピンの各周波数ピンについて、そのピンをF0であると仮定したときの倍音に対応する周波数ピンのパワーを重みつきで足し合わせることで、当該ピンにF0が存在する尤度を計算する。この音高尤度関数の計算は、対数周波数スケールでは以下で定式化される。

【0035】

【数3】

$$H(t, s) = \sum_{n=1}^N h_n P(t, s + 1200 \log_2 n),$$

10

ここで、 t 、 s はそれぞれ時間インデクスと対数周波数 [cents] を表し、 $P(t, s)$ は時間フレーム t 、周波数 s [cents]における入力スペクトログラムの振幅である。 N は足し合わせる倍音数、 h_n は各倍音の重み関数であり、本実施の形態ではそれぞれ15および 0.86^{n-1} とする。人間の聴覚特性の非線形性を考慮するため、SHSを適用する前に、入力スペクトルピンに対してA特性補正をかけるものとする。

【0036】

SHSによる音高尤度関数 $H(t, s)$ から歌声音高 $F(t)$ は以下の式で計算される。

【0037】

【数4】

20

$$F(t) = \arg \max_{c_l^{(t)} \leq s \leq c_h^{(t)}} H(t, s)$$

ここで、 $c_l^{(t)}$ 、 $c_h^{(t)}$ はそれぞれ、時間フレーム t における音高探索周波数範囲の下限と上限 ([cents]) である。

【0038】

図6は、F0軌跡推定部74によって推定された歌声基本周波数F0軌跡から第2のタイプの時間周波数マスク(バイナリマスク)を生成する過程の一例を示している。第2のタイプの時間周波数マスク生成部75は、歌声基本周波数F0軌跡に基づいて、歌声基本周波数F0と倍音周辺以外のスペクトルピンをマスクする機能を有する第2のタイプの時間周波数マスクを生成する(ステップST6)。生成された第2のタイプの時間周波数マスクは第2のタイプの時間周波数マスク記憶部76に記憶される。

30

【0039】

ロバスト主成分分析を用いた従来の歌声分離では、曲の一部しか現れないベースやドラム、メインボーカルと音高をずらして唱和するバックコーラスなども、歌声として分離されてしまう。歌声・伴奏音分離と歌声のF0推定は相互依存性をもっている。つまり、歌声のF0軌跡が与えられていれば、歌声分離に利用することができる一方、歌声が分離されていれば、そのF0軌跡を推定することは比較的容易である。そこでこの相補的な関係を利用した歌声分離のために入力音響信号に対して、統合マスク(第3のタイプの時間周波数マスク)を用いて、精密な歌声分離を行う。そこで、第2のタイプの時間周波数マスク生成部75は、歌声基本周波数F0軌跡を利用して、さらに精度の高い歌声分離を行うために、歌声基本周波数F0軌跡から、基本周波数(F0)と倍音周辺以外のパワーをマスクする調波マスクを第2のタイプの時間周波数として生成する。ここで「歌声基本周波数F0と倍音周辺」とは、歌声基本周波数F0のピークとその倍音のピークを中心として、予め定めた周波数幅に入る周波数である。この周波数幅は、歌声基本周波数F0とその倍音のスペクトルの形状から自動的に定めることもできる。

40

【0040】

50

【数 5】

$$M_h(t, f) = \begin{cases} 1 & \left[\begin{array}{l} H_t^h - \frac{w}{2} < C(f) < H_t^h + \frac{w}{2} \\ H_t^h = F_t + 1200 \log_2 h, 1 \leq h \leq H \end{array} \right. \\ 0 & \text{otherwise} \end{cases}$$

ここで、 F_t は時間フレーム t における F_0 [cents]、 $C(f)$ は周波数ビン f に対応する対数周波数 [cents]、 H は倍音数、 w は各倍音でマスクを取る幅 [cents]を示す。ロバスト主成分分析によるバイナリマスクと調波マスクを用いて、最終的な歌声と伴奏のスペクトログラム $X_s(t, f)$ 、 $X_m(t, f)$ はそれぞれ以下のように得られる。

10

【0041】

【数 6】

$$X_s(t, f) = M_r(t, f)M_h(t, f)X(t, f),$$

$$X_m(t, f) = X(t, f) - X_s(t, f)$$

マスク統合部77は、第1のタイプの時間周波数マスクと第2のタイプの時間周波数マスクとを統合して第3のタイプの時間周波数マスクを時間周波数マスク（統合マスク）として作成する（ステップS T7）。この第3のタイプの時間周波数マスクからなる時間周波数マスク（統合マスク）は、上位概念で言えば、調波構造から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、出現可能性から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン及び出現可能性から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、調波構造から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビンを音楽スペクトログラムからマスクングする機能を有するものである。より具体的に言えば、第3のタイプの時間周波数マスク（統合マスク）は、第2のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、第1のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビン及び第1のタイプの時間周波数マスクでは歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、第2のタイプの時間周波数マスクでは歌声を含むスペクトルビンではないと判断されるスペクトルビンをマスクングする機能を有する。第3のタイプの時間周波数マスクを時間周波数マスク（統合マスク）は、第3のタイプの時間周波数マスク記憶部4に記憶される。

20

30

【0042】

図7は、第1のタイプの時間周波数マスク（ロバスト主成分分析によるバイナリマスク）と第2のタイプの時間周波数マスク（歌声基本周波数 F_0 による調波構造のバイナリマスク）を統合する場合の一例を画像で示している。図8は、図7の画像の理解を高めるために、図7に示した複数の画像の一部をそれぞれ拡大して示す図である。この例では、第1のタイプの時間周波数マスク（ロバスト主成分分析によるバイナリマスク）と第2のタイプの時間周波数マスク（歌声基本周波数 F_0 によるバイナリマスクまたは調波マスク）との統合を、第1のタイプの時間周波数マスクの選択領域と第2のタイプの時間周波数マスクの選択領域との論理積（AND）をとることにより両マスクを統合して第3のタイプの時間周波数マスク（統合バイナリマスク）を得ている。

40

【0043】

マスクング部5は、第3のタイプの時間周波数マスク（統合バイナリマスク）を音楽スペクトログラムに適用して分離用歌声スペクトログラムを生成する（ステップS T8）。このマスクング部5から出力される分離用歌声スペクトログラムを記憶部に記憶しておい

50

てもよいのは勿論である。図9は、マスク部5における処理を画像で示すための図である。また図10は、図9の画像の理解を高めるために、図9に示した複数の画像の一部をそれぞれ拡大して示す図である。図11(A)乃至(D)は、マスク部5によるマスク処理の状況を示す波形図である。なお図11(A)乃至(D)においては、スペクトログラムに含まれる1フレーム分のスペクトルを図示の対象としている。図11(A)は音楽スペクトログラムに含まれる混合音スペクトル $X(f)$ である。そして図11(B)は混合音スペクトル $X(f)$ に対応する第1のタイプの時間周波数マスク(ロバスト主成分分析マスクに含まれる1フレーム分の周波数マスク $M_b(f)$)であり、図11(C)は混合音スペクトル $X(f)$ に対応する第2のタイプの時間周波数マスクに含まれる1フレーム分の周波数マスク[調波マスク $M_h(f)$]である。そして図11(D)は、第1のタイプの時間周波数マスク(ロバスト主成分分析マスクに含まれる1フレーム分の周波数マスク $M_b(f)$)と第2のタイプの時間周波数マスクに含まれる1フレーム分の周波数マスク[調波マスク $M_h(f)$]が統合されて生成された第3のタイプの時間周波数マスクに含まれる1フレーム分の周波数マスク[統合マスク $M_b(f) \cdot M_h(f)$]によってマスクされて得た分離された歌声スペクトル $[X(f) \cdot M_b(f) \cdot M_h(f)]$ である。図11(D)から分かるように、統合マスクを使用してマスクを行うと分離精度が高くなっているのが分かる。

10

【0044】

そして信号分離生成部6は、分離用歌声スペクトログラムに基づいて歌声信号を分離生成する。具体的には、図12に示すように、信号分離生成部6は、マスク部5から出力された分離用歌声スペクトログラムを時間領域に逆変換することにより歌声信号を分離生成する(ステップST9)。

20

【0045】

本実施の形態の効果を確認するために、ロバスト主成分分析により歌声信号を分離した場合と、本実施の形態で歌声信号を分離した場合について、目的音源の歪みで分離精度を判定するNSDR(Normalized Signal-to-Distortion Ratio[dB])で、110曲の音楽音響信号から歌声信号を分離した結果を比較してみた。その結果、歌声の分離精度に関しては、本実施の形態では5.06[dB]、ロバスト主成分分析では2.09[dB]、伴奏の分離精度に関しては、本実施の形態では6.21[dB]、ロバスト主成分分析では1.71[dB]という結果が得られた。歌声分離及び伴奏分離の両方において、本実施の形態のほうが、RPCAよりも精度が高いことが確認された。

30

【0046】

なお上記各構成要件は、1以上のプロセッサとメモリによって実現するのが好ましい。またマスク生成システム4は、本実施の形態の歌声信号分離システムと一緒に構成する必要はない。すなわち第3のタイプの時間周波数マスク(統合マスク)は、歌声信号分離システムとは別に設けられたマスク生成システムによって事前に生成しておいてもよいのは勿論である。

【0047】

上記実施の形態では、2つバイナリマスクの統合に論理積(AND)を用いたが、本発明におけるマスクの統合は、上記実施の形態に限定されるものではない。図13は、マスクの統合の他の例を示す概念図である。この例では、第1のタイプの時間周波数マスク(RPCAによるバイナリマスク)の選択領域と第2のタイプの時間周波数マスク(歌声基本周波数 F_0 による調波構造のバイナリマスク)の選択領域との論理積をとって、仮統合時間周波数マスクを生成する。そしてこの仮統合時間周波数マスクから、歌が無い区間を推定して、推定された時間フレームの全要素を0にすることにより第3のタイプの時間周波数マスク(統合バイナリマスク)とすることができる。この時間周波数マスクでは、歌が無い区間を推定して、推定された時間フレームの全要素を0にするため、さらに分離精度を高めることができる。

40

【0048】

図14は、マスクの統合のさらに他の例を示す概念図である。この例では、第1のタイ

50

ブの時間周波数マスク（ロバスト主成分分析によるバイナリマスク）の選択領域と第2のタイプの時間周波数マスク（歌声基本周波数F0による調波構造のバイナリマスク）の選択領域との論理積をとって、仮統合時間周波数マスクを生成する。そしてこの仮統合時間周波数マスクから、歌が無い区間を推定して、推定された時間フレームの全要素を0にし、且つ第1のタイプの時間周波数マスク（ロバスト主成分分析によるバイナリマスク）から子音を通過させる要素を得て該要素を仮統合時間周波数マスクに反映して、第3のタイプの時間周波数マスク（統合バイナリマスク）とすることができる。このようにするとさらに分離精度を高めることができる。

【産業上の利用可能性】

【0049】

近年、既存楽曲をユーザが自分好みに編集・加工することを可能にする能動的音楽鑑賞システムの研究が盛んである。中でも、混合音中の歌声の編集は最も実現が難しい課題の一つであり、既存の歌声の声質を他の歌唱者の声質に直接変換する技術は提案されているが、歌声がもつ特徴的な音高軌跡、すなわち歌唱表現を編集する技術は実現されていなかったが、本発明によれば、調波構造から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、出現可能性から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビン及び前記出現可能性から判断すると歌声を含むスペクトルビンであると判断されるスペクトルビンであっても、調波構造から判断すると歌声を含むスペクトルビンではないと判断されるスペクトルビンを音楽スペクトログラムからマスクする機能を有する時間周波数マスクを用いることにより、歌声信号と伴奏音信号とを含む音楽音響信号から歌声信号を従来よりも精度よく分離できる。

【符号の説明】

【0050】

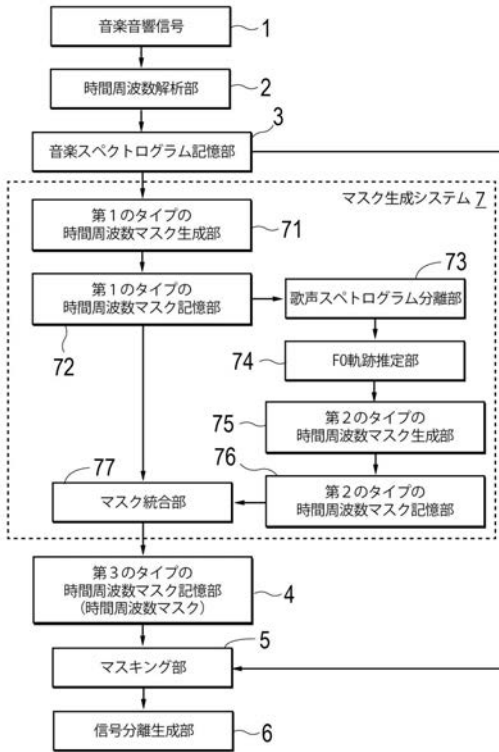
- 1 音楽音響信号
- 2 時間周波数解析部
- 3 音楽スペクトログラム記憶部
- 4 時間周波数マスク記憶部
- 5 マスキング部
- 6 信号分離生成部
- 7 マスク生成システム
- 7 1 時間周波数マスク生成部
- 7 2 時間周波数マスク記憶部
- 7 3 歌声スペクトログラム分離部
- 7 4 F0軌跡推定部
- 7 5 時間周波数マスク生成部
- 7 6 時間周波数マスク記憶部
- 7 7 マスク統合部

10

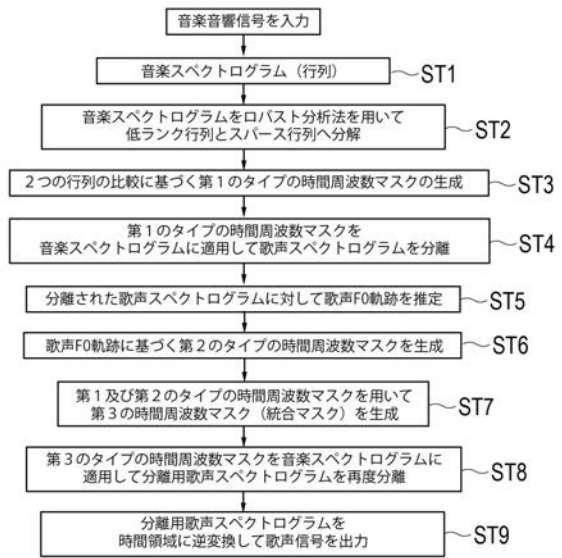
20

30

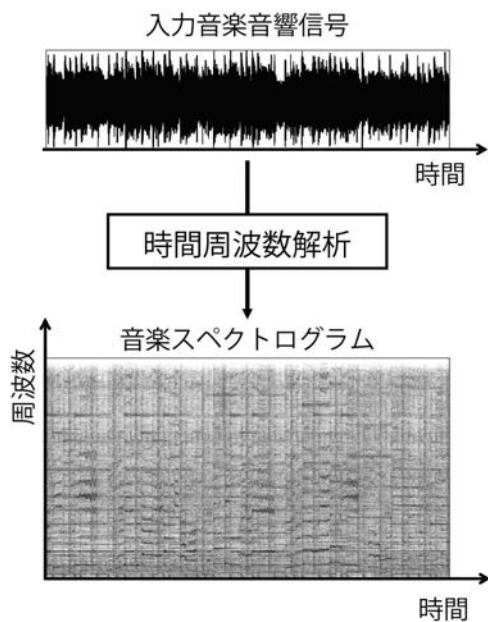
【 図 1 】



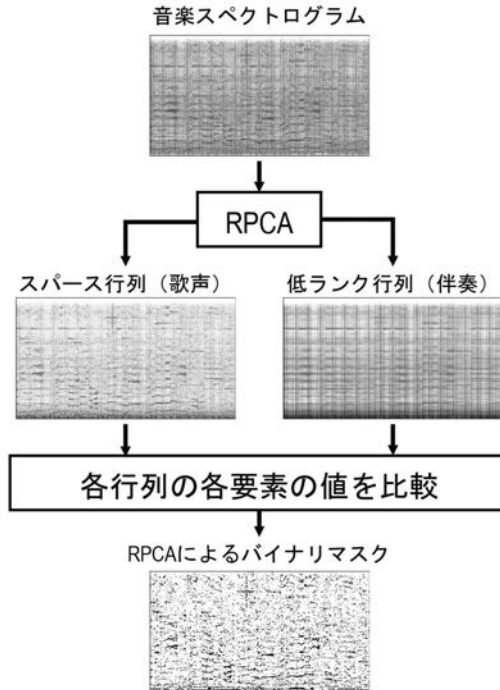
【 図 2 】



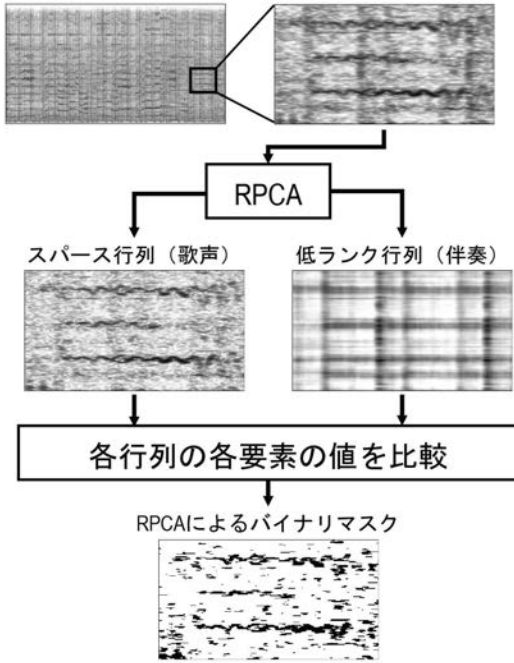
【 図 3 】



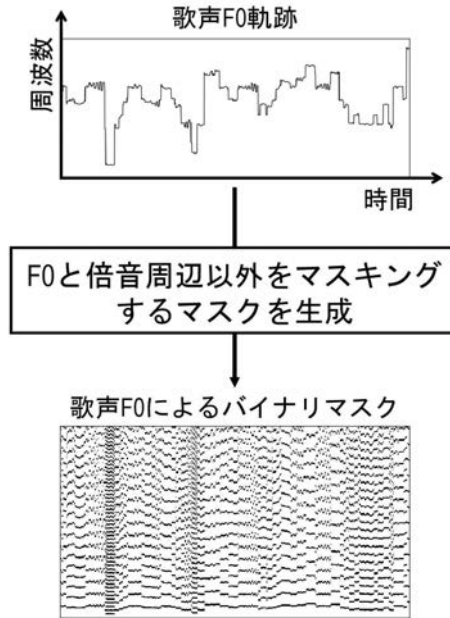
【 図 4 】



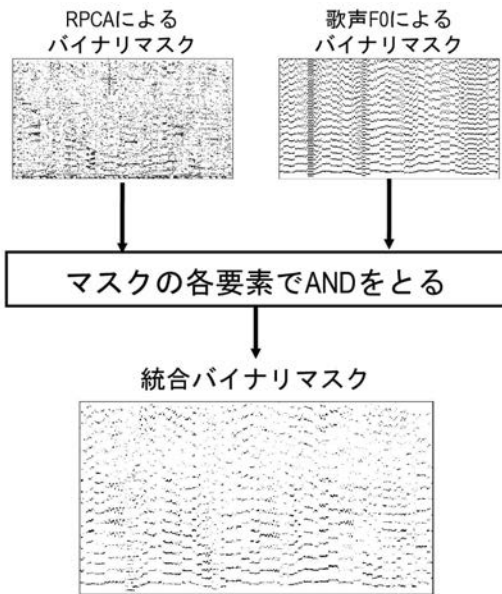
【図5】



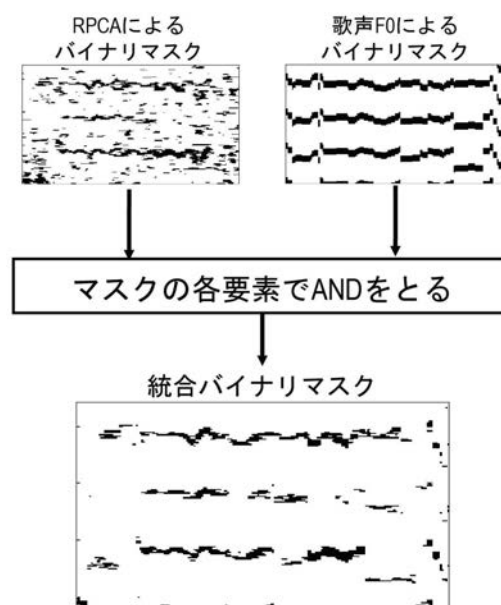
【図6】



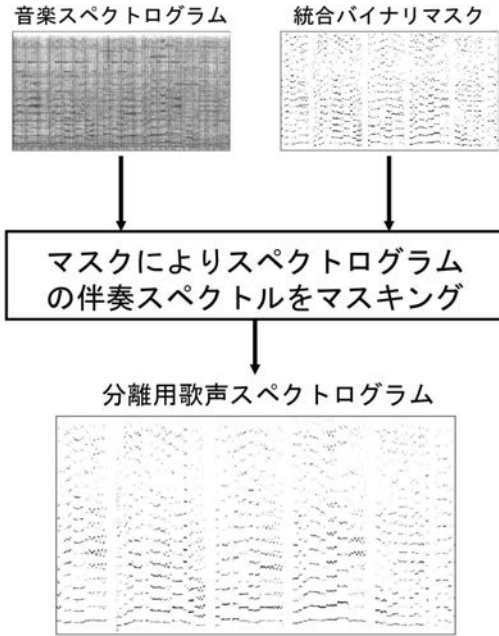
【図7】



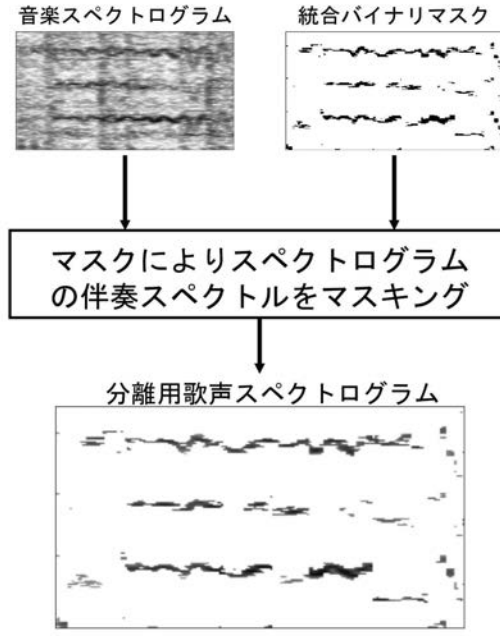
【図8】



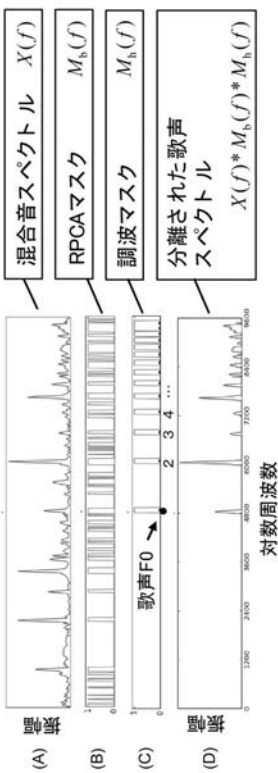
【図 9】



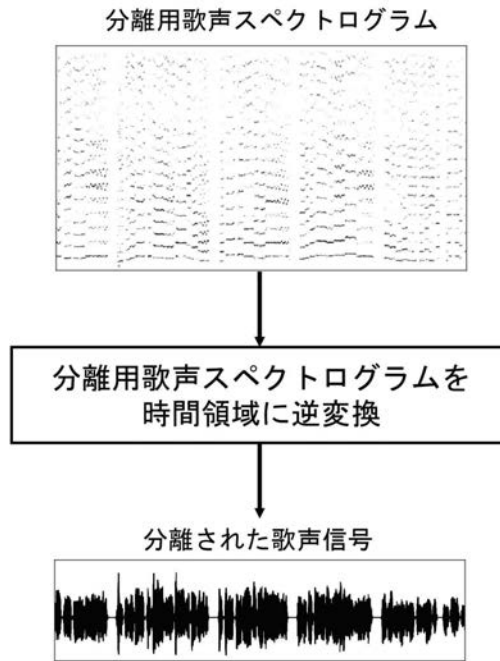
【図 10】



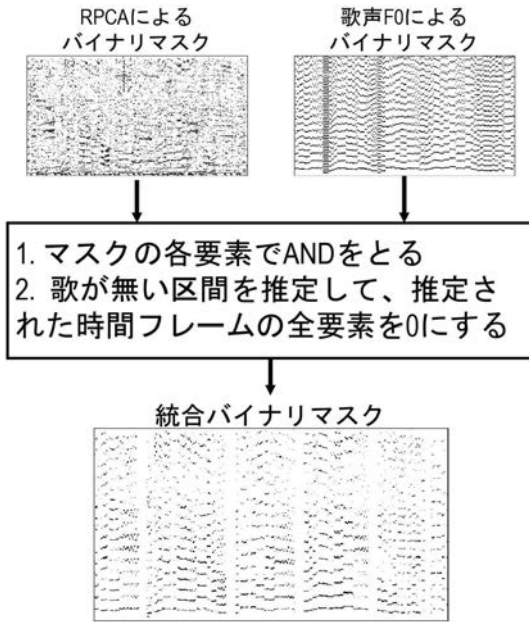
【図 11】



【図 12】



【図 1 3】



【図 1 4】

