

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2015-514239

(P2015-514239A)

(43) 公表日 平成27年5月18日(2015.5.18)

(51) Int.Cl.		F I		テーマコード (参考)
G 1 0 L 15/10 (2006.01)		G 1 0 L 15/10	5 0 0 Z	5 L 0 9 6
G 0 6 T 7/00 (2006.01)		G 0 6 T 7/00	P	
G 0 6 T 7/20 (2006.01)		G 0 6 T 7/20	B	
G 1 0 L 25/51 (2013.01)		G 1 0 L 25/51	4 0 0	
G 1 0 L 15/00 (2013.01)		G 1 0 L 15/00	2 0 0 G	
審査請求 未請求 予備審査請求 有 (全 104 頁) 最終頁に続く				

(21) 出願番号 特願2015-505720 (P2015-505720)
 (86) (22) 出願日 平成25年3月7日 (2013.3.7)
 (85) 翻訳文提出日 平成26年12月3日 (2014.12.3)
 (86) 国際出願番号 PCT/US2013/029558
 (87) 国際公開番号 W02013/154701
 (87) 国際公開日 平成25年10月17日 (2013.10.17)
 (31) 優先権主張番号 61/623, 910
 (32) 優先日 平成24年4月13日 (2012.4.13)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 13/664, 295
 (32) 優先日 平成24年10月30日 (2012.10.30)
 (33) 優先権主張国 米国 (US)

(71) 出願人 595020643
 クゥアルコム・インコーポレイテッド
 QUALCOMM INCORPORATED
 アメリカ合衆国、カリフォルニア州 92
 121-1714、サン・ディエゴ、モア
 ハウス・ドライブ 5775
 (74) 代理人 100108855
 弁理士 蔵田 昌俊
 (74) 代理人 100109830
 弁理士 福原 淑弘
 (74) 代理人 100103034
 弁理士 野河 信久
 (74) 代理人 100075672
 弁理士 峰 隆司

最終頁に続く

(54) 【発明の名称】 マルチモーダル整合方式を使用するオブジェクト認識

(57) 【要約】

シーン中の1つまたは複数のオブジェクトを認識し、位置を特定するための方法、システムおよび製造品が開示される。シーンの画像および/またはビデオがキャプチャされる。シーンにおいて記録されたオーディオを使用して、キャプチャされたシーンのオブジェクト探索が狭められる。たとえば、キャプチャされた画像/ビデオ中の探索エリアを限定するために、音の到来方向(DOA)が判断され、使用され得る。別の例では、記録されたオーディオ中で識別される音のタイプに基づいてキーポイントシグネチャが選択され得る。キーポイントシグネチャは、本システムが認識するように構成された特定のオブジェクトに対応する。次いで、キャプチャされたシーン中で識別されるキーポイントを、選択されたキーポイントシグネチャと比較する、シフト不変特徴変換(SIFT)分析を使用して、シーン中のオブジェクトが認識され得る。

【選択図】 図2

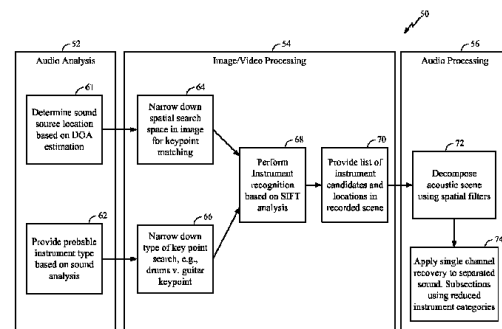


FIG. 2

【特許請求の範囲】**【請求項 1】**

デバイスにおいて、シーン中のオブジェクトを認識する方法であって、
前記シーンにおいて記録されたオーディオに基づいて前記オブジェクトに対応するキーポイントを選択することと、
前記選択されたキーポイントに基づいて前記オブジェクトを識別することと
を備える、方法。

【請求項 2】

前記シーンにおいて記録されたオーディオに基づいて、1つまたは複数のオブジェクトに対応する1つまたは複数のキーポイントシグネチャを選択することと、
前記シーンの画像中の複数のキーポイントを識別することと、
前記オブジェクトを識別するために前記キーポイントを前記キーポイントシグネチャと比較することと
をさらに備える、請求項 1 に記載の方法。

10

【請求項 3】

前記シーンにおいて記録された前記オーディオに基づいてシーン画像の一部分を選択することと、
前記画像の前記一部分内からのみ前記キーポイントを選択することと
をさらに備える、請求項 1 に記載の方法。

20

【請求項 4】

前記シーンにおいて記録された前記オーディオに基づいて前記画像の一部分を選択することが、
前記オーディオからオーディオ到来方向(DOA)を判断することと、
前記オーディオDOAに基づいて前記画像の前記一部分を選択することと
を含む、請求項 3 に記載の方法。

【請求項 5】

前記オーディオDOAを判断することが、
前記シーンに位置する複数のマイクロフォンにおいて前記オーディオを受信し、それによって複数のマイクロフォン信号を生成することと、
前記マイクロフォン信号に基づいて前記オーディオDOAを判断することと
を含む、請求項 4 に記載の方法。

30

【請求項 6】

前記シーンのビデオ記録から複数の局所動きベクトルを計算することと、
前記局所動きベクトルを1つまたは複数のオブジェクトに対応する所定の局所動きベクトルのデータベースと比較することによって、および前記キーポイントを1つまたは複数のキーポイントシグネチャと比較することによって前記オブジェクトを識別することと
をさらに備える、請求項 1 に記載の方法。

【請求項 7】

前記シーンにおいて記録された前記オーディオから複数の音響認識特徴を計算することと、
前記音響認識特徴を1つまたは複数のオブジェクトに対応する所定の音響認識特徴のデータベースと比較することによって、および前記キーポイントを1つまたは複数のキーポイントシグネチャと比較することによって前記オブジェクトを識別することと
をさらに備える、請求項 1 に記載の方法。

40

【請求項 8】

前記音響認識特徴がメル周波数ケプストラム係数を含む、請求項 7 に記載の方法。

【請求項 9】

前記画像中に現れる1つまたは複数のオブジェクトについての範囲情報を判断することと、
前記範囲情報に基づいて前記キーポイントを分析することと

50

をさらに備える、請求項 1 に記載の方法。

【請求項 10】

範囲情報を判断することが、オートフォーカスカメラを使用して範囲情報を判断することと、マルチカメラ画像視差推定を使用して範囲情報を判断することと、上記の任意の最適な組合せとからなるグループから選択される、請求項 9 に記載の方法。

【請求項 11】

シーンにおいて記録されたオーディオに基づいて前記シーン中のオブジェクトに対応するキーポイントを選択するように構成されたキーポイントセクタと、

前記選択されたキーポイントに基づいて前記オブジェクトを識別するように構成された整合デバイスと

を備える、装置。

【請求項 12】

シーンの画像中の複数のキーポイントを識別するように構成されたキーポイント検出器をさらに備え、

前記キーポイントセクタが、前記シーンにおいて記録されたオーディオに基づいて、1 つまたは複数のオブジェクトに対応する 1 つまたは複数のキーポイントシグネチャを選択するように構成され、

前記整合デバイスが、前記シーン中のオブジェクトを識別するために前記キーポイントを前記キーポイントシグネチャと比較するように構成された、

請求項 11 に記載の装置。

【請求項 13】

前記シーンにおいて記録された前記オーディオに基づいて前記シーンの画像の一部を選択するように構成された第 1 のセクタと、

前記画像の前記一部分内からのみ前記キーポイントを選択するように構成された第 2 のセクタと

をさらに備える、請求項 11 に記載の装置。

【請求項 14】

前記第 1 のセクタが、

前記オーディオからオーディオ到来方向 (DOA) を判断するように構成された検出器と、

前記オーディオ DOA に基づいて前記画像の前記一部分を選択するように構成された第 3 のセクタと

を含む、請求項 13 に記載の装置。

【請求項 15】

前記検出器が、

前記オーディオを受信して、複数のマイクロフォン信号を生成するための、前記シーンに位置する複数のマイクロフォンと、

前記マイクロフォン信号に基づいて前記オーディオ DOA を判断するように構成されたオーディオプロセッサと

を含む、請求項 14 に記載の装置。

【請求項 16】

前記シーンのビデオ記録から複数の局所動きベクトルを計算するように構成されたビデオプロセッサ

をさらに備え、

前記整合デバイスが、前記局所動きベクトルを 1 つまたは複数のオブジェクトに対応する所定の局所動きベクトルのデータベースと比較することによって、および前記キーポイントを 1 つまたは複数のキーポイントシグネチャと比較することによって前記オブジェクトを識別するように構成された、

請求項 11 に記載の装置。

【請求項 17】

前記シーンにおいて記録された前記オーディオから複数の音響認識特徴を計算するように構成されたオーディオプロセッサをさらに備え、

前記整合デバイスが、前記音響認識特徴を１つまたは複数のオブジェクトに対応する所定の音響認識特徴のデータベースと比較することによって、および前記キーポイントを１つまたは複数のキーポイントシグネチャと比較することによって前記オブジェクトを識別するように構成された、

請求項 11 に記載の装置。

【請求項 18】

前記音響認識特徴がメル周波数ケプストラム係数を含む、請求項 17 に記載の装置。

10

【請求項 19】

前記画像中に現れる１つまたは複数のオブジェクトについての範囲情報を判断するように構成された範囲検出器と、

前記範囲情報に基づいて前記キーポイントを分析するように構成されたキーポイント検出器と

をさらに備える、請求項 11 に記載の装置。

【請求項 20】

前記範囲検出器が、オートフォーカスカメラと、マルチカメラアレイと、上記の任意の好適な組合せとからなるグループから選択される検出器を含む、請求項 19 に記載の装置。

20

【請求項 21】

シーンにおいて記録されたオーディオに基づいて前記シーン中のオブジェクトに対応するキーポイントを選択するための手段と、

前記選択されたキーポイントに基づいて前記オブジェクトを識別するための手段とを備える、装置。

【請求項 22】

前記シーンにおいて記録されたオーディオに基づいて、１つまたは複数のオブジェクトに対応する１つまたは複数のキーポイントシグネチャを選択するための手段と、

前記シーンの画像中の複数のキーポイントを識別するための手段と、

前記シーン中の前記オブジェクトを識別するために前記キーポイントを前記キーポイントシグネチャと比較するための手段と

30

をさらに備える、請求項 21 に記載の装置。

【請求項 23】

前記シーンにおいて記録された前記オーディオに基づいて前記シーンの画像の一部分を選択するための手段と、

前記画像の前記一部分内からのみ前記キーポイントを選択するための手段と

をさらに備える、請求項 21 に記載の装置。

【請求項 24】

前記シーンにおいて記録された前記オーディオに基づいて前記画像の一部分を選択するための前記手段が、

40

前記オーディオからオーディオ到来方向(DOA)を判断するための手段と、

前記オーディオDOAに基づいて前記画像の前記一部分を選択するための手段とを含む、請求項 23 に記載の装置。

【請求項 25】

前記オーディオDOAを判断するための手段が、

前記シーンに位置する複数のマイクロフォンにおいて前記オーディオを受信し、それによって複数のマイクロフォン信号を生成するための手段と、

前記マイクロフォン信号に基づいて前記オーディオDOAを判断するための手段とを含む、請求項 24 に記載の装置。

【請求項 26】

50

前記シーンのビデオ記録から複数の局所動きベクトルを計算するための手段と、

前記局所動きベクトルを１つまたは複数のオブジェクトに対応する所定の局所動きベクトルのデータベースと比較することによって、および前記キーポイントを１つまたは複数のキーポイントシグネチャと比較することによって前記オブジェクトを識別するための手段と

をさらに備える、請求項２１に記載の装置。

【請求項２７】

前記シーンにおいて記録された前記オーディオから複数の音響認識特徴を計算するための手段と、

前記音響認識特徴を１つまたは複数のオブジェクトに対応する所定の音響認識特徴のデータベースと比較することによって、および前記キーポイントを１つまたは複数のキーポイントシグネチャと比較することによって前記オブジェクトを識別するための手段と
をさらに備える、請求項２１に記載の装置。

10

【請求項２８】

前記音響認識特徴がメル周波数ケプストラム係数を含む、請求項２７に記載の装置。

【請求項２９】

画像中に現れる１つまたは複数のオブジェクトについての範囲情報を判断するための手段と、

前記範囲情報に基づいて前記キーポイントを分析するための手段と

をさらに備える、請求項２１に記載の装置。

20

【請求項３０】

範囲情報を判断するための手段が、オートフォーカスカメラを使用して範囲情報を判断するための手段と、マルチカメラ画像視差推定を使用して範囲情報を判断するための手段と、上記の任意の好適な組合せとからなるグループから選択される、請求項２９に記載の装置。

【請求項３１】

シーンにおいて記録されたオーディオに基づいて前記シーン中のオブジェクトに対応するキーポイントを選択するためのコードと、

前記選択されたキーポイントに基づいて前記オブジェクトを識別するためのコードと
を備える、１つまたは複数のプロセッサによって実行可能な命令のセットを具備するコンピュータ可読媒体。

30

【請求項３２】

前記シーンにおいて記録されたオーディオに基づいて、１つまたは複数のオブジェクトに対応する１つまたは複数のキーポイントシグネチャを選択するためのコードと、

前記シーンの画像中の複数のキーポイントを識別するためのコードと、

前記シーン中の前記オブジェクトを識別するために前記キーポイントを前記キーポイントシグネチャと比較するためのコードと

をさらに備える、請求項３１に記載のコンピュータ可読媒体。

【請求項３３】

前記シーンにおいて記録された前記オーディオに基づいて画像の一部分を選択するためのコードと、

前記画像の前記一部分内からのみ前記キーポイントを選択するためのコードと

をさらに備える、請求項３１に記載のコンピュータ可読媒体。

40

【請求項３４】

前記シーンにおいて記録された前記オーディオに基づいて前記画像の一部分を選択するための前記コードが、

前記オーディオからオーディオ到来方向(DOA)を判断するためのコードと、

前記オーディオDOAに基づいて前記画像の前記一部分を選択するためのコードと
を含む、請求項３３に記載のコンピュータ可読媒体。

【請求項３５】

50

前記オーディオD O Aを判断するためのコードが、

前記シーンに位置する複数のマイクロフォンにおいて前記オーディオを受信し、それによって複数のマイクロフォン信号を生成するためのコードと、

前記マイクロフォン信号に基づいて前記オーディオD O Aを判断するためのコードとを含む、請求項34に記載のコンピュータ可読媒体。

【請求項36】

前記シーンのビデオ記録から複数の局所動きベクトルを計算するためのコードと、

前記局所動きベクトルを1つまたは複数のオブジェクトに対応する所定の局所動きベクトルのデータベースと比較することによって、および前記キーポイントを1つまたは複数のキーポイントシグネチャと比較することによって前記オブジェクトを識別するためのコードと

をさらに備える、請求項31に記載のコンピュータ可読媒体。

【請求項37】

前記シーンにおいて記録された前記オーディオから複数の音響認識特徴を計算するためのコードと、

前記音響認識特徴を1つまたは複数のオブジェクトに対応する所定の音響認識特徴のデータベースと比較することによって、および前記キーポイントを1つまたは複数のキーポイントシグネチャと比較することによって前記オブジェクトを識別するためのコードとをさらに備える、請求項31に記載のコンピュータ可読媒体。

【請求項38】

前記音響認識特徴がメル周波数ケプストラム係数を含む、請求項37に記載のコンピュータ可読媒体。

【請求項39】

画像中に現れる1つまたは複数のオブジェクトについての範囲情報を判断するためのコードと、

前記範囲情報に基づいて前記キーポイントを分析するためのコードとをさらに備える、請求項31に記載のコンピュータ可読媒体。

【請求項40】

範囲情報を判断するためのコードが、オートフォーカスカメラを使用して範囲情報を判断するためのコードと、マルチカメラ画像視差推定を使用して範囲情報を判断するためのコードと、上記の任意の好適な組合せとからなるグループから選択される、請求項39に記載のコンピュータ可読媒体。

【発明の詳細な説明】

【技術分野】

【0001】

米国特許法第119条による優先権の主張

[0001]本特許出願は、本出願の譲受人に譲渡され、参照により本明細書に明確に組み込まれる、2012年4月13日に出願された米国仮出願第61/623,910号の優先権を主張する。

【0002】

[0002]本開示は、一般にオブジェクト認識装置および方法に関する。

【背景技術】

【0003】

[0003]様々な適用例は、視覚シーン中のオブジェクトを識別することが可能である機械またはプロセッサを有することから恩恵を受け得る。コンピュータビジョンの分野は、シーン中のオブジェクトまたは特徴を識別することを可能にする技法および/またはアルゴリズムを提供することを試みており、オブジェクトまたは特徴は、1つまたは複数のキーポイント(keypoint)を識別する記述子によって特徴づけられ得る。SIFT(スケール不変特徴変換(Scale Invariant Feature Transform))など、これらの技法はまた、しばしば、適用例の中でも、オブジェクト認識、オブジェクト検出、画像整合、3次元構造

10

20

30

40

50

構築、ステレオ対応、および／または動き追跡、顔認識に適用される。

【 0 0 0 4 】

[0004]大部分のオブジェクト認識技法は、シーンからキャプチャされた視覚情報、たとえば、ビデオ、画像またはピクチャのみに依拠する。

【 発明の概要 】

【 0 0 0 5 】

[0005]この概要は、すべての企図される実施形態の包括的な概観ではなく、すべての実施形態の主要または重要な要素を識別するものでも、いずれかまたはすべての実施形態の範囲を定めるものでもない。その唯一の目的は、後で提示するより詳細な説明の導入として、1つまたは複数の実施形態のいくつかの概念を簡略化された形式で提示することである。

10

【 0 0 0 6 】

[0006]シーン中の1つまたは複数のオブジェクトを認識し、位置を特定するための改善された技法が開示される。これらの技法は、シーン中のオブジェクトを認識するのを助けるために、シーンにおいて記録されたオーディオならびに視覚情報の使用を取り入れる。これらの技法の一態様によれば、装置は、シーンにおいて記録されたオーディオに基づいてシーン中のオブジェクトに対応するキーポイントを選択するように構成されたキーポイントセレクトと、選択されたキーポイントに基づいてオブジェクトを識別するように構成されたキーポイント整合デバイスとを含む。

【 0 0 0 7 】

20

[0007]さらなる態様によれば、シーン中のオブジェクトを認識する方法は、シーンにおいて記録されたオーディオに基づいてオブジェクトに対応するキーポイントを選択することと、選択されたキーポイントに基づいてオブジェクトを識別することを含む。

【 0 0 0 8 】

[0008]さらなる態様によれば、1つまたは複数のプロセッサによって実行可能な命令のセットを具備するコンピュータ可読媒体が提供される。媒体は、シーンにおいて記録されたオーディオに基づいてシーン中のオブジェクトに対応するキーポイントを選択するためのプログラミングコードと、選択されたキーポイントに基づいてオブジェクトを識別するためのプログラミングコードとを記憶する。

【 0 0 0 9 】

30

[0009]他の態様、特徴および利点は、以下の図および詳細な説明を審査すると当業者に明らかであるかまたは明らかになる。すべてのそのような追加の特徴、態様、および利点は本明細書内に含まれ、添付の特許請求の範囲によって保護されるものである。

【 図面の簡単な説明 】

【 0 0 1 0 】

[0010]図面は例示のためのものにすぎないことを理解されたい。さらに、図中の構成要素は必ずしも一定の縮尺であるとは限らず、代わりに、本明細書で説明する技法およびデバイスの原理を示すことに強調が置かれる。図において、同様の参照番号は、異なる図全体にわたって対応する部分を示す。

【 図 1 】 [0011]例示的な聴覚シーンおよび視聴覚シーン分析システムを示す概念図。

40

【 図 2 】 [0012]図1の視聴覚シーン分析システムを動作させる方法を示すプロセスブロック図。

【 図 3 】 [0013]図1の聴覚シーン分析システムを使用して聴覚シーンを分解し、処理する例示的な方法を示すフローチャート。

【 図 4 】 [0014]聴覚シーン中の音生成オブジェクトを識別する例示的な方法を示すフローチャート。

【 図 5 A 】 [0015]聴覚シーン中の音生成オブジェクトを識別する第2の例示的な方法のフローチャート。

【 図 5 B 】 聴覚シーン中の音生成オブジェクトを識別する第2の例示的な方法のフローチャート。

50

【図 6】[0016]図 1 の聴覚シーン分析システムとともに使用され得る例示的なオブジェクト認識サブシステムのいくつかの構成要素を示すブロック図。

【図 7】[0017]記録された画像上でオブジェクト認識を実行するための機能段階を示すブロック図。

【図 8】[0018]例示的な画像処理段階におけるガウススケール空間発生を示す図。

【図 9】[0019]例示的な画像処理段階における特徴検出を示す図。

【図 10】[0020]例示的な画像処理段階における特徴記述子抽出を示す図。

【図 11】[0021]照度不変特徴 / キーポイント検出のためのスケール空間正規化の一例を示す図。

【図 12】[0022]正規化されたスケール空間差分の 1 つのレベルがどのように取得され得るかの一例を示す図。 10

【図 13】[0023]改善された特徴検出のために正規化されたスケール空間差分を発生するための方法を示す図。

【図 14】[0024]特徴 / キーポイント検出デバイスの一例を示すブロック図。

【図 15】[0025]画像整合デバイスの一例を示すブロック図。

【図 16 A】[0026]一般的構成による装置 A 1 0 0 のブロック図。

【図 16 B】[0027]マイクロフォンアレイ R 1 0 0 と装置 A 1 0 0 のインスタンスとを含むデバイス D 1 0 のブロック図。

【図 16 C】[0028]点音源 j から受信された信号成分の、アレイ R 1 0 0 のマイクロフォン M C 1 0 および M C 2 0 の軸に対する到来方向 θ_j を示す図。 20

【図 17】[0029]装置 A 1 0 0 の実装形態 A 1 1 0 のブロック図。

【図 18 A】[0030]フィルタ更新モジュール U M 1 0 の実装形態 U M 2 0 のブロック図。

【図 18 B】[0031]フィルタ更新モジュール U M 2 0 の実装形態 U M 2 2 のブロック図。

【図 19 A】[0032]カメラ C M 1 0 をもつアレイ R 1 0 0 の 4 マイクロフォン実装形態 R 1 0 4 の構成の一例の上面図。

【図 19 B】[0033]到来方向の推定のための遠距離場モデルを示す図。

【図 20】[0034]装置 A 1 0 0 の実装形態 A 1 2 0 のブロック図。

【図 21】[0035]装置 A 1 2 0 および A 2 0 0 の実装形態 A 2 2 0 のブロック図。

【図 22】[0036]D O A 推定に S R P - P H A T を使用した結果によるヒストグラムの例を示す図。 30

【図 23】[0037]I V A 適応ルール (4 0 ~ 6 0 度の音源分離) を使用して適応された逆混合行列の異なる出力チャネルに関する 4 つのヒストグラムのセットの一例を示す図。

【図 24】[0038]画像またはビデオキャプチャ中にシーン中のオブジェクトの視差を検出するように構成された例示的な画像キャプチャデバイスの図。

【図 25】[0039]図 2 4 のデバイス中に含まれ得る例示的な画像処理システムのブロック図。

【図 26 A】[0040]知覚されたオブジェクト深さと相関させられたオブジェクト視差の例示的な例の図。

【図 26 B】知覚されたオブジェクト深さと相関させられたオブジェクト視差の例示的な例の図。 40

【図 27 A】[0041]一般的構成による方法 M 1 0 0 のフローチャート。

【図 27 B】[0042]方法 M 1 0 0 の実装形態 M 2 0 0 のフローチャート。

【図 27 C】[0043]一般的構成による、オーディオ信号を分解するための装置 M F 1 0 0 のブロック図。

【図 27 D】[0044]別の一般的構成による、オーディオ信号を分解するための装置 A 1 0 0 のブロック図。

【図 28 A】[0045]方法 M 1 0 0 の実装形態 M 3 0 0 のフローチャート。

【図 28 B】[0046]装置 A 1 0 0 の実装形態 A 3 0 0 のブロック図。

【図 28 C】[0047]装置 A 1 0 0 の別の実装形態 A 3 1 0 のブロック図。

【図 29 A】[0048]方法 M 2 0 0 の実装形態 M 4 0 0 のフローチャート。 50

【図 2 9 B】[0049]方法 M 2 0 0 の実装形態 M 5 0 0 のフローチャート。

【図 3 0 A】[0050]方法 M 1 0 0 の実装形態 M 6 0 0 のフローチャート。

【図 3 0 B】[0051]装置 A 1 0 0 の実装形態 A 7 0 0 のブロック図。

【図 3 1】[0052]装置 A 1 0 0 の実装形態 A 8 0 0 のブロック図。

【図 3 2】[0053]モデル $B f = y$ を示す図。

【図 3 3】[0054]図 3 2 のモデルの変形 $B' f = y$ を示す図。

【図 3 4】[0055]複数の音源がアクティブであるシナリオを示す図。

【発明を実施するための形態】

【0011】

[0056]図面を参照し、組み込む以下の詳細な説明は、1つまたは複数の特定の実施形態について説明し、例示する。限定するためではなく、例示し、教示するためだけに提供されるこれらの実施形態について、当業者が特許請求の範囲を实践することを可能にするのに十分詳細に図示し、説明する。したがって、簡潔のために、説明は、当業者に知られているある情報を省略し得る。

10

【0012】

[0057]「例示的」という単語は、本開示全体にわたって、「例、事例、または例示の働きをすること」を意味するために使用する。本明細書で「例示的」と記載されたものはどんなものも、必ずしも他の手法または特徴よりも好ましいまたは有利であると解釈されるべきではない。その文脈によって明確に限定されない限り、「信号」という用語は、本明細書では、ワイヤ、バス、または他の伝送媒体上に表されたメモリロケーション（またはメモリロケーションのセット）の状態を含む、その通常の意味のいずれかを示すために使用される。

20

【0013】

[0058]本明細書で説明するオブジェクト認識技法は多くの異なるシーンに適用され得るが、本明細書で説明する例は、多くの音源、たとえば、ミュージシャン、演奏者、楽器などが単一のシーン中に存在する、音楽シーンに関係する。いくつかのビデオゲーム（たとえば、Guitar Hero（登録商標）、Rock Band（登録商標））およびコンサート音楽シーンは、複数の楽器およびボーカリストが同時にプレイすることを伴い得る。現在の商用ゲームおよび音楽生成システムでは、これらのシナリオから記録されたオーディオを別々に分析し、後処理し、アップミックスすることが可能であるように、これらのシナリオが、連続的にプレイされるか、または近接して配置されたマイクロフォンを用いてプレイされる必要がある。これらの制約は、音楽生成の場合、干渉を制御する能力および/または空間効果を記録する能力を制限し得、ビデオゲームの場合、制限されたユーザエクスペリエンスをもたらし得る。

30

【0014】

[0059]どんな楽器がプレイされているのか、およびどれくらいのミュージシャン/音源がシーン中に存在するのかを狭めるのを助ける、何らかのアプリオリ(a priori)な知識または他の情報が利用可能にされた場合、音楽聴覚シーン分解は大幅に簡略化され得る。

【0015】

[0060]本明細書で開示するオブジェクト認識技法は、複数の音源を有するシーンにおいて記録されたオーディオを分解するための従来の試みの制限の多くを克服する。概して、オブジェクトまたは特徴認識は、特徴識別および/またはオブジェクト認識のために画像中の関心ポイント（キーポイントとも呼ばれる）を識別することおよび/またはそれらのキーポイントの周りの局所的特徴(localized features)を識別することを伴い得る。本明細書で開示するシステムおよび方法では、いくつかの画像ベースの楽器およびオーディオベースのノート/楽器認識技法が組み合わせられる。いくつかの異なるデバイスは、コンピュータビジョンを使用して特徴識別および/またはオブジェクト認識を実行することが可能であり得る。そのようなデバイスの例は、電話ハンドセット（たとえば、セルラーハンドセット）、ビデオ記録が可能なハンドヘルドモバイルデバイス、オーディオおよびビデオコンテンツを記録する個人メディアプレーヤ、携帯情報端末(PDA)または他の

40

50

ハンドヘルドコンピューティングデバイス、ならびにノートブック、ノートブックコンピュータ、ラップトップコンピュータ、タブレットコンピュータ、または他のポータブルコンピューティングデバイス中に実装され得る。さらに、複数の楽器およびボーカリストが同時にプレイすることを伴い得る、ビデオゲーム（たとえば、Guitar Hero（登録商標）、Rock Band（登録商標））およびコンサート音楽シーンを実行することが可能なデバイス。ポータブルコンピューティングデバイスの種類は現在、ラップトップコンピュータ、ノートブックコンピュータ、ウルトラポータブルコンピュータ、タブレットコンピュータ、モバイルインターネットデバイス、スマートブックおよびスマートフォンなどの名称を有するデバイスを含む。

【0016】

10

[0061]第1の方法では、オーディオ知識のみに基づく情報が、シーン中で探索される音源のタイプを定義するのを助け、画像ベースのスケール不変特徴変換（SIFT）探索において考慮されるべきオブジェクト形状のキーポイントシグネチャの探索ユニバースを低減する。さらに、そのようなキーポイント探索は、必ずしも静止画像に制限されずともは限らないが、深さ（範囲）レイヤ探索の必要に応じて単一または複数のカメラを使用する、典型的なミュージシャンの動きパターンのための周囲ビデオフレームの分析をも伴い得る。キーポイント探索は、認識された楽器を関連する尤度で与えるために、マルチモーダルベイズ推定（multi-modal Bayesian estimation）を介して組み込まれる。

【0017】

20

[0062]第2の方法では、画像のいくつかの部分中の安定した楽器キーポイントを計算するために、マルチマイクロフォン音定位情報と楽器形状認識の両方から音源（たとえば、楽器）ロケーションが推定される。この方法は、オブジェクト認識を改善するために第1の方法と組み合わせられ得る。

【0018】

[0063]第3の方法では、第2の方法を使用して判断される情報など、関連するマルチマイクロフォン音源定位情報とともに、音声／オーディオ認識において使用されるメル周波数ケプストラム係数（MFCC：mel-frequency cepstral coefficient）などの音響特徴が、音源認識を行うためにマルチモーダルベイズ推定において直接使用される。第3の方法は、オブジェクト認識を改善するために第1の方法と組み合わせられ得る。

【0019】

30

[0064]上記の方法は、たとえば、シーンから記録されたオーディオを分解することに対するスパース復元分解手法（sparse recovery decomposition approach）の場合、基底関数インベントリ（basis function inventory）のサイズのより改良された定義を可能にし得る。

【0020】

[0065]図1は、例示的な聴覚シーン10と例示的な聴覚シーン分析システム12とを示す概念図である。聴覚シーン分析システム12は、シーン分析システム14と、マイクロフォンのアレイ18と、1つまたは複数のカメラ16とを含む。カメラ16は、シーン10に対して様々なロケーションおよび角度に配置された1つまたは複数の静止画像カメラおよび／または1つまたは複数のビデオカメラを含み得る。

40

【0021】

[0066]シーン分析システム14は、オブジェクト認識サブシステム20と、音響分解サブシステム22とを含む。オブジェクト認識サブシステム20は、本明細書で説明する方法に従って、シーンにおいて記録されたオーディオ、（1つまたは複数の）画像および／またはビデオに基づいてシーン10中の音源を認識し、位置を特定するように構成される。音響分解サブシステム22は、分離されたオーディオが個々に処理され得るように、オブジェクト認識サブシステム20からの情報に基づいて、シーンを別個の音源に分解するように構成される。

【0022】

[0067]図2は、図1の視聴覚シーン分析システム10を動作させる例示的な方法を示す

50

プロセスブロック図 50 である。本プロセスは、シーン中の 1 つまたは複数の楽器の位置を特定し、識別するために、シーンにおいて記録された視聴覚情報を分析するためのステップを示している。本方法は、オーディオ分析ブロック 52 と、画像および / またはビデオ処理ブロック 54 と、オーディオ処理ブロック 56 とを含む。

【0023】

[0068] 開示する方法は、キャプチャされたオーディオおよび / またはビデオ信号を一連のセグメントとして処理し得る。典型的なセグメント長は約 1 ~ 10 秒にわたる。1 つの特定の例では、信号は、約 1 秒の長さをそれぞれ有する一連の重複しないセグメントまたは「フレーム」に分割される。また、そのような方法によって処理されるセグメントは、異なる演算によって処理されるより大きいセグメントのセグメント（すなわち、「サブフレーム」）であり得、またはその逆も同様である。

10

【0024】

[0069] オーディオ分析ブロック 52 は、シーンにおいて記録されたオーディオ情報に基づいてシーン中の（1 つまたは複数の）音源ロケーションを判断するステップを含む（ボックス 61）。オーディオ情報はマイクロフォンアレイ 18 によってキャプチャされ得る。音ロケーションは、シーン中の音源および / または音源について判断された範囲情報から音の推定される到来方向（DOA: direction of arrival）に基づいて判断され得る。音源の DOA は、本明細書において以下で説明するオーディオ DOA 推定技法を使用して推定され得、音源の範囲は、図 18 ~ 図 29 を参照しながら本明細書において以下で説明する範囲発見技法を使用して推定され得る。

20

【0025】

[0070] オーディオ分析ブロック 52 はまた、シーン中の各音源に音源の推定タイプを与えるステップを含む（ボックス 62）。たとえば、楽器について、シーンにおいて記録された音は、その音を生成している楽器の可能性のあるタイプと音源を整合させるために、楽器ノートライブラリを使用して分析され得る。

【0026】

[0071] ボックス 61、62 からの音源ロケーションおよびタイプ推定は、画像 / ビデオ処理ブロック 54 に受け渡され、音源の視覚的識別のために探索を制限するために使用される。ボックス 64 において、推定されたロケーション情報を使用して、キーポイント整合のためにシーンの記録画像中の空間探索空間を狭める。ボックス 66 において、画像キーポイント探索が推定楽器タイプに基づいて狭められる。これらのステップの両方は、シーン中の（1 つまたは複数の）楽器を識別する信頼性を著しく改善し得、また、（1 つまたは複数の）楽器の視覚的認識を行うために必要とされる処理の量を低減し得る。

30

【0027】

[0072] ボックス 68 において、シーン中の（1 つまたは複数の）楽器を識別するために、シーンにおいて記録された画像および / またはビデオデータ上で視覚的オブジェクト認識分析が実行される。この分析は、視覚特徴分析方式、たとえば、シーンのスケール不変特徴変換（SIFT）分析を伴うことができ、分析されるべき画像のキーポイントおよびエリアは、ボックス 61、62 からのオーディオ導出情報に基づいて狭められる。例示的な SIFT 分析方法の詳細については、本明細書において以下で図 7 ~ 図 17 に関して開示する。

40

【0028】

[0073] 視覚特徴分析の結果（ボックス 70）は、シーン中の音源（たとえば、楽器）候補とそれらの対応するロケーションとのリストであり、そのリストはオーディオ処理ブロック 56 に与えられる。

【0029】

[0074] オーディオ処理ブロック 56 は、記録されたオーディオの品質を向上させるために、別個の音源がより良く分離され、識別され、処理され得るように、シーンから記録されたオーディオをさらに分析してオーディオを分解する。ボックス 72 において、画像 / ビデオ処理ブロック 52 からのロケーション情報を使用して、識別された音源ロケーショ

50

ンサブセクタの各々のほうへそれぞれ向けられたマルチマイクロフォンアレイのための空間フィルタを発生する。これは、記録されたオーディオデータ中の音源を分離するのを支援する。ボックス74において、楽器音源の識別を改善するために、シングルチャネル基底関数インベントリベースのスパース復元技法が、分離された音サブセクタの各々に適用される。信号チャネル復元技法は、基底関数インベントリを低減するために楽器カテゴリーノートの低減されたセットを使用することができ、この低減は、画像/ビデオ処理ブロック54によって与えられた楽器候補のリスト(ボックス70)によって誘導される。ボックス70において使用され得る例示的なスパース復元技法については、本明細書において以下で図30~図37に関して説明する。

【0030】

[0075]図3は、図1の聴覚シーン分析システム12を使用して聴覚シーンを分解する例示的な方法を示すフローチャート200である。ステップ202において、システム12がオーディオおよび視覚情報(静止画像および/またはビデオ)を記録する。ステップ204において、オブジェクト認識サブシステム20がシーン10中の音生成オブジェクトのうちの1つまたは複数を識別し、その位置を特定する。ステップ206において、音響分解サブシステム22は音響シーンを別個の音源に分解する。ステップ208において、音響分解サブシステム22は、分離された音に信号チャネル基底関数インベントリベースのスパース復元を適用する。

【0031】

[0076]図4は、聴覚シーン中の音生成オブジェクトを識別する第1の例示的な方法を示すフローチャート300である。この方法はオブジェクト認識サブシステム20によって実行され得る。ステップ302において、キャプチャされた画像中のキーポイントを識別する。ステップ304において、シーンにおいて記録されたオーディオに基づいて、楽器などの音生成オブジェクトに対応する1つまたは複数のキーポイントシグネチャを選択する。ステップ306において、画像中のキーポイントを、選択されたキーポイントシグネチャと比較することによって、シーン中の少なくとも1つのオブジェクトを識別する。

【0032】

[0077]図5A~図5Bに、聴覚シーン中の音生成オブジェクトを識別する第2の例示的な方法のフローチャート400を示す。この方法はオブジェクト認識サブシステム20によって実行され得る。ステップ402において、キャプチャされた画像中のキーポイントを識別する。ステップ404において、識別されたキーポイントから安定したキーポイントを選択する。ステップ406において、シーンから記録されたオーディオに基づいて、シーンの画像中の関心領域(ROI: region of interest)を選択する。ステップ408において、ROI中の安定したキーポイントを選択する。

【0033】

[0078]ステップ410において、シーンのビデオから局所動きベクトル(LMV: local motion vector)を計算する。ステップ412において、ROI中のLMVを選択する。

【0034】

[0079]ステップ414において、シーンにおいて記録されたオーディオに基づいて、楽器などの音生成オブジェクトに対応する1つまたは複数のキーポイントシグネチャを選択する。

【0035】

[0080]ステップ416において、シーンからの記録されたオーディオに基づいてオーディオ信頼性値(CV: confidence value)を計算する。オーディオCVは、オーディオ特徴整合デバイス、たとえば、MFC分類器の出力に基づき得る。オーディオCVはベクトルであり得、ベクトルの各要素は、オブジェクトが特定のタイプのオブジェクト、たとえば、トランペット、ピアノなどである尤度を示す。

【0036】

[0081]ステップ418において、シーンのキャプチャされたデジタル画像に基づいて画

10

20

30

40

50

像信頼性値 (C V) を計算する。画像 C V は、整合デバイス、たとえば、S I F T 整合デバイスの出力に基づき得る。S I F T 整合デバイスは、画像 C V を生成するために、R O I 中の安定したキーポイントを、選択されたキーポイントシグネチャと比較する。画像 C V はベクトルであり得、ベクトルの各要素は、オブジェクトが特定のタイプのオブジェクト、たとえば、トランペット、ピアノなどである尤度を示す。

【0037】

[0082] ステップ 420 において、シーンからの記録されたビデオに基づいてビデオ信頼性値 (C V) を計算する。ビデオ C V は、R O I 中で選択された L M V を比較するヒストグラム整合プロセスの出力に基づき得る。ビデオ C V はベクトルであり得、ベクトルの各要素は、オブジェクトが特定のタイプのオブジェクト、たとえば、トランペット、ピアノなどである尤度を示す。

10

【0038】

[0083] オーディオ C V、画像 C V およびビデオ C V はそれぞれ正規化され得る。

【0039】

[0084] ステップ 422 において、オーディオ C V と画像 C V とビデオ C V とに基づいてシーン中のオブジェクトを識別する。たとえば、最終 C V は、オーディオ C V と画像 C V とビデオ C V との重み付き和として計算され得る。各 C V の重み付け係数は、それぞれの録音モダリティの信号対雑音比 (S N R) に基づくことができ、特に現在の録音フレームの S N R の関数であり得る。

【0040】

20

[0085] モダリティ C V がベクトルである場合、最終 C V もベクトルであり、ベクトルの各要素は、オブジェクトが特定のタイプのオブジェクト、たとえば、トランペット、ピアノなどである尤度を示す。最大尤度を示す要素がオブジェクトを識別する。

【0041】

[0086] 図 6 は、図 1 の聴覚シーン分析システム 12 とともに使用され得る例示的なオブジェクト認識サブシステム 500 のいくつかの構成要素を示すブロック図である。サブシステム 500 は、オーディオプロセッサ 502 と、画像プロセッサ 504 と、ビデオプロセッサ 506 と、S I F T 整合デバイス 532 と、キーポイントシグネチャデータベース (D B) 534 と、音響特徴データベース 536 と、音響特徴整合デバイス 538 と、ヒストグラム整合デバイス 540 と、オブジェクト局所動きベクトル (L M V) ヒストグラムデータベース 542 と、マルチモーダル分類器 544 とを含む。

30

【0042】

[0087] オーディオプロセッサ 502 は、シーンにおいてマイクロフォンアレイ 18 からオーディオ信号を受信し、記録する。画像プロセッサ 504 は、シーンのピクチャを撮っている 1 つまたは複数のカメラ 508 から、シーンの 1 つまたは複数の画像を受信し、記録する。ビデオプロセッサ 506 は、シーンを記録している 1 つまたは複数のビデオカメラ 510 から、ビデオ信号を受信し、記録する。

【0043】

[0088] オーディオプロセッサ 502 は、到来方向 (D O A) 検出器 512 と、関心領域 (R O I) セレクタ 514 と、音分類器 516 と、音響特徴抽出器 518 とを含む。マイクロフォンアレイ 18 から受信されたマイクロフォン信号から、D O A 検出器 512 は、シーン内に位置する音源から出ている音の到来方向を判断する。D O A 検出器 512 の例示的な構成要素および機能については、本明細書において図 18 ~ 図 25 に関して説明する。D O A とアレイの位置とから、シーン中の音源のロケーションの推定が判断され得る。この D O A 情報は R O I セレクタ 514 に受け渡される。R O I セレクタ 514 は、D O A 情報とマイクロフォンアレイ 18 の既知の位置とに基づいて音源のロケーションを推定する。R O I セレクタ 514 は、次いで、ロケーション情報に基づいてシーンの画像の特定の部分を選択する。選択された部分または R O I は、音源を含んでおり、したがって、キーポイント探索と L M V 計算をシーンの一部分のみに制限するために使用され得る。

40

【0044】

50

[0089]音分類器 5 1 6 は、記録されたオーディオの特性に基づいて音源のタイプを分類する。たとえば、音源として楽器のタイプを識別するために、分類器 5 1 6 によって楽器ノートライブラリが使用され得る。

【 0 0 4 5 】

[0090]音分類器 5 1 6 の出力はオーディオ信頼性値であり、それはキーポイントシグネチャデータベース 5 3 4 への入力として与えられる。オーディオ信頼性値に基づいてキーポイントシグネチャデータベース 5 3 4 から 1 つまたは複数のキーポイントシグネチャが選択される。これらの選択されたキーポイントシグネチャは S I F T 整合デバイス 5 3 2 に与えられる。

【 0 0 4 6 】

[0091]音響特徴抽出器 5 1 8 は、M F C C など、マイクロフォン信号から導出された音響特性を計算する。これらの抽出された特徴は音響特徴整合デバイス 5 3 8 に与えられ、音響特徴整合デバイス 5 3 8 は、抽出された特徴を様々なタイプの音源の音響特徴のデータベース 5 3 6 と比較することによって音源を識別する。音響特徴整合デバイスの出力は音響特徴信頼性値であり得、この音響特徴信頼性値は、他の C V について上記で説明したのと同様の要素を有するベクトルであり得る。

【 0 0 4 7 】

[0092]画像プロセッサ 5 0 4 は、キーポイント検出器 5 2 0 と、安定キーポイント検出器 5 2 2 と、R O I キーポイントセクタ 5 2 4 とを含む。キーポイント検出器 5 2 0 は、本明細書で説明する方法を使用して、シーンのキャプチャされたデジタル画像中のキーポイントを判断する。安定キーポイント検出器 5 2 2 は、キーポイント探索を改善し、安定しているそれらの検出されたキーポイントのみを選択する。R O I キーポイントセクタ 5 2 4 は、R O I セクタ 5 1 4 から、キャプチャされた画像中の R O I を識別する座標情報を受信する。この座標情報に基づいて、R O I キーポイントセクタは、画像キーポイント選択を、R O I 内に位置するそれらの安定したキーポイントに狭める。

【 0 0 4 8 】

[0093]R O I 内で検出された安定したキーポイントは、次いで、S I F T 整合デバイス 5 3 2 に与えられる。

【 0 0 4 9 】

[0094]本質的に、S I F T 整合デバイス 5 3 2 は、画像 C V を発生するために、安定した R O I キーポイントを、キーポイントシグネチャデータベース 5 3 4 から取り出されたキーポイントシグネチャと比較する。

【 0 0 5 0 】

[0095]ビデオプロセッサ 5 0 6 は、L M V 計算器 5 2 6 と、R O I L M V セクタ 5 2 8 と、R O I L M V ヒストグラム計算器 5 3 0 とを含む。L M V 計算器 5 2 6 は、(1 つまたは複数の) カメラ 5 1 0 からデジタルビデオ信号を受信し、シーンの所定の録音持続時間について L M V を計算する。L M V は、次いで、R O I L M V セクタ 5 2 8 に受け渡される。R O I L M V セクタ 5 2 8 は、R O I セクタ 5 1 4 から R O I の座標情報を受信し、その座標情報に基づいて R O I 内のそれらの L M V を選択する。

【 0 0 5 1 】

[0096]R O I 内の L M V は、次いで、R O I L M V ヒストグラム計算器 5 3 0 に受け渡され、R O I L M V ヒストグラム計算器 5 3 0 は、R O I から L M V ヒストグラムを計算する。シーンの L M V ヒストグラムは、次いで、ヒストグラム整合デバイス 5 4 0 に受け渡される。ヒストグラム整合デバイス 5 4 0 は、最も近接した整合を見つけるために、シーン L M V ヒストグラムを、オブジェクト L M V ヒストグラムデータベース 5 4 2 に記憶されたオブジェクト L M V ヒストグラムと比較する。ヒストグラム整合デバイス 5 4 0 は、この比較に基づいてビデオ C V を出力する。

【 0 0 5 2 】

[0097]マルチモーダル分類器 5 4 4 は、S I F T 整合デバイス 5 3 2 と、音分類器 5 1 6 と、音響特徴整合デバイス 5 3 8 と、ヒストグラム整合デバイス 5 4 0 との出力に基づ

10

20

30

40

50

いてシーン中のオブジェクトを識別する。マルチモーダル分類器 5 4 4 は、オーディオ C V と画像 C V とビデオ C V と音響特徴 C V との重み付き和であり得る、最終信頼性値ベクトルを計算することによってこれを達成することができる。分類器 5 4 4 は、認識された楽器を関連する尤度で与えるためにベイズ推定を実行し得る。C V の重み付け係数は、図 4 A ~ 図 4 B に関して説明したものと同様であり得る。

【 0 0 5 3 】

[0098]さらに、サブシステム 5 0 0 はまた、シーン中の認識された各オブジェクトについて改善されたオブジェクトロケーションを出力し得る。改善された（１つまたは複数の）オブジェクトロケーションは、マルチモーダル分類器 5 4 4、カメラ 5 0 8 から出力と、オーディオプロセッサ 5 0 2 の R O I セレクタ 5 1 4 からの推定オブジェクトロケーションとに基づくことができる。改善された（１つまたは複数の）オブジェクトロケーションは、関心領域またはオブジェクトロケーションを推定する際のそれらの精度および速度を改善するために D O A 検出器 5 1 2 および / または R O I セレクタ 5 1 4 にフィードバックされ得、たとえば、前のビデオ / 画像フレームにおいて判断された推定 D O A またはオブジェクトロケーションは、オーディオプロセッサ 5 0 2 がその R O I 選択プロセスにおいて使用する初期座標として次のフレームに手渡され得る。

【 0 0 5 4 】

キーポイント選択および S I F T 整合デバイス

[0099]例示的なキーポイント検出器 5 2 0、キーポイントセレクタ 5 2 2 および S I F T 整合デバイス 5 3 2 の動作について以下のように説明する。

【 0 0 5 5 】

[00100]概して、オブジェクトまたは特徴認識は、オブジェクト認識のために画像中の関心ポイント（キーポイントとも呼ばれる）を識別することおよび / またはそれらのキーポイントの周りの局所的特徴を識別することを伴い得る。画像データ中のそのような特徴的な要素を本明細書では「キーポイント」と呼ぶが、本明細書で使用するキーポイントという用語は、個々のピクセル、ピクセルのグループ、分数ピクセル部分、１つまたは複数の記述子、他の画像成分、あるいはそれらの任意の組合せを指し得ることを理解されたい。特徴の高い安定性および再現性を有することは、これらの認識アルゴリズムでは非常に重要である。したがって、キーポイントは、それらが画像スケール変化および / または回転に対して不変であり、ひずみ、視点の変化、および / または雑音および照度の変化の実質的な範囲にわたってロバスト（robust）な整合を与えるように、選択および / または処理され得る。さらに、オブジェクト認識などのタスクに好適であるように、特徴記述子は、好ましくは、単一の特徴が、複数のターゲット画像からの特徴の大規模データベースに対して高い確率で正しく整合され得るという意味において特徴的であり得る。

【 0 0 5 6 】

[00101]画像中のキーポイントが検出され、位置を特定された後、それらのキーポイントは、様々な記述子を使用することによって識別または記述され得る。たとえば、記述子は、画像特性の中でも、形状、色、テクスチャ、回転、および / または動きなど、画像中のコンテンツの視覚特徴を表し得る。キーポイントに対応し、記述子によって表される個々の特徴は、次いで、既知のオブジェクトからの特徴のデータベースに整合される。

【 0 0 5 7 】

[00102]画像のためのキーポイントを識別し、選択することの一部として、選択されたいくつかのポイントは、精度または信頼性の欠如により廃棄される必要があり得る。たとえば、いくつかの初期に検出されたキーポイントは、エッジ沿いの不十分なコントラストおよび / または不十分な定位を理由に拒否され得る。そのような拒否は、照度、雑音および配向の変動に対してキーポイント安定性を増加させる際に重要である。また、特徴整合の再現性を減少させ得る、誤ったキーポイント拒否を最小限に抑えることが重要である。

【 0 0 5 8 】

[00103]概して、画像中の照度は、空間的に変動する関数によって表され得ることを認識されたい。したがって、照度の影響（たとえば、シェーディング、明るい画像、暗い画

10

20

30

40

50

像など)は、照度関数を排除する正規化プロセスによって特徴/キーポイント検出のために無効にされ得る。たとえば、画像は、画像の平滑化スケール空間(smoothened scale space) L を発生するために平滑化ファクタの範囲をもつ関数 G (すなわち、カーネルまたはフィルタ)を使用して、画像を漸進的(progressively)にぼかすことによって処理され得る。次いで、平滑化スケール空間レベルの隣接するペア間の差分($L_i - L_{i-1}$)を取ることによって、画像のためのスケール空間差分 D が取得され得る。次いで、スケール空間レベルの特定の差分 D_i を取得するために使用されるスケール空間レベル L_i のうち最も平滑なスケール空間レベルと同程度に平滑であるかそれよりも平滑である平滑化スケール空間レベル L_k でスケール空間レベルの各差分 D_i を除算することによって、スケール空間 L の差分の正規化が達成される。

10

【0059】

[00104]図7は、記録された画像上でオブジェクト認識を実行するための機能段階を示すブロック図である。画像キャプチャ段階702において、関心画像708(すなわち、記録画像)がキャプチャされ得る。画像708は、デジタルキャプチャ画像を取得するために、1つまたは複数の画像センサーおよび/またはアナログデジタル変換器を含み得る、画像キャプチャデバイスによってキャプチャされ得る。画像センサー(たとえば、電荷結合デバイス(CCD)、相補型金属半導体(CMOS))は光を電子に変換し得る。電子はアナログ信号を形成し得、次いで、そのアナログ信号は、アナログデジタル変換器によってデジタル値に変換される。このようにして、画像 $I(x, y)$ を、たとえば、対応する色、照度、および/または他の特性をもつ複数のピクセルとして定義し得るデジタルフォーマットで画像708はキャプチャされ得る。

20

【0060】

[00105]画像処理段階704において、キャプチャされた画像708は、次いで、対応するスケール空間710(たとえば、ガウススケール空間)を発生し、特徴検出712を実行し、特徴記述子抽出716を実行することによって処理される。特徴検出712は、キャプチャされた画像708について高度に特徴的なキーポイントおよび/または幾何学的関心キーポイントを識別し得、それらのキーポイントは、その後、特徴記述子抽出716において複数の記述子を取得するために使用され得る。画像比較段階706において、これらの記述子は、既知の記述子のデータベースとの(たとえば、キーポイントおよび/またはキーポイントの他の特性あるいはキーポイントを囲むパッチを比較することによる)特徴整合722を実行するために使用される。次いで、特徴整合が正しいことを確認するために、キーポイント整合に対する幾何学的検証または一貫性検査724が実行されて、整合結果726が与えられる。このようにして、記録画像が、ターゲット画像のデータベースと比較されおよび/またはそれから識別され得る。

30

【0061】

[00106]画像中の照度の変化は、画像のための特徴/キーポイント認識の安定性および/または再現性に有害な影響を及ぼし得ることが観測されている。すなわち、画像中の局所および/または大域(global)照度変化は、画像のための特徴/キーポイントの検出に影響を及ぼすことがある。たとえば、特徴/キーポイントの数および/またはロケーションが、画像中の照度(たとえば、シェーディング、コントラストなど)に応じて変化し得る。したがって、画像中の特徴/キーポイント検出より前に、局所および/または大域照度変化の影響を実質的になくすかまたは最小限に抑えることが有益であろう。

40

【0062】

[00107]これを行うための1つの方法は、特徴/キーポイント検出を開始するより前に、局所および/または大域照度変化を除去または補償するように画像自体を処理することであり得る。しかしながら、そのようなプロセスは計算集約的であり得る。さらに、画像中に局所および/または大域照度変化が存在するかどうかを判断することがしばしば困難である。そのようなプロセスは、データベース中の画像にも適用されなければならないであろう。照度変化を補正するためにターゲット画像とデータベース画像の両方を最初に処理することなしには、特徴/キーポイント整合は成功しないことがある。しかし、照度が

50

特定の画像にどのように影響を及ぼし得るかの事前知識なしには、このプロセスは自動的に実装することがかなり困難である。

【 0 0 6 3 】

[00108]したがって、実質的な処理オーバーヘッドなしに実行され得る代替形態が必要とされる。一例によれば、特徴検出の目的での画像上の（一様あるいは非一様な）照度の影響は、スケール空間差分に特徴ノキーポイント検出が実行されるより前にスケール空間差分を正規化することによって、なくされるかまたは低減され得る。この正規化プロセスは、すでに利用可能である平滑化スケール空間を使用して実行され、したがって、追加の計算が最小限に抑えられ得る。

【 0 0 6 4 】

[00109]一例によれば、スケール空間正規化器 7 1 4 は、照度変化が画像中のキーポイントノ特徴認識に及ぼす影響を低減するかまたはなくすために、スケール空間発生 7 1 0 の一部として実装され得る。

【 0 0 6 5 】

[00110]図 8 に、例示的な画像処理段階 7 0 4 におけるガウススケール空間発生を示す。画像中の特徴検出を実行するために、スケール不変特徴変換（SIFT）など、いくつかのアルゴリズムが開発されている。画像中の特定のオブジェクトの検出への第 1 のステップは、その局所特徴に基づいて記録されたオブジェクトを分類することである。その目的は、たとえば、照度、画像雑音、回転、スケーリング、およびノまたは視点の小さい変化に対して不変およびノまたはロバストである特徴を識別し、選択することである。すなわち、クエリ画像と比較ターゲット画像との間に照度、画像雑音、回転、スケール、およびノまたは視点の差があるにもかかわらず、これらの 2 つの画像間の整合が発見されなければならない。これを行うための 1 つの方法は、高度に特徴的な特徴（たとえば、画像中の特徴的なポイント、ピクセル、およびノまたは領域）を識別するために画像のパッチ上の極値検出（たとえば、極大値または極小値）を実行することである。

【 0 0 6 6 】

[00111]SIFT は、照度、画像雑音、回転、スケーリングの変化、およびノまたは視点の小さい変化に対して適度に不変である局所特徴を検出し、抽出するための 1 つの手法である。SIFT の画像処理段階 7 0 4 は、（a）スケール空間極値検出、（b）キーポイント定位、（c）配向割当て、およびノまたは（d）キーポイント記述子の発生を含み得る。特に、高速ロバスト特徴（SURF：Speed Up Robust Features）、勾配位置および配向ヒストグラム（LOH：Gradient Location and Orientation Histogram）、局所エネルギーベース形状ヒストグラム（LESH：Local Energy based Shape Histogram）、勾配の圧縮ヒストグラム（CHoG：Compressed Histogram of Gradients）を含む、特徴検出と、後続の特徴記述子発生とのための代替アルゴリズムは、本明細書で説明する特徴からも恩恵を受け得ることが明らかなはずである。

【 0 0 6 7 】

[00112]ガウススケール空間発生 7 1 0 において、デジタル画像 $I(x, y)$ 7 0 8 は漸進的にガウスぼかし（すなわち、平滑化）されて、ガウスピラミッド 7 5 2 が構成される。ガウスぼかし（平滑化）は、概して、元の画像 $I(x, y)$ をスケール c_s におけるガウスぼかしノ平滑化関数 $G(x, y, c_s)$ で畳み込み、したがって、ガウスぼかしノ平滑化関数 $L(x, y, c_s)$ が $L(x, y, c_s) = G(x, y, c_s) * I(x, y)$ として定義されることを伴う。ここで、 G はガウスクアーネルであり、 c_s は、画像 $I(x, y)$ をぼかすために使用されるガウス関数の標準偏差を示す。乗数 c が変化すると、 $(c_0 < c_1 < c_2 < c_3 < c_4)$ 、標準偏差 c_s は変化し、漸進的なぼかしノ平滑化が得られる。シグマ s は、ベーススケール変数（たとえば、ガウスクアーネルの幅）である。高いスケール（すなわち、低い解像度）ほど、より低いスケール（すなわち、より高い解像度）よりもぼかされるノ平滑化される。したがって、広いスケールレベル（すなわち、低い解像度）ほど、画像はより平滑になる（よりぼかされる）。

【 0 0 6 8 】

[00113]ぼけた画像 L を生成するために初期画像 $I(x, y)$ がガウシアン G で増分的に畳み込まれるとき、ぼけた画像 L は、スケール空間において定数ファクタ c だけ分離される。ガウスぼかしされた（平滑化された）画像 L の数が増加し、ガウスピラミッド 7 5 2 のために与えられる近似が連続空間に近づくにつれて、これらの 2 つのスケールも 1 つのスケールに近づく。一例では、畳み込まれた画像 L はオクターブによってグループ化され得、1 オクターブは、標準偏差 s の値の倍化に対応し得る。その上、乗数 c （たとえば、 $c_0 < c_1 < c_2 < c_3 < c_4 \dots$ ）の値は、固定数の畳み込まれた画像 L がオクターブごとに取得されるように選択される。スケーリングの各オクターブは、明示的な画像サイズ変更に対応する。したがって、元の画像 $I(x, y)$ が漸進的ぼかし / 平滑化関数によってぼかされる / 平滑化されるにつれて、ピクセルの数は漸進的に低減される。本明細書では説明のためにガウス平滑化関数を使用した、他のタイプの平滑化カーネル / 関数が採用され得ることに留意されたい。

10

【0069】

[00114]ガウスピラミッド 7 5 2 中の任意の 2 つの連続するガウスぼかし画像の差分を計算することによって、ガウス差分（D o G : difference of Gaussian）ピラミッド 7 5 4 が構成される。D o G 空間 7 5 4 において、 $D(x, y, a) = L(x, y, c_n s) - L(x, y, c_{n-1} s)$ である。D o G 画像 $D(x, y, s)$ は、スケール $c_n s$ および $c_{n-1} s$ における 2 つの隣接するガウスぼかし画像 L 間の差分である。 $D(x, y, s)$ のスケールは、 $c_n s$ と $c_{n-1} s$ との間のどこかにある。D o G 画像 D が、オクターブごとに隣接するガウスぼかし画像 L から取得され得る。各オクターブ後に、ガウス画像が 2 分の 1 にダウンサンプリングされ、次いでこのプロセスが繰り返される。このようにして、画像は、並進、回転、スケール、および / または他の画像パラメータおよび / またはひずみに対してロバストまたは不変である局所特徴に変換され得る。

20

【0070】

[00115]記録画像の D o G 空間 7 5 4 は、発生されると、関心特徴を識別する（たとえば、画像中の高度に特徴的なポイントを識別する）ための極値検出のために利用され得る。これらの高度に特徴的なポイントは、本明細書ではキーポイントと呼ばれる。これらのキーポイントは、各キーポイントを囲むパッチまたは局所領域の特性によって識別され得る。記述子が、キーポイントおよびその対応するパッチごとに生成され得、それは、クエリ画像と記憶されたターゲット画像との間のキーポイントの比較のために使用され得る。「特徴」は、記述子（すなわち、キーポイントおよびその対応するパッチ）を指し得る。特徴（すなわち、キーポイントおよび対応するパッチ）のグループはクラスタと呼ばれることがある。

30

【0071】

[00116]図 9 に、例示的な画像処理段階 7 0 4 における特徴検出を示す。特徴検出 7 1 2 において、D o G 空間 7 5 4 を使用して画像 $I(x, y)$ のキーポイントを識別し得る。特徴検出 7 1 2 は、画像中の特定のサンプルポイントがピクセルの周りの局所領域またはパッチが、（幾何学的に言って）潜在的に関心のあるパッチであるかどうかを判断しようとする。

【0072】

[00117]概して、D o G 空間 7 5 4 中に極大値および / または極小値が識別され、これらの極大値および極小値のロケーションが D o G 空間 7 5 4 中のキーポイントロケーションとして使用される。図 9 に示す例では、キーポイント 7 6 0 はパッチ 7 5 8 で識別されている。極大値および極小値を発見すること（局所的極値検出としても知られる）は、D o G 空間 7 5 4 中の各ピクセル（たとえば、キーポイント 7 6 0 に対するピクセル）を、その 8 つの隣接するピクセルと、同じスケールで比較し、ならびに（隣接するパッチ 7 5 6 および 7 6 2 中の）9 つの隣接するピクセルと、キーポイント 8 0 8 の両側に隣接するスケールの各々で比較し、合計 26 個のピクセル（ $9 \times 2 + 8 = 26$ ）に対して比較することによって達成され得る。ここで、パッチは 3×3 ピクセル領域として定義される。概して、キーポイント 7 5 8 に対するピクセル値が、パッチ 7 5 8、7 5 6、および 7 6

40

50

0 中のすべての 2 6 個の比較されたピクセルの間で最大値または最小値である場合、それがキーポイントとして選択される。キーポイントは、それらのロケーションがより正確に識別されるようにさらに処理され得、低コントラストキーポイントおよびエッジキーポイントなど、キーポイントのうちのいくつかが廃棄され得る。

【 0 0 7 3 】

[00118] 図 1 0 に、例示的な画像処理段階 7 0 4 における特徴記述子抽出を示す。概して、特徴（たとえば、キーポイントおよびその対応するパッチ）は記述子によって表され得、記述子は、（クエリ画像からの）特徴と、ターゲット画像のデータベースに記憶された特徴との効率的な比較を可能にする。特徴記述子抽出 7 1 6 の一例では、各キーポイントは、局所画像勾配の方向に基づいて、1 つまたは複数の配向、または方向を割り当てられ得る。局所画像特性に基づいて各キーポイントに一貫した配向を割り当てることによって、キーポイント記述子は、この配向に対して表され、したがって、画像回転に対する不変性を達成することができる。ガウスぼかし画像 L 中でおよび / またはキーポイントスケールにおいて、キーポイント 7 6 0 の周りの隣接する領域中のピクセルごとに大きさおよび方向の計算が実行され得る。 (x, y) に位置するキーポイント 7 6 0 に対する勾配の大きさは $m(x, y)$ として表され得、 (x, y) におけるキーポイントに対する勾配の配向または方向は $\Gamma(x, y)$ として表され得る。キーポイントのスケールを使用して、すべての計算がスケール不変方式で実行されるように、キーポイント 7 6 0 のスケールに最も近いスケールで、ガウス平滑化された画像 L を選択する。各画像サンプル $L(x, y)$ について、このスケールで、勾配の大きさ $m(x, y)$ と配向 $\Gamma(x, y)$ とが、ピクセル差分を使用して計算される。たとえば、大きさ $m(x, y)$ は次のように計算され得る。

【 数 1 】

$$m(x, y) = \sqrt{(L(x+1, y) - L(x-1, y))^2 + (L(x, y+1) - L(x, y-1))^2} \quad (\text{式1})$$

【 0 0 7 4 】

方向または配向 $\Gamma(x, y)$ は次のように計算され得る。

【 数 2 】

$$\Gamma(x, y) = \arctan \left[\frac{(L(x, y+1) - L(x, y-1))}{(L(x+1, y) - L(x-1, y))} \right] \quad (\text{式2})$$

【 0 0 7 5 】

ここで、 $L(x, y)$ は、キーポイントのスケールでもあるスケール s における、ガウスぼかし画像 $L(x, y, s)$ のサンプルである。

【 0 0 7 6 】

[00119] キーポイント 7 6 0 に対する勾配は、D o G 空間中のキーポイントの平面より上に、より高いスケールで存在するガウスピラミッド中の平面に対して、またはキーポイントより下に、より低いスケールで存在するガウスピラミッドの平面中のいずれかで、一貫して計算され得る。どちらにしても、各キーポイントについて、勾配は、キーポイントを囲む矩形エリア（たとえば、パッチ）中ですべて 1 つの同じスケールで計算される。その上、画像信号の周波数は、ガウスぼかし画像のスケールに反映される。しかし、S I F T は、単にパッチ（たとえば、矩形エリア）中のすべてのピクセルにおいて勾配値を使用する。パッチがキーポイントの周りで定義され、サブブロックがブロック内で定義され、

サンプルがサブブロック内で定義され、この構成は、キーポイントのスケールが異なるときでさえ、すべてのキーポイントに対して同じままである。したがって、画像信号の周波数は、同じオクターブ中のガウス平滑化フィルタの連続適用とともに変化する一方で、異なるスケールにおいて識別されたキーポイントは、スケールで表される画像信号の周波数の変化にかかわらず、同じサンプル数でサンプリングされ得る。

【0077】

[00120]キーポイント配向を特徴づけるために、(SIFTでは)キーポイント760の近傍において(キーポイントのスケールに最も近接したスケールにおけるガウス画像を使用して)勾配配向のベクトルが生成され得る。しかしながら、キーポイント配向は、たとえば、勾配の圧縮ヒストグラム(CHoG)を使用することによって、勾配配向ヒストグラム(図10参照)によっても表され得る。隣接する各ピクセルの寄与は、勾配の大きさとガウス窓とによって重み付けされ得る。ヒストグラムのピークは支配的な配向に対応する。キーポイントのすべての特性はキーポイント配向に対して測定され得、これにより、回転に対する不変性が与えられる。

【0078】

[00121]一例では、各ブロックについてガウス重み付け勾配の分布が計算され得、各ブロックは、2サブブロック×2サブブロックで合計4サブブロックである。ガウス重み付け勾配の分布を計算するために、いくつかのピンをもつ配向ヒストグラムが形成され、各ピンはキーポイントの周りのエリアの部分のカバーする。たとえば、配向ヒストグラムは36個のピンを有し得、各ピンは配向の360度範囲のうちの10度をカバーする。代替的に、ヒストグラムは8つのピンを有し得、各ピンは360度範囲のうちの45度をカバーする。本明細書で説明するヒストグラムコーディング技法は、任意の数のピンのヒストグラムに適用可能であることが明らかかなはずである。ヒストグラムを最終的に生成する他の技法も使用され得ることに留意されたい。

【0079】

[00122]勾配分布および配向ヒストグラムは様々な方法で取得され得る。たとえば、2次元勾配分布(d_x, d_y) (たとえば、ブロック806)が1次元分布(たとえば、ヒストグラム814)に変換される。キーポイント760は、キーポイント760を囲むパッチ806(セルまたは領域とも呼ばれる)の中心に位置する。各レベルのピラミッドについて事前計算された勾配が、各サンプルロケーション808において小さい矢として示されている。図示のように、サンプル808の4×4領域はサブブロック810を形成し、サブブロックの2×2領域はブロック806を形成する。ブロック806は記述子窓と呼ばれることもある。ガウス重み付け関数は、円802で示され、各サンプルポイント808の大きさに重みを割り当てるために使用される。円形窓802中の重みは平滑に低下する。ガウス窓802の目的は、窓の位置の小さな変化によって記述子が突然変化することを回避し、記述子の中心から遠い勾配にあまり重点を与えないことである。2×2サブブロックから配向ヒストグラム812の2×2=4アレイが取得され、ヒストグラムの各ピン中に8つの配向があり、それにより(2×2)×8=32次元の特徴記述子ベクトルが得られる。たとえば、配向ヒストグラム813および815は、サブブロック810に対する勾配分布に対応し得る。しかしながら、各ヒストグラム中に8つの配向をもつヒストグラム(8ピンヒストグラム)の4×4アレイを使用して、それにより各キーポイントについて(4×4)×8=128次元の特徴記述子ベクトルが得られると、より良好な結果が与えられ得る。勾配分布を取得するために、他のタイプの(たとえば、異なるボロノイセル構造を用いた)量子化ピンコンストラクションも使用され得ることに留意されたい。

【0080】

[00123]本明細書で使用するヒストグラムは、ピンとして知られている様々な独立したカテゴリーに分類される観測、サンプル、または出現(たとえば、勾配)の数を計数するマッピング k_i である。ヒストグラムのグラフは、ヒストグラムを表すための1つの方法にすぎない。したがって、 k が観測、サンプル、または出現の総数であり、 m がピンの総

10

20

30

40

50

数である場合、ヒストグラム k_i における周波数は下記の条件を満たす。

【数 3】

$$n = \sum_{i=1}^m k_i \quad (\text{式3})$$

【0081】

ただし、 \sum は総和演算子である。

【0082】

[00124] キーポイントに対する特徴記述子ベクトルを取得するために、サブブロックからのヒストグラムは連結され得る。16個のサブブロックからの8ピンヒストグラム中の勾配が使用される場合、128次元の特徴記述子ベクトルが得られ得る。

10

【0083】

[00125] このようにして、記述子は、識別されたキーポイントごとに取得され得、そのような記述子は、ロケーション (x, y) と、配向と、ガウス重み付け勾配の分布の記述子とによって特徴づけられ得る。画像は、1つまたは複数のキーポイント記述子 (画像記述子とも呼ばれる) によって特徴づけられ得ることに留意されたい。さらに、記述子はまた、ロケーション情報 (たとえば、キーポイントの座標)、スケール (たとえば、キーポイントが検出されたガウススケール)、およびクラスタ識別子などの他の情報などを含み得る。

20

【0084】

[00126] ガウス差分空間 754 中で演算することによって、画像のルミナンスのいかなるレベルシフト (ルミナンスへの空間的に一様な加法的バイアス) も完全に無視される。しかし、ルミナンスのスケールシフトは、キーポイントが判定され、最終的に選択または拒否される方法に影響を及ぼす。これは、一様な乗法的ルミナンスファクタ、ならびに空間的に変動する乗法的ルミナンスファクタの両方に当てはまる。キーポイント検出とまさに同程度に重要であるのが、画像内のその定位である。オブジェクトは、その特徴の幾何学的コンテンツと、それらの空間的相互関係とによってカテゴリー分類される。したがって、キーポイントが検出された場合でも、その定位はルミナンススケール変化に関して不変の方法で計算されるべきであるように、キーポイントの計算されたロケーション

30

【0085】

[00127] したがって、キーポイントを識別し、記述子を生成するより前に、キーポイントが検出されたスケール空間から照度の影響を低減、除去、および / またはフィルタ処理するために、1つの特徴がガウス差分空間 754 を正規化することを行う。

【0086】

スケール空間正規化の例示的な差分

[00128] 図 11 に、照度不変特徴 (illumination invariant feature) / キーポイント検出のためのスケール空間正規化の一例を示す。画像 $I(x, y)$ 822 は、平滑化されたスケール空間ピラミッド 826 を発生するために、異なるスケール c_i において平滑化カーネル $G(x, y, c_i)$ 824 で畳み込まれ得、ただし、 i は 0 と n との間の整数である。平滑化カーネルはガウスカーネルおよび / または他のタイプの平滑化関数であり得ることに留意されたい。スケール空間差分 828 を取得するために、平滑化されたスケール空間ピラミッド 826 の 2 つの隣接するスケール空間の間の差分が取られ得る。

40

【0087】

[00129] 最初に、スケール空間差分 828 の各レベルは、画像 $I(x, y)$ 822 で畳み込まれた異なるスケールにおける平滑化カーネル 824 の差分 (たとえば、 $G(x, y, c_{j+1}) - G(x, y, c_j)$) として定義され得ることがわかる。これは、2つの対応する平滑化スケール空間差分 (たとえば、 $L(x, y, c_{j+1}) - L(x, y, c_j)$) に等しい。したがって、2つの平滑化スケール空間差分は次のように表され得る。

50

【数 4】

$$D(x, y, \sigma) = (G(x, y, c_{j+1}\sigma) - G(x, y, c_j\sigma)) * I(x, y) = L(x, y, c_{j+1}\sigma) - L(x, y, c_j\sigma) \quad (\text{式4})$$

【0088】

[00130]また、照度がスケーリング関数 $S(x, y)$ として表される場合、2つの平滑化スケール空間差分に対する照度変化は次のように表され得ることがわかる。

10

【数 5】

$$D(x, y, \sigma) = (G(x, y, c_{j+1}\sigma) - G(x, y, c_j\sigma)) * (I(x, y)S(x, y)) \quad (\text{式5})$$

【0089】

ここで、一般的な場合、照度スケーリング関数 $S(x, y)$ は、空間的に変動するか、または空間的に一定であり得る。

20

【0090】

[00131]しかしながら、照度スケーリング関数 $S(x, y)$ を取得するために照度をランタイムでモデル化することは実際のおよび/または実現可能でない。したがって、本明細書では、特徴選択および/またはブルーニングが一般に実行される特徴空間（たとえば、DOG空間828）からの照度によってバイアスされない、基礎をなす特徴（たとえば、キーポイント）を引き出す代替手法が開示される。この代替手法によれば、画像 $I(x, y)$ 822 のルミナンス分布は、画像 $I(x, y)$ 822 から抽出されたスケール空間情報を利用することによって正規化される。照度に関する事前情報は必要とされない。この方法は、何らかの大きい計算および処理を導入することなしに、異なる照度変化にわたって一貫したレベルで安定した特徴を選定することと再現性を高めることを可能にする。

30

【0091】

[00132]これを行うために、特徴検出が行われるスケール空間差分828は、より広いスケール空間によって正規化され得る。この手法は次式によって定義され得る。

【数 6】

$$D'(x, y, \sigma) = \left[\frac{[G(x, y, c_{j+1}\sigma) - G(x, y, c_j\sigma)] * [I(x, y)S(x, y)]}{G(x, y, c_{j+1+h}\sigma) * [I(x, y)S(x, y)]} \right] \quad (\text{式6})$$

40

【0092】

ただし、

第1のガウス平滑化カーネル $G(x, y, c_{j+1})$ は第2のガウス平滑化カーネル $G(x, y, c_j)$ よりも広く（すなわち、スケール c_{j+1} はスケール c_j よりも広く、ただし、 j は0と n との間の正の整数であり）、

$I(x, y)$ は、処理されている画像またはその派生物（たとえば、画像の反射特性）であり、

$S(x, y)$ は照度スケーリング関数であり、

50

$G(x, y, c_{j+1+h})$ は、第 2 の平滑化カーネル $G(x, y, c_{j+1})$ と同程度に
 広いまたはそれよりも広いスケール空間を有する第 3 の平滑化カーネルであり、ただし
 、 h は 0 と n との間の正の整数である。スケール空間差分 8 2 8 のあらゆるレベル上でこ
 のプロセスを繰り返すことによって、正規化されたスケール空間 8 3 0 の差分が発生され
 得る。たとえば、 $G(x, y, c_{j+1})$ と $G(x, y, c_j)$ とによって定義される差
 分スケール空間では、正規化関数は $G(x, y, c_{j+1})$ またはそれより高い任意のも
 の（すなわち、 $G(x, y, c_{j+2})$ 、 $G(x, y, c_{j+3})$ 、 \dots ）であり得る。
 正規化関数は、差分スケール空間中で使用される両方の平滑化カーネルよりも大きい必要
 はなく、それは平滑器である必要のみがある。別の例では、正規化関数は、使用される第
 1 の平滑化カーネルと第 2 の平滑化カーネルとの和（すなわち、 $G(x, y, c_{j+1})$
 $+ G(x, y, c_j)$ ）であり得、したがって、

【数 7】

$$D'(x, y, \sigma) = \left[\frac{[G(x, y, c_{j+1}\sigma) - G(x, y, c_j\sigma)] * [I(x, y)S(x, y)]}{[G(x, y, c_{j+1}\sigma) + G(x, y, c_j\sigma)] * [I(x, y)S(x, y)]} \right] \quad (式7)$$

【0093】

[00133] 式 6 は次のようにも表され得ることに留意されたい。

【数 8】

$$D'(x, y, \sigma) = \left[\frac{[L(x, y, c_{j+1}\sigma) - L(x, y, c_j\sigma)] * S(x, y)}{L(x, y, c_{j+1+h}\sigma) * S(x, y)} \right] \quad (式8)$$

【0094】

照度スケーリング関数 $S(x, y)$ は（式 6、式 7 および式 8 の）分子と分母の両方に現
 れるので、そのスケーリングの影響は相殺される。すなわち、照度スケーリング関数 S
 (x, y) は正規化のために使用される平滑化画像 $L(x, y, c_{j+1+h}) * S(x, y)$ 中に存在するので、それは、スケール空間差分 $[L(x, y, c_{j+1}) - L(x, y, c_j)] * S(x, y)$ における照度スケーリング関数 $S(x, y)$ の影響を完全
 にまたは実質的に相殺する。前述のように、 $L(x, y, c_{j+1+h})$ は、 $L(x, y, c_{j+1})$ またはより高いスケール画像（すなわち、 $L(x, y, c_{j+2})$ 、 $L(x, y, c_{j+3})$ 、 \dots ）に等しくなり得る。このようにして、分母中の画像コンテンツは、それがごくわずかな空間アーティファクトしか導入しない程度まで平滑化される。

【0095】

[00134] スケール空間差分を正規化する際に、正規化する平滑化画像 $L(x, y, c_{j+1+h})$ は、（キーポイント / 特徴を識別する）局所的極値位置をシフトしないように、特徴空間（すなわち、スケール空間差分）をあまりに多く変化させないように選択されなければならない。すなわち、スケール不変特徴を達成するためにはスケール空間差分が最良であることが知られているので、スケール空間差分の密接な近似が正規化後に保持されなければならない。この目的で、平滑化画像 $L(x, y, c_{j+1+h})$ は、高周波数成分が平均されるようにそのスケールレベルが十分に平滑であるように選択される。すなわち、平滑化画像 $L(x, y, c_{j+1+h})$ が十分に平坦である場合、スケール空間 L の差分 $(x, y, c_{j+1}) - L(x, y, c_j)$ の形状はほとんど変化しない（すなわち、特徴 / キーポイントの位置は変化しない）。一実施形態では、正規化されている差分スケールレベルを取得するために使用されるスケールレベルに近接した（それと同じであるかまたはその次に最も高い）スケールレベルにおける正規化関数を選択することは、多すぎる雑音を導入することを回避するので、好適であり得ることに留意されたい。たとえば、

$G(x, y, c_{j+1})$ と $G(x, y, c_j)$ とによって定義される差分スケール空間のために $G(x, y, c_{j+1})$ のような平滑スケールを選ぶことによって、スケール空間中のその特定のレベルについて典型的な局所不規則性が維持され得る。

【0096】

[00135] 前記のように、画像中で検出される特徴の数は、画像の乗法的ルミナンススケール変化によって大幅に影響を受け得る（たとえば、低減され得る）。ルミナンスによって生じるスケーリングは、幾何学的変換がなくても最終の特徴空間中のコンテンツを大幅に低減する、画像上のマスクのように働く傾向がある。したがって、式6および式7の適用によって達成される正規化により、照度変化にかかわらず幾何学的有意性が「等しい」特徴が検出され、それによって再現性が増加することが保証される。

10

【0097】

[00136] 図12に、正規化されたスケール空間差分の1つのレベルがどのように取得され得るかの一例を示す。ここで、画像 $I(x, y)$ 852は、第1の平滑化スケール空間画像 $L(x, y, c_j)$ 858を取得するために、第1の平滑化カーネル $G(x, y, c_j)$ 854で畳み込まれ得る。画像 $I(x, y)$ 852はまた、第2の平滑化スケール空間画像 $L(x, y, c_{j+1})$ 860を取得するために、第2の平滑化カーネル $G(x, y, c_{j+1})$ 856で畳み込まれ得る。第2の平滑化画像860と第1の平滑化画像858との間の差分が取られて、スケール空間レベルの差分 $D_j(x, y,)$ 862が取得され得る。このスケール空間レベルの差分 $D_j(x, y,)$ 862は、より高いスケール平滑化カーネル $G(x, y, c_{j+1+h})$ 866または平滑化スケール空間画像 $L(x, y, c_{j+1+h})$ 868に基づいて（すなわち、式6および/または式7に従って）正規化されて、正規化スケール空間レベル $D'_j(x, y,)$ 864が取得され得る。このプロセスは、（スケーリングファクタ c_j によって設定される）異なる幅の異なる平滑化カーネルを画像 $I(x, y)$ に適用することによって繰り返され、それにより平滑化されたスケール空間ピラミッドが構築され得る。スケール空間差分（たとえば、図11中の828）は、平滑化されたスケール空間ピラミッド（たとえば、図11中の826）の隣接するレベル間の差分を取ることによって構築され得る。正規化されたスケール空間差分（たとえば、図11中の830）は、式6および/または式7に従って発生され得る。

20

【0098】

30

[00137] 図13に、照度の変化に対して耐性がある改善された特徴検出のために正規化されたスケール空間差分を発生するための方法を示す。902において、 $i = 0 \sim n$ について、平滑化されたスケール空間ピラミッドを構成する複数の平滑化画像 $L(x, y, c_i)$ を取得するために、（ $i = 0 \sim n$ について、異なる c_i によって設定される）異なるスケーリング幅の平滑化カーネル $G(x, y, c_i)$ で画像 $I(x, y)$ を畳み込む。画像 $I(x, y)$ は、照度関数 $S(x, y)$ によって完全にまたはピクセルごとに変更されているベース画像 $I_0(x, y)$ によって特徴づけられ得る。一例では、平滑化カーネル $G(x, y, c_i)$ は、平滑化されたスケール空間ピラミッドがガウススケール空間ピラミッドであるように、ガウスカーネルであり得る。

【0099】

40

[00138] 次に、904において、 $j = 0 \sim n - 1$ について、平滑化されたスケール空間ピラミッドにわたって平滑化画像の隣接するペアの差分 $L(x, y, c_{j+1}) - L(x, y, c_j)$ を取ることによって、スケール空間差分 $D_j(x, y,)$ を発生する。このプロセスは、複数のレベルを有するスケール空間差分を取得するために、隣接する平滑化画像の複数のセットについて繰り返される。第2の平滑化画像 $L(x, y, c_{j+1})$ を取得するために使用される第2の平滑化カーネル $G(x, y, c_{j+1})$ は、第1の平滑化画像 $L(x, y, c_j)$ を取得するために使用される第1の平滑化カーネル $G(x, y, c_j)$ よりも広くなり得ることに留意されたい。

【0100】

[00139] 次いで906において、 $j = 0 \sim n - 1$ について、スケール空間の各差分 $D_j(x, y,)$

50

x, y, \quad レベルを対応する平滑化画像 $L(x, y, c_{j+1+h})$ で除算することによって正規化されたスケール空間差分 $D'_j(x, y, \quad)$ を発生し、各平滑化画像 $L(x, y, c_{j+1+h})$ は、画像 $L(x, y, c_{j+1})$ および $L(x, y, c_j)$ の2つの異なる平滑化バージョンのうちの平滑なほうと同程度に平滑であるかまたはそれよりも平滑である。すなわち、正規化する平滑化画像 $L(x, y, c_{j+1+h})$ は、画像 $L(x, y, c_{j+1})$ および $L(x, y, c_j)$ の2つの異なる平滑化バージョンのためのスケール（たとえば、平滑化カーネル）のうちの大きいほうに等しいかまたはそれよりも広いスケール（たとえば、平滑化カーネル）を有し得る。

【0101】

[00140] 次いで 908 において、 $j = 0 \sim n - 1$ について、正規化されたスケール空間差分 $D'_j(x, y, \quad)$ を使用して画像 $I(x, y)$ の特徴を識別する。たとえば、特徴がその周りで定義され得るキーポイントとして局所的極値（すなわち、極小値または極大値）が識別され得る。次いで 910 において、識別された特徴に基づいて画像 $I(x, y)$ のための記述子を発生する。

【0102】

[00141] 図 11、図 12、および図 13 に示す方法は、画像の照度に関する事前情報を必要としない。この方法は、何らかの大きい（有意な）計算および処理を導入することなしに、異なる照度変化にわたって一貫したレベルで画像中の安定した特徴を選定することと再現性を高めることを可能にする。すなわち、平滑化スケール空間は、スケール空間差分 $D_j(x, y, \quad)$ を正規化するために使用される平滑化画像 $L(x, y, c_{j+1+h})$ をすでに含むので、正規化のために除算演算の他に追加の処理は必要とされない。

【0103】

[00142] さらに、特徴が選択される信頼性を適応させることによって、特徴が検出されるスケール（たとえば、平滑化レベル）に従ってより安定した特徴が取得され得る。すなわち、より高いスケールは、概して、より平滑な（すなわち、よりぼかされた）バージョンの画像を含み、そのようなスケールにおいて検出されたキーポイント／特徴は、より高い程度の信頼性を有する。

【0104】

[00143] 図 14 は、照度不変特徴検出デバイスの一例を示すブロック図である。特徴検出デバイス 1200 は、デジタルクエリ画像 1202 を受信または取得し得る。次いで、スケール空間発生器 1204（たとえば、ガウススケール空間発生器）が、クエリ画像 1202 を異なるスケール幅の複数の異なる平滑化カーネル 1203 で畳み込んで、スケール空間を発生し得る。スケール空間は、異なるスケール幅に平滑化された画像の複数の平滑化バージョンを備え得る。次いで、スケール空間差分発生器 1206 が、スケール空間からスケール空間差分を発生する。次いで、スケール空間差分正規化器 1208 が、たとえば、スケール空間レベルの各差分を対応する平滑化画像で除算することによって、スケール空間差分を正規化し、そのような平滑化画像は、除算されるスケール空間差分を発生するために使用される平滑化画像のうちの大きいほうと同程度に広いかまたはそれよりも広いスケールを有する。次いで、キーポイント発生器 1210 が、正規化されたスケール空間差分中のキーポイントを識別または検出する。これは、たとえば、正規化されたスケール空間差分のピクセルの間で局所的極値（すなわち、極大値または極小値）を見つけることによって行われ得る。特徴発生器 1212 が、次いで、たとえば、識別されたキーポイントの周りの局所ピクセルを特徴づけることによって、特徴を発生し得る。キーポイント発生器 1210 と特徴発生器 1212 との機能は特徴検出器によって実行され得ることに留意されたい。次いで、特徴記述子発生器 1214 が、各特徴について記述子を発生して、クエリ画像を識別するように働くことができる複数の画像記述子 1216 を与える。図 14 に示す機能は、別個の回路によってあるいは 1 つまたは複数のプロセッサによって実行され得る。

【0105】

[00144] 図 15 は、特徴検出のために正規化されたスケール空間差分を使用する画像整

合デバイスの一例を示すブロック図である。画像整合デバイス 1300 は、通信インターフェース 1304、画像キャプチャデバイス 1306、および / または記憶デバイス 1308 に結合された、処理回路 1302 を含み得る。通信インターフェース 1304 は、ワイヤード / ワイヤレスネットワーク上で通信し、画像および / または 1 つまたは複数の画像のための特徴記述子を受信するように適合され得る。画像キャプチャデバイス 1306 は、たとえば、クエリ画像をキャプチャすることができるデジタルカメラであり得る。処理回路 1302 は、画像から特徴を抽出する画像処理回路 1314 と、クエリ画像をターゲット画像のデータベース 1310 におよび / またはクエリ画像記述子を記述子データベース 1312 に整合させるために、抽出された特徴を使用する画像整合回路 1316 とを含み得る。例示的な一実装形態によれば、画像整合アプリケーションが、クエリ画像を画像データベース中の 1 つまたは複数の画像に整合させることを試みる。画像データベースは、データベース 1310 に記憶された 1 つまたは複数の画像に関連する何百万もの特徴記述子を含み得る。

10

20

30

40

50

【0106】

[00145] 画像処理回路 1314 は、ガウススケール空間発生器 1322、スケール空間差分発生器 1324、スケール空間差分正規化器 1326、キーポイント検出器 1328、特徴発生器 1330、および / または特徴記述子発生器 1332 を含む、特徴識別回路 1320 を含み得る。ガウススケール空間発生器 1322 は、たとえば、図 8 および図 11 に示すように、複数の異なるスケール空間を発生するために画像をぼかし関数（たとえば、平滑化カーネル）で畳み込むように働き得る。次いで、スケール空間差分発生器 1324 がスケール空間からスケール空間差分を発生する。次いで、スケール空間差分正規化器 1326 が、たとえば、スケール空間レベルの各差分を対応する平滑化画像で除算することによって、スケール空間差分を正規化し、そのような平滑化画像は、（図 12 に示した）除算されるスケール空間差分を発生するために使用される平滑化画像のいずれよりも広い。次いで、キーポイント発生器 1328 が、正規化されたスケール空間差分中のキーポイントを識別または検出する。これは、たとえば、正規化されたスケール空間差分のピクセルの間で局所的極値（すなわち、極大値または極小値）を見つけることによって行われ得る。特徴発生器 1330 が、次いで、たとえば、（図 9 に示した）識別されたキーポイントの周りの局所ピクセルを特徴づけることによって、特徴を発生し得る。次いで、特徴記述子発生器 1332 が、各特徴について記述子を発生して、（図 10 に示した）クエリ画像を識別するように働くことができる複数の画像記述子を与える。

【0107】

[00146] 次いで、画像整合回路 1316 が、特徴記述子に基づいてクエリ画像を画像データベース 1310 中の画像に整合させることを試み得る。整合結果は、（たとえば、画像または特徴記述子を送るモバイルデバイスに）通信インターフェースを介して与えられ得る。

【0108】

[00147] いくつかの実装形態では、クエリ画像のためのキーポイントに関連する特徴記述子のセットは画像整合デバイスによって受信され得ることに留意されたい。この状況では、クエリ画像は、（記述子を取得するために）すでに処理されている。したがって、画像処理回路 1314 は、画像整合デバイス 1300 からバイパスされるかまたは除去され得る。

【0109】

DOA 検出器およびオーディオシーン分解

[00148] 本明細書で開示するシステムおよび方法のいくつかの構成では、例示的な DOA 検出器 512 の機能と、空間フィルタ 72 を使用してオーディオシーンを分解するプロセスとは、このセクションにおいて説明する技法を使用して達成され得る。

【0110】

[00149] 遠距離場オーディオ処理（たとえば、オーディオ音源強調）の適用は、1 つまたは複数の音源が録音デバイスから比較的大きい距離（たとえば 2 メートル以上の距離）

に位置するときに生じ得る。

【 0 1 1 1 】

[00150]遠距離場使用事例の第 1 の例では、いくつかの異なる音源を含む音響シーンの記録を分解して、個別の音源のうちの 1 つまたは複数からそれぞれの音響成分を取得する。たとえば、異なる音源（たとえば、異なる音声および / または楽器）からの音が分離されるように、生の音楽演奏を記録することが望ましいことがある。別のそのような例では、「ロックバンド」タイプのビデオゲームなどのビデオゲームの 2 人以上の異なるプレーヤからの音声入力（たとえば、命令および / または歌唱）を区別することが望ましいことがある。

【 0 1 1 2 】

[00151]遠距離場使用事例の第 2 の例では、マルチマイクロフォンデバイスを使用して、（「ズームインマイクロフォン（zoom-in microphone）」とも呼ばれる）ビューの音場を狭めることによって遠距離場オーディオ音源強調を実行する。カメラを通じてシーンを見ているユーザは、カメラのレンズのズーム機能を使用して、たとえば、個々の話者または他の音源に対するビューの視界を選択的にズームし得る。相補的音響「ズームイン」効果をもたらすために、視覚的ズーム動作と同期して、記録される音響領域も被選択音源に狭められるように、カメラを実装することが望ましいことがある。

【 0 1 1 3 】

[00152]特定の遠くの音源から到来する音響成分を区別することは、単にビームパターンを特定の方向に狭めることではない。ビームパターンの空間幅は、フィルタのサイズを増加させることによって（たとえば、初期係数値のより長いセットを使用してビームパターンを定義することによって）狭められ得るが、音源の単一の到来方向にのみ依存すると、実際にはフィルタが音源エネルギーの大部分を逃すことになり得る。残響などの影響により、たとえば、音源信号は通常、異なる周波数においてやや異なる方向から到来し、結果的に、遠くの音源の到来方向は一般にはっきりしない。したがって、信号のエネルギーは、特定の方向に集中するのではなく、角度範囲にわたって拡散することがあり、特定の音源の到来角を、単一の方向におけるピークとしてではなく周波数範囲にわたる重心として特徴づけることがより有用であり得る。

【 0 1 1 4 】

[00153]フィルタのビームパターンが、単一の方向（たとえば、任意の 1 つの周波数における最大エネルギーによって示される方向）だけでなく、異なる周波数における方向の集中の幅をカバーすることが望ましいことがある。たとえば、ビームが、様々な対応する周波数において、そのような集中の幅内で、わずかに異なる方法に向くことを可能にすることが望ましいことがある。

【 0 1 1 5 】

[00154]適応ビームフォーミングアルゴリズムを使用して、1 つの周波数における特定の方向での最大応答と、別の周波数における異なる方向での最大応答とを有するフィルタを取得し得る。適応ビームフォーマーは一般に、正確な音声アクティビティ検出に依存するが、これは遠距離場の話者にとって達成するのが難しい。そのようなアルゴリズムはまた、所望の音源からの信号と干渉音源からの信号とが類似のスペクトルを有するとき（たとえば、2 つの音源の両方が人々の話声であるとき）パフォーマンスが芳しくないことがある。また、適応ビームフォーマーの代替として、ブラインド音源分離（BSS: blind source separation）ソリューションを使用して、1 つの周波数における特定の方向での最大応答と、別の周波数における異なる方向での最大応答とを有するフィルタを取得し得る。しかしながら、そのようなアルゴリズムは、遅い収束、極小値への収束、および / またはスケーリングのあいまいさを示すことがある。

【 0 1 1 6 】

[00155]良好な初期条件を提供するデータ独立型、開ループ手法（たとえば、MVDR ビームフォーマー）を、音声アクティビティ検出器を使用せずに出力間の相関を最小限に抑える閉ループ方法（たとえば、BSS）と組み合わせて、それによって改良されたロバ

10

20

30

40

50

ストな分離ソリューションを提供することが望ましいことがある。BSS方法は経時的に適応を実行するので、残響環境においてもロバストなソリューションを生成することが期待され得る。

【0117】

[00156]ヌルビームを使用してフィルタを初期化する既存のBSS初期化手法とは対照的に、本明細書で説明するソリューションは、音源ビームを使用してフィルタを初期化し、指定の音源方向に集中する。そのような初期化なしに、BSS方法がリアルタイムで有用なソリューションに適応するのを期待することは現実的でないことがある。

【0118】

[00157]図16Aに、フィルタバンクBK10と、フィルタ配向モジュールOM10と、フィルタ更新モジュールUM10とを含み、マルチチャネル信号（この例では、入力チャネルMCS10-1およびMCS10-2）を受信するように構成された、装置A100のブロック図を示す。フィルタバンクBK10は、マルチチャネル信号に基づく第1の信号に複数の第1の係数を適用して、第1の出力信号OS10-1を生成するように構成される。フィルタバンクBK10はまた、マルチチャネル信号に基づく第2の信号に複数の第2の係数を適用して、第2の出力信号OS10-2を生成するように構成される。フィルタ配向モジュールOM10は、第1の音源方向DA10に基づく複数の第1の係数の値の初期セットCV10を生成し、第1の音源方向DA10とは異なる第2の音源方向DA20に基づく複数の第2の係数の値の初期セットCV20を生成するように構成される。フィルタ更新モジュールUM10は、第1および第2の出力信号からの情報に基づいて、複数の第1および第2の係数の値の初期セットを更新して、値の対応する更新されたセットUV10およびUV20を生成するように構成される。

【0119】

[00158]音源方向DA10およびDA20の各々が、入力チャネルMCS10-1およびMCS10-2を生成するマイクロフォンアレイに対する（たとえば、アレイのマイクロフォンの軸に対する）対応する音源の推定方向を示すことが望ましいことがある。図16Bに、マイクロフォンアレイR100と、アレイから（たとえば、入力チャネルMCS10-1およびMCS10-2を含む）マルチチャネル信号MCS10を受信するように構成された装置A100のインスタンスとを含む、デバイスD10のブロック図を示す。アレイR100は、図1のアレイ18と、図1のシステム14中の装置A100との中に含まれ得る。

【0120】

[00159]図16Cに、点音源jから受信された信号成分の、アレイR100のマイクロフォンMC10およびMC20の軸に対する到来方向 θ_j を示す。アレイの軸は、マイクロフォンの音響的に敏感な面の中心を通る線として定義される。この例では、標示dは、マイクロフォンMC10とマイクロフォンMC20との間の距離を示す。

【0121】

[00160]フィルタ配向モジュールOM10は、ビームフォーミングアルゴリズムを実行して、それぞれの音源方向DA10、DA20におけるビームを記述した係数値の初期セットCV10、CV20を発生するように実装され得る。ビームフォーミングアルゴリズムの例としては、DSB（遅延和ビームフォーマー（delay-and-sum beamformer））、LCMV（線形制約最小分散（linear constraint minimum variance））、およびMVD R（最小分散無ひずみ応答（minimum variance distortionless response））がある。一例では、フィルタ配向モジュールOM10は、次のようなデータ独立式に従って、各フィルタが他の音源方向においてゼロ応答（またはヌルビーム）を有するように、ビームフォーマーの $N \times M$ 係数行列Wを計算するように実装される。

【0122】

$$W(\quad) = D^H(\quad, \quad) [D(\quad, \quad) D^H(\quad, \quad) + r(\quad) \times I]^{-1}$$

ただし、 $r(\quad)$ は、非反転性（noninvertibility）を補償するための正則化項である。

別の例では、フィルタ配向モジュールOM10は、次のような式に従って、MVD Rビー

ムフォーマーの $N \times M$ 係数行列 W を計算するように実装される。

【数 9】

$$W = \frac{\Phi^{-1}D(\omega)}{D^H(\omega)\Phi^{-1}D(\omega)}. \quad (1)$$

【0123】

これらの例では、 N は出力チャネルの数を示し、 M は入力チャネルの数（たとえば、マイクロフォンの数）を示し、 Φ は雑音の正規化クロスパワースペクトル密度行列を示し、 D （ ω ）は、（指向性行列とも呼ばれる） $M \times N$ アレイマニホールド行列を示し、上付き文字 H は共役転置関数を示す。通常、 M は N 以上である。

【0124】

[00161] 係数行列 W の各行は、フィルタバンク $BK10$ の対応するフィルタの係数の初期値を定義する。一例では、係数行列 W の第 1 の行は初期値 $CV10$ を定義し、係数行列 W の第 2 の行は初期値 $CV20$ を定義する。別の例では、係数行列 W の第 1 の行は初期値 $CV20$ を定義し、係数行列 W の第 2 の行は初期値 $CV10$ を定義する。

【0125】

[00162] 行列 D の各列 j は、次の式として表され得る周波数 ω_j にわたる遠距離場音源 j の指向性ベクトル（または「ステアリングベクトル」）である。 $D_{mj}(\omega_j) = \exp(-i \times \cos(\theta_j) \times \cos(\phi_j) \times \omega_j / c)$ この式において、 i は虚数を示し、 c は媒質中の音の伝播速度（たとえば、空中では 340 m/s ）を示し、 θ_j は、図 16C に示される入射到来角としてマイクロフォンアレイの軸に対する音源 j の方向（たとえば、 $j = 1$ の方向 $DA10$ および $j = 2$ の方向 $DA20$ ）を示し、 $\cos(\phi_j)$ は、 M 個のマイクロフォンのアレイにおける m 番目のマイクロフォンの空間座標を示す。均一なマイクロフォン間隔 d をもつマイクロフォンの線形アレイの場合、ファクタ $\cos(\phi_j)$ は $(m - 1)d$ として表され得る。

【0126】

[00163] 拡散雑音場の場合、行列 Γ は次のようなコヒーレンス関数 Γ_{ij} を使用して置き換えられ得る。

【数 10】

$$\Gamma_{ij} = \begin{cases} \text{sinc}\left(\frac{\omega d_{ij}}{c}\right), & i \neq j, \\ 1, & i = j \end{cases}$$

【0127】

ただし、 d_{ij} はマイクロフォン i とマイクロフォン j との間の距離を示す。さらなる一例では、行列 Γ は、 $(\Gamma + (\sigma^2 / N) I)$ に置き換えられ、ただし、 σ^2 は、（たとえば安定性に関する）対角ローディングファクタである。

【0128】

[00164] 一般に、フィルタバンク $BK10$ の出力チャネルの数 N は、入力チャネルの数 M 以下である。図 16A は、 N の値が 2 である（すなわち、2 つの出力チャネル $OS10 - 1$ および $OS10 - 2$ をもつ）装置 $A100$ の実装形態を示しているが、 N および M は 2 よりも大きい値（たとえば、3、4、またはそれ以上）を有し得ることを理解されたい。そのような一般的な場合、フィルタバンク $BK10$ は、 N 個のフィルタを含むように実装され、フィルタ配向モジュール $OM10$ は、これらのフィルタのために初期係数値の N 個の対応するセットを生成するように実装され、これらの原理のそのような拡張は、明確に企図され、本明細書によって開示される。

【0129】

10

20

30

40

50

[00165]たとえば、図17に、NとMの両方の値が4である装置A100の実装形態A110のブロック図を示す。装置A110は、4つのフィルタを含むフィルタバンクBK10の実装形態BK12を含み、各フィルタは、入力チャネルMCS10-1、MCS10-2、MCS10-3、およびMCS10-4の各々をフィルタ処理して、出力信号（またはチャネル）OS10-1、OS10-2、OS10-3、およびOS10-4のうちの対応する1つを生成するように構成される。装置A100はまた、フィルタバンクBK12のフィルタのために係数値の初期セットCV10、CV20、CV30、およびCV40を生成するように構成された、フィルタ配向モジュールOM10の実装形態OM12と、係数値の初期セットを適応させて、値の対応する更新されたセットUV10、UV20、UV30、およびUV40を生成するように構成された、フィルタ適応モジュールAM10の実装形態AM12とを含む。

10

【0130】

[00166]（「ビームパターン」とも呼ばれる）周波数ビン対入射角に関するフィルタバンクBK10のフィルタの初期応答は、MVD Rビームフォーミングアルゴリズム（たとえば、上の式（1））に従ってフィルタ配向モジュールOM10によって発生されたフィルタの係数値によって判断される。この応答は、入射角0（たとえば、マイクロフォンアレイの軸の方向）を中心として対称的であり得る。初期条件の異なるセット（たとえば、所望の音源からの音と、干渉音源からの音との推定到来方向の異なるセット）の下でのこのビームパターンの変形が有され得る。

20

【0131】

[00167]特定の適用例に適していると考えられる指向性とサイドローブ発生との間の折衷に従って選択されたビームフォーマー設計に従って係数値CV10およびCV20を生成するようにフィルタ配向モジュールOM10を実装することが望ましいことがある。上記の例は、周波数領域ビームフォーマー設計について説明しているが、時間領域ビームフォーマー設計に従って係数値のセットを生成するように構成されたフィルタ配向モジュールOM10の代替実装形態も、明確に企図され、本明細書によって開示される。

【0132】

[00168]フィルタ配向モジュールOM10は、（たとえば、上記で説明したようにビームフォーミングアルゴリズムを実行することによって）係数値CV10およびCV20を発生させるように、またはストレージから係数値CV10およびCV20を取り出すように実装され得る。たとえば、フィルタ配向モジュールOM10は、音源方向（たとえば、DA10およびDA20）に従って値（たとえば、ビーム）の事前計算されたセットの中から選択することによって、係数値の初期セットを生成するように実装され得る。そのような係数値の事前計算されたセットをオフラインで計算して、対応する所望の解像度における方向および/または周波数の所望の範囲をカバーし得る（たとえば、0、20、または30度から150、160、または180度までの範囲における、5度、10度、または20度の各間隔についての係数値の異なるセット）。

30

【0133】

[00169]フィルタ配向モジュールOM10によって生成される初期係数値（たとえば、CV10およびCV20）は、音源信号の間に所望のレベルの分離をもたらすようにフィルタバンクBK10を構成するには十分でないことがある。これらの初期値が基づく推定音源方向（たとえば、方向DA10およびDA20）が完全に正確であったとしても、フィルタを一定の方向にステアリングするだけでは、アレイから遠く離れた音源間の最良の分離、または特定の離れた音源への最良の集中は実現しないことがある。

40

【0134】

[00170]フィルタ更新モジュールUM10は、第1および第2の出力信号OS10-1およびOS10-2からの情報に基づいて、第1および第2の係数の初期値CV10およびCV20を更新して、値の対応する更新されたセットUV10およびUV20を生成するように構成される。たとえば、フィルタ更新モジュールUM10は、これらの初期係数値によって記述されるビームパターンを適応させるために適応BSSアルゴリズムを実行

50

するように実装され得る。

【 0 1 3 5 】

[00171] B S S 方法は、 $Y_j(\omega, l) = W(\omega) X_j(\omega, l)$ などの式に従って様々な音源から、統計的に独立した信号成分を分離し、ただし、 X_j は周波数領域における入力（混合）信号の j 番目のチャンネルを示し、 Y_j は周波数領域における出力（分離）信号の j 番目のチャンネルを示し、 ω は周波数ビンインデックスを示し、 l は時間フレームインデックスを示し、 W はフィルタ係数行列を示す。概して、B S S 方法は、次のような式による逆混合行列 W の経時的適応として記述され得る。

【 数 1 1 】

$$W_{l+r}(\omega) = W_l(\omega) + \mu [I - \{\Phi(Y(\omega, l))Y(\omega, l)^H\}] W_l(\omega), \quad (2)$$

10

【 0 1 3 6 】

ただし、 r は適応間隔（または更新レート）パラメータを示し、 μ は適応速度（または学習レート）ファクタを示し、 I は恒等行列を示し、上付き文字 H は共役転置関数を示し、

Φ はアクティブ化関数（activation function）を示し、括弧 $\langle \cdot \rangle$ は（たとえば、フレーム l から $l + L - 1$ にわたるものであって、 L は一般に r 以下である）時間平均化演算を示す。一例では、 μ の値は $0 \sim 1$ である。式（2）はB S S 学習ルールまたはB S S 適応ルールとも呼ばれる。アクティブ化関数 Φ は一般に、所望の信号の累積密度関数に近似するように選択され得る非線形有界関数である。そのような方法において使用されるアクティブ化関数 Φ の例としては、双曲正接関数（hyperbolic tangent function）、シグモイド関数（sigmoid function）、および符号関数（sign function）がある。

20

【 0 1 3 7 】

[00172] フィルタ更新モジュール U M 1 0 は、本明細書で説明するB S S 方法に従ってフィルタ配向モジュール O M 1 0 によって生成された係数値（たとえば、C V 1 0 および C V 2 0）を適応させるように実装され得る。そのような場合、出力信号 O S 1 0 - 1 および O S 1 0 - 2 は、周波数領域信号 Y のチャンネル（たとえば、それぞれ第 1 のチャンネルおよび第 2 のチャンネル）であり、係数値 C V 1 0 および C V 2 0 は、逆混合行列 W の対応する行（たとえば、それぞれ第 1 の行および第 2 の行）の初期値であり、適応された値は、適応後の逆混合行列 W の対応する行（たとえば、それぞれ第 1 の行および第 2 の行）によって定義される。

30

【 0 1 3 8 】

[00173] 周波数領域における適応のためのフィルタ更新モジュール U M 1 0 の典型的な実装形態では、逆混合行列 W は有限インパルス応答（F I R）多項式行列である。そのような行列は、要素として F I R フィルタの周波数変換（たとえば、離散フーリエ変換）を有する。時間領域における適応のためのフィルタ更新モジュール U M 1 0 の典型的な実装形態では、逆混合行列 W は F I R 行列である。そのような行列は要素として F I R フィルタを有する。そのような場合、係数値の各初期セット（たとえば、C V 1 0 および C V 2 0）は、一般に複数のフィルタを記述することになることを理解されよう。たとえば、係数値の各初期セットは、逆混合行列 W の対応する行の要素ごとにフィルタを記述し得る。周波数領域実装形態の場合、係数値の各初期セットは、マルチチャンネル信号の周波数ビンごとに、逆混合行列 W の対応する行の各要素のフィルタの変換を記述し得る。

40

【 0 1 3 9 】

[00174] B S S 学習ルールは、一般に、出力信号間の相関を減らすように設計される。たとえば、B S S 学習ルールは、出力信号間の相互情報量を最小限に抑えるように、出力信号の統計的独立性を高めるように、または出力信号のエントロピーを最大にするように選択され得る。一例では、フィルタ更新モジュール U M 1 0 は、独立成分分析（I C A : independent component analysis）として知られているB S S 方法を実行するように実装される。そのような場合、フィルタ更新モジュール U M 1 0 は、上記で説明したアクティ

50

ブ化関数、または、たとえば、次のようなアクティブ化関数を使用するように構成され得る。

【数 1 2】

$$\Phi(Y_j(\omega, l)) = Y_j(\omega, l) / |Y_j(\omega, l)|$$

【0 1 4 0】

周知のICA実装形態の例としては、Infomax、FastICA (www-dot-cis-dot-hut-dot-fi/projects/ica/fasticaでオンライン入手可能)、およびJADE (固有行列の結合近似対角化 (Joint Approximate Diagonalization of Eigenmatrices)) がある。

10

【0 1 4 1】

[00175] スケーリングおよび周波数置換は、BSSにおいて一般に遭遇される2つのあいまいさである。フィルタ配向モジュールOM10によって生成される初期ビームは置換されないが、そのようなあいまいさは、ICAの場合に適応中に生じ得る。置換されない解を維持するために、代わりに、周波数ビン間の予想される依存性をモデル化するソースプライアを使用する複素ICAの一変形である独立ベクトル解析 (IVA) を使用するようにフィルタ更新モジュールUM10を構成することが望ましいことがある。この方法では、アクティブ化関数は、たとえば、次の式などの多変量アクティブ化関数である。

20

【数 1 3】

$$\Phi(Y_j(\omega, l)) = Y_j(\omega, l) / (\sum_{\omega} |Y_j(\omega, l)|^p)^{1/p}$$

【0 1 4 2】

ただし、pは1以上の整数値 (たとえば、1、2、または3) を有する。この関数において、分母の項は、すべての周波数ビンにわたる分離された音源スペクトルに関係する。この場合、置換のあいまいさは解決される。

30

【0 1 4 3】

[00176] 得られた適応係数値によって定義されるビームパターンは、直線ではなく畳み込まれているように見え得る。そのようなパターンは、遠くの音源の分離には一般に不十分である初期係数値CV10およびCV20によって定義されるビームパターンよりも良好な分離をもたらすと予想され得る。たとえば、10 ~ 12 dB から 18 ~ 20 dB への干渉消去の増加が観測されている。適応係数値によって表されるソリューションはまた、マイクロフォン応答 (たとえば、利得および/または位相応答) の不整合に対し、開ループビームフォーミングソリューションよりもロバストであると予想され得る。

【0 1 4 4】

40

[00177] 上記の例は、周波数領域におけるフィルタ適応について説明しているが、時間領域における係数値のセットを更新するように構成されたフィルタ更新モジュールUM10の代替実装形態も、明確に企図され、本明細書によって開示される。時間領域BSS方法は、置換のあいまいさの影響を受けないが、一般に、周波数領域BSS方法よりも長いフィルタの使用を伴い、実際には扱いにくいことがある。

【0 1 4 5】

[00178] BSS方法を使用して適応されたフィルタは概して、良好な分離を達成するが、そのようなアルゴリズムも、特に音源が遠くにある場合に、分離信号にさらなる残響をもたらす傾向がある。特定の到来方向において単位利得を強制する幾何学的制約を追加することによって、適応BSSソリューションの空間応答を制御することが望ましいことが

50

ある。ただし、上述のように、単一の到来方向に対してフィルタ応答を調整するのは、残響環境では不十分であり得る。その上、B S S 適応において（ヌルビーム方向とは反対の）ビーム方向を強制しようとすると、問題が生じかねない。

【 0 1 4 6 】

[00179]フィルタ更新モジュール U M 1 0 は、方向に対する値の適応されたセットの判断された応答に基づいて、複数の第 1 の係数の値の適応されたセットと複数の第 2 の係数の値の適応されたセットとのうちの少なくとも 1 つを調整するように構成される。この判断された応答は、指定の特性を有する応答に基づき、異なる周波数では異なる値を有し得る。一例では、判断された応答は、最大応答である（たとえば、指定の特性は最大値である）。調整されるべき係数のセット j ごとに、また調整されるべき範囲内の各周波数において、たとえば、この最大応答 $R_j(\omega)$ は、次のような式に従って、その周波数における適応されたセットの複数の応答のうちの最大値として表され得る。

10

【 数 1 4 】

$$R_j(\omega) = \max_{\theta \in [-\pi, \pi]} |W_{j1}(\omega)D_{\theta 1}(\omega) + W_{j2}(\omega)D_{\theta 2}(\omega) + \dots + W_{jM}(\omega)D_{\theta M}(\omega)|, (3)$$

【 0 1 4 7 】

ただし、W は適応された値の行列（たとえば、F I R 多項式行列）であり、 W_{jm} は、行 j および列 m における行列 W の要素を示し、列ベクトル D () の各要素 m は、次の式で表され得る距離 の遠距離場音源から受信される信号に関する周波数 における位相遅延を示す。

20

【 0 1 4 8 】

$$D_m(\omega) = \exp(-i \times \cos(\theta) \times \text{pos}(m) \times \omega / c)$$

別の例では、判断された応答は、最小応答（たとえば、各周波数における適応されたセットの複数の応答の中の最小値）である。

【 0 1 4 9 】

[00180]一例では、式 (3) は、範囲 [- , +] において の 6 4 個の均一に離間した値について評価される。他の例では、式 (3) は、 の異なる数の値（たとえば、1 6 個または 3 2 個の均一に離間した値、5 度または 1 0 度の増分における値など）について、不均一な間隔で（たとえば、横方向の範囲にわたって、縦方向における範囲よりも大きい解像度で、またはその逆）、および / または異なる関心領域（たとえば、[- , 0]、[- / 2 , + / 2]、[- , + / 2]）にわたって評価され得る。均一なマイクロフォン間隔 d をもつマイクロフォンの線形アレイの場合、係数 $\text{pos}(m)$ は $(m - 1) d$ として表され得、したがって、ベクトル D () の各要素 m は次のように表され得る。 $D_m(\omega) = \exp(-i \times \cos(\theta) \times (m - 1) d \times \omega / c)$ 式 (3) が最大値を有する方向 の値は、周波数 の値が異なる場合には異なると予想され得る。音源方向（たとえば、D A 1 0 および / または D A 2 0 ）は、式 (3) が評価される の値の中に含まれ得、または、代替的に、それらの値とは別個であり得る（たとえば、音源方向が、式 (3) が評価される の値の隣接するものの間の角度を示す場合）ことに留意されたい。

30

40

【 0 1 5 0 】

[00181]図 1 8 A に、フィルタ更新モジュール U M 1 0 の実装形態 U M 2 0 のブロック図を示す。フィルタ更新モジュール U M 1 0 は、出力信号 O S 1 0 - 1 および O S 1 0 - 2 からの情報に基づいて係数値 C V 1 0 および係数値 C V 2 0 を適応させて、値の対応する適応されたセット A V 1 0 および A V 2 0 を生成するように構成された適応モジュール A P M 1 0 を含む。たとえば、適応モジュール A P M 1 0 は、本明細書で説明する B S S 方法のいずれか（たとえば、I C A、I V A）を実行するように実装され得る。

【 0 1 5 1 】

50

[00182]フィルタ更新モジュールUM20はまた、(たとえば、上記の式(3)による)方向に対する値の適応されたセットAV10の最大応答に基づいて適応された値AV10を調整して、値の更新されたセットUV10を生成するように構成された調整モジュールAJM10を含む。この場合、フィルタ更新モジュールUM20は、更新された値UV20としてそのような調整をせずに適応された値AV20を生成するように構成される。(本明細書で開示する構成の範囲はまた、係数値CV20が適応も調整もされないという点で、装置A100とは異なる装置を含むことに留意されたい。そのような構成は、たとえば、信号が残響をほとんどまたはまったく伴わずに直接経路を介して対応する音源から到来する状況において使用され得る。)

[00183]調整モジュールAJM10は、値の適応されたセットを、方向に対する各周波数における所望の利得応答(たとえば、最大の単位利得応答)を有するようにセットを正規化することによって調整するように実装され得る。そのような場合、調整モジュールAJM10は、係数値の適応されたセットj(たとえば、適応された値AV10)の各値を、セットの最大応答 R_j ()で除算して、係数値の対応する更新されたセット(たとえば、更新された値UV10)を取得するように実装され得る。

【0152】

[00184]所望の利得応答が単位利得応答以外である場合、調整モジュールAJM10は、適応された値および/または正規化された値に利得係数を適用することを調整演算が含むように実装され得、ここで、利得係数値の値は周波数とともに変化して、所望の利得応答を記述する(たとえば、音源のピッチ周波数のハーモニックを選好し、および/または干渉物によって支配され得る1つまたは複数の周波数を減衰させる)。判断された応答が最小応答である場合、調整モジュールAJM10は、(たとえば、各周波数の)最小応答を減算することによって、または方向に対する各周波数における所望の利得応答(たとえば、最小のゼロの利得応答)を有するようにセットを再マッピングすることによって、適応されたセットを調整するように実装され得る。

【0153】

[00185]係数値のセットのうちの2つ以上について、また場合によってはすべてについて(たとえば、少なくとも、定位された音源に関連しているフィルタについて)そのような正規化を実行するように調整モジュールAJM10を実装することが望ましいことがある。図18Bに、調整モジュールAJM10の実装形態AJM12を含むフィルタ更新モジュールUM20の実装形態UM22のブロック図を示し、AJM12はまた、方向に対する値の適応されたセットAV20の最大応答に基づいて、適応された値AV20を調整して、値の更新されたセットUV20を生成するように構成される。

【0154】

[00186]そのようなそれぞれの調整は、追加の適応フィルタに(たとえば、適応行列Wの他の行に)同じ方法で拡張され得ることを理解されたい。たとえば、図17に示したフィルタ更新モジュールUM12は、係数値の4つのセットCV10、CV20、CV30、およびCV40を適応させて、値の4つの対応する適応されたセットを生成するように構成された適応モジュールAPM10の一実装形態と、値の対応する適応されたセットの最大応答に基づいて、値の更新されたセットUV30およびUV40の一方または両方の各々を生成するように構成された調整モジュールAJM12の一実装形態とを含むように、フィルタ更新モジュール22の一実装形態として構成され得る。

【0155】

[00187]従来のオーディオ処理ソリューションは、雑音基準の計算と、計算された雑音基準を適用する後処理ステップとを含み得る。本明細書で説明する適応ソリューションは、後処理への依存を弱め、フィルタ適応への依存を強めて、干渉する点音源を除去することによって干渉消去と残響除去とを改善するように実装され得る。残響は、周波数とともに変化する利得応答を有する伝達関数(たとえば、室内応答伝達関数)として考えられ得、減衰する周波数成分もあれば、増幅する周波数成分もある。たとえば、室内のジオメトリは、様々な周波数における信号の相対強度に影響を与えることがあり、いくつかの周波

10

20

30

40

50

数が支配的になり得る。ある周波数から別の周波数に変化する方向において（すなわち、各周波数における主要ビームの方向において）所望の利得応答を有するようにフィルタを抑制することによって、本明細書で説明する正規化演算は、異なる周波数において空間中で信号のエネルギーが拡散される度合いの差異を補償することによって、信号を残響除去するのを助け得る。

【 0 1 5 6 】

[00188]最良の分離および残響除去の結果を達成するために、一部の到来角範囲内で音源から到来するエネルギーを通過させ、他の角度で干渉音源から到来するエネルギーをブロックする空間応答を有するように、フィルタバンク B K 1 0 のフィルタを構成することが望ましいことがある。本明細書で説明するように、B S S 適応を使用して、フィルタが初期解の近傍でより良い解を見つけることを可能にするように、フィルタ更新モジュール U M 1 0 を構成することが望ましいことがある。ただし、所望の音源に向けられた主要ビームを維持する制約なしに、フィルタ適応は、同様の方向からの干渉音源が（たとえば、干渉音源からのエネルギーを除去する広いヌルビームを作ることによって）主要ビームを損なうのを許容し得る。

10

【 0 1 5 7 】

[00189]フィルタ更新モジュール U M 1 0 は、制約付き B S S を介して適応ヌルビームフォーミングを使用して、音源定位解からの大きい逸脱を防ぐ一方、小さい定位誤差を訂正することができるように構成され得る。しかしながら、フィルタが異なる音源に方向を変えるのを防ぐフィルタ更新ルールに関する空間制約を課することが望ましいこともある。たとえば、フィルタを適応させるプロセスが、干渉音源の到来方向にヌル制約を含めることが望ましいことがある。そのような制約は、ビームパターンが低周波数において当該干渉方向にその配向を変えるのを防ぐことが望ましいことがある。

20

【 0 1 5 8 】

[00190] B S S 逆混合行列の一部のみを適応させるようにフィルタ更新モジュール U M 1 0 を実装する（たとえば、適応モジュール A P M 1 0 を実装する）ことが望ましいことがある。たとえば、フィルタバンク B K 1 0 のフィルタのうちの 1 つまたは複数を固定することが望ましいことがある。そのような制約は、（たとえば、上記の式（2）に示した）フィルタ適応プロセスが係数行列 W の対応する行を変えるのを防止することによって実装され得る。

30

【 0 1 5 9 】

[00191]一例では、そのような制約は、固定されるべき各フィルタに対応する（たとえば、フィルタ配向モジュール O M 1 0 によって生成された）係数値の初期セットを維持するために、適応プロセスの開始時から適用される。そのような実装形態は、たとえば、静止した干渉物にビームパターンが向けられているフィルタにとって適切であり得る。別の例では、そのような制約は、係数値の適応されたセットのさらなる適応を防止するために（たとえば、フィルタが収束したことが検出されたときに）後で適用される。そのような実装形態は、たとえば、安定した残響環境における静止した干渉物にビームパターンが向けられているフィルタにとって適切であり得る。フィルタ係数値の正規化されたセットが固定されると、セットが固定されている間は調整モジュール A J M 1 0 はそれらの値の調整を実行する必要がないが、調整モジュール A J M 1 0 は係数値の他のセットを（たとえば、調整モジュール A J M 1 0 によるそれらの適応に応答して）調整し続け得ることに留意されたい。

40

【 0 1 6 0 】

[00192]代替または追加として、周波数範囲の一部分のみでフィルタのうちの 1 つまたは複数を適応させるようにフィルタ更新モジュール U M 1 0 を実装する（たとえば、適応モジュール A P M 1 0 を実装する）ことが望ましいことがある。フィルタのそのような固定化は、当該範囲から外れた周波数に（たとえば、上記の式（2）中の の値に）対応するフィルタ係数値を適応させないことによって達成され得る。

【 0 1 6 1 】

50

[00193] 有用な情報を含んでいる周波数範囲でのみ、フィルタのうちの1つまたは複数（場合によってはすべて）の各々を適応させ、別の周波数範囲ではフィルタを固定することが望ましいことがある。適応されるべき周波数範囲は、マイクロフォンアレイから話者までの予想される距離、マイクロフォン間の距離（例：たとえば空間エイリアシングにより、空間フィルタ処理がいずれにせよ失敗する周波数でフィルタを適応させるのを回避するため）、部屋のジオメトリ、および/または室内のデバイスの配置などのファクタに基づき得る。たとえば、入力信号は、特定の周波数範囲（たとえば、高周波数範囲）にわたって、その範囲で正しいBSS学習をサポートするのに十分な情報を含んでいないことがある。そのような場合、適応なしにこの範囲で初期の（または場合によっては直近の）フィルタ係数値を使用し続けることが望ましいことがある。

10

【0162】

[00194] 音源がアレイから3～4メートル以上離れているとき、一般的に、音源によって放出される高周波エネルギーで、マイクロフォンに到達するものはほとんどない。そのような場合、フィルタ適応を適切にサポートする情報は、高周波数範囲ではほとんど得られないことがあるので、高周波数でフィルタを固定し、低周波数でのみそれらを適応させることが望ましいことがある。

【0163】

[00195] 追加または代替として、どの周波数を適応させるべきかの決定は、周波数帯域において現在利用可能なエネルギーの量、および/またはマイクロフォンアレイから現在の話者までの推定距離などのファクタに従って、実行時間中に変わり得、フィルタごとに異なり得る。たとえば、ある時間には最高2kHz（あるいは3kHzまたは5kHz）の周波数でフィルタを適応させ、別の時間には最高4kHz（あるいは5kHz、8kHz、または10kHz）の周波数でフィルタを適応させることが望ましいことがある。特定の周波数のために固定され、すでに調整されている（たとえば、正規化されている）フィルタ係数値を調整モジュールAJM10が調整する必要はないが、調整モジュールAJM10は他の周波数で係数値を（たとえば、適応モジュールAPM10によるそれらの適応に応答して）調整し続け得ることに留意されたい。

20

【0164】

[00196] フィルタバンクBK10は、更新された係数値（たとえば、UV10およびUV20）をマルチチャネル信号の対応するチャネルに適用する。更新された係数値は、（たとえば、調整モジュールAJM10による）本明細書で説明する調整後の（たとえば、適応モジュールAPM10によって適応された）逆混合行列Wの対応する行の値であるが、そのような値が本明細書で説明するように固定されている場合は除く。係数値の各更新されたセットは一般に、複数のフィルタを記述することになる。たとえば、係数値の各更新されたセットは、逆混合行列Wの対応する行の要素ごとにフィルタを記述し得る。

30

【0165】

[00197] 概して、各推定音源方向（たとえば、DA10および/またはDA20）は、測定、計算、予測、予想、および/または選択され得、所望の音源、干渉音源、または反射からの音の到来方向を示し得る。フィルタ配向モジュールOM10は、別のモジュールまたはデバイスから（たとえば、音源定位モジュールから）推定音源方向を受信するように構成され得る。そのようなモジュールまたはデバイスは、（たとえば、顔および/または動き検出を実行することによる）カメラからの画像情報および/または超音波反射からの測距情報に基づいて推定音源方向を生成するように構成され得る。そのようなモジュールまたはデバイスはまた、音源の数を推定するように、および/または動いている1つまたは複数の音源を追跡するように構成され得る。図19Aに、そのような画像情報をキャプチャするために使用され得るカメラCM10をもつアレイR100の4マイクロフォン実装形態R104の構成の一例の上面図を示す。

40

【0166】

[00198] 代替的に、装置A100は、マルチチャネル信号MCS10内の情報および/またはフィルタバンクBK10によって生成される出力信号内の情報に基づいて、推定音

50

源方向（たとえば、D A 1 0 および D A 2 0）を計算するように構成された方向推定モジュール D M 1 0 を含むように実装され得る。そのような場合、方向推定モジュール D M 1 0 はまた、上記で説明したように画像情報および / または測距情報に基づいて推定音源方向を計算するように実装され得る。たとえば、方向推定モジュール D M 1 0 は、マルチチャネル信号 M C S 1 0 に適用される、一般化相互相関（G C C : generalized cross-correlation）アルゴリズム、またはビームフォーマーアルゴリズムを使用して音源 D O A を推定するように実装され得る。

【 0 1 6 7 】

[00199] 図 2 0 に、マルチチャネル信号 M C S 1 0 内の情報に基づいて推定音源方向 D A 1 0 および D A 2 0 を計算するように構成された方向推定モジュール D M 1 0 のインスタンスを含む装置 A 1 0 0 の実装形態 A 1 2 0 のブロック図を示す。この場合、方向推定モジュール D M 1 0 およびフィルタバンク B K 1 0 は、同じ領域中で動作する（たとえば、周波数領域信号としてマルチチャネル信号 M C S 1 0 を受信し、処理する）ように実装される。図 2 1 に、装置 A 1 2 0 および A 2 0 0 の実装形態 A 2 2 0 のブロック図を示し、ここでは、方向推定モジュール D M 1 0 は、変換モジュール X M 2 0 から周波数領域においてマルチチャネル信号 M C S 1 0 からの情報を受信するように構成される。

【 0 1 6 8 】

[00200] 一例では、方向推定モジュール D M 1 0 は、位相変換を使用したステアード応答パワー（S R P - P H A T : steered response power using the phase transform）アルゴリズムを使用して、マルチチャネル信号 M C S 1 0 内の情報に基づいて推定音源方向を計算するように実装される。S R P - P H A T アルゴリズムは、最尤音源定位から得られるものであり、出力信号の相関が最大となる時間遅延を判断する。相互相関は、各ピンにおいて電力によって正規化され、それにより、より良いロバストネスが与えられる。残響環境では、S R P - P H A T は、競合する音源定位方法よりも良い結果をもたらすことが予想され得る。

【 0 1 6 9 】

[00201] S R P - P H A T アルゴリズムは、周波数領域における受信信号ベクトル X（すなわち、マルチチャネル信号 M C S 1 0）

$$X(\omega) = [X_1(\omega), \dots, X_p(\omega)]^T = S(\omega)G(\omega) + S(\omega)H(\omega) + N(\omega)$$

で表され得、

ただし、S は音源信号ベクトルを示し、利得行列 G、室内伝達関数ベクトル H、および雑音ベクトル N は次のように表され得る。

【数 1 5】

$$X(\omega) = [X_1(\omega), \dots, X_p(\omega)]^T,$$

$$G(\omega) = [\alpha_1(\omega)e^{-j\omega\tau_1}, \dots, \alpha_p(\omega)e^{-j\omega\tau_p}]^T,$$

$$H(\omega) = [H_1(\omega), \dots, H_p(\omega)]^T,$$

$$N(\omega) = [N_1(\omega), \dots, N_p(\omega)]^T.$$

【 0 1 7 0 】

これらの式において、P はセンサーの数（すなわち、入力チャネルの数）を示し、 τ_i は利得ファクタを示し、 τ_i は音源からの伝搬時間を示す。

【 0 1 7 1 】

[00202] この例では、複合雑音ベクトル $N^c(\omega) = S(\omega)H(\omega) + N(\omega)$ は、以下のゼロ平均、周波数独立、結合ガウス分布（zero-mean, frequency-independent, joint

t Gaussian distribution) を有すると仮定され得る。

【数 1 6】

$$p(N^c(\omega)) = \rho \exp \left\{ -\frac{1}{2} [N^c(\omega)]^H Q^{-1}(\omega) N^c(\omega) \right\},$$

【0 1 7 2】

ただし、 Q () は共分散行列であり、 ρ は定数である。音源方向は、次の式を最大化することによって推定され得る。

【数 1 7】

$$J_z = \int_{\omega} \frac{[G^H(\omega)Q^{-1}(\omega)X(\omega)]^H G^H(\omega)Q^{-1}(\omega)X(\omega)}{G^H(\omega)Q^{-1}(\omega)G(\omega)} d\omega.$$

【0 1 7 3】

N () = 0 であるとの仮定の下で、この式は次のように書き直され得る。

【数 1 8】

$$J_z = \frac{1}{\gamma P} \int \left| \sum_{i=1}^P \frac{X_i(\omega) e^{j\omega\tau_i}}{|X_i(\omega)|} \right|^2 d\omega, \quad (4)$$

【0 1 7 4】

ただし、 $0 < \gamma < 1$ は設計定数であり、式 (4) の右辺を最大化する時間遅延 τ_i は音源の到来方向を示す。

【0 1 7 5】

[00203] 図 2 2 に、周波数 f の範囲にわたる異なる 2 音源シナリオの DOA 推定に S R P - P H A T のそのような実装形態を使用した結果によるプロットの例を示す。これらのプロットでは、 y 軸は

【数 1 9】

$$\left| \sum_{i=1}^P \frac{X_i(\omega) e^{j\omega\tau_i}}{|X_i(\omega)|} \right|^2$$

【0 1 7 6】

の値を示し、 x 軸は、アレイ軸に対する推定音源到来方向 θ_i (

【数 2 0】

$$\theta_i (= \cos^{-1}(\tau_i c / d))$$

【0 1 7 7】

) を示す。各プロットにおいて、各線は範囲内の異なる周波数に対応し、各プロットはマイクロフォンアレイの縦方向を中心として対称的である (すなわち、 $\theta_i = 0$)。左上のプロットは、アレイから 4 メートルの距離にある 2 つの音源のヒストグラムを示している。右上のプロットは、アレイから 4 メートルの距離にある 2 つの近接した音源のヒストグラ

10

20

30

40

50

ムを示している。左下のプロットは、アレイから 2 . 5 メートルの距離にある 2 つの音源のヒストグラムを示している。右下のプロットは、アレイから 2 . 5 メートルの距離にある 2 つの近接した音源のヒストグラムを示している。これらのプロットの各々は、推定音源方向を、全周波数にわたる単一のピークとしてではなく、重心によって特徴づけられ得る角度範囲として示すことがわかるであろう。

【 0 1 7 8 】

[00204]別の例では、方向推定モジュール D M 1 0 は、ブラインド音源分離 (B S S) アルゴリズムを使用して、マルチチャネル信号 M C S 1 0 内の情報に基づいて推定音源方向を計算するように実装される。B S S 方法は、干渉音源からのエネルギーを除去する信頼できるヌルビームを発生する傾向があり、これらのヌルビームの方向は、対応する音源の到来方向を示すために使用され得る。方向推定モジュール D M 1 0 のそのような実装形態は、次のような式に従って、マイクロフォン j および j ' のアレイの軸に対する周波数 f における音源 i の到来方向 (D O A) を計算するように実装され得る。

10

$$\hat{\theta}_{i,jj'}(f) = \cos^{-1} \left(\frac{\arg \left([W^{-1}]_{ji} / [W^{-1}]_{ji'} \right)}{2\pi f c^{-1} \|p_j - p_{j'}\|} \right), \quad (5)$$

【 0 1 7 9 】

ただし、W は逆混合行列を示し、 p_j および $p_{j'}$ は、それぞれマイクロフォン j および j ' の空間的座標を示す。この場合、本明細書で説明するようにフィルタ更新モジュール U M 1 0 によって更新されるフィルタとは別個に方向推定モジュール D M 1 0 の B S S フィルタ (たとえば、逆混合行列 W) を実装することが望ましいことがある。

20

【 0 1 8 0 】

[00205]図 2 3 に、4 つのヒストグラムのセットの一例を示し、各ヒストグラムは、4 行逆混合行列 W の対応するインスタンスの (アレイ軸に対する) 各入射角に式 (5) がマッピングする周波数ビンの数を示し、ただし、W は、マルチチャネル信号 M C S 1 0 内の情報に基づいており、本明細書で説明する I V A 適応ルールに従って方向推定モジュール D M 1 0 の一実装形態によって計算される。この例では、入力マルチチャネル信号は、約 4 0 ~ 6 0 度の角度だけ分離された 2 つのアクティブな音源からのエネルギーを含んでいる。左上のプロットは、(音源 1 の方法を示す) I V A 出力 1 のヒストグラムを示しており、右上のプロットは、(音源 2 の方法を示す) I V A 出力 2 のヒストグラムを示している。これらのプロットの各々は、推定音源方向を、全周波数にわたる単一のピークとしてではなく、重心によって特徴づけられ得る角度範囲として示すことがわかるであろう。下のプロットは、I V A 出力 3 および 4 のヒストグラムを示しており、これらは、両方の音源からのエネルギーをブロックし、残響からのエネルギーを含んでいる。

30

【 0 1 8 1 】

[00206]別の例では、方向推定モジュール D M 1 0 は、複数の異なる周波数成分の各々についてマルチチャネル信号 M C S 1 0 のチャネル間の位相差に基づいて推定音源方向を計算するように実装される。(たとえば、図 1 9 B に示された平面波面の仮定が有効になるように) 遠距離場に点音源が 1 つあり、残響がない理想的な場合、位相差と周波数との比は周波数に対して一定である。図 1 5 B に示されたモデルを参照すると、方向推定モジュール D M 1 0 のそのような実装形態は、量

40

【 数 2 2 】

$$\frac{c \Delta \varphi_i}{d 2 \pi f_i}$$

【 0 1 8 2 】

50

の（アークコサインとも呼ばれる）逆コサインとして音源方向 θ_i を計算するように構成され得、ただし、 c は音速（約 340 m / 秒）を示し、 d はマイクロフォン間の距離を示し、 ϕ_i は 2 つのマイクロフォンチャンネルの対応する位相推定間のラジアン差分を示し、 f_i は、位相推定が対応する周波数成分（たとえば、対応する FFT サンプルの周波数、あるいは対応するサブバンドの中心周波数またはエッジ周波数）である。

【0183】

画像中のオブジェクト深さ判断

[00207] 以下で、画像からオブジェクト深さ情報を判断するための例示的な構成について説明する。第 1 の構成では、画像中のオブジェクトの推定深さを判断するために、マルチカメラ画像視差技法が使用される。第 2 の構成では、画像シーン中のオブジェクト範囲を推定するために単一カメラ自動フォーカス技法が使用され得る。SIFT キーポイント探索は、推定キーポイント深さ情報を含むことによってよりロバストにされ得る。

【0184】

[00208] 図 24 は、画像またはビデオキャプチャ中にシーン中のオブジェクトの視差を検出するように構成された画像キャプチャデバイス 1350 の特定の構成の図である。画像キャプチャデバイス 1350 は、画像処理モジュール 1356 に結合された画像センサーペア 1352 を含む。画像処理モジュール 1356 は外部メモリ 1362 に結合される。画像処理モジュール 1356 は、同期およびインターフェースモジュール 1354 と、画像処理機能モジュール 1358 と、視差検出モジュール 1342 と、符号化モジュール 1360 とを含む。

【0185】

[00209] 画像センサーペア 1352 は、画像データ 1370 を画像処理モジュール 1356 に与えるように構成される。単一のシーンに対応する第 1 の画像と第 2 の画像とを使用してオブジェクト深さ判断が実行され得る。第 1 の画像は、第 1 のセンサー（たとえば、右センサー）によるシーンの第 1 の画像キャプチャに対応し得、第 2 の画像は、第 2 のセンサー（たとえば、左センサー）によるシーンの第 2 の画像キャプチャに対応し得、第 2 の画像キャプチャは、図 24 に示すセンサーペア 1352 などによって、第 1 の画像キャプチャと実質的に同時である。

【0186】

[00210] 同期およびインターフェースモジュール 1354 は、データ 1372 を画像処理機能モジュール 1358 に与えるように構成される。画像処理機能モジュール 1358 は、処理された画像データ 1380 を視差検出モジュール 1342 に与えるように構成される。符号化モジュール 1360 は、画像 / ビデオデータ 1382 を受信し、オブジェクト深さデータで符号化された画像 / ビデオデータ 1384 を発生するように構成される。

【0187】

[00211] 視差検出モジュール 1342 は、画像センサーペア 1352 によってキャプチャされたシーン内のオブジェクトに対応する視差値を判断するように構成され得る。特定の構成では、視差検出モジュール 1342 は、シーン固有オブジェクト検出またはキーポイント検出および視差判断機能を組み込む。

【0188】

[00212] 画像センサーペア 1352 は、代表的な図では、右センサー（すなわち、閲覧者の右眼によって知覚されるシーンに関連する画像をキャプチャする第 1 のセンサー）と、左センサー（すなわち、閲覧者の左眼によって知覚されるシーンに関連する画像をキャプチャする第 2 のセンサー）とを含むセンサーのペアとして示されている。画像データ 1370 は、左センサーによって生成された左画像データと、右センサーによって生成された右画像データとを含む。各センサーは、水平方向に延在する感光性構成要素の行と、垂直方向に延在する感光性構成要素の列とを有するものとして示されている。左センサーと右センサーは、水平方向に沿って互いに距離 d において実質的に位置合わせされる。本明細書で使用する画像データ内の「水平」方向は、右画像データ中のオブジェクトのロケーションと、左画像データ中の同じオブジェクトのロケーションとの間の変位の方向である

。

【 0 1 8 9 】

[00213]図 2 5 は、図 2 4 のシステム中に含まれ得る画像処理システム 1 4 4 0 の特定の実施形態の図である。処理システム 1 4 4 0 は、入力画像データ 1 4 0 4 を受信し、出力画像データ 1 4 2 8 を発生するように構成される。処理システム 1 4 4 0 は、較正入力 1 4 5 0 を介して受信されるカメラ較正パラメータ 1 4 0 6 に応答し得る。

【 0 1 9 0 】

[00214]画像処理システム 1 4 4 0 は、微細ジオメトリ補正モジュール 1 4 1 0 と、キーポイント検出モジュール 1 4 1 2 と、キーポイント整合モジュール 1 4 1 4 と、深さ計算モジュール 1 4 1 6 とを含む。

10

【 0 1 9 1 】

[00215]ジオメトリ補正モジュール 1 4 1 0 は、データ経路 1 4 7 0 を介して入力画像データ 1 4 0 4 を受信し、補正された画像データ 1 4 5 4 を発生するように構成される。ジオメトリ補正モジュール 1 4 1 0 は、カメラ較正パラメータ 1 4 0 6 からのデータを使用し得、入力画像データ 1 4 0 4 を調整して、画像データ 1 4 0 4 のレンダリングに悪影響を及ぼし得る不整合、収差、または他の較正状態について訂正し得る。例示のために、ジオメトリ補正モジュール 1 4 1 0 は、較正パラメータ 1 4 0 6 について調整するために、任意のグリッド上で画像データ 1 4 0 4 のリサンプリングを効果的に実行し得る。

【 0 1 9 2 】

[00216]処理システム 1 4 4 0 がコンピューティングデバイス中に実装され得る構成では、カメラ較正パラメータ 1 4 0 6 は、画像 / ビデオデータファイルのヘッダ中でなど、入力画像データ 1 4 0 4 とともに受信され得る。処理システム 1 4 4 0 が図 2 4 の画像キャプチャデバイス 1 3 5 0 などの画像キャプチャデバイス中に実装される構成では、カメラ較正パラメータ 1 4 0 6 は、画像キャプチャデバイスの画像センサーペアに対応し得、微細ジオメトリ補正モジュール 1 4 1 0 にとってアクセス可能なメモリに記憶され得る。

20

【 0 1 9 3 】

[00217]キーポイント検出モジュール 1 4 1 2 は、補正された画像データ 1 4 5 4 を受信し、キーポイントロケーションデータ 1 4 5 6 を発生するように構成される。キーポイント検出モジュール 1 4 1 2 は、補正された画像データ 1 4 5 4 中の特徴的なポイントを識別するように構成される。たとえば、特徴的なポイントは、シーン中のオブジェクトの垂直エッジ、または水平方向において高周波成分を有するそのシーンの他のポイントに対応し得る。画像データ中のそのような特徴的な要素を本明細書では「キーポイント」または「オブジェクト」と呼ぶが、そのような識別された要素は、個々のピクセル、ピクセルのグループ、分数ピクセル部分、他の画像成分、またはそれらの任意の組合せに対応し得ることを理解されたい。たとえば、キーポイントは、受信された画像データのサブサンプリングされたルーマ成分をもつピクセルに対応し得、垂直エッジ検出フィルタを使用して検出され得る。

30

【 0 1 9 4 】

[00218]キーポイント整合モジュール 1 4 1 4 は、キーポイントロケーションデータ 1 4 5 4 を受信し、識別されたキーポイントに対応する視差データ 1 4 5 8 を発生するように構成される。キーポイント整合モジュール 1 4 1 4 は、探索範囲内でキーポイントの周りを探索し、視差ベクトルの信頼性測度を生成するように構成され得る。

40

【 0 1 9 5 】

[00219]深さ計算モジュール 1 4 1 6 は、視差データ 1 4 5 8 を受信し、センサー 1 3 5 2 からのキーポイントの推定距離を示すレンジデータ 1 4 6 0 を発生するように構成される。

【 0 1 9 6 】

[00220]処理システム 1 4 4 0 の動作中に、レンジ評価プロセスが実行される。画像データ 1 4 0 4 をキャプチャした 2 つのセンサー間の相対位置を推定し、補正するように設計された較正手順が、オフラインで（たとえば、デバイスのエンドユーザへの配信より前

50

に) 実行され得るが、ジオメトリ補正は画像データ 1 4 0 4 のフレームごとに実行され得る。

【0 1 9 7】

[00221] 処理は、(たとえば、キーポイント検出モジュール 1 4 1 2 において) キーポイント検出を続ける。視差を確実に推定するために使用され得る画像のオブジェクトまたはピクセル(キーポイント)のセットが選択される。推定視差における高い信頼性が達成され得るが、シーン中のすべての領域またはオブジェクトが使用されるとは限らない。キーポイントのセットの選択は、適切な(1つまたは複数の)解像度を生成するために、画像サブサンプリングを含み得る。(たとえば、垂直方向の特徴に対応する水平周波数のみを探すために) 画像高域フィルタを適用し、その後、フィルタを適用することによって発生した結果の平方値または絶対値を取り得る。所定のしきい値を超える結果は、潜在的キーポイントとして識別され得る。一部の局所近傍内の最良のキーポイント(たとえば、所定の領域内にあるすべてのキーポイントの最大フィルタ結果に対応するキーポイント)を選択するために、潜在的キーポイントに対してキーポイントブルーニングプロセスが実行され得る。

10

【0 1 9 8】

[00222] 検出されたキーポイントを使用して、(たとえば、キーポイント整合モジュール 1 4 1 4 において) キーポイント整合が実行され得る。第 1 の画像(たとえば左画像または右画像)中のキーポイントと、第 2 の画像(たとえば左画像および右画像のうちの他方)中の対応するエリアとの間の対応が判断され得る。信頼性推定値が生成され得、それは、キーポイント選択とともに視差推定精度を著しく改善し得る。左画像中のキーポイントと右画像中のキーポイントとの間の整合がどれくらい近接しているかの判断を可能にするために、整合は、正規化された相互共分散(cross-covariance)を使用して実行され得る。信頼性測度はこの正規化された相互共分散に基づき得る。特定の実施形態では、第 1 の画像中のキーポイントに対応する第 2 の画像中のキーポイントの位置を特定するための探索範囲は、センサー較正のための画像補正がすでに実行されているので、水平のみであり、探索範囲は、第 1 の画像中のキーポイントの周りの一定の範囲のみをカバーするように調整される。視差値はこれらの比較から計算される。

20

【0 1 9 9】

[00223] 図 2 6 A および図 2 6 B は、知覚されたオブジェクト深さと相関させられたオブジェクト視差の例示的な実施形態の図である。オブジェクト深さ判断は、異なる画像を各眼 1 5 0 4、1 5 0 6 にダイレクトすることに依拠する。目的は、オブジェクト視差(水平シフト)が深さと相関させられるように、左および右(L/R)画像から深さの錯覚を再生成することである。図 2 6 A は、ディスプレイ表面 1 5 2 4 を越えて知覚されるオブジェクト 1 5 3 0 に対応する正の視差 1 5 5 0 を示している。視差 1 5 5 0 は、左画像中のオブジェクトのロケーション 1 5 2 0 と、右画像中のオブジェクトのロケーション 1 5 2 2 との間の距離を示す。観察者は、左画像中のオブジェクト 1 5 3 0 の画像と、右画像中のオブジェクト 1 5 3 0 の画像とを融合させて、左眼 1 5 0 4 の見通し線 1 5 6 0 と、右眼 1 5 0 6 の見通し線 1 5 6 2 との交点においてオブジェクト 1 5 3 0 を知覚することになる。

30

40

【0 2 0 0】

[00224] 図 2 6 B は、ディスプレイ表面 1 5 2 4 の前で知覚されるオブジェクト 1 5 3 0 に対応する負の視差 1 5 5 0 を示している。視差 1 5 5 0 は、左画像中のオブジェクトのロケーション 1 5 2 0 と、右画像中のオブジェクトのロケーション 1 5 2 2 との間の距離を示す。観察者は、左画像中のオブジェクト 1 5 3 0 の画像と、右画像中のオブジェクト 1 5 3 0 の画像とを融合させて、左眼 1 5 0 4 の見通し線 1 5 6 0 と、右眼 1 5 0 6 の見通し線 1 5 6 2 との交点において、ディスプレイ表面 1 5 3 4 の前でオブジェクト 1 5 3 0 を知覚することになる。

【0 2 0 1】

[00225] 2 つの眼から見えるオブジェクト変位は、視覚野によって深さとして解釈され

50

る。2つのキャプチャされた画像間の視差はシーンに依存することになる。シーン深さを感知することを使用すると、画像中のキーポイント探索を特定の深さでのまたはその近くでのオブジェクトのみに狭めることができ、したがって、オブジェクト認識の信頼性を高めることができる。

【0202】

[00226] 深さ計算モジュール602によって実行されるシーンレンジ推定は、左画像と右画像との間のスパーズ動きベクトル推定として一般化され得る。シーンレンジ評価プロセスはキー（特徴的な）ポイント識別を含むことができる。水平シフトのみが存在する（および測定される）ので、垂直変化は必要とされない。水平変化（何らかの垂直成分をもつエッジ）が使用される。いくつかの構成では、キーポイントは異なる解像度で検出され得る。オブジェクトレンジ推定プロセスはまた、キーポイント整合を含むことができる。光源レベル非依存になるために、およびロバストな視差信頼性メトリックを生成するために、キーポイント整合は、正規化された相互共分散を使用して実行され得る。その結果、キーポイントを異なる解像度で整合させることは不要になり得る。

【0203】

オーディオシーン分解

[00227] 音響分解サブシステム22は、シーンから記録されたオーディオ信号を分解するために、このセクションで説明する技法を採用することができる。本明細書で開示するものは、楽音（note）のペンデンス（pendency）にわたる楽音のスペクトルの変化に関係する情報を含む基底関数インベントリと、スパーズ復元技法とを使用する、オーディオ信号の分解である。そのような分解は、信号の分析、符号化、再生、および/または合成をサポートするために使用され得る。本明細書では、調波楽器（すなわち、非打楽器）および打楽器からの混合音を含むオーディオ信号の定量分析の例を示す。

【0204】

[00228] 開示する技法は、キャプチャされたオーディオ信号を一連のセグメントとして処理するように構成され得る。典型的なセグメント長は約5または10ミリ秒から約40または50ミリ秒にわたり、セグメントは、重複しても（たとえば、隣接するセグメントが25%または50%だけ重複する）、重複しなくてもよい。1つの特定の例では、信号は、10ミリ秒の長さをそれぞれ有する一連の重複しないセグメントまたは「フレーム」に分割される。また、そのような方法によって処理されるセグメントは、異なる演算によって処理されるより大きいセグメントのセグメント（すなわち、「サブフレーム」）であり得、またはその逆も同様である。

【0205】

[00229] 2つ以上の楽器および/またはボーカル信号の混合から個々のノート/ピッチプロファイルを抽出するために音楽シーンを分解することが望ましいことがある。潜在的な使用事例としては、複数のマイクロフォンを用いてコンサート/ビデオゲームシーンをテープに記録すること、空間/スパーズ復元処理を用いて楽器とボーカルとを分解すること、ピッチ/ノートプロファイルを抽出すること、補正ピッチ/ノートプロファイルを用いて個々の音源を部分的にまたは完全にアップミックスすることがある。そのような動作は、音楽アプリケーション（たとえば、QualcommのQUSICアプリケーション、Rock BandまたはGuitar Heroなどのビデオゲーム）の機能をマルチプレーヤ/シンガーシナリオに拡張するために使用され得る。

【0206】

[00230]（たとえば、図34に示すように）同時に2人以上のボーカリストがアクティブであり、および/または複数の楽器がプレイされるシナリオを音楽アプリケーションが処理することを可能にすることが望ましいことがある。そのような機能は、現実的な音楽テープ記録シナリオ（マルチピッチシーン）をサポートするために望ましいことがある。ユーザは、各音源を別々に編集および再合成する能力を希望し得るが、サウンドトラックを生成することは、それらの音源を同時に記録することを伴い得る。

【0207】

10

20

30

40

50

[00231]本開示では、複数の音源が同時にアクティブになり得る音楽アプリケーションのための使用事例を可能にするために使用され得る方法について説明する。そのような方法は、基底関数インベントリベースのスパース復元（たとえば、スパース分解）技法を使用してオーディオ混合信号を分析するように構成され得る。

【0208】

[00232]基底関数のセットについて（たとえば、効率的なスパース復元アルゴリズムを使用して）アクティブ化係数の最もスパースなベクトルを見つけることによって混合信号スペクトルを音源成分に分解することが望ましいことがある。基底関数のセットは、図2の画像/ビデオ処理ブロック54によってシーン中に存在すると示された特定のタイプの楽器に減少させられ得る。アクティブ化係数ベクトルを（たとえば、基底関数のセットととも）使用して、混合信号を再構成するか、または混合信号の（たとえば、1つまたは複数の選択された楽器からの）選択された部分を再構成し得る。また、（たとえば、大きさおよび時間サポートに従って）スパース係数ベクトルを後処理することが望ましいことがある。

10

【0209】

[00233]図27Aに、オーディオ信号を分解する方法M100のフローチャートを示す。方法M100は、オーディオ信号のフレームからの情報に基づいて、周波数範囲にわたる対応する信号表現を計算するタスクT100を含む。方法M100は、タスクT100によって計算された信号表現と、複数の基底関数とに基づいて、アクティブ化係数のベクトルを計算するタスクT200をも含み、アクティブ化係数の各々は、複数の基底関数のうちの異なる1つに対応する。

20

【0210】

[00234]タスクT100は、信号表現を周波数領域ベクトルとして計算するように実装され得る。そのようなベクトルの各要素は、メルまたはバーク尺度（mel or Bark scale）に従って取得され得る、サブバンドのセットのうち対応する1つのサブバンドのエネルギーを示し得る。しかしながら、そのようなベクトルは、一般に、高速フーリエ変換（FFT）、または短時間フーリエ変換（STFT）など、離散フーリエ変換（DFT）を使用して計算される。そのようなベクトルは、たとえば、64、128、256、512、または1024ピンの長さを有し得る。一例では、オーディオ信号は、8kHzのサンプリングレートを有し、0～4kHz帯域は、長さ32ミリ秒の各フレームについて256ピンの周波数領域ベクトルによって表される。別の例では、信号表現は、オーディオ信号の重複セグメントにわたる修正離散コサイン変換（MDCT）を使用して計算される。

30

【0211】

[00235]さらなる一例では、タスクT100は、フレームの短期電力スペクトルを表すケプストラム係数（たとえば、メル周波数ケプストラム係数またはMFCC）のベクトルとして信号表現を計算するように実装される。この場合、タスクT100は、フレームのDFT周波数領域ベクトルの大きさにメル尺度フィルタバンクを適用することと、フィルタ出力の対数をとることと、対数値のDCTをとることとによって、そのようなベクトルを計算するように実装され得る。そのような手順は、たとえば、「STQ: DSR - Front-end feature extraction algorithm; compression algorithm」と題する、ETS IドキュメントES 201 108に記載されているオーロラ規格（欧州通信規格協会、2000年）において説明されている。

40

【0212】

[00236]楽器は、一般に、明確な音色を有する。楽器の音色は、そのスペクトルエンベロープ（たとえば、周波数範囲にわたるエネルギーの分布）によって記述され得るので、異なる楽器の音色の範囲は、個々の楽器のスペクトルエンベロープを符号化する基底関数のインベントリを使用してモデル化され得る。

【0213】

[00237]各基底関数は、周波数範囲にわたる対応する信号表現を備える。これらの信号表現の各々は、タスクT100によって計算された信号表現と同じ形態を有することが望

50

ましいことがある。たとえば、各基底関数は、長さ 64、128、256、512、または 1024 ピンの周波数領域ベクトルであり得る。代替的に、各基底関数は、MFCC のベクトルなどのケプストラム領域ベクトルであり得る。さらなる一例では、各基底関数はウェーブレット領域ベクトルである。

【0214】

[00238]基底関数インベントリ A は、各楽器 n (たとえば、ピアノ、フルート、ギター、ドラムなど) の基底関数のセット A_n を含み得る。たとえば、楽器の音色は、概して、各楽器 n の基底関数のセット A_n が、一般に、楽器ごとに異なり得るある所望のピッチ範囲にわたる各ピッチについて少なくとも 1 つの基底関数を含むようなピッチ従属である。たとえば、半音階スケールにチューニングされた楽器に対応する基底関数のセットは、10
オクターブ当たり 12 ピッチの各々の異なる基底関数を含み得る。ピアノの基底関数のセットは、ピアノの各キーについて異なる基底関数を含み、合計で 88 個の基底関数を含み得る。別の例では、各楽器の基底関数のセットは、5 オクターブ (たとえば、56 ピッチ) または 6 オクターブ (たとえば、67 ピッチ) など、所望のピッチ範囲内の各ピッチについて異なる基底関数を含む。基底関数のこれらのセット A_n は独立であり得、または 2 つ以上のセットが 1 つまたは複数の基底関数を共有し得る。

【0215】

[00239]セットの各基底関数は、楽器の音色を異なる対応するピッチで符号化し得る。音楽信号のコンテキストでは、人間ボイスは、インベントリが 1 つまたは複数の人間ボイスモデルの各々の基底関数のセットを含み得るような楽器と見なされ得る。20

【0216】

[00240]基底関数のインベントリは、アドホック記録された個々の楽器記録から学習された一般的な楽器ピッチデータベースに基づき得、および / または (たとえば、独立成分分析 (ICA)、期待値最大化 (EM: expectation-maximization) などの分離方式を使用した) 混合の分離されたストリームに基づき得る。

【0217】

[00241]オーディオを処理するための基底関数のセットの選択は、図 2 の画像 / ビデオ処理ブロック 54 によって与えられる楽器候補のリストに基づき得る。たとえば、基底関数のセットは、画像 / ビデオ処理ブロック 54 のオブジェクト認識プロセスによってシーン中で識別される楽器のみに制限され得る。30

【0218】

[00242]タスク T100 によって計算された信号表現と、インベントリ A からの基底関数の複数 B とに基づいて、タスク T200 はアクティブ化係数のベクトルを計算する。このベクトルの各係数は、基底関数の複数 B のうちの異なる 1 つに対応する。たとえば、タスク T200 は、基底関数の複数 B に従って、ベクトルが信号表現のための最も有望なモデルを示すように、ベクトルを計算するように構成され得る。図 32 に、そのようなモデル $Bf = y$ を示し、ここで、基底関数の複数 B は、B 個の列が個々の基底関数であり、f が基底関数アクティブ化係数の列ベクトルであり、y が、記録された混合信号のフレーム (たとえば、スペクトログラム周波数ベクトルの形態の 5、10、または 20 ミリ秒フレーム) の列ベクトルであるような行列である。40

【0219】

[00243]タスク T200 は、線形プログラミング問題を解くことによって、オーディオ信号の各フレームのアクティブ化係数ベクトルを復元するように構成され得る。そのような問題を解くために使用され得る方法の例としては、非負値行列因子分解 (NNMF: nonnegative matrix factorization) がある。NNMF に基づくシングルチャネル基準方法は、(たとえば、以下で説明するように) 期待値最大化 (EM) 更新ルールを使用して、基底関数とアクティブ化係数とを同時に計算するように構成され得る。

【0220】

[00244]既知または部分的に既知の基底関数空間における最もスパースなアクティブ化係数ベクトルを見つけることによって、オーディオ混合信号を (1 つまたは複数の人間ボ

10

20

30

40

50

イスを含み得る)個々の楽器に分解することが望ましいことがある。たとえば、タスク T 200 は、既知の楽器基底関数のセットを使用して、(たとえば、効率的なスパース復元アルゴリズムを使用して)基底関数インベントリにおける最もスパースなアクティブ化係数ベクトルを見つけることによって、入力信号表現を音源成分(たとえば、1つまたは複数の個々の楽器)に分解するように構成され得る。

【0221】

[00245]劣決定系(underdetermined system)の連立一次方程式(すなわち、式よりも多い未知数を有する系)の最小 L1 ノルム解は、しばしばその系の最もスパースな解でもあることが知られている。L1 ノルムの最小化によるスパース復元は、以下のように実行され得る。

10

【0222】

[00246]ターゲットベクトル f_0 は、 $K < N$ 個の非 0 成分を有する長さ N のスパースベクトルであり(すなわち、「 K スパース」であり)、射影行列(すなわち、基底関数行列) A は、サイズ約 K のセットについてインコヒーレント(ランダム様)であると仮定する。信号 $y = A f_0$ であることがわかる。次いで、 $A f = y z$ を条件として

【数 23】

$$\min \|f\|_1$$

20

【0223】

を解くこと(ただし、 $\|f\|_1$ は

【数 24】

$$\sum_{i=1}^n |f_i|$$

【0224】

30

として定義される)により、 f_0 が正確に復元される。その上、扱いやすいプログラムを解くことによって、 $M = K \cdot \log N$ 個のインコヒーレント測定値から f_0 を復元することができる。測定値の数 M は、アクティブな成分の数にほぼ等しい。

【0225】

[00247] 1つの手法は、圧縮感知(compressive sensing)からのスパース復元アルゴリズムを使用することである。(「圧縮感知(compressed sensing)」とも呼ばれる)圧縮感知の一例では、信号復元 $x = y$ であり、 y は、長さ M の観測信号ベクトルであり、 x は、 y の凝縮(condensed)表現である、 $K < N$ 個の非 0 成分を有する長さ N のスパースベクトル(すなわち、「 K スパースモデル」)であり、 A は、サイズ $M \times N$ のランダム射影行列である。ランダム射影 A はフルランクでないが、それは、高確率でスパース/圧縮可能信号モデルについて可逆である(すなわち、それは不良設定逆問題(ill-posed inverse problem)を解く)。

40

【0226】

[00248]アクティブ化係数ベクトル f は、対応する基底関数セット A_n のアクティブ化係数を含む各楽器 n のサブベクトル f_n を含むと見なされ得る。これらの楽器固有のアクティブ化サブベクトルは独立して(たとえば、後処理動作において)処理され得る。たとえば、1つまたは複数のスパースリティ制約(たとえば、ベクトル要素の少なくとも半分が 0 であること、楽器固有のサブベクトル中の非 0 要素の数が最大値を超えないことなど)を強制することが望ましいことがある。アクティブ化係数ベクトルの処理は、各フレームについて各非 0 アクティブ化係数のインデックス番号を符号化すること、各非 0 アクティブ

50

化係数のインデックスと値とを符号化すること、またはスパースベクトル全体を符号化することを含み得る。そのような情報は、示されたアクティブな基底関数を使用して混合信号を再生するため、または混合信号の特定の部分のみ（たとえば、特定の楽器によってプレイされるノートのみ）を再生するために、（たとえば、別の時間および/またはロケーションにおいて）使用され得る。

【0227】

[00249]楽器によって生成されるオーディオ信号は、ノートと呼ばれる一連のイベントとしてモデル化され得る。ノートをプレイする調波楽器の音は、たとえば、（アタックとも呼ばれる）オンセット段階、（サステーンとも呼ばれる）定常段階、および（リリースとも呼ばれる）オフセット段階の、時間的に異なる領域に分割され得る。ノートの時間エンベロープの別の記述（ADSR）は、アタックとサステーンとの間の追加のディケイ（decay）段階を含む。このコンテキストでは、ノートの持続時間は、アタック段階の開始からリリース段階の終了まで（または、同じ弦上の別のノートの開始など、そのノートを終了する別のイベントまで）の間隔として定義され得る。ノートは単一のピッチを有すると仮定されるが、インベントリは、単一のアタックと（たとえば、ビブラートまたはポルタメントなどのピッチベンディング効果によって生成される）複数のピッチとを有するノートをモデル化するようにも実装され得る。いくつかの楽器（たとえば、ピアノ、ギター、またはハープ）は、コードと呼ばれるイベントにおいて一度に2つ以上のノートを生成し得る。

【0228】

[00250]異なる楽器によって生成されるノートはサステーン段階中に同様の音色を有し得るので、そのような期間中にどの楽器がプレイしているかを識別することは困難であり得る。しかしながら、ノートの音色は、段階ごとに変化することが予想され得る。たとえば、アクティブな楽器を識別することは、サステーン段階中よりもアタックまたはリリース段階中に容易であり得る。

【0229】

[00251]アクティブ化係数ベクトルが適切な基底関数を示す可能性を高めるために、基底関数間の差分を最大にすることが望ましいことがある。たとえば、基底関数が時間に対するノートのスペクトルの変化に関係する情報を含むことが望ましいことがある。

【0230】

[00252]時間に対する音色の変化に基づいて基底関数を選択することが望ましいことがある。そのような手法は、ノートの音色のそのような時間領域展開に関係する情報を基底関数インベントリに符号化することを含み得る。たとえば、特定の楽器 n の基底関数のセット A_n は、2つ以上の対応する信号表現の各々がノートの展開における異なる時間（たとえば、アタック段階の時間、サステーン段階の時間、およびリリース段階の時間）に対応するように、各ピッチにおいてこれらの信号表現を含み得る。これらの基底関数は、ノートをプレイする楽器の記録の対応するフレームから抽出され得る。

【0231】

[00253]図27Cに、一般的構成による、オーディオ信号を分解するための装置MF100のブロック図を示す。装置MF100は、（たとえば、タスクT100に関して本明細書で説明したように）オーディオ信号のフレームからの情報に基づいて、周波数範囲にわたる対応する信号表現を計算するための手段F100を含む。装置MF100は、（たとえば、タスクT200に関して本明細書で説明したように）手段F100によって計算された信号表現と、複数の基底関数とに基づいて、アクティブ化係数の各々が複数の基底関数のうちの異なる1つに対応する、アクティブ化係数のベクトルを計算するための手段F200をも含む。

【0232】

[00254]図27Dに、変換モジュール2100と係数ベクトル計算器2200とを含む、別の一般的構成による、オーディオ信号を分解するための装置A100のブロック図を示す。変換モジュール2100は、（たとえば、タスクT100に関して本明細書で説明

したように)オーディオ信号のフレームからの情報に基づいて、周波数範囲にわたる対応する信号表現を計算するように構成される。係数ベクトル計算器 2 2 0 0 は、(たとえば、タスク T 2 0 0 に関して本明細書で説明したように)変換モジュール 2 1 0 0 によって計算された信号表現と、複数の基底関数とに基づいて、アクティブ化係数の各々が複数の基底関数のうちの異なる 1 つに対応する、アクティブ化係数のベクトルを計算するように構成される。

【0 2 3 3】

[00255]図 2 7 B に、基底関数インベントリが各ピッチにおける各楽器のための複数の信号表現を含む、方法 M 1 0 0 の実装形態 M 2 0 0 のフローチャートを示す。これらの複数の信号表現の各々は、周波数範囲にわたるエネルギーの複数の異なる分布(たとえば、複数の異なる音色)を記述する。インベントリはまた、異なる時間関係モダリティのために異なる複数の信号表現を含むように構成され得る。1 つのそのような例では、インベントリは、各ピッチにおける弾かれた(bowed)弦の複数の信号表現と、各ピッチにおけるはじかれた(plucked)(たとえば、ピッツィカート)の弦の異なる複数の信号表現とを含む。

10

【0 2 3 4】

[00256]方法 M 2 0 0 は、タスク T 1 0 0 の複数のインスタンス(この例では、タスク T 1 0 0 A および T 1 0 0 B)を含み、各インスタンスは、オーディオ信号の対応する異なるフレームからの情報に基づいて、周波数範囲にわたる対応する信号表現を計算する。様々な信号表現は連結され得、同様に、各基底関数は複数の信号表現の連結であり得る。この例では、タスク T 2 0 0 は、混合フレームの連結を各ピッチにおける信号表現の連結に整合させる。図 3 3 に、混合信号 y のフレーム $p 1$ とフレーム $p 2$ とが整合のために連結された、図 3 2 のモデル $B f = y$ の変形 $B' f = y$ の一例を示す。

20

【0 2 3 5】

[00257]インベントリは、各ピッチにおける複数の信号表現がトレーニング信号の連続するフレームからとられるように構築され得る。他の実装形態では、各ピッチにおける複数の信号表現が、時間的により大きい窓にわたる(たとえば、連続するフレームではなく時間的に分離されたフレームを含む)ことが望ましいことがある。たとえば、各ピッチにおける複数の信号表現が、アタック段階と、サステーン段階と、リリース段階との中の少なくとも 2 つからの信号表現を含むことが望ましいことがある。ノートの時間領域展開に関するより多くの情報を含むことによって、異なるノートの基底関数のセット間の差分が増加され得る。

30

【0 2 3 6】

[00258]図 2 8 A に、セグメントの高周波を強調するタスク T 3 0 0 を含む方法 M 1 0 0 の実装形態 M 3 0 0 のフローチャートを示す。この例では、タスク T 1 0 0 は、事前強調の後にセグメントの信号表現を計算するように構成される。図 2 9 A に、タスク T 3 0 0 の複数のインスタンス T 3 0 0 A、T 3 0 0 B を含む、方法 M 2 0 0 の実装形態 M 4 0 0 のフローチャートを示す。一例では、事前強調タスク T 3 0 0 は、2 0 0 H z を上回るエネルギーと総エネルギーとの比を増加させる。

【0 2 3 7】

40

[00259]図 2 8 B に、変換モジュール 2 1 0 0 の上流でオーディオ信号に対して高周波強調を実行するように構成された事前強調フィルタ 2 3 0 0 (たとえば、1 次高域フィルタなどの高域フィルタ)を含む装置 A 1 0 0 の実装形態 A 3 0 0 のブロック図を示す。図 2 8 C に、事前強調フィルタ 2 3 0 0 が変換係数に対して高周波事前強調を実行するように構成された、装置 A 1 0 0 の別の実装形態 A 3 1 0 のブロック図を示す。これらの場合、また、基底関数の複数 B に対して高周波事前強調(たとえば、高域フィルタ処理)を実行することが望ましいことがある。

【0 2 3 8】

[00260]楽音は、ビブラートおよび/またはトレモロなどのカラーレーション効果を含み得る。ビブラートは、一般に、4 または 5 から 7、8、1 0、または 1 2 ヘルツまでの

50

範囲内にある変調レートをもつ周波数変調である。ビブラートによるピッチ変化は、シンガーの場合には、0.6から2半音の間で変動し得、管弦楽器の場合には、概して+/-0.5半音よりも少ない（たとえば、弦楽器の場合には、0.2から0.35半音の間である）。トレモロは、一般に同様の変調レートを有する振幅変調である。

【0239】

[00261]基底関数インベントリにおいてそのような効果をモデル化することは困難であり得る。そのような効果の存在を検出することが望ましいことがある。たとえば、ビブラートの存在は、4~8Hzの範囲内の周波数領域ピークによって示され得る。また、そのような特性は、再生中に効果を復元するために使用され得るので、検出された効果のレベルの測度を（たとえば、このピークのエネルギーとして）記録することが望ましいことがある。トレモロの検出および定量化では、同様の処理が時間領域において実行され得る。効果が検出され、場合によっては定量化された後、ビブラートの場合には時間に対して周波数を平滑化することによって、またはトレモロの場合には時間に対して振幅を平滑化することによって変調を除去することが望ましいことがある。

【0240】

[00262]図30Bに、変調レベル計算器（MLC: modulation level calculator）を含む装置A100の実装形態A700のブロック図を示す。MLCは、上記で説明したように、オーディオ信号のセグメント中の検出された変調の測度（たとえば、時間または周波数領域における検出された変調ピークのエネルギー）を計算し、場合によっては記録するように構成される。

【0241】

[00263]本開示では、複数の音源が同時にアクティブになり得る音楽アプリケーションのための使用事例を可能にするために使用され得る方法について説明する。そのような事例では、可能な場合、アクティブ化係数ベクトルを計算する前に音源を分離することが望ましいことがある。この目的を達成するために、マルチチャネル技法とシングルチャネル技法との組合せが提案される。

【0242】

[00264]図29Bに、信号を空間クラスタに分離するタスクT500を含む方法M100の実装形態M500のフローチャートを示す。タスクT500は、音源をできる限り多くの空間クラスタに隔離するように構成され得る。一例では、タスクT500は、記録された音響シナリオをできる限り多くの空間クラスタに分離するために、マルチマイクロフォン処理を使用する。そのような処理は、マイクロフォン信号間の利得差および/または位相差に基づき得、そのような差は、周波数帯域全体にわたって評価されるか、あるいは複数の異なる周波数サブバンドまたは周波数ビンの各々において評価され得る。

【0243】

[00265]空間分離方法のみでは、所望の分離レベルを達成するには不十分であり得る。たとえば、いくつかの音源は、マイクロフォンアレイに対して近接しすぎているか、または場合によっては準最適に構成されることがある（たとえば、複数のバイオリン奏者および/または調波楽器が1つのコーナーに位置し得、打楽器奏者が通常後方に位置する）。典型的な音楽バンドシナリオでは、音源は（たとえば、図34に示すように）互いに近接して位置するかまたは他の音源のさらに後ろに位置し得るので、空間情報のみを使用して、バンドに対して同じ概略的な方向にあるマイクロフォンのアレイによってキャプチャされた信号を処理すると、音源のすべてを互いから区別することができないことがある。タスクT100およびT200は、（たとえば、図34に示すように）個々の楽器を分離するために、本明細書で説明するシングルチャネル基底関数インベントリベースのスパース復元（たとえば、スパース分解）技法を使用して個々の空間クラスタを分析し得る。

【0244】

[00266]計算しやすさのために、基底関数の複数Bは、基底関数のインベントリAよりもかなり小さいことが望ましいことがある。大きいインベントリから開始して、所与の分離タスクのためのインベントリを狭めることが望ましいことがある。基底関数Bのセット

の選択は、記録されたシーン中の楽器の視覚的認識に基づいて低減され得る。たとえば、B個の基底関数は、図2の画像/ビデオ処理ブロック54によって与えられる楽器候補のリストに対応するもの、または図6のシステム500によって識別されるものに制限され得る。

【0245】

[00267]別の例では、そのような低減はまた、セグメントが打楽器からの音を含むのか調波楽器からの音を含むのかを判断することと、整合のためにインベントリから基底関数の適切な複数Bを選択することによって実行され得る。打楽器は、調波音の場合の水平線とは反対にインパルス様のスペクトログラム（たとえば、垂直線）を有する傾向がある。

10

【0246】

[00268]調波楽器は、一般に、ある基本ピッチおよび関連する音色と、この調波パターンの対応する高周波拡張とによって、スペクトログラムにおいて特徴づけられ得る。したがって、別の例では、スペクトルの高周波レプリカは、低周波スペクトルに基づいて予測され得るので、これらのスペクトルのより低いオクターブのみを分析することによって計算タスクを低減することが望ましいことがある。整合の後に、アクティブな基底関数を高周波に外挿し、混合信号から減算して、残差信号を取得し得、残差信号は、符号化されおよび/またはさらに分解され得る。

【0247】

[00269]そのような低減はまた、グラフィカルユーザインターフェースにおけるユーザ選択を通して実行され、ならびに/あるいは、第1のスパース復元ラン（sparse recovery run）または最尤適合に基づく、可能性が最も高い楽器および/またはピッチの事前分類によって実行され得る。たとえば、スパース復元演算の第1のランを実行して、復元されたスパース係数の第1のセットを取得し得、この第1のセットに基づいて、適用可能なノート基底関数がスパース復元演算の別のランのために狭められ得る。

20

【0248】

[00270]1つの低減手法は、いくつかのピッチ間隔においてスパースシテスコアを測定することによって、いくつかの楽器ノートの存在を検出することを含む。そのような手法は、初期ピッチ推定値に基づいて、1つまたは複数の基底関数のスペクトル形状を改善することと、方法M100において、改善された基底関数を複数Bとして使用することとを含む得る。

30

【0249】

[00271]低減手法は、対応する基底関数に射影された音楽信号のスパースシテスコアを測定することによってピッチを識別するように構成され得る。最良のピッチスコアが与えられれば、基底関数の振幅形状は、楽器ノートを識別するために最適化され得る。アクティブな基底関数の低減されたセットは、次いで、方法M100において複数Bとして使用され得る。

【0250】

[00272]図30Aに、基底関数の第1ランインベントリ低減を含む方法M100の実装形態M600のフローチャートを示す。方法M600は、（たとえば、メルまたはバーク尺度の場合のように、隣接する要素間の周波数距離が周波数とともに増加する）非線形周波数領域におけるセグメントの信号表現を計算するタスクT600を含む。一例では、タスクT600は、定Q変換（constant-Q transform）を使用して非線形信号表現を計算するように構成される。方法M600はまた、非線形信号表現と、複数の同様に非線形の基底関数とに基づいて、アクティブ化係数の第2のベクトルを計算するタスクT700を含む。第2のアクティブ化係数ベクトルからの（たとえば、アクティブなピッチ範囲を示し得るアクティブ化された基底関数の識別情報からの）情報に基づいて、タスクT800は、タスクT200において使用する基底関数の複数Bを選択する。また、方法M200、M300、およびM400は、そのようなタスクT600、T700、およびT800を含むように実装され得ることに明確に留意されたい。

40

50

【 0 2 5 1 】

[00273]図 3 1 に、基底関数のより大きいセットから（たとえば、インベントリから）複数の基底関数を選択するように構成されたインベントリ低減モジュール（IRM: inventory reduction module）を含む装置 A 1 0 0 の実装形態 A 8 0 0 のブロック図を示す。モジュールIRMは、（たとえば、定Q変換に従って）非線形周波数領域におけるセグメントの信号表現を計算するように構成された第2の変換モジュール2 1 1 0を含む。モジュールIRMは、本明細書で説明するように、非線形周波数領域における計算された信号表現と、第2の複数の基底関数とに基づいて、アクティブ化係数の第2のベクトルを計算するように構成された第2の係数ベクトル計算器2 2 1 0をも含む。モジュールIRMは、本明細書で説明するように、第2のアクティブ化係数ベクトルからの情報に基づいて、基底関数のインベントリの中から複数の基底関数を選択するように構成された基底関数セクタをも含む。

10

【 0 2 5 2 】

[00274]方法M 1 0 0 は、オンセット検出（たとえば、楽音のオンセットを検出すること）と、調波楽器スパース係数を改善するための後処理とを含むことが望ましいことがある。アクティブ化係数ベクトル f は、楽器固有の基底関数セット B_n のアクティブ化係数を含む、各楽器 n の対応するサブベクトル f_n を含むと見なされ得、これらのサブベクトルは独立して処理され得る。

【 0 2 5 3 】

[00275]一般的なオンセット検出方法はスペクトルの大きさ（たとえば、エネルギー差）に基づき得る。たとえば、そのような方法は、スペクトルエネルギーおよび/またはピークスロープに基づいてピークを見つけることを含み得る。

20

【 0 2 5 4 】

[00276]また、各個々の楽器のオンセットを検出することが望ましいことがある。たとえば、調波楽器の中のオンセット検出の方法は、時間的な対応する係数差に基づき得る。1つのそのような例では、調波楽器 n のオンセット検出は、現在のフレームの楽器 n の係数ベクトル（サブベクトル f_n ）の最大大きさの要素のインデックスが、前のフレームの楽器 n の係数ベクトルの最大大きさの要素のインデックスに等しくない場合にトリガされる。そのような動作は、各楽器について反復され得る。

【 0 2 5 5 】

[00277]調波楽器のスパース係数ベクトルの後処理を実行することが望ましいことがある。たとえば、調波楽器では、大きい大きさを有し、および/または指定された基準を満たす（たとえば、十分に鋭い）アタックプロファイルを有する、対応するサブベクトルの係数を保持すること、ならびに/あるいは残差係数を除去（たとえば、ゼロアウト）することが望ましいことがある。

30

【 0 2 5 6 】

[00278]各調波楽器について、支配的な大きさと許容できるアタック時間とを有する係数が保持され、残差係数がゼロ化されるように、（たとえば、オンセット検出が示されたときに）各オンセットフレームにおいて係数ベクトルを後処理することが望ましいことがある。アタック時間は、時間に対する平均大きさなどの基準に従って評価され得る。1つのそのような例では、係数の現在の平均値が係数の過去の平均値よりも小さい場合（たとえば、フレーム $(t - 5)$ からフレーム $(t + 4)$ までなど、現在の窓にわたる係数の値の和が、フレーム $(t - 15)$ からフレーム $(t - 6)$ までなど、過去の窓にわたる係数の値の和よりも小さい場合）、現在のフレーム t の楽器の各係数はゼロアウトされる（すなわち、アタック時間は許容できない）。各オンセットフレームにおける調波楽器のための係数ベクトルのそのような後処理は、最大大きさをもつ係数を保持し、他の係数をゼロアウトすることをも含み得る。各非オンセットフレームにおける各調波楽器について、前のフレーム中の値が0でなかった係数のみを保持し、ベクトルの他の係数をゼロアウトするように係数ベクトルを後処理することが望ましいことがある。

40

【 0 2 5 7 】

50

[00279] 上述のように、E M アルゴリズムは、初期基底関数行列を発生するため、および / または (たとえば、アクティブ化係数ベクトルに基づいて) 基底関数行列を更新するために使用され得る。E M 手法のための更新ルールの例について次に説明する。スペクトログラム $V_{f,t}$ が与えられれば、各時間フレームについてスペクトル基底ベクトル $P(f|z)$ と重みベクトル $P_t(z)$ とを推定することが望まれる。これらの分布から行列分解が得られる。

【0258】

[00280] E M アルゴリズムは、以下のように適用される。最初に、重みベクトル $P_t(z)$ とスペクトル基底ベクトル $P(f|z)$ とをランダムに初期化する。次いで、収束するまで後続のステップ間を反復する。1) 予想 (E) ステップ - スペクトル基底ベクトル $P(f|z)$ と重みベクトル $P_t(z)$ とが与えられれば、後の分布 $P_t(z|f)$ を推定する。この推定は、以下のように表され得る。

10

【数25】

$$P_t(z|f) = \frac{P_t(f|z)P(z)}{\sum_z P_t(f|z)P(z)}$$

【0259】

[00281] 2) 最大化 (M) ステップ - 後の分布 $P_t(z|f)$ が与えられれば、重みベクトル $P_t(z)$ とスペクトル基底ベクトル $P(f|z)$ とを推定する。重みベクトルの推定は、以下のように表され得る。

20

【数26】

$$P_t(z) = \frac{\sum_f V_{f,t} P_t(z|f)}{\sum_z \sum_f V_{f,t} P_t(z|f)}$$

【0260】

スペクトル基底ベクトルの推定は、以下のように表され得る。

【数27】

30

$$P(f|z) = \frac{\sum_t V_{f,t} P_t(z|f)}{\sum_t \sum_f V_{f,t} P_t(z|f)}$$

【0261】

[00282] 本明細書で開示するシステムおよび方法は、コンピュータ、ゲーミングコンソール、またはセルラーフォン、携帯情報端末 (PDA)、スマートフォンなどのハンドヘルドデバイスなどを含む、任意の好適な視聴覚システム中に含まれ得る。本明細書で説明した構成要素の主な機能は、概してデジタル処理領域において実装される。しかしながら、これらの構成要素は、代替的に、好適なアナログ構成要素を使用するアナログ領域において実装されるか、またはアナログ電子構成要素とデジタル電子構成要素との任意の好適な組合せにおいて実装され得る。

40

【0262】

[00283] 音響信号を受信するように構成された2つ以上のマイクロフォンのアレイと、1つまたは複数のカメラとを有するポータブル視聴覚感知デバイス内で、本明細書で説明したシステムおよび方法を実装することが望ましいことがある。そのようなアレイを含むように実装され得、オーディオ記録および / または音声通信適用例のために使用され得るポータブルオーディオ感知デバイスの例としては、電話ハンドセット (たとえば、セルラー電話ハンドセット)、ハンドヘルドオーディオおよび / またはビデオレコーダ、携帯情報端末 (PDA) または他のハンドヘルドコンピューティングデバイス、およびノートブ

50

ックコンピュータ、ラップトップコンピュータ、ネットブックコンピュータ、タブレットコンピュータ、あるいは他のポータブルコンピューティングデバイスがある。ポータブルコンピューティングデバイスの種類は現在、ラップトップコンピュータ、ノートブックコンピュータ、ネットブックコンピュータ、ウルトラポータブルコンピュータ、タブレットコンピュータ、モバイルインターネットデバイス、スマートブック、およびスマートフォンなどの名称を有するデバイスを含む。そのようなデバイスは、ディスプレイスクリーンを含む上部パネルと、キーボードを含み得る下部パネルとを有し得、それらの2つのパネルは、クラムシェルまたは他のヒンジ結合関係で接続され得る。そのようなデバイスは、上面上にタッチスクリーンディスプレイを含むタブレットコンピュータとして同様に実装され得る。そのような方法を実行するように構築され得、オーディオ記録および/または音声通信適用例のために使用され得るオーディオ感知デバイスの他の例としては、テレビジョンディスプレイ、セットトップボックス、ならびにオーディオおよび/またはビデオ会議デバイスがある。

10

【0263】

[00284]本明細書で開示するシステムおよび方法は、リアルタイムで視聴覚情報を処理するように、ならびに以前に記録された視聴覚情報を処理するように実装され得る。

【0264】

[00285]本明細書で説明したシステム、装置、デバイスおよびそれらのそれぞれの構成要素の機能、ならびに方法ステップおよびモジュールは、ハードウェアで実装されるか、ハードウェアによって実行されるソフトウェア/ファームウェアで実装されるか、またはそれらの任意の好適な組合せで実装され得る。ソフトウェア/ファームウェアは、マイクロプロセッサ、DSP、埋込みコントローラまたは知的財産(IP: intellectual property)コアなど、1つまたは複数のデジタル回路によって実行可能な命令のセット(たとえば、プログラミングコードセグメント)を有するプログラムであり得る。ソフトウェア/ファームウェアで実装される場合、機能は、命令またはコードとして1つまたは複数のコンピュータ可読媒体上に記憶され得る。コンピュータ可読媒体はコンピュータ記憶媒体を含み得る。記憶媒体は、コンピュータによってアクセスされ得る任意の利用可能な媒体であり得る。限定ではなく例として、そのようなコンピュータ可読媒体は、RAM、ROM、EEPROM(登録商標)、CD-ROMまたは他の光ディスクストレージ、磁気ディスクストレージまたは他の磁気ストレージデバイス、あるいは命令またはデータ構造の形態の所望のプログラムコードを搬送または記憶するために使用され得、コンピュータによってアクセスされ得る、任意の他の媒体を備えることができる。本明細書で使用するディスク(disk)およびディスク(disc)は、コンパクトディスク(disc)(CD)、レーザディスク(登録商標)(disc)、光ディスク(disc)、デジタル多用途ディスク(disc)(DVD)、フロッピー(登録商標)ディスク(disk)およびblue-ray(登録商標)ディスク(disc)を含み、ディスク(disk)は、通常、データを磁氣的に再生し、ディスク(disc)は、データをレーザで光学的に再生する。上記の組合せもコンピュータ可読媒体の範囲内に含まれるべきである。

20

30

【0265】

[00286]視聴覚シーン分析システムおよび方法のいくつかの例が開示された。これらのシステムおよび方法は例であり、可能な結合は本明細書で説明したものに限定されない。その上、これらの例に対する様々な変更が可能であり、本明細書で提示した原理は他のシステムにも同様に適用され得る。たとえば、本明細書で開示する原理は、パーソナルコンピュータ、エンターテインメントカウンセラー、ビデオゲームなどのデバイスに適用され得る。さらに、様々な構成要素および/または方法ステップ/ブロックは、特許請求の範囲から逸脱することなく、明確に開示したものの以外の構成で実装され得る。

40

【0266】

[00287]したがって、これらの教示に鑑みて、他の実施形態および変更形態は当業者に容易に行われる。したがって、以下の特許請求の範囲は、上記の明細書および添付の図面とともに閲覧されたとき、すべてのそのような実施形態および変更形態を包含するもので

50

ある。

【図 1】

図 1

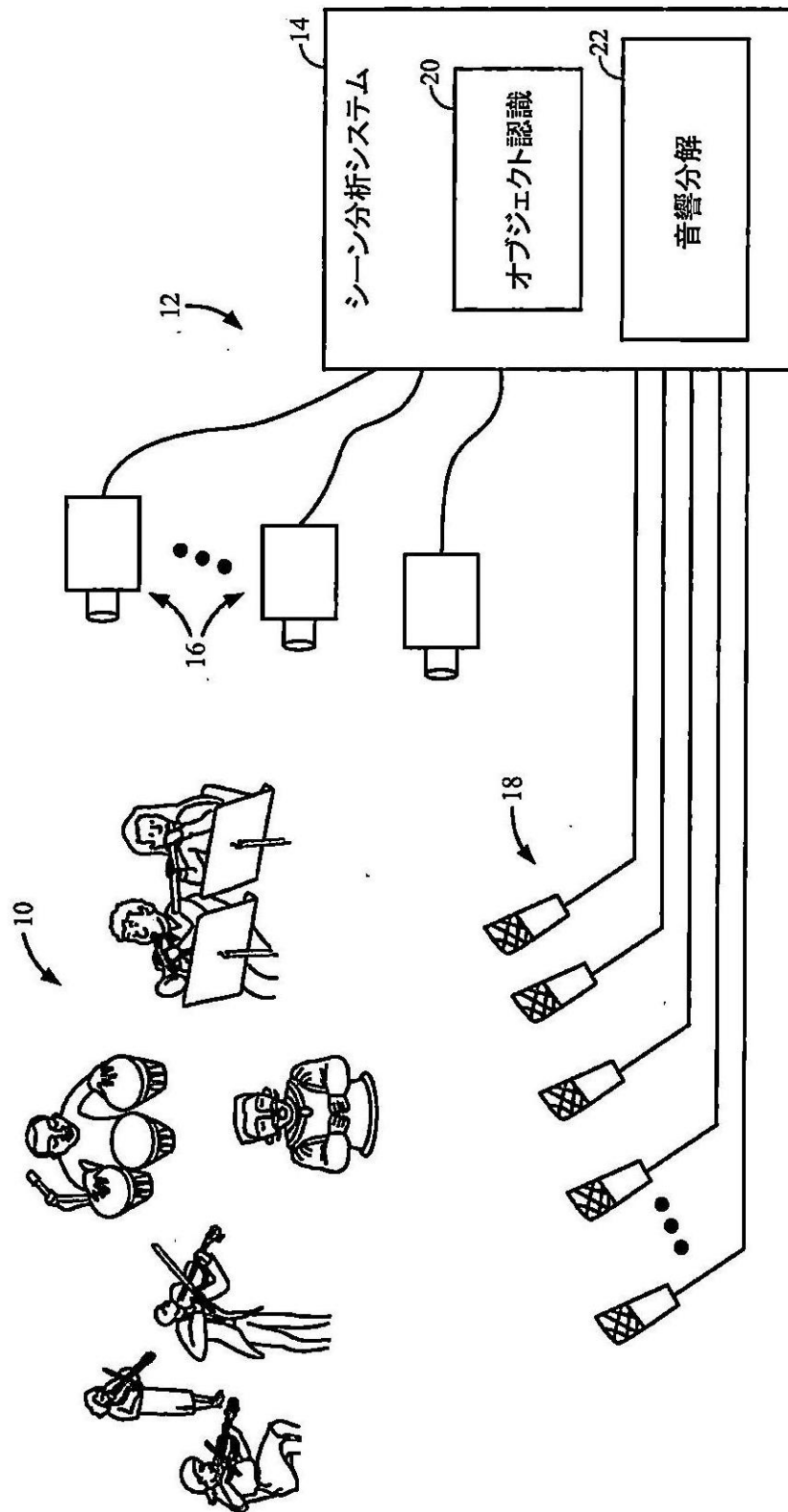


FIG. 1

【図 2】

図 2

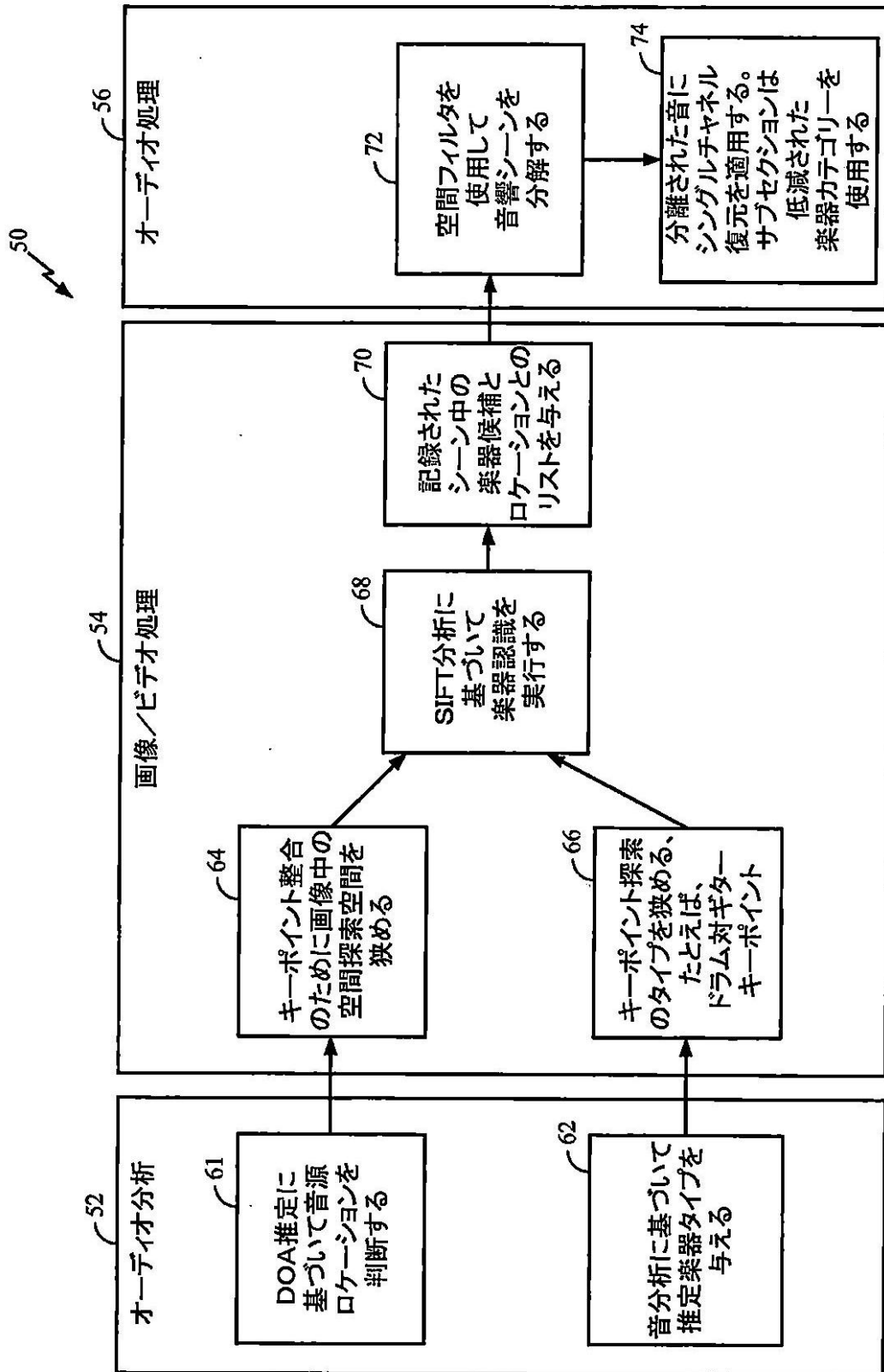


FIG. 2

【図 3】

図 3

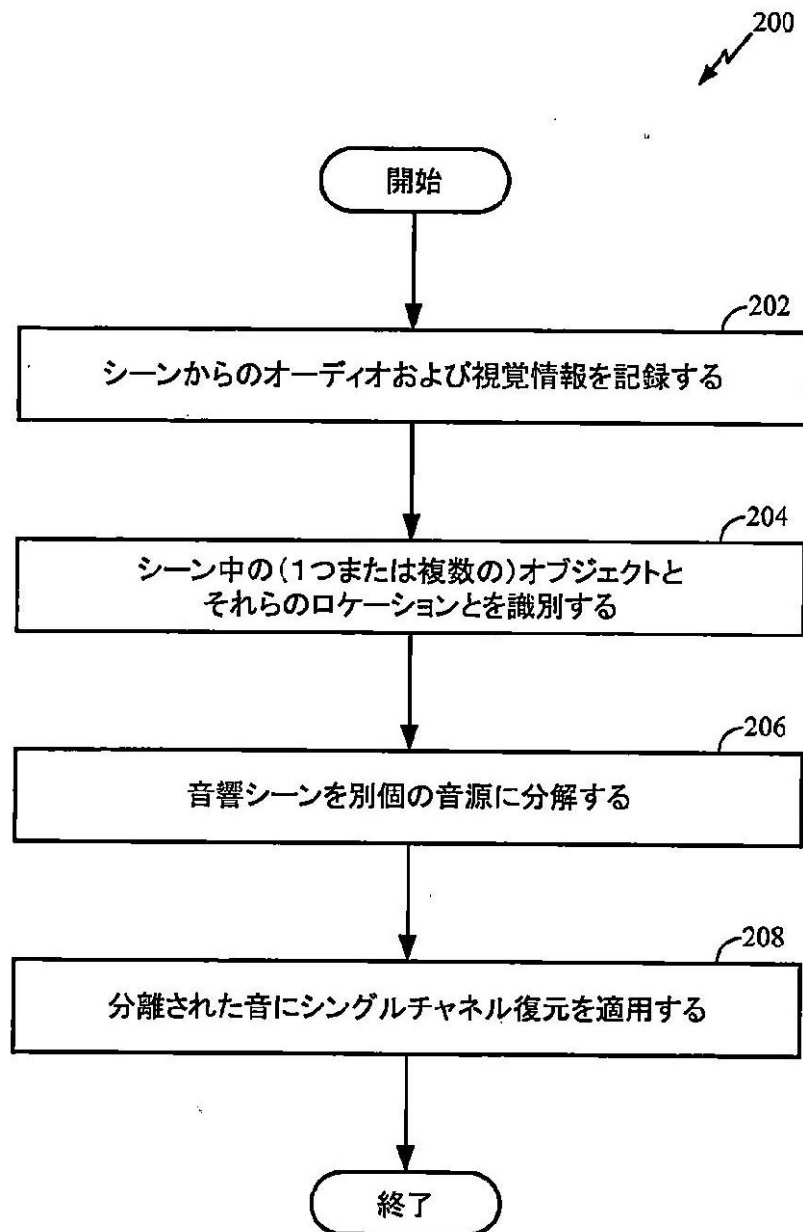


FIG. 3

【図 4】

図 4

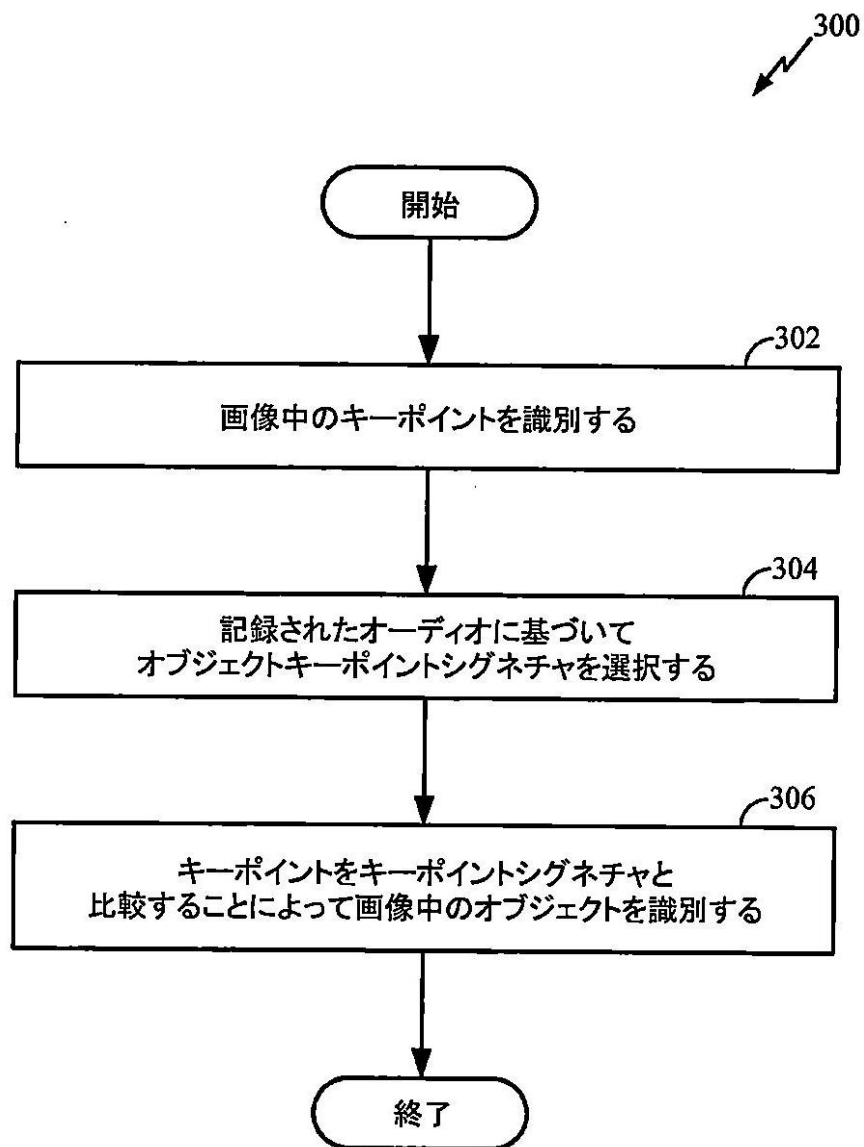


FIG. 4

【図 5 A】

図 5A

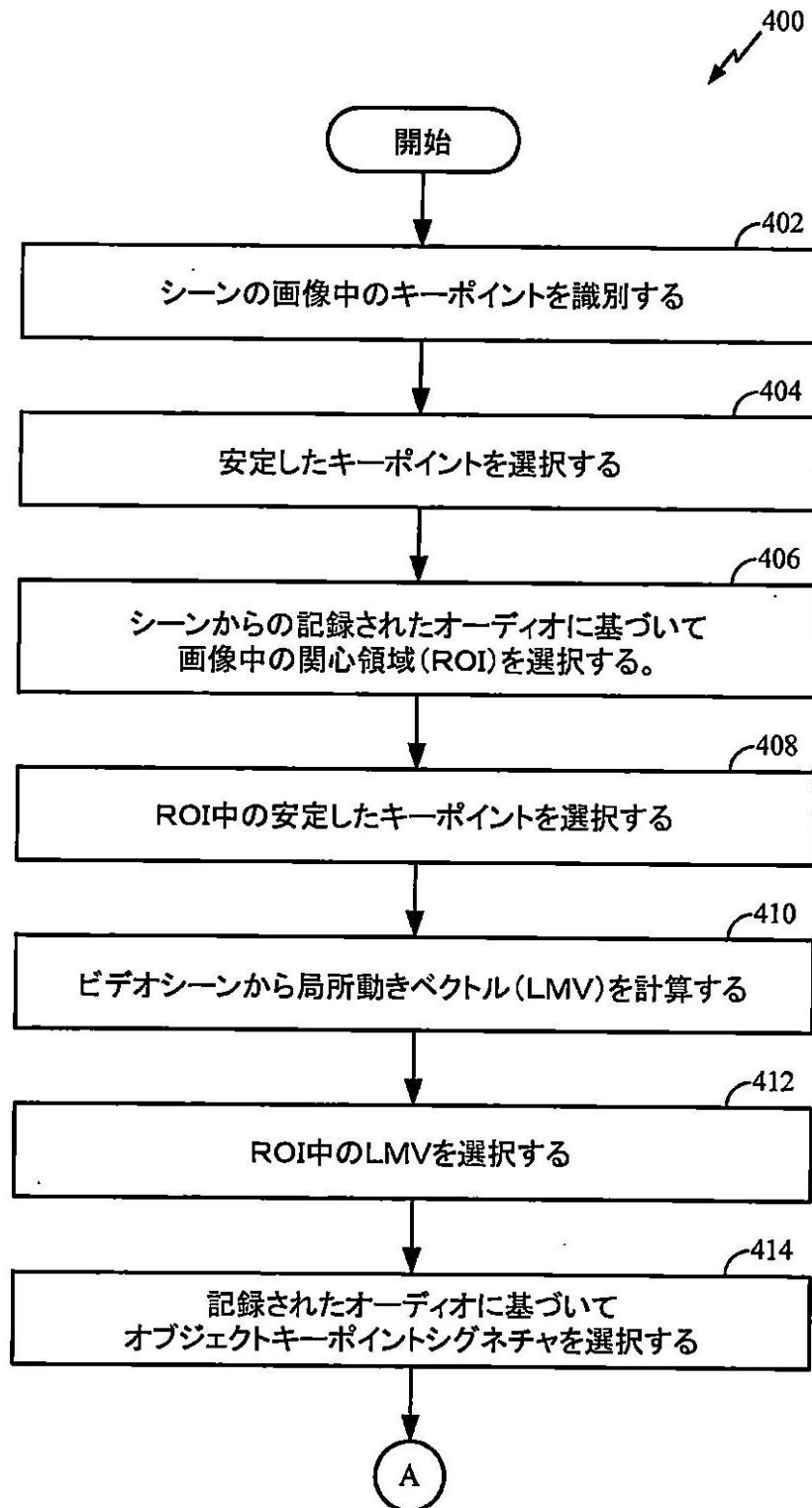


FIG. 5A

【図 5 B】

図 5B

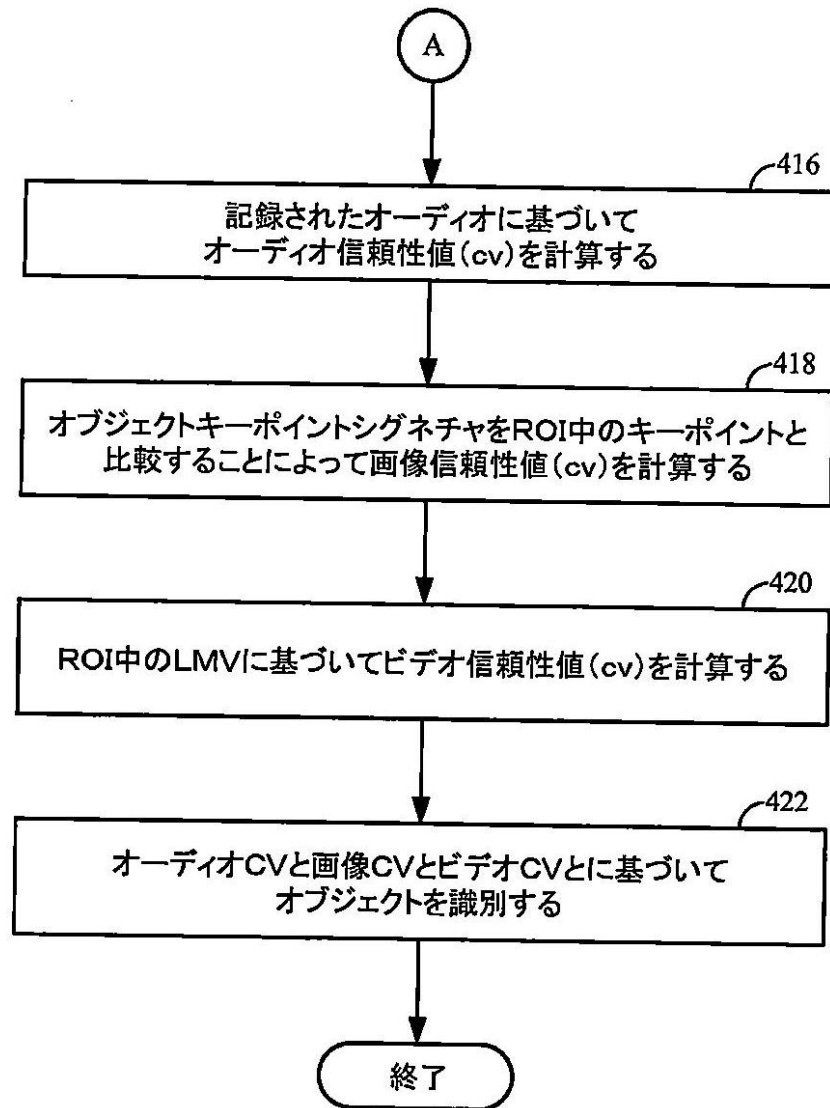


FIG. 5B

【図6】

図 6

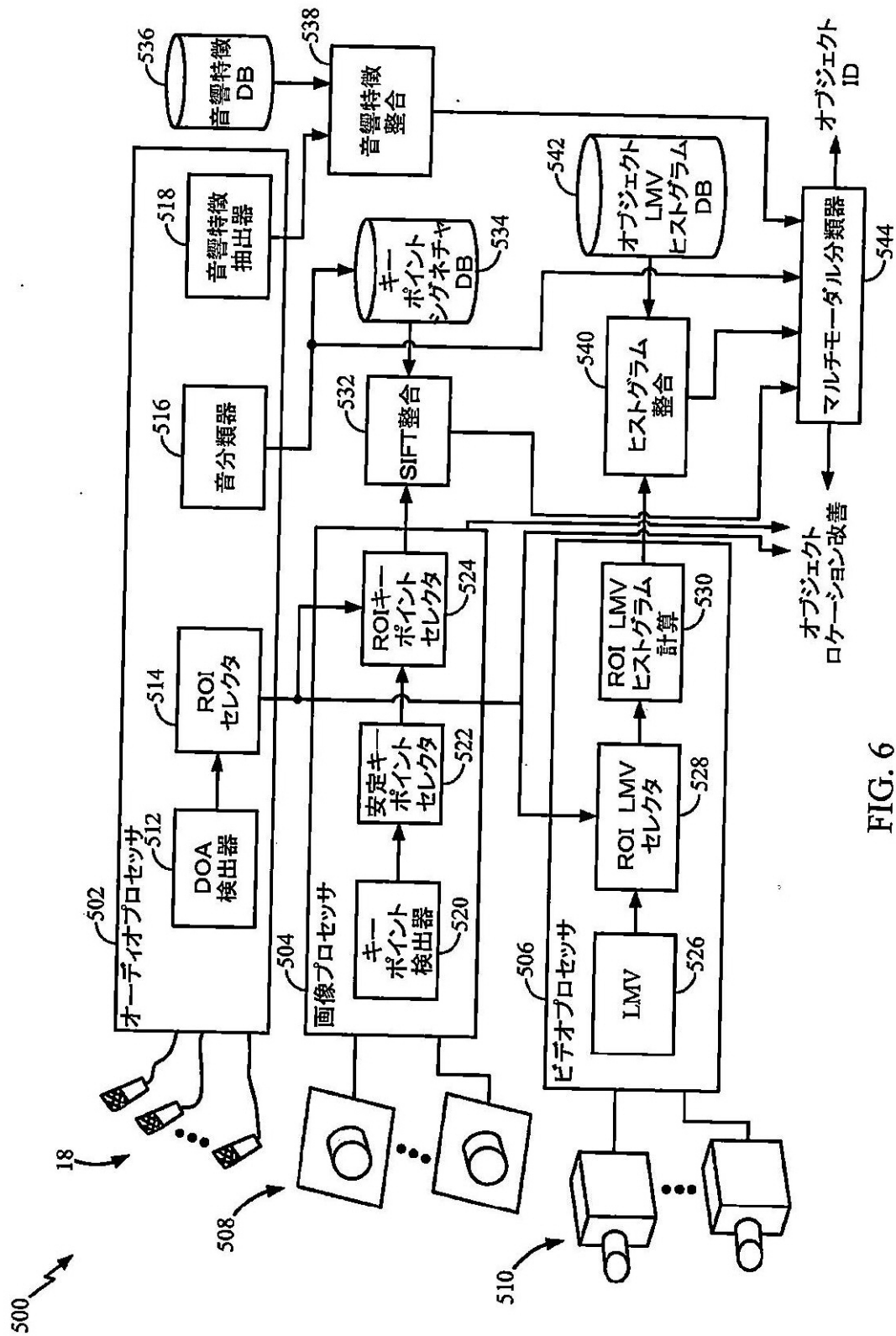


FIG. 6

【図 7】

図 7

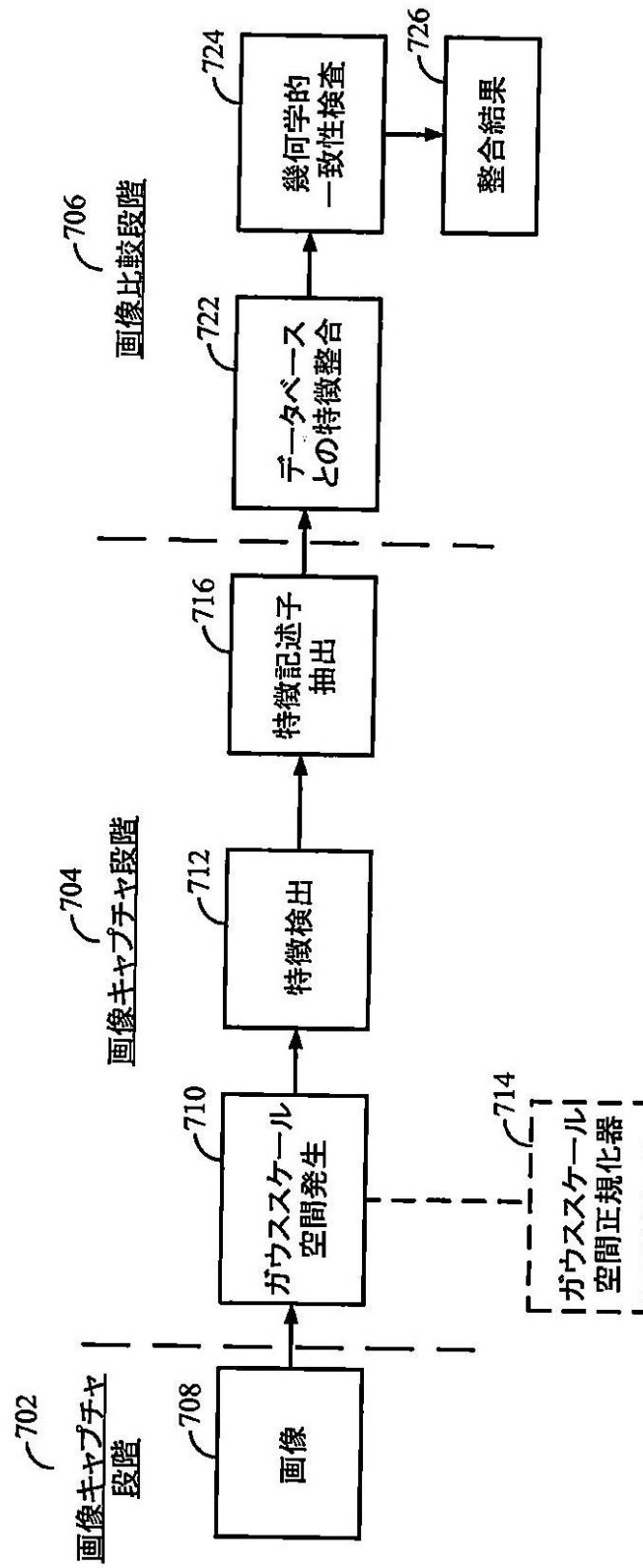


FIG. 7

【図 8】

図 8

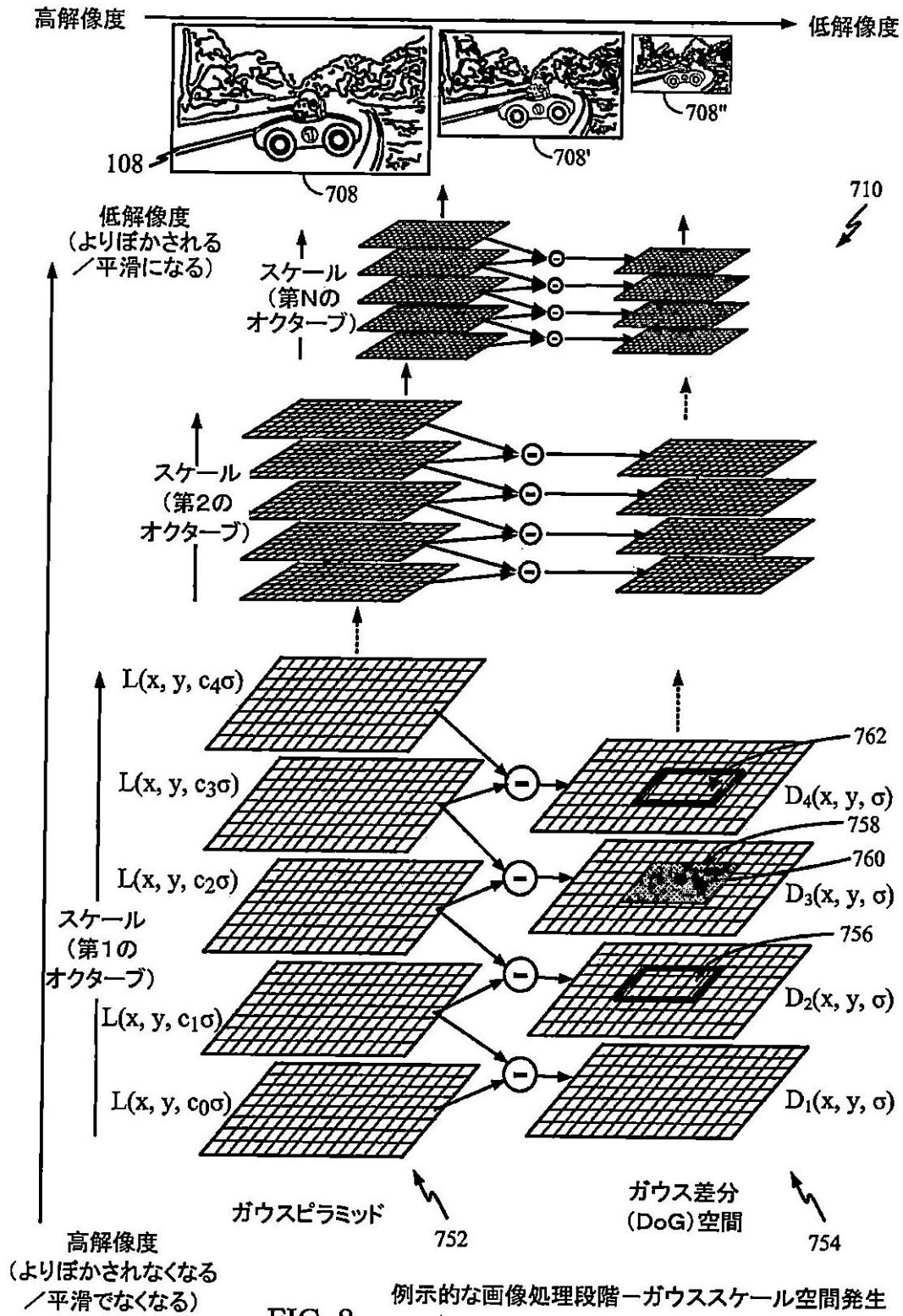


FIG. 8

【図 9】

図 9

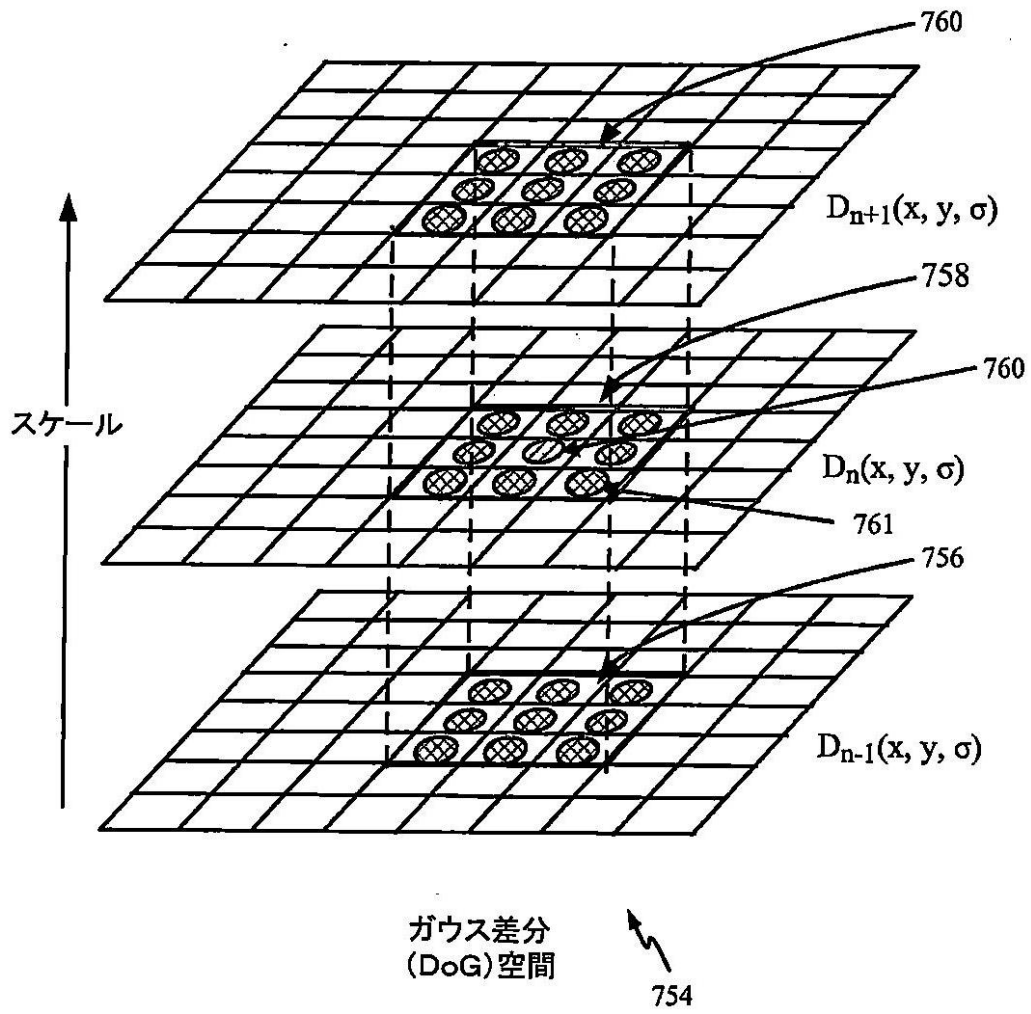


FIG. 9

【図 10】

図 10

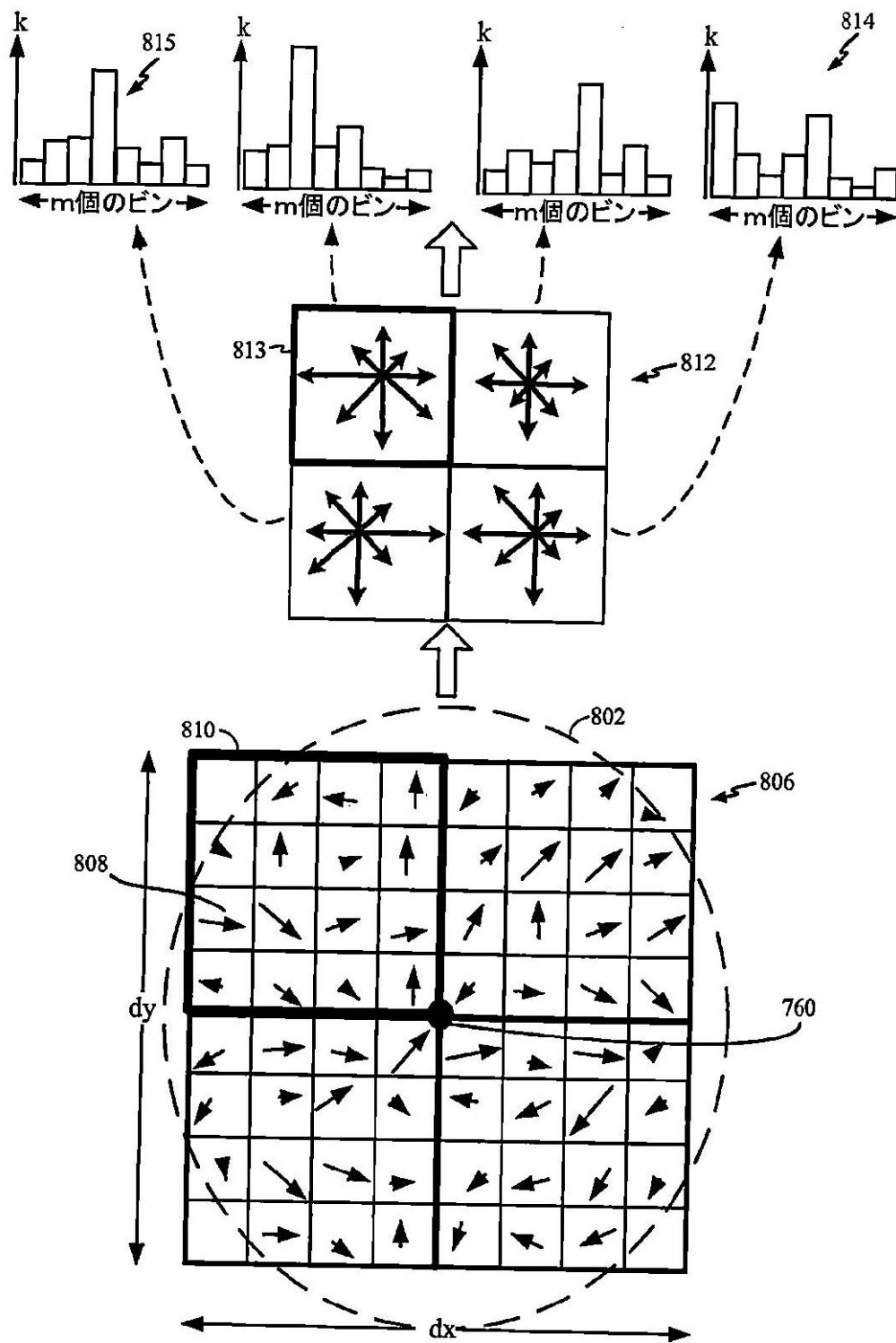
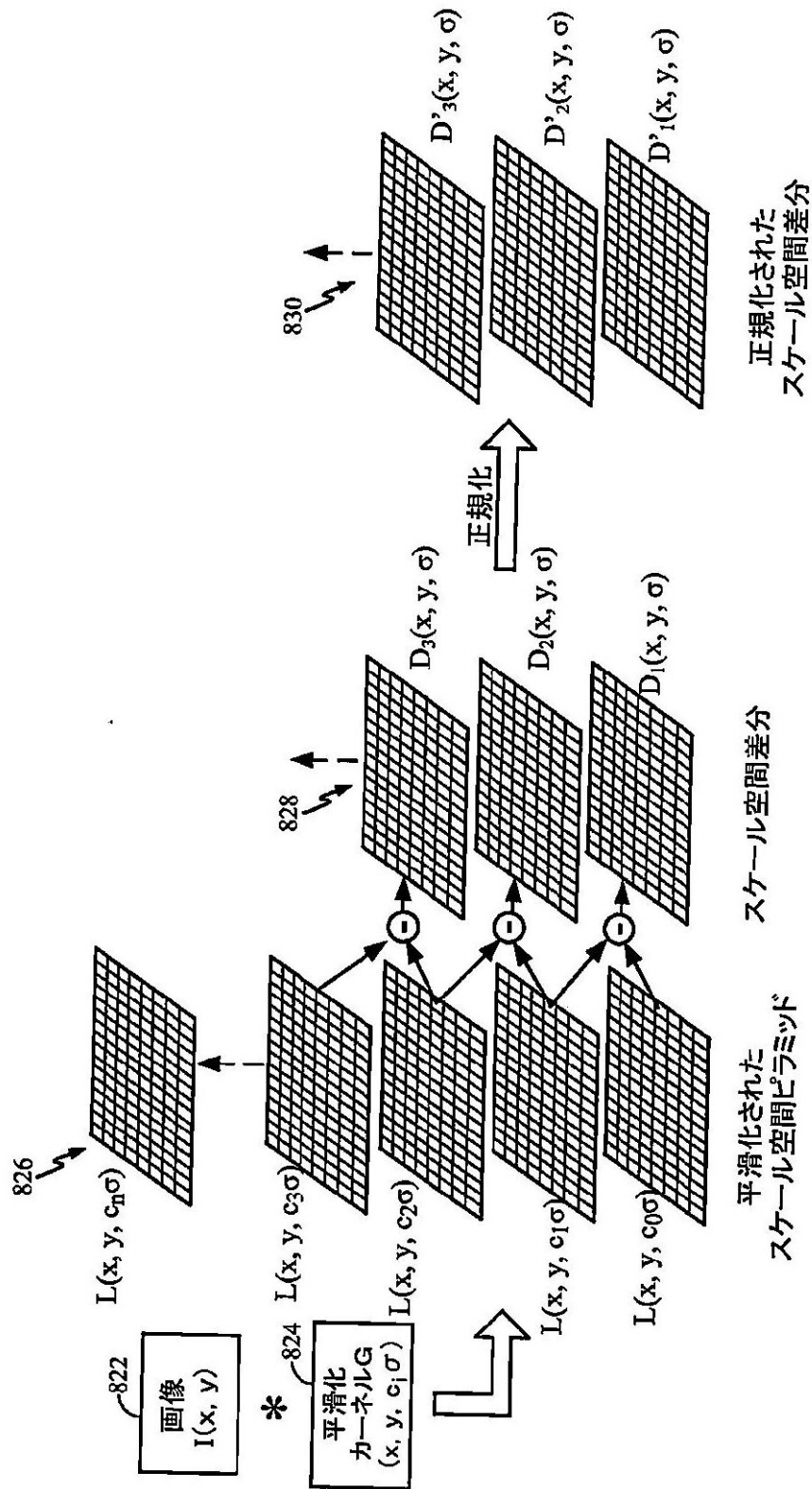


FIG. 10

【図 11】

図 11

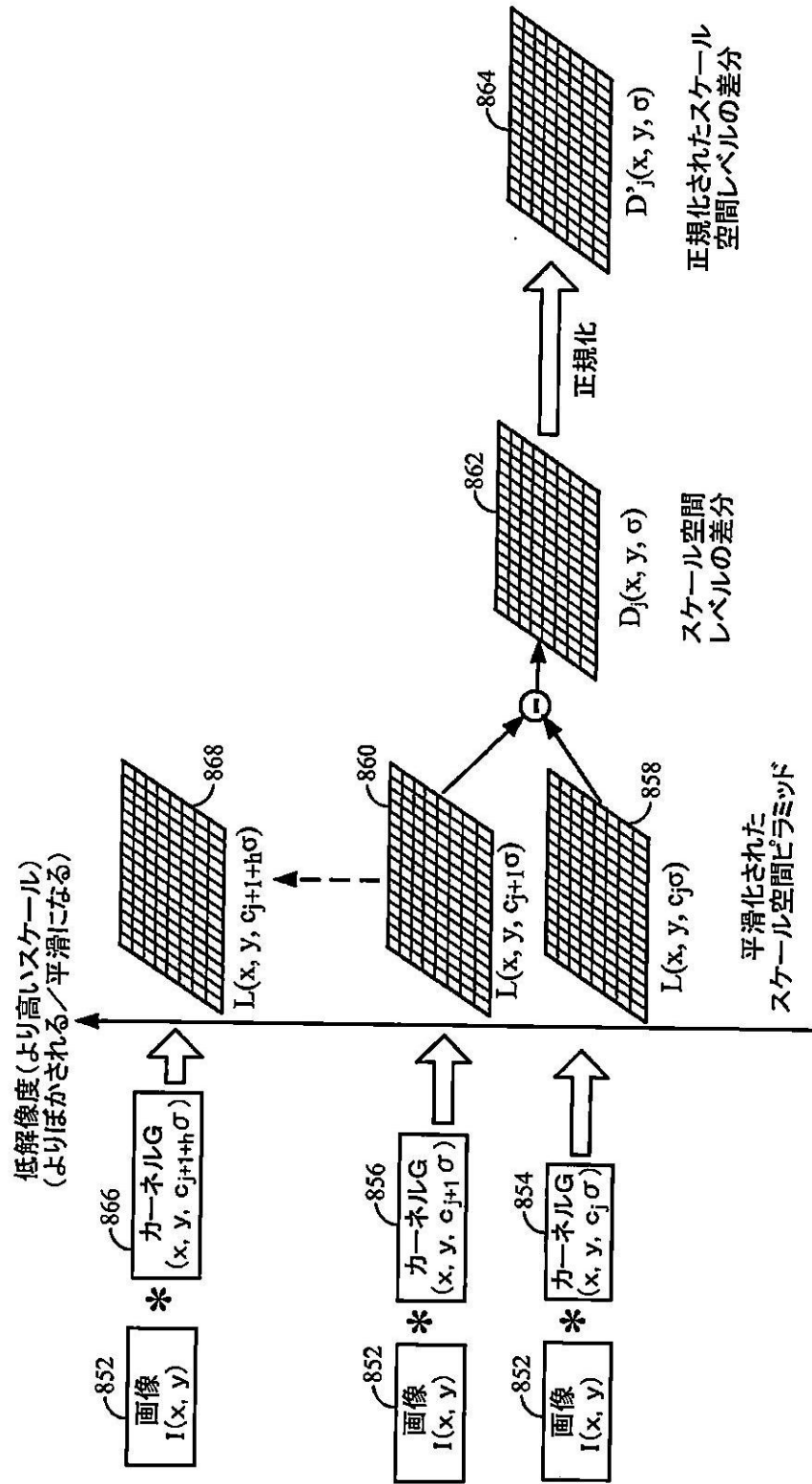


$$D'_j(x, y, \sigma) = \left[\frac{(G(x, y, c_{j+1}\sigma) - G(x, y, c_j\sigma)) * [I(x, y)S(x, y)]}{G(x, y, c_{j+1}\sigma) * [I(x, y)S(x, y)]} \right]$$

FIG. 11

【図 12】

図 12



$$D'_j(x, y, \sigma) = \left[\frac{(G(x, y, c_{j+1}\sigma) - G(x, y, c_{j+1+h}\sigma)) * [I(x, y)S(x, y)]}{G(x, y, c_{j+1+h}\sigma) * [I(x, y)S(x, y)]} \right]$$

高解像度(より低いスケール)
(よりぼかされなくなる/平滑でなくなる)

FIG. 12

【図 13】

図 13

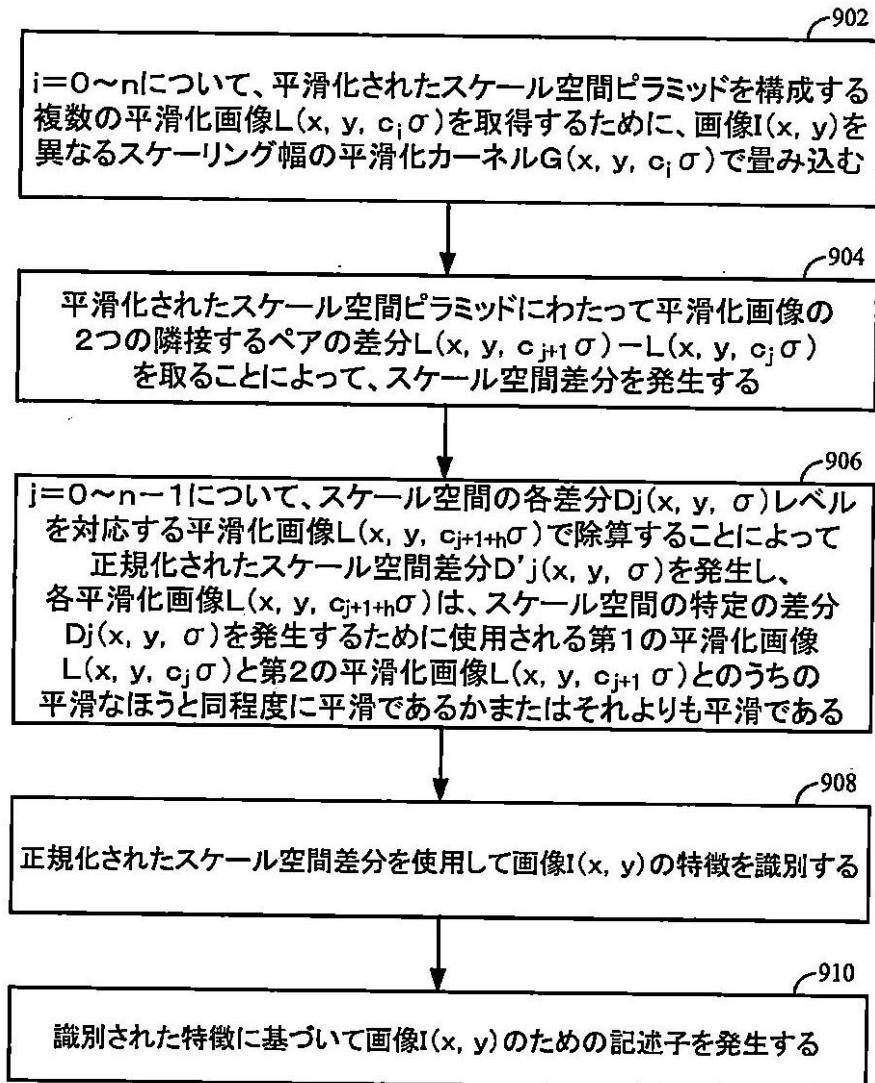


FIG. 13

【 図 1 4 】

図 14

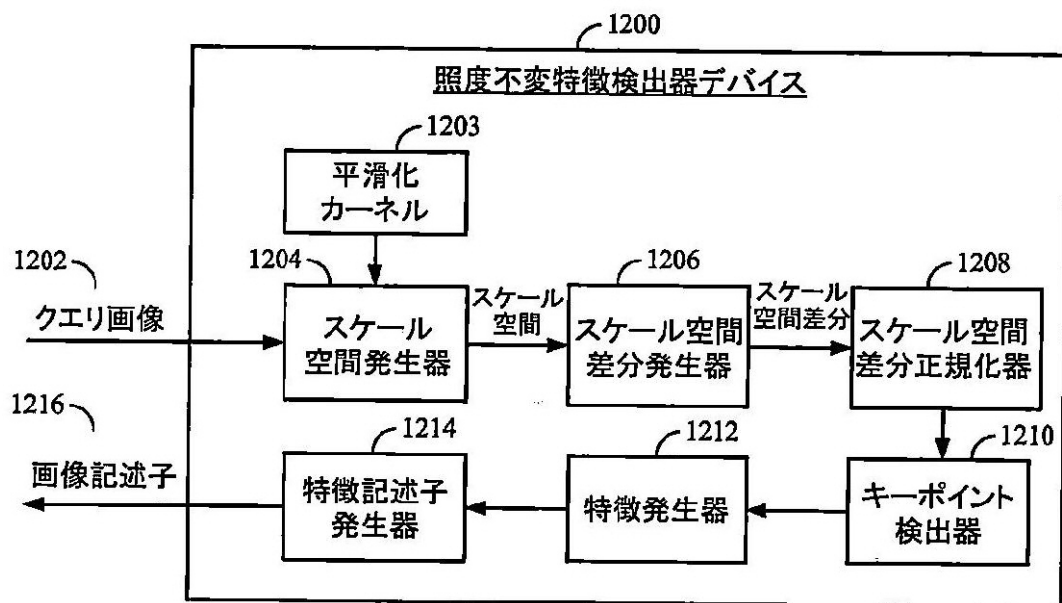


FIG. 14

【図 15】

図 15

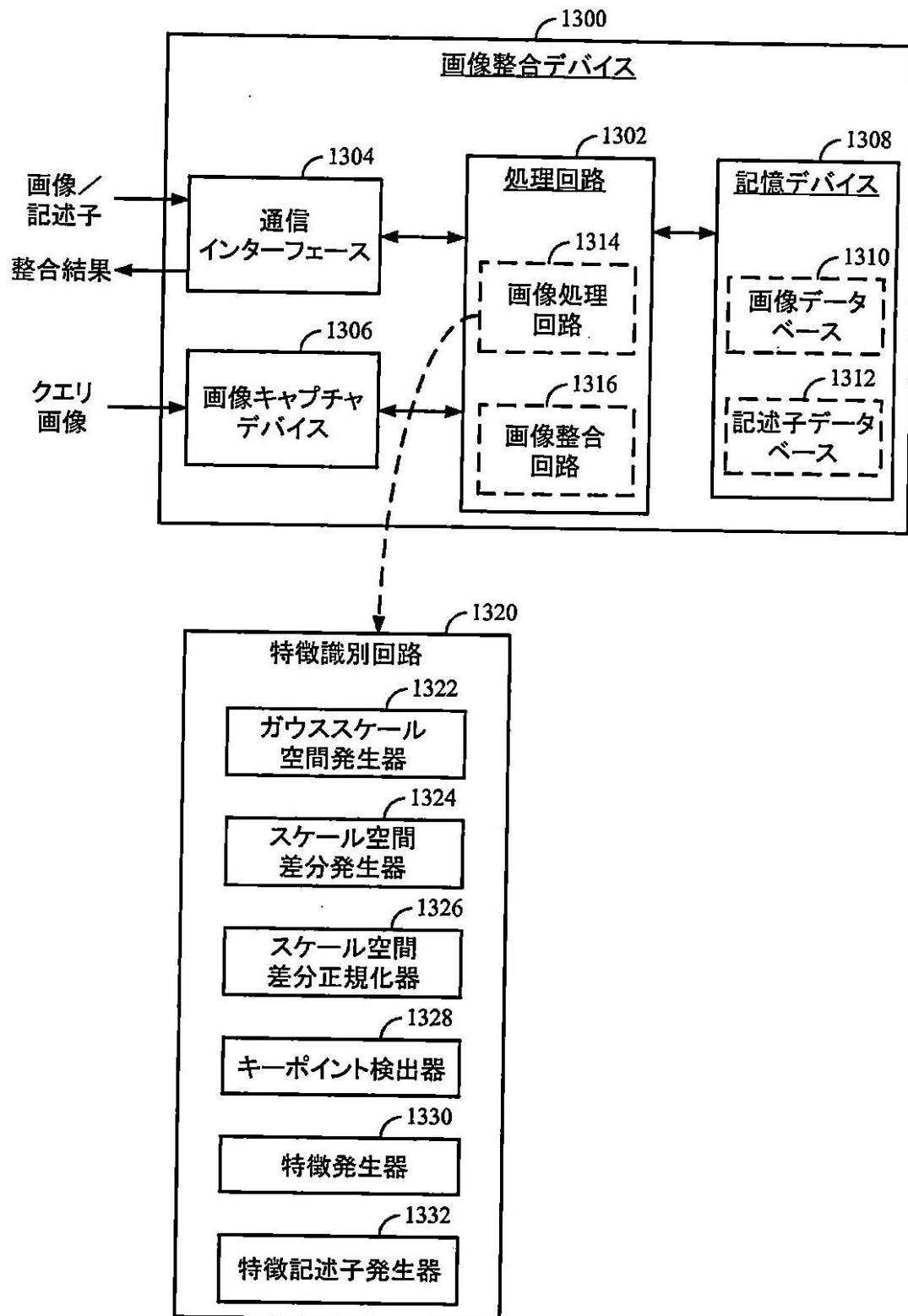
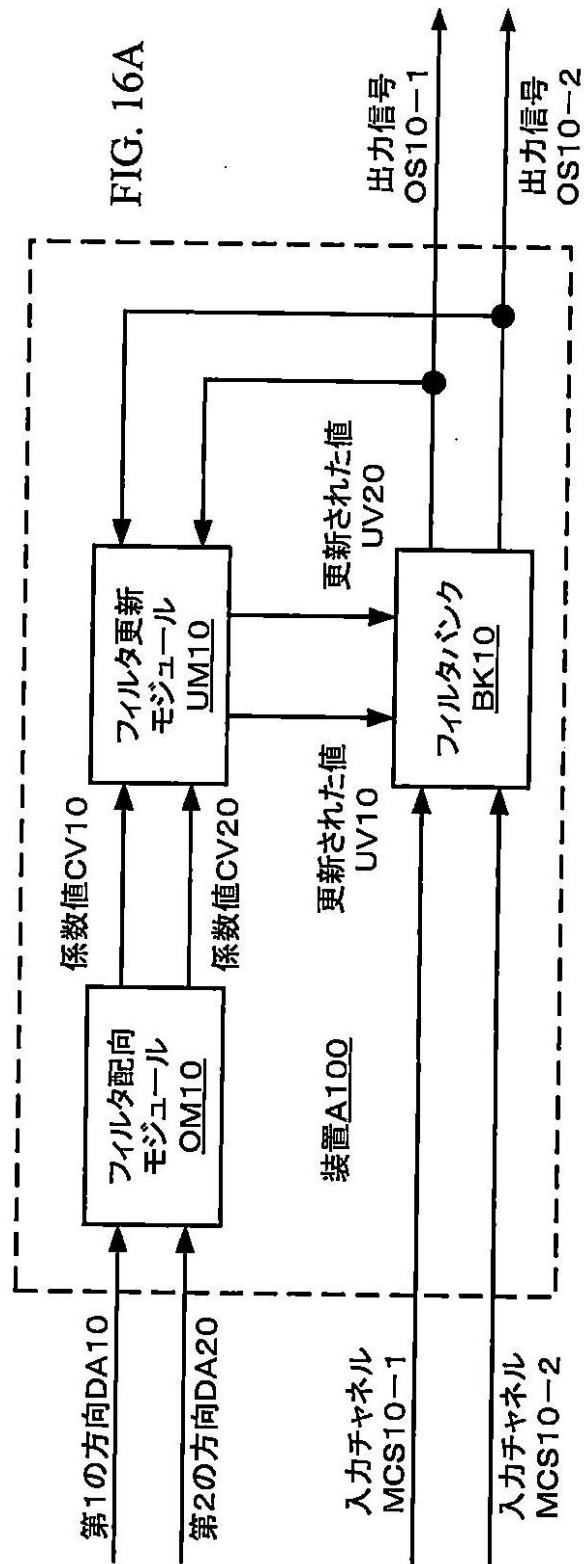


FIG. 15

【図 16 A】

図 16A



【図 16 B】

図 16B

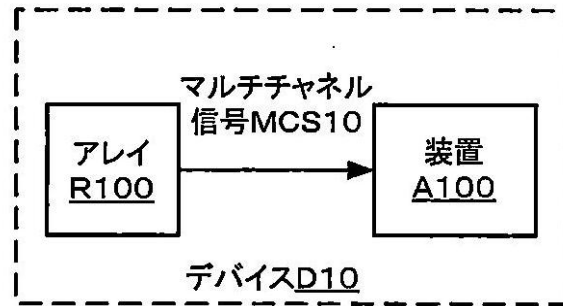


FIG. 16B

【図 16 C】

図 16C

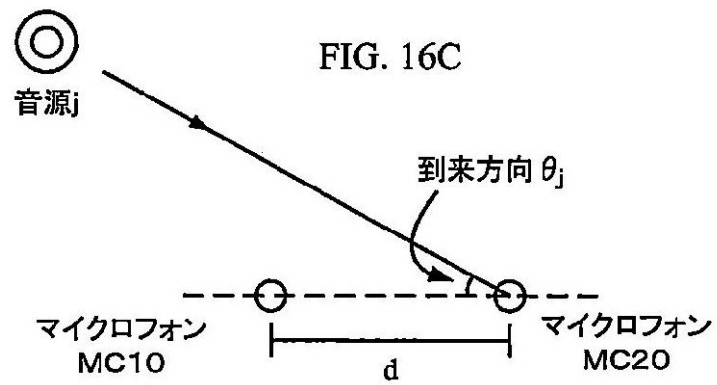


FIG. 16C

【図 17】

図 17

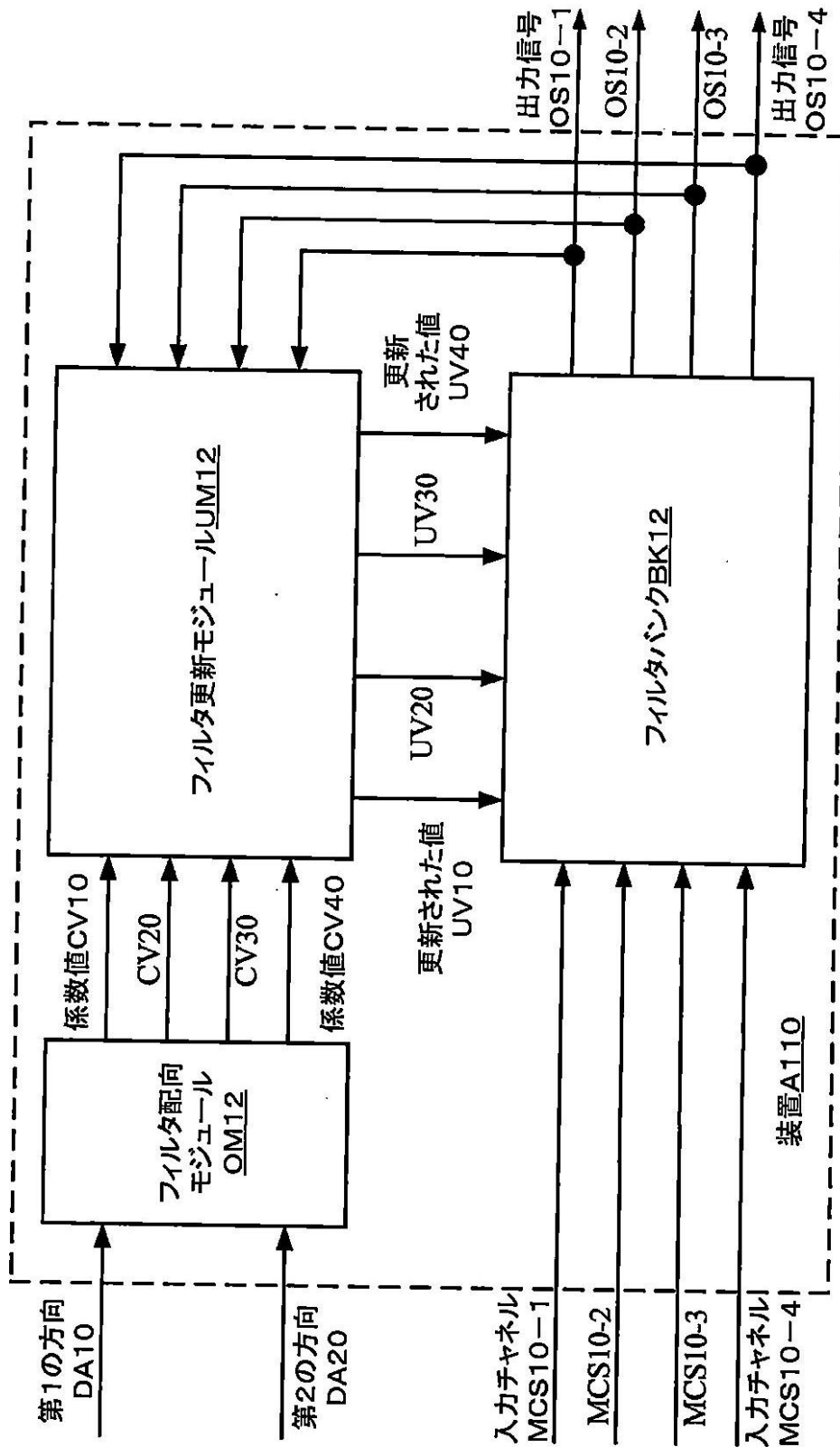


FIG. 17

【図 18 A】

図 18A

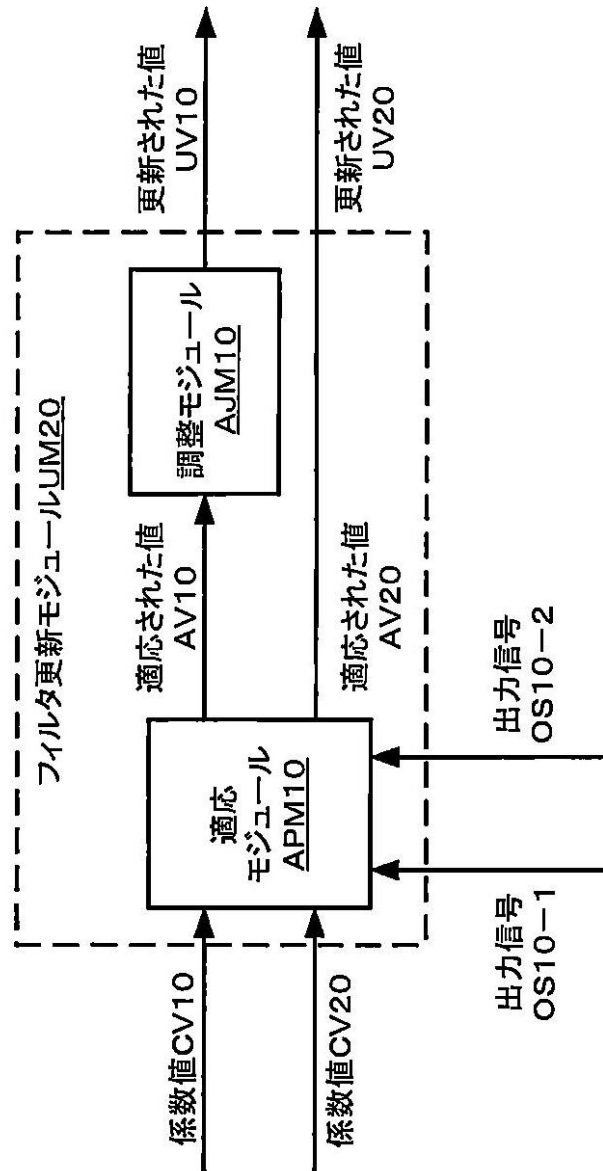


FIG. 18A

【図 18 B】

図 18B

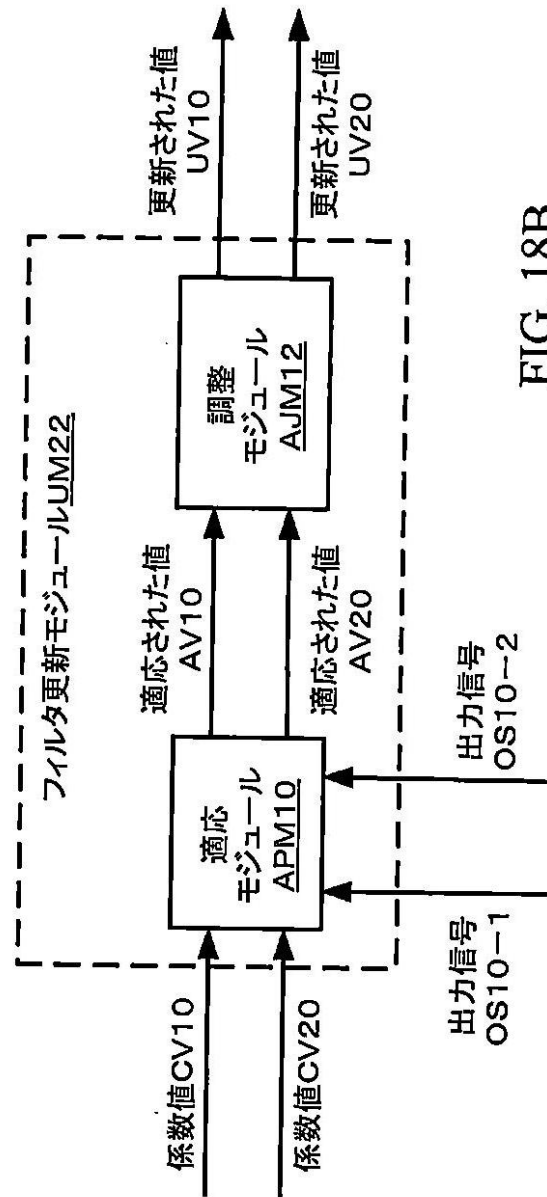


FIG. 18B

【図 19 A】

図 19A



FIG. 19A

【図 19 B】

図 19B

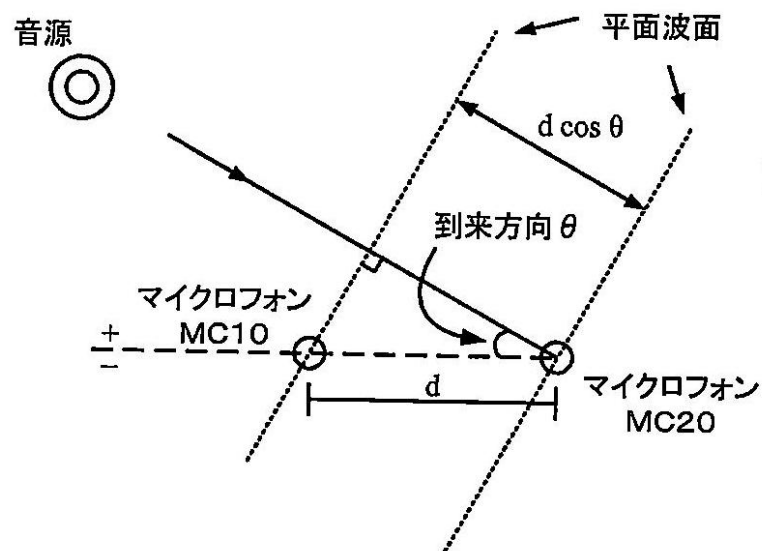


FIG. 19B

【図 20】

図 20

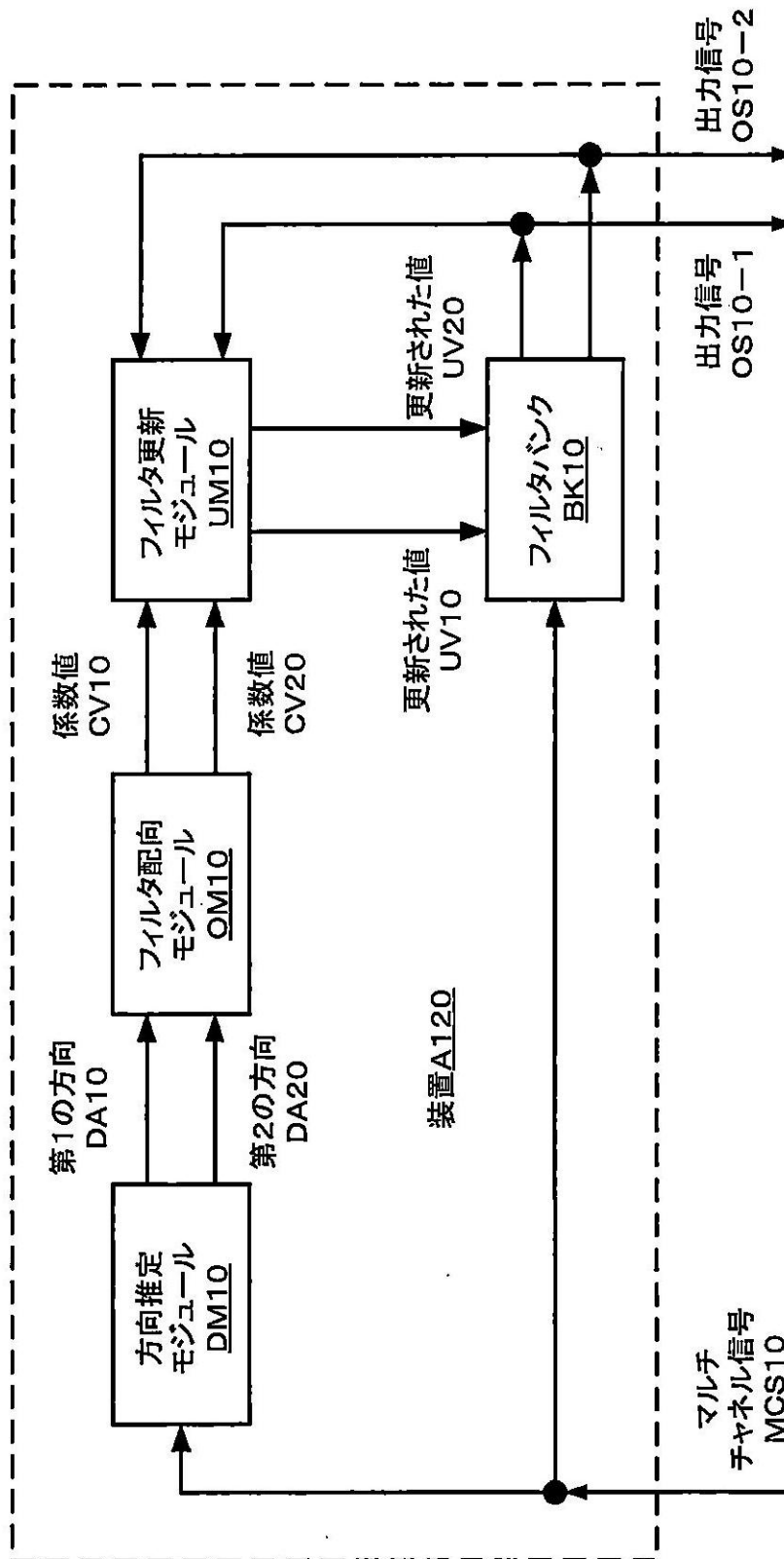


FIG. 20

【図 2 1】

図 21

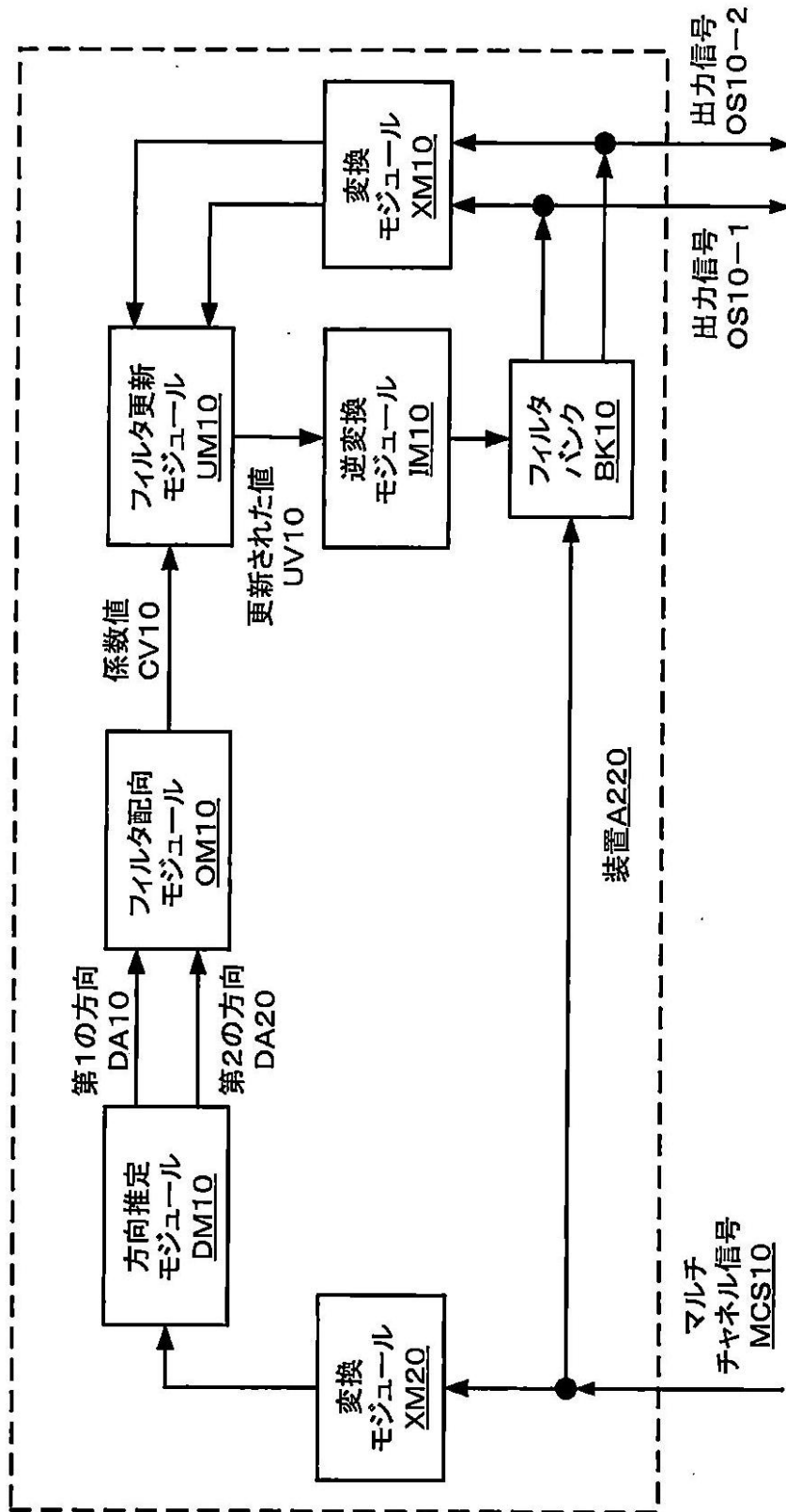


FIG. 21

【図 22】

図 22

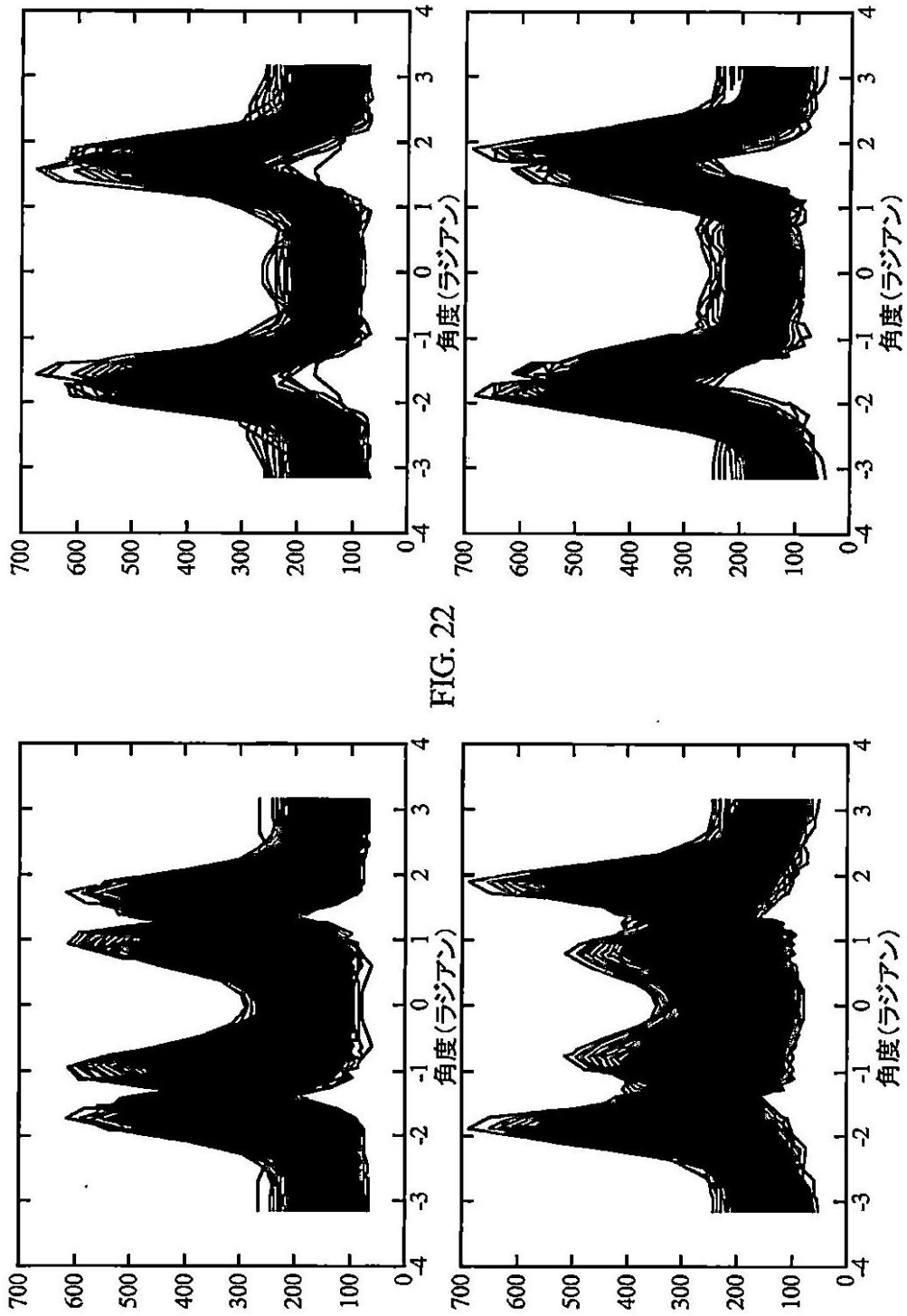
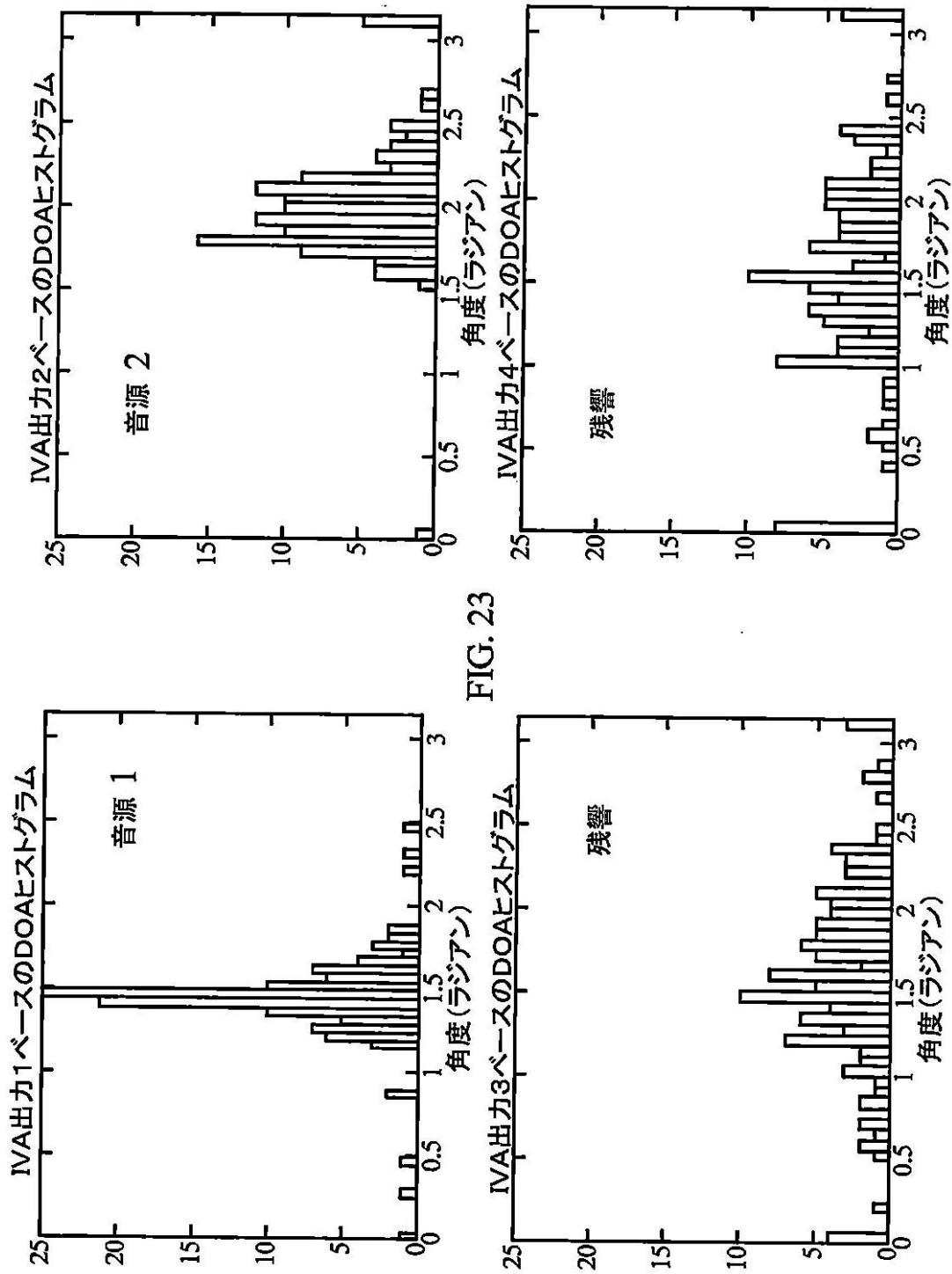


FIG. 22

【図 23】

図 23



【図24】

図24

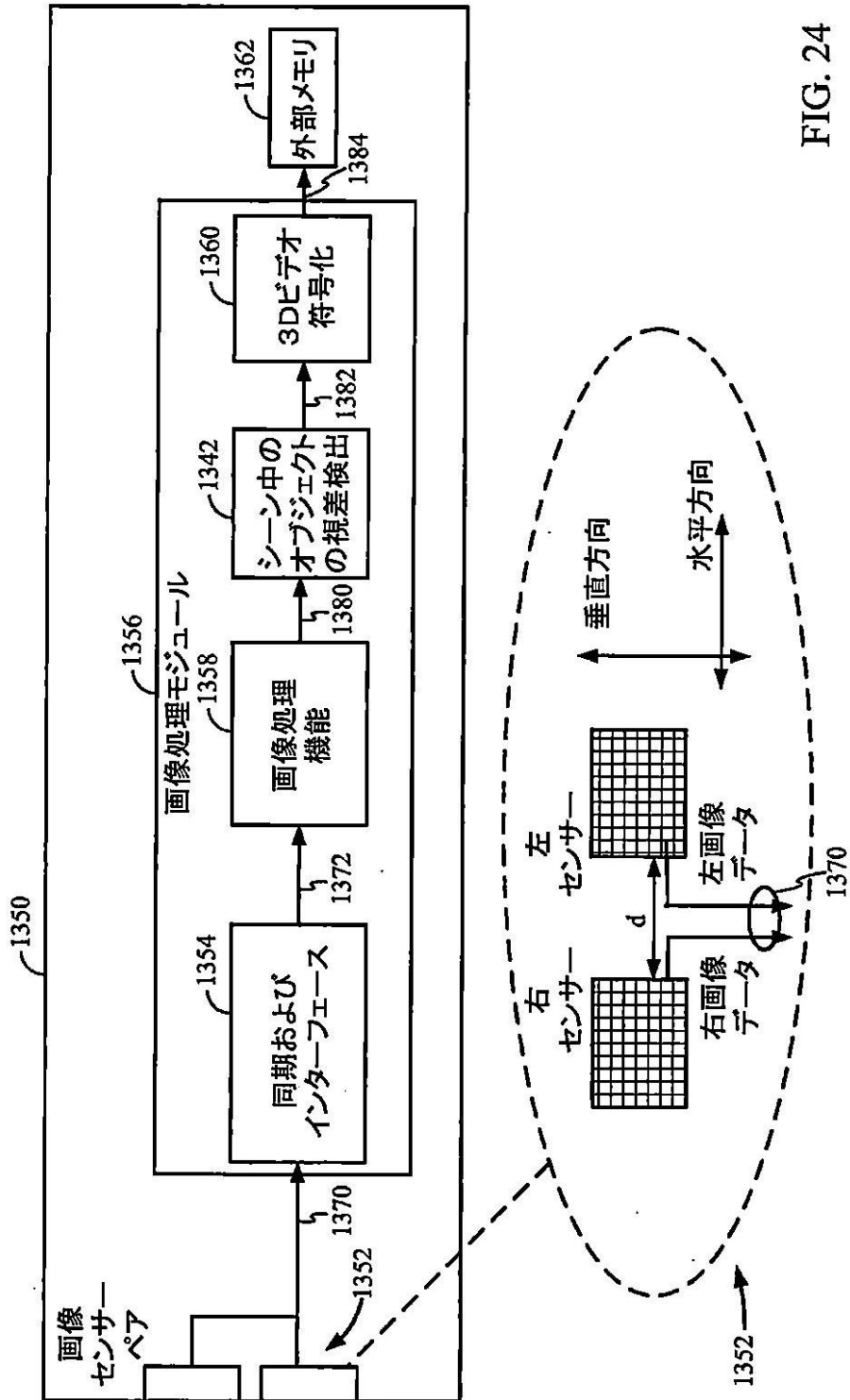


FIG. 24

【図 25】

図 25

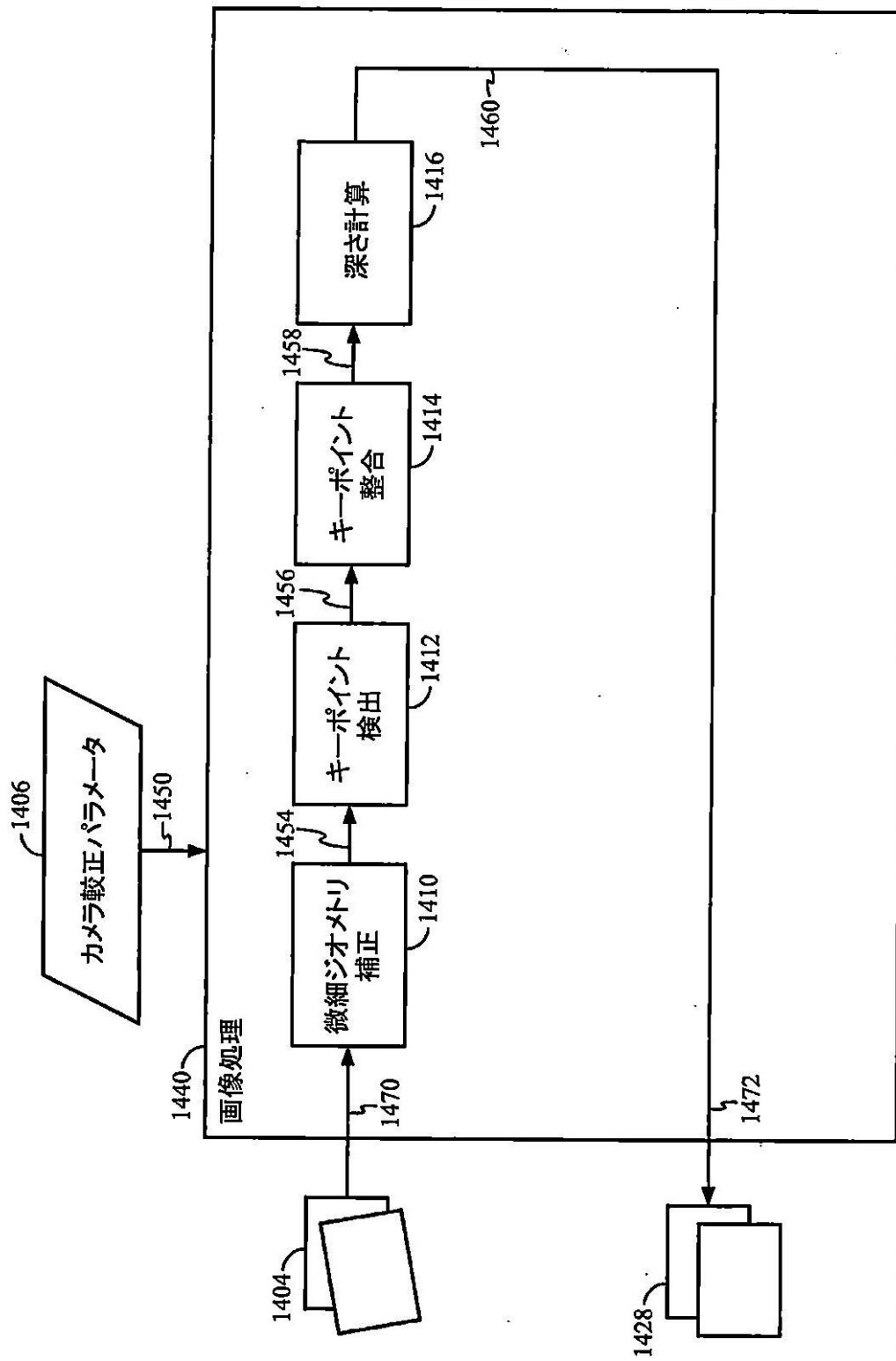


FIG. 25

【 図 2 6 A 】

図 26A

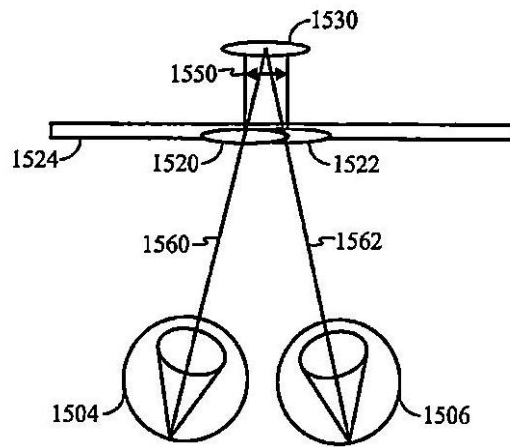


FIG. 26A

【 図 2 6 B 】

図 26B

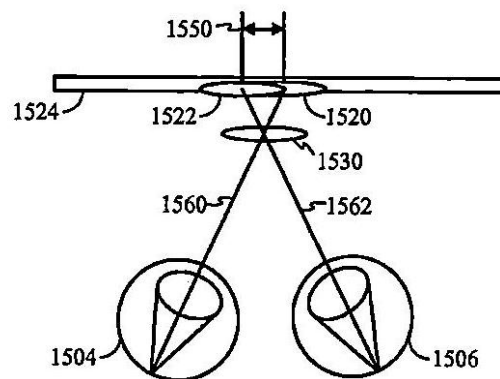


FIG. 26B

【図 27 A】

図 27A

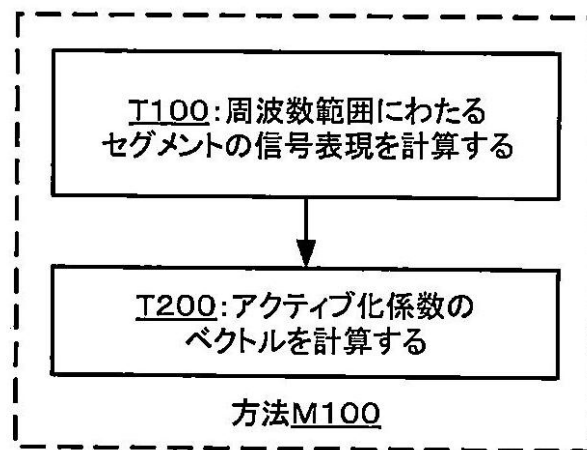


FIG. 27A

【図 27 B】

図 27B

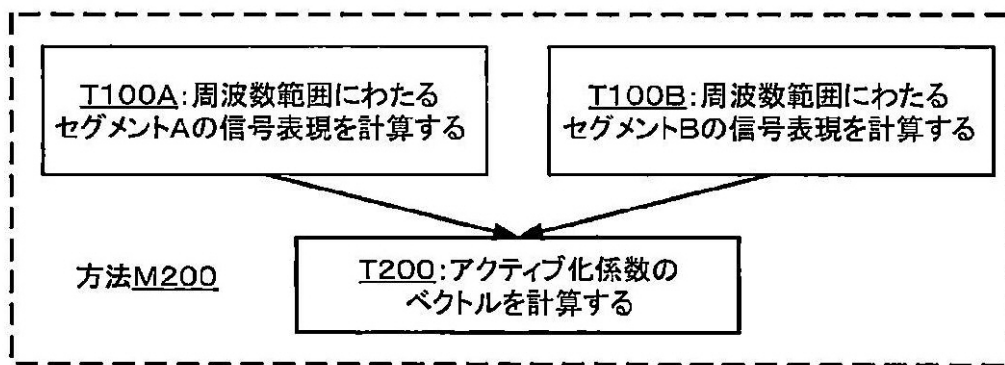


FIG. 27B

【図 27 C】

図 27C

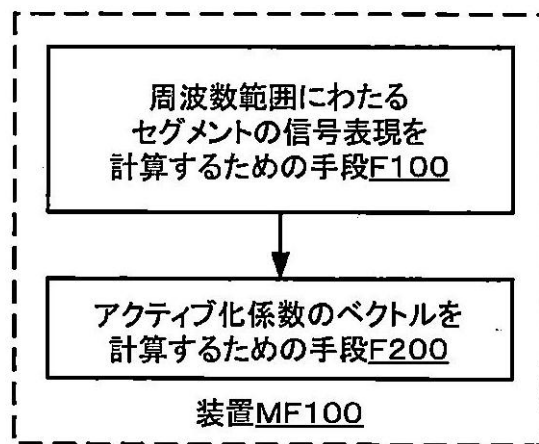


FIG. 27C

【図 27 D】

図 27D

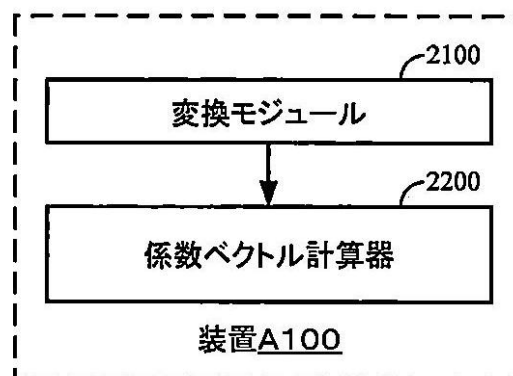


FIG. 27D

【図 28 A】

図 28A

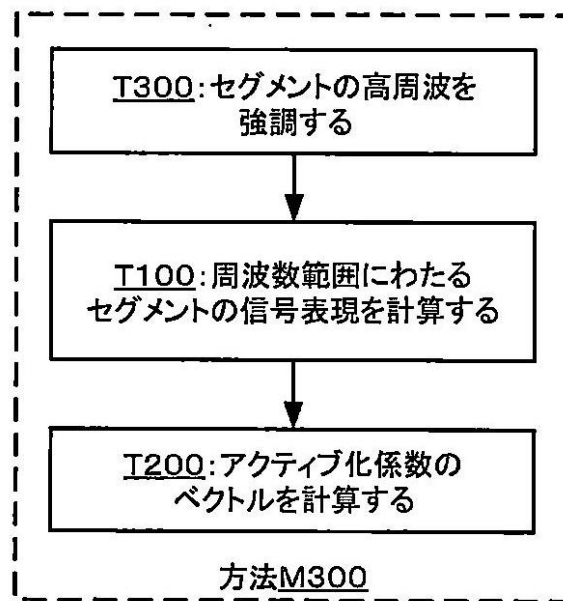
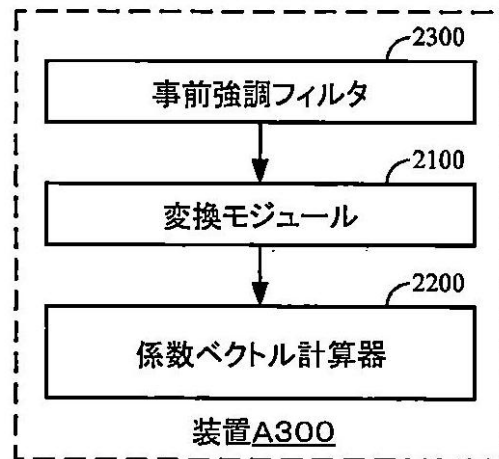


FIG. 28A

【図 28 B】

図 28B

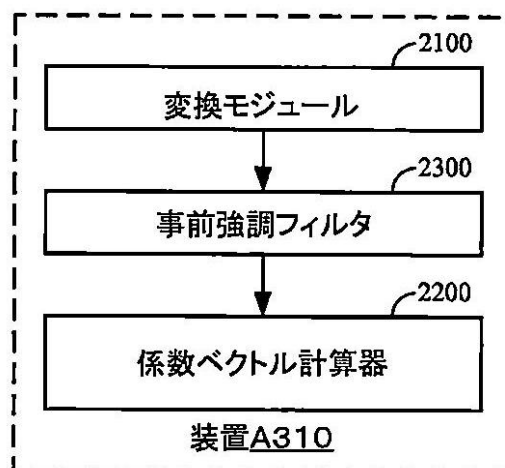
FIG. 28B



【図 28 C】

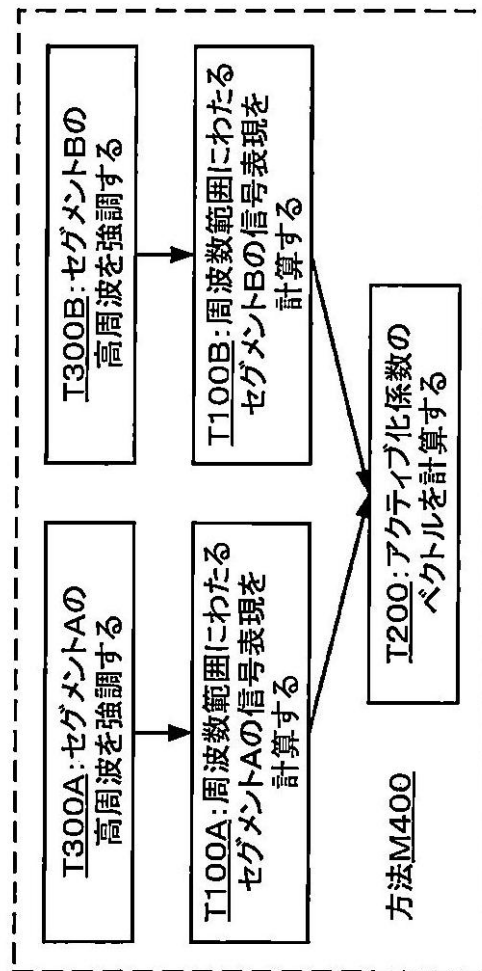
図 28C

FIG. 28C



【図 29 A】

図 29A



【図 29B】

図 29B

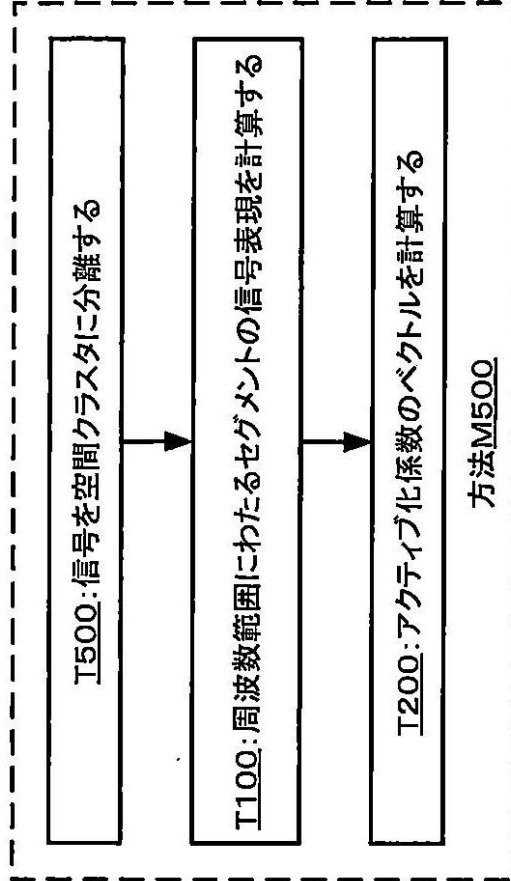


FIG. 29B

【図 30 A】

図 30A

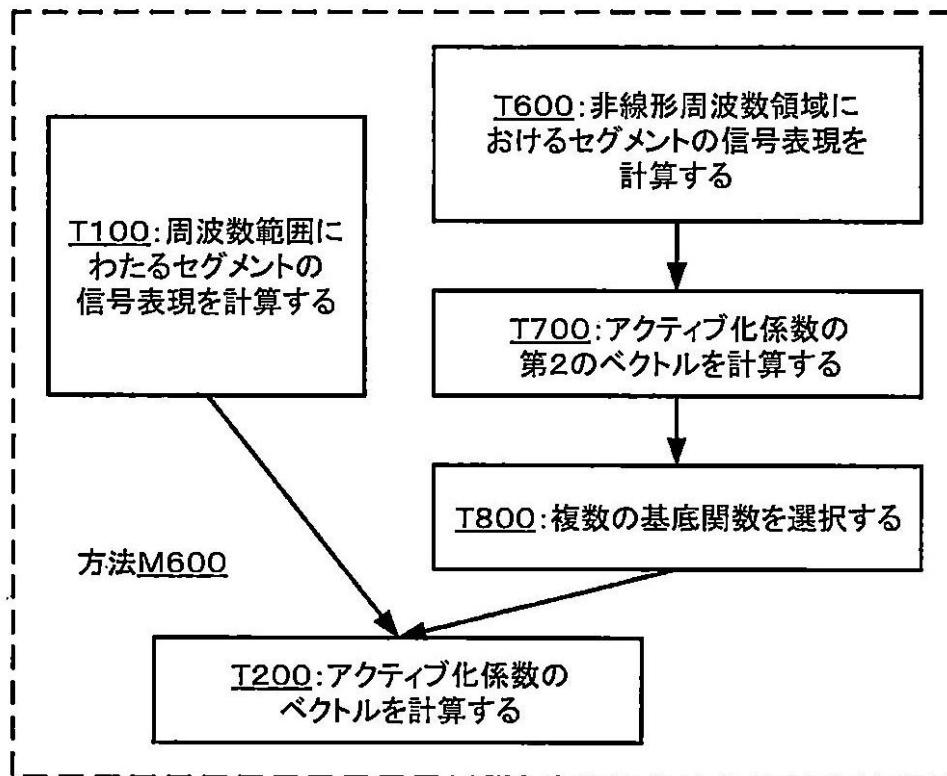


FIG. 30A

【図 30 B】

図 30B

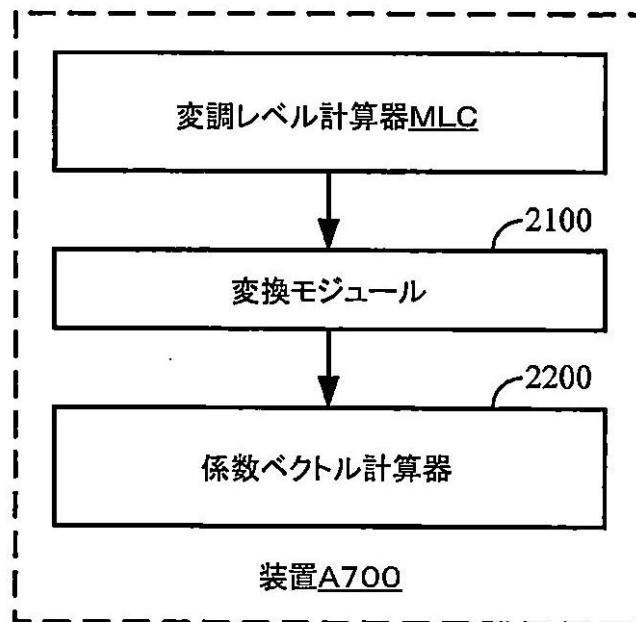


FIG. 30B

【図 3 1】

図 31

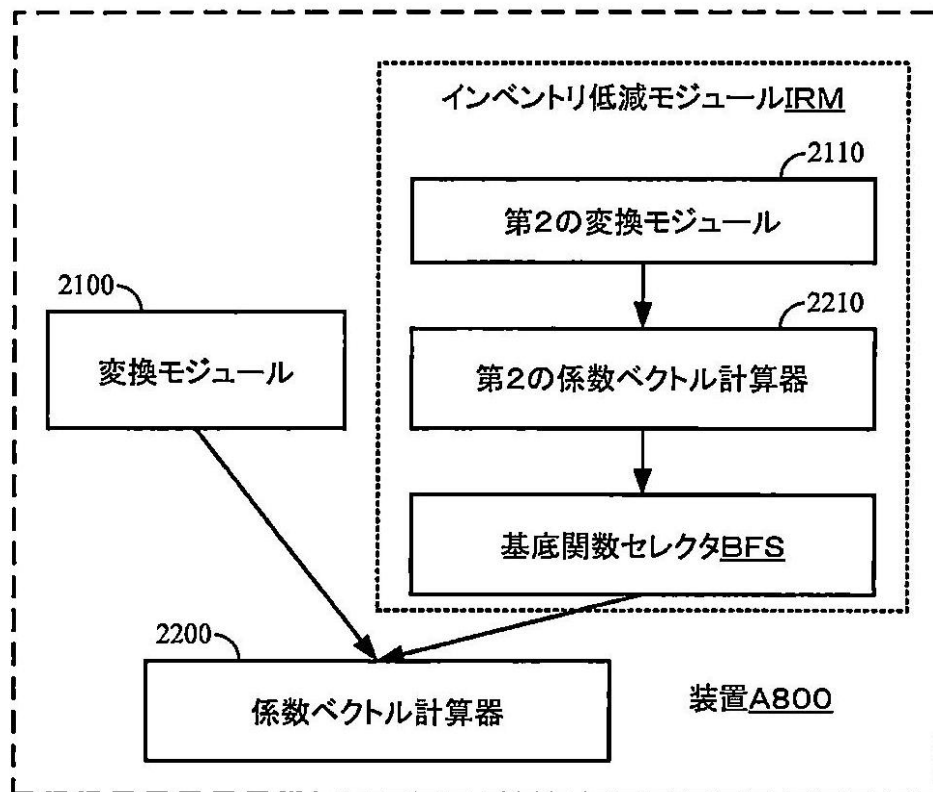


FIG. 31

【図 3 2】

図 32

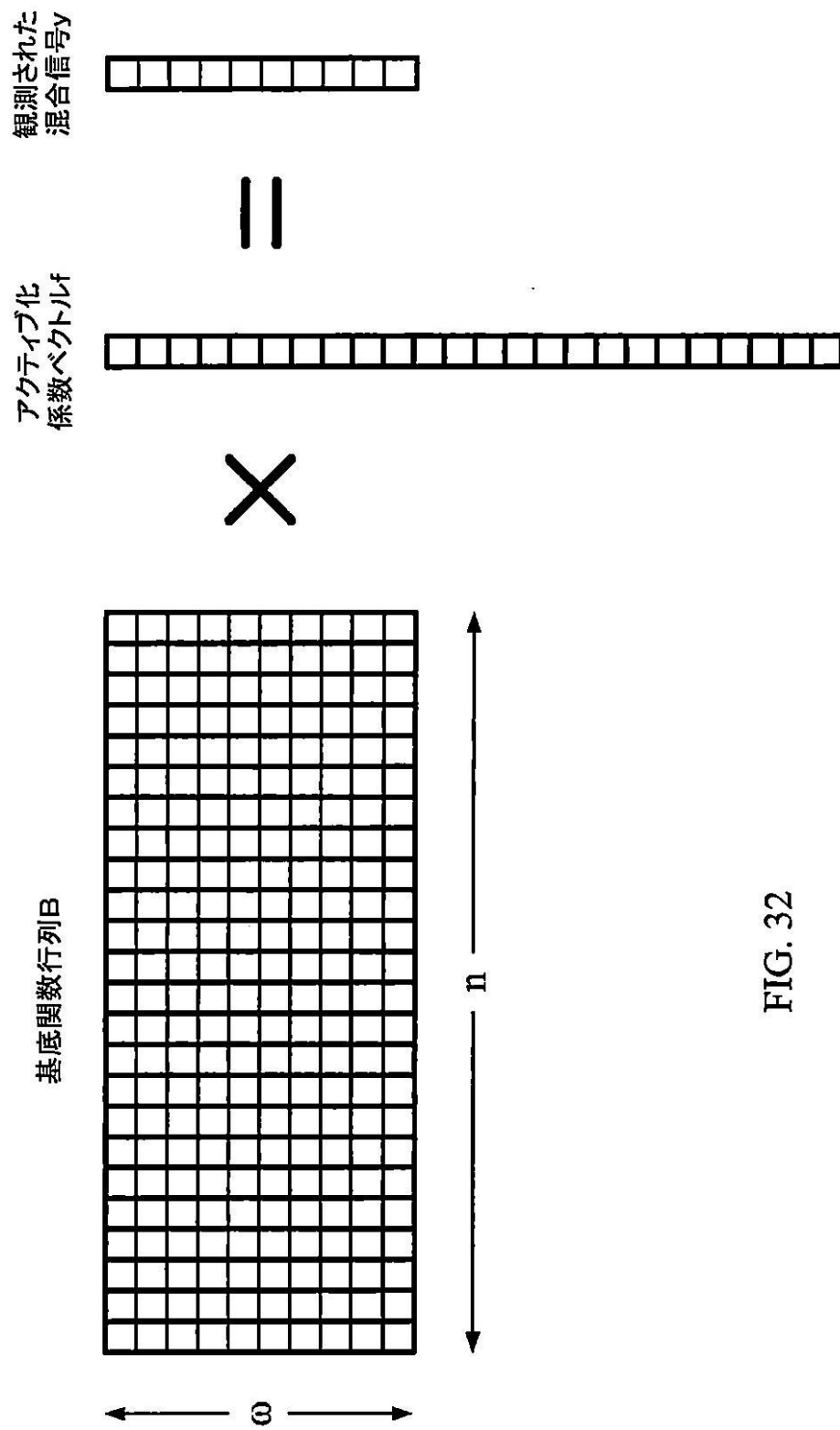


FIG. 32

【図 33】

図 33

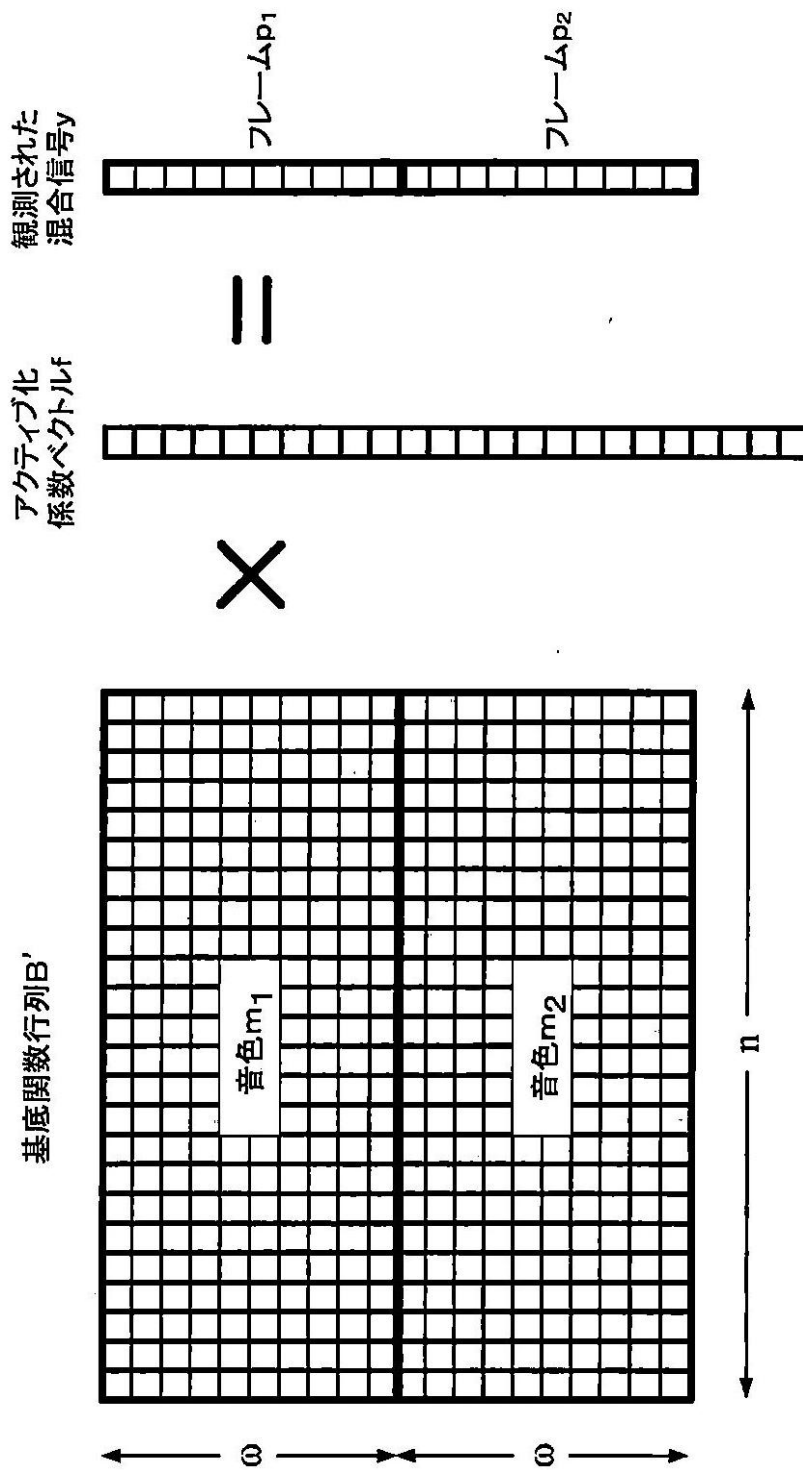


FIG. 33

【 図 3 4 】

図 34

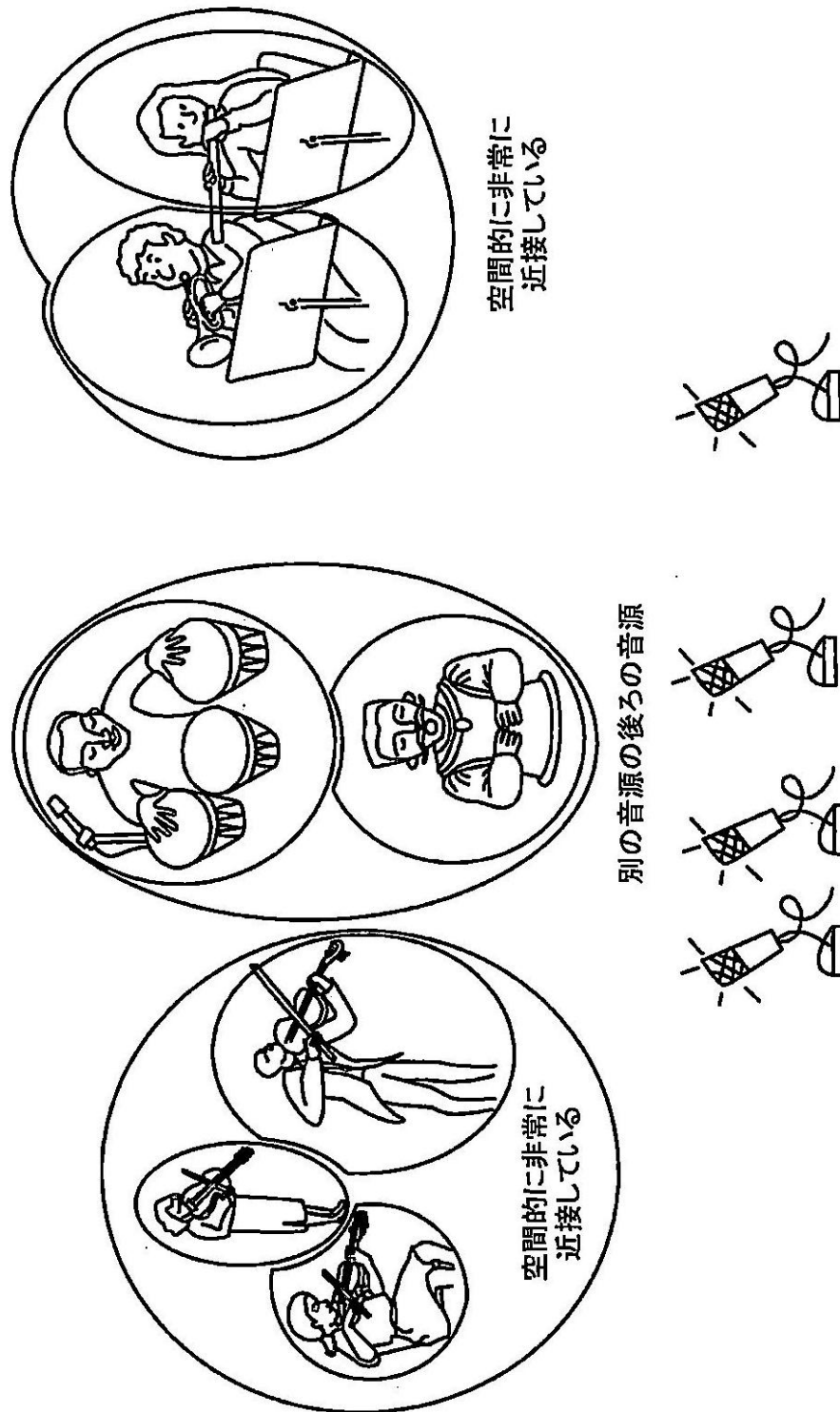


FIG. 34

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No

PCT/US2013/029558

A. CLASSIFICATION OF SUBJECT MATTER

INV. G06K9/62 G06K9/00 G06T7/20 H04R3/00
ADD.

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

G06K G06T H04R

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

EPO-Internal, WPI Data

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	STAFFAN EKVALL ET AL: "Integrating Active Mobile Robot Object Recognition and SLAM in Natural Environments", INTELLIGENT ROBOTS AND SYSTEMS, 2006 IEEE/RSJ INTERNATIONAL CONFERENCE ON, IEEE, PI, 1 October 2006 (2006-10-01), pages 5792-5797, XP031006097, ISBN: 978-1-4244-0258-8 sect III, III.B, IV, V.B ----- -/--	1-40

☒ Further documents are listed in the continuation of Box C.☒ See patent family annex.

* Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&" document member of the same patent family

Date of the actual completion of the international search

24 June 2013

Date of mailing of the international search report

01/07/2013

Name and mailing address of the ISA/

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040,
Fax: (+31-70) 340-3016

Authorized officer

Darolti, Cristina

INTERNATIONAL SEARCH REPORT

International application No

PCT/US2013/029558

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	VERMAAK J ET AL: "Sequential Monte Carlo fusion of sound and vision for speaker tracking", PROCEEDINGS OF THE EIGHT IEEE INTERNATIONAL CONFERENCE ON COMPUTER VISION. (ICCV). VANCOUVER, BRITISH COLUMBIA, CANADA, JULY 7 - 14, 2001; [INTERNATIONAL CONFERENCE ON COMPUTER VISION], LOS ALAMITOS, CA : IEEE COMP. SOC, US, vol. 1, 7 July 2001 (2001-07-07), pages 741-746, XP010554056, DOI: 10.1109/ICCV.2001.937600 ISBN: 978-0-7695-1143-6 abstract -----	1-40
A	STROBEL ET AL: "Joint Audio-Video Object Localization and Tracking", IEEE SIGNAL PROCESSING MAGAZINE, IEEE SERVICE CENTER, PISCATAWAY, NJ, US, 1 January 2001 (2001-01-01), pages 22-31, XP002358730, ISSN: 1053-5888, DOI: 10.1109/79.911196 sect. Object Localization; page 24 -----	1-40
A	LO D ET AL: "Multimodal talker localization in video conferencing environments", HAPTIC, AUDIO AND VISUAL ENVIRONMENTS AND THEIR APPLICATIONS, 2004. HA VE 2004. PROCEEDINGS. THE 3RD IEEE INTERNATIONAL WORKSHOP ON OTTAWA, ONT., CANADA 2-3 OCT. 2004, PISCATAWAY, NJ, USA, IEEE, US, 2 October 2004 (2004-10-02), pages 195-200, XP010765318, DOI: 10.1109/HAVE.2004.1391905 ISBN: 978-0-7803-8817-8 sect. III; figures 1-4 -----	1-40
X	US 2003/103647 A1 (RUI YONG [US] ET AL) 5 June 2003 (2003-06-05) paragraphs [0071], [0073], [0105], [0106], [0163] -----	1-5, 11-15, 21-25, 31-35
X	EP 1 643 769 A1 (SAMSUNG ELECTRONICS CO LTD [KR]) 5 April 2006 (2006-04-05) paragraph [0054] - paragraph [0057]; figures 4A-C ----- -/--	1-5, 11-15, 21-25, 31-35

INTERNATIONAL SEARCH REPORT

International application No

PCT/US2013/029558

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>SOCIAL ROBOT ED - ANONYMOUS: "Attending to Learn and Learning to Attend for a", HUMANOID ROBOTS, 2006 6TH IEEE-RAS INTERNATIONAL CONFERENCE ON, IEEE, PI, 1 December 2006 (2006-12-01), pages 618-623, XP031053086, ISBN: 978-1-4244-0199-4</p> <p>page 68; figure 4.1</p> <p>& Lijin Aryananda: "A Few Days of A Robot's Life in the Human's World: Toward Incremental Individual Recognition", PhD Thesis, 3 April 2007 (2007-04-03), XP055060888, Retrieved from the Internet: URL:http://dspace.mit.edu/handle/1721.1/37144 [retrieved on 2013-04-24] sect.II; figure 1</p>	1-40
X	<p>-----</p> <p>BENJAMIN FRANSEN ET AL: "Using vision, acoustics, and natural language for disambiguation", HUMAN-ROBOT INTERACTION (HRI), 2007 2ND ACM/IEEE INTERNATIONAL CONFERENCE ON, IEEE, 9 March 2007 (2007-03-09), pages 73-80, XP032211851, ISBN: 978-1-59593-617-2</p> <p>sect.1, 4.9, 4.10, 5.2, 5.3, 6.1, 8, 9; figures 1, 6,7</p>	1-5, 7-15, 17-25, 27-35, 37-40
X	<p>-----</p> <p>GERALD FRIEDLAND ET AL: "Visual speaker localization aided by acoustic models", PROCEEDINGS OF THE SEVENTEEN ACM INTERNATIONAL CONFERENCE ON MULTIMEDIA, MM '09, 1 January 2009 (2009-01-01), page 195, XP055060737, New York, New York, USA DOI: 10.1145/1631272.1631301 ISBN: 978-1-60-558608-3</p> <p>sect.4.3, 4.4, 5, 6, 7; figure 6</p>	1-3,7-10
A	<p>-----</p> <p>Ming-Yu Chen ET AL: "MoSIFT: Recognizing Human Actions in Surveillance Videos", CMU-CS Report, 24 September 2009 (2009-09-24), XP055067839, Retrieved from the Internet: URL:http://www.cs.cmu.edu/~mychen/publication/ChenMoSIFTCMU09.pdf [retrieved on 2013-06-21] sect.1 par.4-5, sect.2, 2.2; figure 2</p> <p style="text-align: center;">-/--</p>	1-40

INTERNATIONAL SEARCH REPORT

International application No

PCT/US2013/029558

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	<p>& Ming-Yu Chen: "Long Term Activity Analysis in Surveillance Video Archives", CMU PhD Thesis, 12 September 2010 (2010-09-12), XP055067853, Retrieved from the Internet: URL:http://www.lti.cs.cmu.edu/research/thesis/2010/mingyu_chen.pdf [retrieved on 2013-06-21] chapter 1, e.g. sect.1.2, 1.3, 1.6.3, 1.6.4</p> <p>-----</p> <p>Elie El-Khoury: "Unsupervised video indexing based on audiovisual characterization of persons", These Doctorat Uni Toulouse, 3 June 2010 (2010-06-03), XP55067931, Retrieved from the Internet: URL:http://thesesups.ups-tlse.fr/1025/1/El-Khoury_Elie.pdf [retrieved on 2013-06-24] sect.7</p> <p>-----</p>	1-3,7-10
X	<p>INOUE N ET AL: "High-Level Feature Extraction Using SIFT GMMs and Audio Models", PATTERN RECOGNITION (ICPR), 2010 20TH INTERNATIONAL CONFERENCE ON, IEEE, PISCATAWAY, NJ, USA, 23 August 2010 (2010-08-23), pages 3220-3223, XP031772106, ISBN: 978-1-4244-7542-1 sect.I, II.A, III.A</p> <p>-----</p>	1-3,7-10
X	<p>TALANTZIS F ET AL: "Audioâ Visual Active Speaker Tracking in Cluttered Indoors Environments", IEEE TRANSACTIONS ON SYSTEMS, MAN AND CYBERNETICS. PART B:CYBERNETICS, IEEE SERVICE CENTER, PISCATAWAY, NJ, US, vol. 38, no. 3, 1 June 2008 (2008-06-01), pages 799-807, XP011344943, ISSN: 1083-4419, DOI: 10.1109/TSMCB.2008.922063 sect.II, II.C, III.B, IV</p> <p>-----</p>	1-10

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US2013/029558**Box No. II Observations where certain claims were found unsearchable (Continuation of item 2 of first sheet)**

This international search report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☐ Claims Nos.:
because they relate to subject matter not required to be searched by this Authority, namely:

2. ☐ Claims Nos.:
because they relate to parts of the international application that do not comply with the prescribed requirements to such an extent that no meaningful international search can be carried out, specifically:

3. ☐ Claims Nos.:
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

Box No. III Observations where unity of invention is lacking (Continuation of item 3 of first sheet)

This International Searching Authority found multiple inventions in this international application, as follows:

see additional sheet

1. ☒ As all required additional search fees were timely paid by the applicant, this international search report covers all searchable claims.

2. ☐ As all searchable claims could be searched without effort justifying an additional fees, this Authority did not invite payment of additional fees.

3. ☐ As only some of the required additional search fees were timely paid by the applicant, this international search report covers only those claims for which fees were paid, specifically claims Nos.:

4. ☐ No required additional search fees were timely paid by the applicant. Consequently, this international search report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

Remark on Protest

- ☐ The additional search fees were accompanied by the applicant's protest and, where applicable, the payment of a protest fee.
- ☐ The additional search fees were accompanied by the applicant's protest but the applicable protest fee was not paid within the time limit specified in the invitation.
- ☒ No protest accompanied the payment of additional search fees.

International Application No. PCT/ US2013/ 029558

FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210

This International Searching Authority found multiple (groups of) inventions in this international application, as follows:

1. claims: 1-5, 11-15, 21-25, 31-35

Claims 3-5 referring to the features of determining a DOA using an array of microphones to select a portion of the scene based on audio.

2. claims: 6, 16, 26, 36

Claim 6 referring to the feature of computing local motion vectors in a video and using them together with the keypoints for recognition

3. claims: 7-10, 17-20, 27-30, 37-40

1) Claims 7-8 referring to the features of computing MFCC acoustic recognition features using them together with the keypoints for recognition.

2) Claims 9-10 referring to the features of determining range information (by multicamera disparity) and analysing the keypoints based on range.

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2013/029558

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 2003103647	A1	05-06-2003	AT 397354 T 15-06-2008
			AT 551676 T 15-04-2012
			CN 1423487 A 11-06-2003
			CN 101093541 A 26-12-2007
			EP 1330128 A2 23-07-2003
			EP 1838104 A2 26-09-2007
			EP 1942679 A2 09-07-2008
			EP 1944975 A2 16-07-2008
			JP 4142420 B2 03-09-2008
			JP 4536789 B2 01-09-2010
			JP 4607984 B2 05-01-2011
			JP 4642093 B2 02-03-2011
			JP 2003216951 A 31-07-2003
			JP 2008204479 A 04-09-2008
			JP 2008243214 A 09-10-2008
			JP 2008243215 A 09-10-2008
			KR 20030045624 A 11-06-2003
			TW 1222031 B 11-10-2004
			US 2003103647 A1 05-06-2003
			US 2005129278 A1 16-06-2005
			US 2005147278 A1 07-07-2005
			US 2005188013 A1 25-08-2005
			US 2005210103 A1 22-09-2005

EP 1643769	A1	05-04-2006	NONE

フロントページの続き

(51)Int.Cl.		F I		テーマコード(参考)
G 1 0 L 15/28 (2013.01)		G 1 0 L 15/28	4 0 0	
G 1 0 L 15/24 (2013.01)		G 1 0 L 15/24	Z	

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC

(74)代理人 100153051
弁理士 河野 直樹

(74)代理人 100140176
弁理士 砂川 克

(74)代理人 100158805
弁理士 井関 守三

(74)代理人 100179062
弁理士 井上 正

(74)代理人 100124394
弁理士 佐藤 立志

(74)代理人 100112807
弁理士 岡田 貴志

(74)代理人 100111073
弁理士 堀内 美保子

(72)発明者 ビッサー、エリック
アメリカ合衆国、カリフォルニア州 9 2 1 2 1、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 ワン、ヒイン
アメリカ合衆国、カリフォルニア州 9 2 1 2 1、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 シディクイ、ハシブ・エー .
アメリカ合衆国、カリフォルニア州 9 2 1 2 1、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 キム、レ - ホン
アメリカ合衆国、カリフォルニア州 9 2 1 2 1、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

F ターム(参考) 5L096 AA09 CA05 FA35 FA67 HA03 HA04 HA07 JA16