

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6654195号
(P6654195)

(45) 発行日 令和2年2月26日 (2020.2.26)

(24) 登録日 令和2年1月31日 (2020.1.31)

(51) Int. Cl.	F I
G 1 0 L 21/028 (2013.01)	G 1 0 L 21/028 C
G 1 0 L 21/0308 (2013.01)	G 1 0 L 21/0308 A
G 1 0 L 19/008 (2013.01)	G 1 0 L 19/008 1 0 0
	G 1 0 L 19/008 2 0 0

請求項の数 11 (全 23 頁)

(21) 出願番号	特願2017-533302 (P2017-533302)	(73) 特許権者	507236292
(86) (22) 出願日	平成27年12月18日 (2015.12.18)		ドルビー ラボラトリーズ ライセンシン
(65) 公表番号	特表2018-503864 (P2018-503864A)		グ コーポレイション
(43) 公表日	平成30年2月8日 (2018.2.8)		アメリカ合衆国 9 4 1 0 3 カリフォル
(86) 国際出願番号	PCT/US2015/066798		ニア州 サンフランシスコ マーケット
(87) 国際公開番号	W02016/106145		ストリート 1 2 7 5
(87) 国際公開日	平成28年6月30日 (2016.6.30)	(74) 代理人	100107766
審査請求日	平成30年11月30日 (2018.11.30)		弁理士 伊東 忠重
(31) 優先権主張番号	201410814973.9	(74) 代理人	100070150
(32) 優先日	平成26年12月22日 (2014.12.22)		弁理士 伊東 忠彦
(33) 優先権主張国・地域又は機関	中国 (CN)	(74) 代理人	100091214
(31) 優先権主張番号	62/108,254		弁理士 大貫 進介
(32) 優先日	平成27年1月27日 (2015.1.27)		
(33) 優先権主張国・地域又は機関	米国 (US)		

最終頁に続く

(54) 【発明の名称】 オーディオ・コンテンツからの投影ベースのオーディオ・オブジェクト抽出

(57) 【特許請求の範囲】

【請求項 1】

複数のチャンネルのオーディオ信号によって表現されているオーディオ・コンテンツからのオーディオ・オブジェクト抽出のための方法であって：

前記複数のチャンネルの相関に基づいて前記複数のチャンネルをクラスタリングし、それによりチャンネル・グループを得る段階と；

前記チャンネル・グループから第一のチャンネルおよび第二のチャンネルを選択する段階と；

前記第一のチャンネルについての第一の部分集合および前記第二のチャンネルについての第二の部分集合を含む投影空間の第一の集合を同定する段階と；

前記第一および第二のチャンネルの間の相関の第一の集合を決定する段階であって、前記第一の集合の相関のそれぞれは前記第一の部分集合の投影空間の一つおよび前記第二の部分集合の投影空間の一つに対応する、段階と；

前記第一の集合の相関のうちで最大の値をもつ前記第一の集合の相関のうちの第一の相関を同定する段階と；

少なくとも部分的には前記第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づいて、前記第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出する段階とを含む、方法。

【請求項 2】

クラスタリングにおいて、前記複数のチャンネルの一对のチャンネルの間の相関が：

10

20

投影空間の第二の集合であって、その対のチャンネルの一方についての第三の部分集合およびその対のチャンネルの他方についての第四の部分集合を含むものを同定する段階と；

その対のチャンネルの間の相関の第二の集合を決定する段階であって、前記第二の集合の相関のそれぞれは前記第三の部分集合の投影空間の一つおよび前記第四の部分集合の投影空間の一つに対応する、段階と；

前記第二の集合の相関のうちの一つを、その対のチャンネルの間の相関として選択する段階であって、選択される相関は第二のあらかじめ定義された閾値より大きい、段階とを実行することによって決定される、
請求項 1 記載の方法。

【請求項 3】

前記チャンネル・グループから前記第一および第二のチャンネルを選択することは：

前記第二のチャンネルのオーディオ信号が前記チャンネル・グループにおいて最大のエネルギーをもつように、前記チャンネル・グループから前記第二のチャンネルを選択することを含む、

請求項 1 記載の方法。

【請求項 4】

前記チャンネル・グループから前記第一および第二のチャンネルを選択することはさらに：

前記第一および第二のチャンネルの間の相関が第三のあらかじめ定義された閾値より大きいように前記チャンネル・グループから前記第一のチャンネルを選択し；

前記第二のチャンネルのオーディオ信号から、少なくとも部分的には前記第一の相関と前記第一の相関に対応する前記第二の部分集合からの投影空間とに基づいてオーディオ・オブジェクトを抽出することを含む、

請求項 3 記載の方法。

【請求項 5】

前記オブジェクト抽出の検証および調整を：

少なくとも部分的には前記第一および第二のチャンネルのオーディオ信号からの抽出されたオーディオ・オブジェクトに基づいて、マルチチャンネル・オブジェクトを生成し；

生成されたマルチチャンネル・オブジェクトをモノ表現にダウンミックスし；

抽出されたオブジェクトのものとマルチチャンネル表現と前記モノ表現との間のダウンミックス類似性を決定し；

少なくとも部分的には前記第一の集合の相関のうちの前記第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づき、前記ダウンミックス類似性にさらに基づいて、前記第一のチャンネルのオーディオ信号から前記オーディオ・オブジェクトを抽出することによって行なうことを含む、

請求項 1 記載の方法。

【請求項 6】

前記オブジェクト抽出の検証および調整を：

少なくとも部分的には前記第一および第二のチャンネルのオーディオ信号からの抽出されたオーディオ・オブジェクトに基づいて、マルチチャンネル・オブジェクトを生成し；

生成されたマルチチャンネル・オブジェクトをモノ表現にダウンミックスし；

少なくとも部分的には推定位置に基づいて前記モノ表現をプリ・レンダリングし；

抽出されたオブジェクトのものとマルチチャンネル表現と前記プリ・レンダリングされたモノ表現との間の、エネルギー分布に関するプリ・レンダリング類似性を決定し；

少なくとも部分的には前記第一の集合の相関のうちの前記第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づき、前記プリ・レンダリング類似性にさらに基づいて、前記第一のチャンネルのオーディオ信号から前記オーディオ・オブジェクトを抽出することによって行なうことを含む、

請求項 1 記載の方法。

【請求項 7】

前記オーディオ・コンテンツは、フル帯域オーディオ信号のフレームを、周波数領域お

10

20

30

40

50

よび時間領域の少なくとも一方において分割することによって得られる一つまたは複数のオーディオ・ブロックを含む、請求項 1 ないし 6 のうちいずれか一項記載の方法。

【請求項 8】

複数のチャンネルのオーディオ信号によって表現されているオーディオ・コンテンツからのオーディオ・オブジェクト抽出のためのシステムであって：

前記複数のチャンネルの相関に基づいて前記複数のチャンネルをクラスタリングすることによって得られるチャンネル・グループから、第一のチャンネルおよび第二のチャンネルを選択するよう構成された選択ユニットと；

前記第一のチャンネルについての第一の部分集合および前記第二のチャンネルについての第二の部分集合を含む投影空間の第一の集合を同定するよう構成された同定ユニットと；

前記第一および第二のチャンネルの間の相関の第一の集合を決定するよう構成された決定ユニットであって、前記第一の集合の相関のそれぞれは前記第一の部分集合の投影空間の一つおよび前記第二の部分集合の投影空間の一つに対応する、決定ユニットと；

前記第一の集合の相関のうちで最大の値をもつ前記第一の集合の相関のうちの第一の相関を同定し、少なくとも部分的には前記第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づいて、前記第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出するよう構成された抽出ユニットとを有する、システム。

【請求項 9】

少なくとも部分的には前記第一および第二のチャンネルのオーディオ信号からの抽出されたオーディオ・オブジェクトに基づいて、マルチチャンネル・オブジェクトを生成するよう構成された生成ユニットと；

生成されたマルチチャンネル・オブジェクトをモノ表現にダウンミックスするよう構成されたダウンミックス・ユニットと；

抽出されたオブジェクトのもとのマルチチャンネル表現と前記モノ表現との間のダウンミックス類似性を決定するよう構成された類似性決定ユニットとをさらに有しており、

前記抽出ユニットは、少なくとも部分的には前記第一の集合の相関のうちの前記第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づき、前記ダウンミックス類似性にさらに基づいて前記第一のチャンネルのオーディオ信号から前記オーディオ・オブジェクトを再抽出するよう構成されている、

請求項 8 記載のシステム。

【請求項 10】

少なくとも部分的には前記第一および第二のチャンネルのオーディオ信号からの抽出されたオーディオ・オブジェクトに基づいて、マルチチャンネル・オブジェクトを生成するよう構成された生成ユニットと；

生成されたマルチチャンネル・オブジェクトをモノ表現にダウンミックスするよう構成されたダウンミックス・ユニットと；

少なくとも部分的には推定位置に基づいて前記モノ表現をプリ・レンダリングするよう構成されたプリ・レンダリング・ユニットと；

抽出されたオブジェクトのもとのマルチチャンネル表現と前記プリ・レンダリングされたモノ表現との間の、エネルギー分布に関するプリ・レンダリング類似性を決定するよう構成された類似性決定ユニットとをさらに有しており、

前記抽出ユニットは、少なくとも部分的には前記第一の集合の相関のうちの前記第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づき、前記プリ・レンダリング類似性にさらに基づいて前記第一のチャンネルのオーディオ信号から前記オーディオ・オブジェクトを再抽出するよう構成されている、請求項 8 記載のシステム。

【請求項 11】

コンピュータに請求項 1 ないし 7 のうちいずれか一項記載の方法を実行させるためのコンピュータ・プログラム。

10

20

30

40

50

【発明の詳細な説明】

【技術分野】

【0001】

関連出願への相互参照

本願は2014年12月22日に出願された中国優先権出願第201410814973.9号および2015年1月27日に出願された米国仮特許出願第62/108,254号の優先権を主張するものである。これらの内容はここに参照によってその全体において組み込まれる。

【0002】

技術

本稿に開示される例示的实施形態は概括的にはオーディオ・コンテンツ処理に関し、より詳細にはオーディオ・コンテンツからのオーディオ・オブジェクト抽出のための方法およびシステムに関する。

【背景技術】

【0003】

伝統的に、オーディオ・コンテンツはチャンネル・ベースのフォーマットで作成され、記憶される。チャンネル・ベースのフォーマットでは、オーディオ・コンテンツは通例、チャンネルの媒体によって表現され、記憶され、伝達され、頒布される。本稿での用法では、用語「オーディオ・チャンネル」または「チャンネル」は、通例あらかじめ定義された物理的位置をもつオーディオ・コンテンツをいう。たとえば、ステレオ、サラウンド5.1、サラウンド7.1などはみな、オーディオ・コンテンツのためのチャンネル・ベースのフォーマットである。各チャンネルは固定位置の物理的スピーカーに対応する。マルチチャンネル・コンテンツが再生されるとき、複数のスピーカーがライブで、没入的な音場を聴取者のまわりに作り出す。近年、いくつかの通常のマルチチャンネル・システムは、チャンネルおよびオーディオ・オブジェクトの両方を含む新しいフォーマットをサポートするよう拡張された。本稿での用法では、用語「オーディオ・オブジェクト」または「オブジェクト」は、定義された継続時間にわたって音場において存在する個別のオーディオ要素をいう。たとえば、オーディオ・オブジェクトはダイアログ、発砲音、雷鳴などを表わしてもよい。これらのオブジェクトは通例、その所望されるサウンド効果を作り出すためにミキサーによって使われる。たとえば、ダイアログは通例、中央前方に定位され、雷の音は通例頭上から発する。人間によるオブジェクトの位置の知覚は、同じオブジェクトのオーディオ信号を再生している複数のスピーカーの発射の結果である。たとえば、オブジェクトが前方左スピーカーおよび前方右スピーカーによって同様のエネルギー・レベルで再生されるときは、人は中央前方からのファントムを知覚する。

【0004】

上述したように、コンテンツがチャンネル・ベースのフォーマットで作成されるとき、それは通例、ミキサーによって特定の再生セッティングのために知覚経験が最適化されることを意味する。しかしながら、異なる再生セッティングによって再生されるときは、その聴取経験は、再生セッティング間の不一致のため劣化することがある。劣化の例は、オブジェクトの位置が変わることがあるということである。このように、チャンネル・ベースのフォーマットは、多様なスピーカー再生構成に適応するには非効率的である。非効率性のもう一つの側面は、バイノーラル・レンダリングにある。ここでは、チャンネル・ベースのフォーマットはスピーカー位置に固有の限られた数の頭部伝達関数(HRTF)を使うことができるだけであり、他の位置については、HRTFの補間が使われ、バイノーラル聴取経験を劣化させる。

【0005】

この問題に対処する一つの潜在的な方法は、チャンネル・ベースの表現から、位置とモノのクリーンな波形とを含むもとの源(またはオブジェクト)を復元し、それらの位置をメタデータとして、スピーカー再生装置のパン・アルゴリズムを制御するために使い、それによりオンザフライでオブジェクトを再レンダリングし、もとの音像と同様の音像を作り出すことである。バイノーラル・レンダリング・セッティングについては(限られた数の

HRTFを使う代わりに)、聴取経験をさらに向上させるために最も適切な諸HRTFを選ぶために位置を使うことができる。

【0006】

しかしながら、メタデータを用いてレンダリングされるべきチャンネル・ベースの表現におけるオブジェクトは常にクリーンであるとは限らない。かかるオブジェクトは、いくつかのチャンネル内で他のオブジェクトと同時に混合されることがある。たとえば、芸術的意図を実施するために、ミキサーは二つのオブジェクトを聴取者の前方に同時に置いて、一方は中央と前方左の間に、他方は中央と前方右の間の何らかの位置に現われるようにすることがある。これは中央前方チャンネルに二つのオブジェクトを含ませることができる。源分離技法が使われない場合には、これら二つのオブジェクトは一つのオブジェクトとみなされ、これはそれらの位置推定を誤らせることになる。

10

【0007】

このように、クリーンなオブジェクトを得てその位置を推定するためには、オブジェクトをそのマルチチャンネル混合から分離してクリーンなマルチチャンネルまたはモノ表現を生成する源分離技法が必要とされる。上述した例において、単一のマルチチャンネル入力が源分離コンポーネントによって、たとえばそれぞれ一つのクリーンなオブジェクトを含んでいるだけの二つのマルチチャンネルまたはモノ出力に分割されることが所望される。

【発明の概要】

【発明が解決しようとする課題】

【0008】

20

上記および他の潜在的な問題に対処するために、本稿に開示される例示的实施形態は、オーディオ・コンテンツからオーディオ・オブジェクトを抽出するための方法およびシステムを提案する。

【課題を解決するための手段】

【0009】

ある側面では、例示的实施形態は、複数のチャンネルのオーディオ信号によって表現されているオーディオ・コンテンツからのオーディオ・オブジェクト抽出のための方法を提供する。本方法は、前記複数のチャンネルの第一のチャンネルについての第一の部分集合および第二のチャンネルについての第二の部分集合を含む投影空間の第一の集合を同定することを含む。本方法はさらに、前記第一および第二のチャンネルの間の相関の第一の集合を決定することを含み、前記第一の集合の相関のそれぞれは前記第一の部分集合の投影空間の一つおよび前記第二の部分集合の投影空間の一つに対応する。本方法はまた、少なくとも部分的には相関の前記第一の集合のうちの第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づいて、前記第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出することを含み、前記第一の相関は第一のあらかじめ定義された閾値より大きい。これに関する諸実施形態はさらに、対応するコンピュータ・プログラム・プロダクトを含む。

30

【0010】

別の側面では、例示的实施形態は、複数のチャンネルのオーディオ信号によって表現されているオーディオ・コンテンツからのオーディオ・オブジェクト抽出のためのシステムを提供する。本システムは、前記複数のチャンネルの第一のチャンネルについての第一の部分集合および第二のチャンネルについての第二の部分集合を含む投影空間の第一の集合を同定するよう構成された同定ユニットを含む。本システムはさらに、前記第一および第二のチャンネルの間の相関の第一の集合を決定するよう構成された決定ユニットを含み、前記第一の集合の相関のそれぞれは前記第一の部分集合の投影空間の一つおよび前記第二の部分集合の投影空間の一つに対応する。本システムはまた、少なくとも部分的には相関の前記第一の集合のうちの第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づいて、前記第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出するよう構成された抽出ユニットも含み、前記第一の相関は第一のあらかじめ定義された閾値より大きい。

40

50

【 0 0 1 1 】

以下の記述を通じて、本稿に開示される例示的实施形態によれば、オーディオ・オブジェクトが複数チャンネルに基づくオーディオ・コンテンツのオーディオ信号のそれぞれから分離されることが理解されるであろう。このようにして、オーディオ・コンテンツ入力、その聴取経験を劣化させることなく、多様な再生構成に適応することが可能となる。例示的实施形態によって達成される他の利点は以下の記述を通じて明白となるであろう。

【 図面の簡単な説明 】

【 0 0 1 2 】

付属の図面を参照しての以下の詳細な説明を通じて、例示的实施形態の上記および他の目的、特徴および利点がよりわかりやすくなるであろう。図面においては、いくつかの例示的实施形態が限定ではなく例として示される。

【 図 1 】 複数のチャンネルに基づくフォーマットのオーディオ信号のセグメントの例を示す図である。

【 図 2 】 例示的实施形態に基づくオーディオ・コンテンツからのオーディオ・オブジェクト抽出のための方法のフローチャートである。

【 図 3 】 ある例示的实施形態に基づくオーディオ・コンテンツからのオーディオ・オブジェクト抽出のためのシステム 3 0 0 のブロック図である。

【 図 4 】 例示的实施形態を実装するのに好適な例示的なコンピュータ・システムのブロック図である。 図面を通じて、同じまたは対応する参照記号は同じまたは対応する部分を指す。

【 発明を実施するための形態 】

【 0 0 1 3 】

例示的实施形態の原理についてこれから図面に示されるさまざまな例示的实施形態を参照して述べる。これらの実施形態の描出は、単に当業者が例示的实施形態をよりよく理解し、さらに実装することができるようにするためだけであり、いかなる仕方であれ本稿に開示される例示的实施形態の範囲を限定することは意図されていないことは理解しておくべきである。また、「第一」、「第二」などの用語は異なる対象を指示するために使われているのであり、対象の序列についていかなる限定をも示唆するものではない。

【 0 0 1 4 】

上述したように、レガシーのチャンネル・ベースのオーディオ・コンテンツは多様な再生セッティングに適応するには不十分である。特に、再生セッティングがミキサーの構成と一致しない場合、該再生セッティングによって表現される聴取経験は劣化する。さらに、芸術的意図を保存しつつ該再生セッティングでオーディオ・コンテンツを表現することも、オブジェクト分離技法にとっての課題となる。

【 0 0 1 5 】

したがって、チャンネル・ベースのオーディオ・コンテンツからできるだけクリーンなオーディオ・オブジェクトを抽出することが所望される。図 1 は、複数のチャンネルに基づくフォーマットのオーディオ信号のセグメントの例を示している。図 1 に示されるように、オーディオ信号のセグメント 1 0 0 は時間および周波数領域において表現されている。横軸によって表現される時間領域では、オーディオ信号のセグメント 1 0 0 は時間軸 T に沿っていくつかのフレームを含みうる。フレームはたとえば、t1 から t2 の時間長さであってもよい。オブジェクト抽出におけるその後の計算およびプロセスの便利のため、オーディオ信号のフレームは時間軸に沿って複数の部分にさらに分割されてもよい（図 1 の破線によって示される）。他方、垂直軸によって表現される周波数領域では、オーディオ信号のセグメント 1 0 0 は、フル帯域信号を表現しており、これもオブジェクト抽出におけるその後の計算およびプロセスの便利のため、周波数領域に沿って複数のサブバンドに分割されることができる。多くの利用可能なスペクトル変換技法がサブバンド分割において適用されうる。たとえば、高速フーリエ変換（FFT）または複素直交ミラー・フィルター（CQMF）がある。人間の聴覚系の特性を考えると、周波数領域における分割は均一でなくても

10

20

30

40

50

よく、低周波数部分ではより細かく、高周波数部分ではより粗くてもよい。図1に示されるように、オーディオ信号100は複数のチャンネル、たとえばチャンネルC1ないしC5に關係している。換言すれば、入力オーディオ信号100は複数のオーディオ信号成分を含み、そのそれぞれはチャンネルC1ないしC5の一つに対応する。したがって、ここでのオーディオ・コンテンツは、複数のチャンネルに基づいて、フル帯域オーディオ信号またはサブバンド・オーディオ信号のチャンネルでありうる。セグメントは、限定なしに、フレーム、フレームの一部、二つ以上のフレームでありうる。いくつかの例示的实施形態では、オーディオ・コンテンツは、フル帯域オーディオ信号のフレームを、周波数領域および時間領域の少なくとも一方において分割することによって得られる一つまたは複数のオーディオ・ブロックを含んでいてもよい。例示的实施形態によれば、オブジェクト抽出があるオーディオ・ブロック、たとえばブロックB1に対して実行されることが所望される場合、その上側の近傍のm個のブロック（単数または複数）および下側の近傍のm個のブロック（単数または複数）が典型的には考慮に入れられる。いくつかの実施形態では、mは1に設定されてもよい。この時点で、それぞれチャンネルC1ないしC5に基づくブロックB0ないしB2と一緒に考慮に入れられ、その全体が、処理されるべきオーディオ・コンテンツをなす。

【0016】

図2は、例示的实施形態に基づくオーディオ・コンテンツからのオーディオ・オブジェクト抽出のための方法200のフローチャートを示している。上述したように、オーディオ・コンテンツは複数のチャンネルのオーディオ信号によって表現されている。

【0017】

図のように、段階S201において、前記複数のチャンネルの第一のチャンネルについての第一の部分集合および第二のチャンネルについての第二の部分集合を含む投影空間の第一の集合が同定される。いくつかの実施形態では、第一および第二のチャンネルは、前記複数のチャンネルのうちの任意のチャンネルであってもよいが、他の例示的实施形態では後述するいくつかの基準に基づいて選択されてもよい。

【0018】

既知のように、チャンネルのオーディオ信号は、個別の成分を得るためにさまざまな空間に投影されることがある。限定ではなく例解のために、オーディオ・コンテンツの第一および第二のチャンネルのオーディオ信号表現についてそれぞれ行列 $X \in \mathbb{R}^{d \times n}$ および $Y \in \mathbb{R}^{k \times n}$ が生成されるとする。ここで、dおよびkはそれぞれのオーディオ信号に含まれる周波数軸に沿ったサブバンド分割（単数または複数）の数を表わし（典型的には $d=k$ ）、nはオーディオ信号における時間軸に沿って分割された部分の数を表わす。すなわち、XおよびYはそれぞれ第一および第二のチャンネルからのオーディオ・コンテンツのオーディオ信号を表わす。すると、XおよびYをそれぞれの投影空間に投影するために投影ベクトル x および y が使用できる。ここで、 $x \in \mathbb{R}^d$ であり、 $y \in \mathbb{R}^k$ である。換言すれば、 $x^T X$ および $y^T Y$ はXおよびYについてのそれぞれの投影空間において投影された成分を表わしうる。ここで、 x^T および y^T はそれぞれ x および y の転置である。さらに、複数の x について、各 x を使ってXを投影することによって得られる対応する複数の空間がある。これら複数の空間の集合は、限定なしに、区別の簡単のために、段階S201の第一のチャンネルについての第一の部分集合と称される。同様に、複数の y について、各 y を使ってYを投影することによって得られる対応する複数の空間がある。これら複数の空間の集合は、限定なしに、区別の簡単のために、段階S201の第二のチャンネルについての第二の部分集合と称される。いくつかの例では、第一の部分集合と第二の部分集合の和集合は投影空間の前記第一の集合をなす。

【0019】

投影空間の第一の集合は通常、第一および第二のチャンネルについて複数の投影空間を含むが、ただ一つの空間を含んでいてもよいことを注意しておくべきである。この場合、Xについての投影空間およびYについての投影空間は同じものである。例示的实施形態の範囲はこの点で限定されない。

【0020】

10

20

30

40

50

次いで、方法は段階S202に進む。ここで、第一のチャネルと第二のチャネルの間の相関の第一の集合が決定される。第一の集合の相関のそれぞれは前記第一の部分集合の投影空間の一つおよび前記第二の部分集合の投影空間の一つに対応する。

【 0 0 2 1 】

投影空間の第一の部分集合および投影空間の第二の部分集合を含む投影空間の前記第一の集合が同定されたのち、オブジェクト抽出を容易にするためにいくつかの基準に基づいて、一対の投影空間が、投影空間の第一の部分集合および投影空間の第二の部分集合からそれぞれ選ばれることができる。例示的实施形態によれば、具体的に、XおよびYの両方に共通のオブジェクトが存在しているが、他の源またはノイズによって汚染されていて、該共通のオブジェクトがXまたはYからより容易に分離されるXおよびYについてのそれぞれの投影空間を見出すことが所望される。

10

【 0 0 2 2 】

例示的实施形態によれば、投影空間の各対について相関が計算される。投影空間の対の一方は第一の部分集合から選ばれ、投影空間の対の他方は第二の部分集合から選ばれる。それにより相関の集合（つまり、段階S202の相関の第一の集合）が形成される。たとえば、 ω_x および ω_y に関するXとYの間の相関は次のように計算されうる：

【 数 1 】

$$\rho = \frac{\omega_x^T X Y^T \omega_y}{\sqrt{(\omega_x^T X X^T \omega_x)(\omega_y^T Y Y^T \omega_y)}}$$

20

(1)

ここで、 ω_x および ω_y の意味は上記と同じままであり、 $\omega_x \in \mathbb{R}^d$ であり、 $\omega_y \in \mathbb{R}^k$ である。

【 0 0 2 3 】

引き続き図2を参照するに、段階S203において、段階S203では、少なくとも部分的には相関の前記第一の集合のうちの第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づいて、前記第一のチャネルのオーディオ信号からオーディオ・オブジェクトが抽出される。ここで、前記第一の相関は第一のあらかじめ定義された閾値より大きい。

30

【 0 0 2 4 】

例示的实施形態によれば、第一のあらかじめ定義された閾値は所望に応じて任意の時点で設定され、調整されることができる。ある例示的实施形態では、第一のあらかじめ定義された閾値は、単に最大相関より小さいが、相関の前記第一の集合における他の相関よりは大きいものとして設定されることができる。この場合、段階S203における目的は、最大の ρ を見出し、それによりさらにオブジェクト抽出のための ω_x および ω_y を同定することである。よって、段階S203では、次が意図される：

【 数 2 】

$$\max_{\omega_x, \omega_y} \omega_x^T X Y^T \omega_y$$

40

(2)

$\omega_x^T X X^T \omega_x = 1, \omega_y^T Y Y^T \omega_y = 1$ の条件の下で

ここで、 X^T 、 Y^T 、 ω_x^T 、 ω_y^T はX、Y、 ω_x 、 ω_y のそれぞれの転置である。 $Y Y^T$ が非特異であれば、次の最適化問題を解くことによって ω_x が得られることが示せる：

【数 3】

$$\max_{\omega_x} \omega_x^T XY^T (YY^T)^{-1} YX^T \omega_x \quad (3)$$

 $\omega_x^T XX^T \omega_x = 1$ の条件の下で

換言すれば、上記の式は、次の一般化された固有値問題の一番上の諸固有値に対応する諸固有ベクトルを見出そうとするものである：

【数 4】

$$XY^T (YY^T)^{-1} YX^T \omega_x = \eta XX^T \omega_x \quad (4)$$

ここで、 η は固有ベクトル ω_x に対応する固有値を表わす。

【0025】

上述したように、いくつかの例示的实施形態によれば、典型的には正規直交性の制約条件のもとで、複数の投影ベクトル ω_x および ω_y があることがある。その場合、これらの複数の投影ベクトルが次の最適化問題を解くことによって同時に計算できる：

【数 5】

$$\max_{W_x} \text{trace}(W_x^T XY^T (YY^T)^{-1} YX^T W_x) \quad (5)$$

 $W_x^T XX^T W_x = I$ の条件の下で

ここで、 $W_x \in \mathbb{R}^{d \times l}$ は投影行列を表わし、 $W_x = [\omega_x^1, \dots, \omega_x^l]$ であり、 l は投影ベクトルの数を表わし、 I は恒等行列を表わす。

【0026】

まとめると、第一および第二のチャネルのオーディオ入力について、オブジェクト抽出の準備のために W_x 、 W_y およびそれらの間の対応する諸相関 R が決定される。ここで、 $W_x = [\omega_x^1, \dots, \omega_x^l]$ であり、 $W_y = [\omega_y^1, \dots, \omega_y^l]$ であり、 ω_x^i または ω_y^i のそれぞれは列ベクトルを表わし、それが投影空間の基底として使用できる。 R は対角線上に 0 でない要素（すなわち）をもつだけの正方相関行列を表わす。 R における i 番目の 0 でない対角要素 r_{ii} について、それは $\omega_x^{iT} X$ と $\omega_y^{iT} Y$ との間の類似性スコアを測る。 $\omega_x^{iT} X$ または $\omega_y^{iT} Y$ のそれぞれが n 次元ベクトルを表わすことを注意しておくべきである。ここで、 n はオーディオ信号のセグメント内の部分の数である。こうして、この測度は、オーディオ・ブロックに基づくオーディオ・コンテンツにおける類似性を反映する。上述したように、 X および Y を X および Y の成分が両者の間の高い相関を示すそれぞれの投影空間に投影することによって、 X と Y の間の高い類似性が観察でき、 X と Y の間の共通のオブジェクトがしかるべく抽出される。

【0027】

たとえば、 i 番目の投影空間について、オブジェクト X_i^* が X から次式により復元される：

【数 6】

$$X_i^* = \omega_x^i \omega_x^{iT} X \quad (6)$$

次いで、（前記第一の部分集合からの l 個の投影空間に対応する） l 個の投影ベクトルが

10

20

30

40

50

らなる W_x について、 X^* が次の代替的な式において計算されうる：

【数 7】

$$F = W_x H W_x^T \quad (7)$$

$$X^* = F X \quad (8)$$

ここで、 H は重み付け対角行列を表わし、0でない要素はその対角線上にあり、その非対角要素はみな0である。 H の導入は、 X^* の復元への諸投影ベクトルの寄与を区別することに有益である。具体的には、ある対の投影空間について、 X と Y がより類似しているほど、 H はより高くなる。結果として、諸 X および諸 Y はそれぞれ前記ある対の投影空間において抽出されることができる。

10

【0 0 2 8】

本稿に開示される例示的实施形態によれば、 H の対角線上の値を決定するための一つの潜在的なアプローチは、それらを相関行列 R に依存して設定することである。上述したように、 R の対角要素は、 W （たとえば W_x または W_y ）の列ベクトルによって構築される投影空間にマッピングされる一対のチャンネルの間の類似性を反映する。このように、より高い類似性スコアは、同じオブジェクトが存在し、これらの空間から復元できる、より大きな確率を示す。よって、高い類似性スコアをもつ空間から「より多くの」オブジェクトを抽出することが合理的である。すなわち、 H は R の適切な関数によって制御できる。つまり、

20

$$H = f(R) \quad (9)$$

ここで、関数 f はいかなる関数でもよく、その値は入力値の増大とともに減少しない。たとえば、 H は、対角要素の和が1に等しい、規格化された R であることができる。

【0 0 2 9】

上述したように、第一および第二のチャンネルは前記複数のチャンネルのうちの任意のチャンネルでありうる。すなわち、段階S203における第一のチャンネルのオーディオ信号からのオブジェクト抽出は第二のチャンネルに関して実行されるように示されているが、実質的に前記複数のチャンネルからのどのチャンネルに関して実行されてもよい。さらに、段階S203では第一のチャンネルのオーディオ信号についてのオーディオ・オブジェクト抽出について述べられているが、同様の動作は第二のチャンネルのオーディオ信号についてのオーディオ・オブジェクト抽出を実行するために第二のチャンネルに適用されてもよい。すなわち、第二のチャンネルのオーディオ信号についてのオブジェクト抽出が、第一のチャンネルまたは前記複数のチャンネルからの他の任意のチャンネルに関して実行されてもよい。これについては簡潔のためここでは詳述しない。例示的实施形態の範囲はこれに関して限定されない。

30

【0 0 3 0】

あるいはまた、いくつかの例示的实施形態では、第一および第二のチャンネルはいくつかの基準に基づいて選択されてもよい。たとえば、両方のチャンネルは、前記複数のチャンネルをその相関に基づいてクラスタリングすることによって得られるあるチャンネル・グループから選択されてもよい。いくつかの例示的实施形態では、前記複数のチャンネルのうちの一対のチャンネルの間の相関は、本稿では、該一対のチャンネルの間の一般的な相関をいう。たとえば、前記複数のチャンネルのうちの一対のチャンネルの間のこの相関は、以下の諸段階によって得られてもよい。

40

【0 0 3 1】

第一に、その対のチャンネルについての投影空間の第二の集合であって、その対のチャンネルの一方についての第三の部分集合およびその対のチャンネルの他方についての第四の部分集合を含むものが同定される。例として、この段階は、段階S201と同様の仕方で実装されることができ、ここでは詳述しない。投影空間の前記第二の集合は、投影空間の前記第一の集合とは異なってもよいことを注意しておく。ただし、場合によってはそれらは同じであってもよい。

【0 0 3 2】

50

次いで、その対のチャンネルの間の相関の第二の集合が決定される。ここで、前記第二の集合の相関のそれぞれは前記第二の部分集合の投影空間の一つおよび前記第四の部分集合の投影空間の一つに対応する。また、この段階は、段階S202と同様の仕方で実装されてもよい。たとえば、その対のチャンネルのそれぞれのオーディオ信号からそれぞれ生成される行列XおよびYについて、前記第二の集合の相関のそれぞれを計算するために式(1)を使う。ここでもまた、相関の前記第一の集合および相関の前記第二の集合は通常は、異なる対のチャンネルについては異なる。

【0033】

次に、前記第二の集合の相関のうちの一つがその対のチャンネルの間の相関として選択される。ここで、選択された相関は第二のあらかじめ定義された閾値より大きい。この選択段階は、段階S203における第一の相関の選択と同様の仕方で実装されてもよく、ここでは詳述しない。たとえば、式(2)～(5)により実装されてもよい。第二のあらかじめ定義された閾値も所望に応じて任意の時点で設定および調整されうる。ある例示的实施形態では、第二のあらかじめ定義された閾値は、単に相関の第二の集合における最大相関より小さいが、他の相関よりは大きいものとして設定されてもよい。この場合、この段階は、相関の前記第二の集合における最大の相関を、その対のチャンネルの間の相関として選択する。

【0034】

前記複数のチャンネルの相関が計算された後、いくつかの例示的实施形態によれば、あらかじめ定義された閾値より大きな相関をもつチャンネルどうしが一つのグループにクラスタリングされることができる。このあらかじめ定義された閾値は、クラスター間の最小の許容される相対類似性スコアとして解釈でき、時間に対して一定の値に設定されることができる。結果として、一つのグループにクラスタリングされるチャンネルどうしは高いグループ内類似性を示す；一方、異なるグループにクラスタリングされるチャンネルどうしは低いグループ間類似性を示す。したがって、一つのグループからのチャンネルのオーディオ信号は通例共通のオブジェクトをもち、該共通のオブジェクトの関係した成分（つまり段階S203のオーディオ・オブジェクト）は段階S201～S203により各チャンネルについて抽出でき、それによりマルチチャンネル・オブジェクトを生成する。これについては後述する。いくつかの例示的实施形態では、チャンネル・グループの数は、クラスタリング手順が終了するときに自動的に決定される。前記複数のチャンネルのうちのチャンネルが互いに似ているまたは前記複数のチャンネルの各対の間の相関がみな前記あらかじめ定義された閾値より大きい場合には、前記複数のチャンネルは単一のグループとみなされてもよいことを注意しておくべきである。

【0035】

いくつかの例示的实施形態によれば、前記複数のチャンネルの相関に基づいて前記複数のチャンネルをクラスタリングすることは、次の手順によって実行されてもよい。

【0036】

初期化：あらかじめ定義された閾値を設定し、対ごとの類似性行列Sを計算する。ここで、項目 S_{ij} はi番目とj番目のチャンネルの間の類似性を表わす。各チャンネルをクラスターとして初期化する。すなわち、 C_1, \dots, C_T であり、Tはチャンネル数を表わす。

【0037】

ループ

各クラスターについてのクラスター内類似性スコアを、クラスター内のチャンネルの対ごとの類似性スコアを平均することによって計算する。すなわち、

【数8】

$$s_{\text{intra}}(m) = \frac{\sum_{i \in C_m} \sum_{j \in C_m} s_{ij}}{N_m}$$

ここで、 N_m はm番目のクラスターの対の数である。

各クラスター対について絶対的なクラスター間類似性スコアを、それぞれの自クラスター内にあるチャンネルの対ごとの類似性スコアを平均することによって計算する。すなわち、

【数 9】

$$s_{inter}(m, n) = \frac{\sum_{i \in C_m} \sum_{j \in C_n} s_{ij}}{N_{mn}}$$

10

ここで、 N_{mn} はm番目とn番目のクラスターの間の対の数を表わす。

各クラスター対についての相対的なクラスター間類似性スコアを、絶対的なクラスター間スコアを二つのクラスター内類似性スコアの平均で割ることによって計算する。すなわち、

【数 10】

$$s_{rela}(m, n) = \frac{s_{inter}(m, n)}{0.5 \times (s_{intera}(m) + s_{intera}(n))}$$

最大の相対クラスター間類似性スコアをもつクラスター対を見出す。最大スコアが前記あらかじめ定義された閾値未満である場合には、ループを終了し；それ以外の場合には、これら二つのクラスターを一つのクラスターに併合する。

20

【0038】

終了。

【0039】

いくつかの例示的实施形態によれば、前記第一のチャンネルが三つ以上のチャンネルを含むグループに属する場合、前記第二のチャンネルについては複数の候補がある。q個のチャンネルからなるチャンネル・グループ $[l_1, \dots, l_{i-1}, l_i, l_{i+1}, \dots, l_q]$ が同定されるとする。 l_i 番目のチャンネルについては、 l_i 番目のチャンネルのオーディオ・オブジェクト抽出のためにq-1個の候補Wがある。すなわち

30

【数 11】

$$W_{l_i}^{(l_1, l_i)}, \dots, W_{l_i}^{(l_{i-1}, l_i)}, W_{l_i}^{(l_{i+1}, l_i)}, \dots, W_{l_i}^{(l_q, l_i)}$$

このように、これらの候補からWを選択するための基準が必要である。

【0040】

上述したように、いくつかの例示的实施形態では、第二のチャンネルはチャンネルのうちの任意のまたはランダムなチャンネルであってもよい。さもなければ、他のいくつかの例示的实施形態では、チャンネル・グループからの第二のチャンネルの選択は、第二のチャンネルのオーディオ信号がチャンネル・グループにおいて最大のエネルギーをもつように実行されてもよい。換言すれば、最も優勢なチャンネルが第二のチャンネルとして選択されることができる。このように、第一のチャンネルについておよびグループ内の他のチャンネルについてのオブジェクト抽出はみな該第二のチャンネル（つまり、最も優勢なチャンネル）に関して実行されてもよい。

40

【0041】

上記のように、第二のチャンネルのオーディオ信号についてのオブジェクト抽出は、第一のチャンネルまたは前記複数のチャンネルからの他の任意のチャンネルに関して実行されてもよい。代替として、いくつかの例示的实施形態によれば、限定なしに、第二のチャンネルがチャンネル・グループにおいて最大のエネルギーをもつ場合、第二のチャンネルのオーディオ信号についてのオブジェクト抽出のためには、単に前記第一のチャンネルを選ぶ代わりに、参

50

照チャンネルを選択することが可能である。たとえば、第二のチャンネルとの相関が第三のあらかじめ定義された閾値より大きいチャンネルが参照チャンネルとして選択されることができ。第三のあらかじめ定義された閾値は、所望に応じて任意の時点において設定され、調整されることができる。ある例示的实施形態では、第三のあらかじめ定義された閾値は、単にチャンネル・グループ内の最大相関よりは小さいが他の相関よりは大きいものとして設定されてもよい。この場合、第二のチャンネルに最も相関しているチャンネルが参照チャンネルとして選択される。方法200の段階S201ないしS203は、第二のチャンネルのオーディオ信号のオーディオ・オブジェクト抽出のために、第二のチャンネルおよび参照チャンネルに適用されることができる。

【0042】

10

いくつかの例示的实施形態において前記第一のチャンネルが、第一と第二のチャンネルの間の相関が第三のあらかじめ定義された閾値より大きいように選択される場合、それがこの場合の参照チャンネルであることができる。したがって、方法200の段階S203において得られるところでは、少なくとも部分的には前記第一の相関と、前記第一の相関に対応する前記第二の部分集合からの投影空間とに基づいて、前記第二のチャンネルのオーディオ信号からオーディオ・オブジェクトが抽出されてもよい。

【0043】

図2に関して上記で示したように、オーディオ・オブジェクトは各チャンネルについてさまざまな投影空間において抽出できる。したがって、いくつかの例では、一つのチャンネル・グループからの諸チャンネルのオーディオ信号から抽出されたオーディオ・オブジェクトに基づいて、マルチチャンネル・オブジェクトが生成されることができる。いくつかのさらなる実施形態によれば、オブジェクト抽出を検証し、調整するために、「ソフト・ゲーティング」手順を導入することが有益でありうる。

20

【0044】

特に、場合によってはある型のオブジェクトの再生がもとの表現への忠誠から逸脱しうるリスクを軽減するために「ソフト・ゲーティング」手順が導入される。「ソフト・ゲーティング」手順を導入するために、たとえば、利得ベクトル g_b が下記のように決定されることができる。

【0045】

第一に、マルチチャンネル・オブジェクトが少なくとも部分的には、第一および第二のチャンネルのオーディオ信号からの抽出されたオーディオ・オブジェクトに基づいて生成される。例示的实施形態によれば、一般に、マルチチャンネル・オブジェクトは、一つのチャンネル・グループからのチャンネルのオーディオ信号から抽出されるオーディオ・オブジェクトに基づいて生成されてもよい。

30

【0046】

第二に、生成されたマルチチャンネル・オブジェクトは、当技術分野で既知の任意の方法を使ってモノ表現にダウンミックスされることができる。抽出されたオブジェクトのモノ表現とものマルチチャンネル表現との間のダウンミックス類似性が次に決定される。たとえば、ダウンミックス類似性は次式のように計算できる。

【0047】

40

【数12】

$$S_b^{(i)} = \left| \frac{\text{Re}(\sum_t M_i(b,t) \times X_i(b,t)^*)}{\sqrt{\sum_t \|M_i(b,t)\|^2} \sqrt{\sum_t \|X_i(b,t)\|^2}} \right|$$

(10)

ここで、 $X_i(b,t)$ はi番目のチャンネルの表現であり、 $M_i(b,t)$ はダウンミックスされたモノ

50

表現であり、 $X_i(b, t)^*$ は $X_i(b, t)$ の共役であり、 $|| \quad ||$ は複素数の絶対値であり、 $\text{Re}()$ の作用は実部を意味する。 b および t はそれぞれサブバンド・インデックスおよび時間部分インデックス、つまり周波数領域および時間領域でのそれぞれのインデックスを表わす。モノ表現ともとの表現との間の全体的なダウンミックス類似性は次式：

【数 1 3】

$$s_b = \frac{1}{T} \sum_{i=1}^T s_b^{(i)} \quad (11)$$

10

のように、あるいは

【数 1 4】

$$s_b = \max_i s_b^{(i)} \quad (12)$$

を介して計算できる。ダウンミックス類似性 s_b によって制御される利得値 g_b 、つまり $g_b^{(1)}$ は

【数 1 5】

$$g_b^{(1)} = f(s_b) \quad (13)$$

20

によって表現できる。

【0 0 4 8】

関数 $f(x)$ が x に関する単調増加関数であることが理解される。 f の定義の一例は次式で書ける。

【0 0 4 9】

【数 1 6】

$$f(x) = \frac{1}{1 + \exp(a \times (x - b))}$$

30

(14)

a の値を負に設定することにより、 $f(x)$ は x に関する単調増加関数になる。

【0 0 5 0】

いくつかの例示的实施形態によれば、計算された利得値は、 X のオブジェクト抽出に影響する重みとして、式(6)または(7)に適用されうる。すなわち、段階S203において第一のチャンネルのオーディオ信号から、また第一のチャンネルが属するチャンネル・グループ内の他の任意のチャンネルのオーディオから、オーディオ・オブジェクトを抽出する際、前記第一の相関または前記対応する投影空間に加えて、式(10)～(12)を介して計算されるダウンミックス類似性も、考慮されるべき因子である。換言すれば、段階S203において第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出することはさらに、ダウンミックス類似性に基づいてオーディオ・オブジェクトを抽出することを含む。したがって、式(6)は次のように変形できる。

40

【0 0 5 1】

【数 1 7】

$$X_i^* = \omega_x^i g_b^{(1)} \omega_x^{iT} X \quad (15)$$

50

式(7)は
【数 1 8】

$$F' = g_b^{(1)} F \quad (16)$$

に変形され、式(8)は
【 0 0 5 2 】
【数 1 9】

$$X^* = F' X \quad (17)$$

に変形される。

10

【 0 0 5 3 】

ダウンミックス類似性 s_b によって制御される利得値 $g_b^{(1)}$ に加えてまたはその代わりに、諸例示の実施形態によれば、利得値 g_b は、次の諸段階によって決定されることもできる：少なくとも部分的には第一および第二のチャンネルのオーディオ信号からの抽出されたオーディオ・オブジェクトに基づいてマルチチャンネル・オブジェクトが生成された後、生成されたマルチチャンネル・オブジェクトがモノ表現にダウンミックスされる。次いで、モノ表現は少なくとも部分的には関係したメタデータ、たとえば推定位置に基づいてプリ・レンダリングされて、「新たな」（つまり、もとのマルチチャンネル・オブジェクトとは異なる）マルチチャンネル・オーディオ信号表現を生成することができる。その後、エネルギー分布に関する、抽出されたオブジェクトのプリ・レンダリングされた表現（つまり、新たなマルチチャンネル・オーディオ信号表現）ともとのマルチチャンネル表現との間のプリ・レンダリングされた類似性が決定される。

20

【 0 0 5 4 】

いくつかの例示の実装では、このプリ・レンダリング類似性は、もとのマルチチャンネル・オブジェクトとモノ・オブジェクトのプリ・レンダリングから帰結したもののエネルギー分布の間の不一致によって反映されることが可能である。すなわち、不一致が大きいほど、プリ・レンダリング類似性は小さくなる。したがって、不一致を測る好適なメトリックが、

【数 2 0】

$$d_b = \sum_{i=1}^T |e_b^i - e_b^{*i}|$$

30

(18)

として、あるいは

【数 2 1】

$$d_b = \max_i |e_b^i - e_b^{*i}|$$

40

(19)

として適切に設計できる。ここで、 e_b^i および e_b^{*i} はそれぞれ、レガシー・コンテンツと、推定されたメタデータと一緒にモノ・オブジェクトをレンダラーを使ってプリ・レンダリングすることから帰結したコンテンツとの、規格化されたエネルギー分布を表わす。 b および i はそれぞれサブバンド・インデックスおよびチャンネル・インデックス、つまり周波数領域およびチャンネル領域でのそれぞれのインデックスを表わす。レンダラーを用いたプリ・レンダリングのチャンネル構成はレガシー・コンテンツのチャンネル構成と同じであることを注意しておくべきである。たとえば、サラウンド5.1レガシー・コンテンツについては、プリ・レンダリングのチャンネル構成もサラウンド5.1であるはずである。規格化さ

50

れたエネルギー分布は

【数 2 2】

$$e_b^i = \frac{E_b^i}{\sum_{i=1}^T E_b^i} \quad (20)$$

を介して計算できる。ここで、 E_b^i は i 番目のチャンネルについての b 番目のサブバンド・エネルギーを表わす。

【0 0 5 5】

したがって、 d_b によって制御される利得値（つまり $g_b^{(2)}$ ）は

10

【数 2 3】

$$g_b^{(2)} = f(d_b) \quad (21)$$

と表現できる。ここで、 $f(d_b)$ は d_b に関する単調減少関数である。

【0 0 5 6】

いくつかの例示的实施形態では、この利得値 $g_b^{(2)}$ も X のオブジェクト抽出に影響する重みとして式(6)または(7)に適用されてもよい。すなわち、段階S203において第一のチャンネルのオーディオ信号から（また第一のチャンネルが属するチャンネル・グループ内の他の任意のチャンネルのオーディオから）オーディオ・オブジェクトを抽出する際、前記第一の相関または前記対応する投影空間に加えて、あるいは前記第一の相関、前記対応する投影空間および前記ダウンミックス類似性に加えて、プリ・レンダリング類似性を反映する、式(18)～(19)を介して計算される不一致も、考慮されるべき因子である。すなわち、段階S203において第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出することはさらに、プリ・レンダリング類似性に基づいてオーディオ・オブジェクトを抽出することを含む。すると、式(6)は次のように変形できる。

20

【0 0 5 7】

【数 2 4】

$$X_i^* = \omega_x^i g_b^{(2)} \omega_x^{iT} X \quad (22)$$

30

または

$$X_i^* = \omega_x^i g_b^{(2)} g_b^{(1)} \omega_x^{iT} X \quad (23)$$

式(7)は

【数 2 5】

$$F'' = g_b^{(2)} F \quad (24)$$

または

$$F'' = g_b^{(1)} \times g_b^{(2)} F \quad (25)$$

40

に変形され、式(8)は

【数 2 6】

$$X^* = F'' X \quad (26)$$

に変形される。

【0 0 5 8】

それぞれダウンミックス類似性およびプリ・レンダリング類似性に関連する利得ベクトル $g_b^{(1)}$ および $g_b^{(2)}$ の少なくとも一方の導入は、抽出されたオブジェクトの再生がもとの表現への忠誠から逸脱するかどうか、および抽出されたオブジェクトの再生が芸術的意図

50

を保存するかどうかを検証しうる。逸脱は、もしあれば、少なくとも、たとえばモノ表現とものマルチチャンネル表現との間の音色の不一致があることを示すことができる。したがって、変形された式(15)、(17)、(22)、(23)、(26)は $g_b^{(1)}$ または $g_b^{(2)}$ の因子を導入することによる逸脱を減らすことができる。

【0059】

図3は、ある例示的实施形態に基づく、オーディオ・コンテンツからのオーディオ・オブジェクト抽出のためのシステム300のブロック図である。ここで、オーディオ・コンテンツは複数のチャンネルのオーディオ信号によって表現されている。図のように、システム300は、前記複数のチャンネルの第一のチャンネルについての第一の部分集合および前記複数のチャンネルの第二のチャンネルについての第二の部分集合を含む投影空間の第一の集合を同定するよう構成された同定ユニット301を含む。システム300はさらに、前記第一および第二のチャンネルの間の相関の第一の集合を決定するよう構成された決定ユニット302を含み、前記第一の集合の相関のそれぞれは前記第一の部分集合の投影空間の一つおよび前記第二の部分集合の投影空間の一つに対応する。システム300はまた、少なくとも部分的には相関の前記第一の集合のうちの第一の相関と、前記第一の相関に対応する前記第一の部分集合からの投影空間とに基づいて、前記第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出するよう構成された抽出ユニット303をも含み、前記第一の相関は第一のあらかじめ定義された閾値より大きい。

【0060】

いくつかの実施形態では、システム300はさらに、チャンネル・グループから前記第一および第二のチャンネルを選択するよう構成された選択ユニットを有していてもよい。前記チャンネル・グループは、前記複数のチャンネルの相関に基づいて前記複数のチャンネルをクラスタリングすることによって得られる。

【0061】

いくつかの実施形態では、前記複数のチャンネルの一对のチャンネルの間の相関が：投影空間の第二の集合であって、その対のチャンネルの一方についての第三の部分集合およびその対のチャンネルの他方についての第四の部分集合を含むものを同定し；その対のチャンネルの間の相関の第二の集合を決定し、ここで、前記第二の集合の相関のそれぞれは前記第三の部分集合の投影空間の一つおよび前記第三の部分集合の投影空間の一つに対応し；前記第二の集合の相関のうちの一つを、その対のチャンネルの間の相関として選択することによって決定される。ここで、選択された相関は第二のあらかじめ定義された閾値より大きい。

【0062】

いくつかの実施形態では、前記チャンネル・グループからの前記第一および第二のチャンネルの選択は、前記第二のチャンネルのオーディオ信号が前記チャンネル・グループにおいて最大のエネルギーをもつように、前記チャンネル・グループから前記第二のチャンネルを選択することを含む。

【0063】

いくつかの実施形態では、前記チャンネル・グループからの前記第一および第二のチャンネルの選択はさらに、前記第一および第二のチャンネルの間の相関が第三のあらかじめ定義された閾値より大きいように前記チャンネル・グループから前記第一のチャンネルを選択し；前記第二のチャンネルのオーディオ信号から、少なくとも部分的には前記第一の相関と前記第一の相関に対応する前記第二の部分集合からの投影空間とに基づいてオーディオ・オブジェクトを抽出することを含んでいてもよい。

【0064】

いくつかの実施形態では、システム300はさらに、少なくとも部分的には前記第一および第二のチャンネルのオーディオ信号からの抽出されたオーディオ・オブジェクトに基づいて、マルチチャンネル・オブジェクトを生成するよう構成された生成ユニットと；生成されたマルチチャンネル・オブジェクトをモノ表現にダウンミックスするよう構成されたダウンミックス・ユニットと；抽出されたオブジェクトのものとマルチチャンネル表現と前記モノ表現との間のダウンミックス類似性を決定するよう構成された類似性決定ユニットとを

有する。ここで、前記第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出することは、前記ダウンミックス類似性にさらに基づいてオーディオ・オブジェクトを抽出することを含む。

【0065】

いくつかの代替的实施形態では、システム300はさらに、少なくとも部分的には前記第一および第二のチャンネルのオーディオ信号からの抽出されたオーディオ・オブジェクトに基づいて、マルチチャンネル・オブジェクトを生成するよう構成された生成ユニットと；生成されたマルチチャンネル・オブジェクトをモノ表現にダウンミックスするよう構成されたダウンミックス・ユニットと；少なくとも部分的には推定位置に基づいて前記モノ表現をプリ・レンダリングするよう構成されたプリ・レンダリング・ユニットと；抽出されたオブジェクトのものとマルチチャンネル表現と前記プリ・レンダリングされたモノ表現との間の、エネルギー分布に関するプリ・レンダリング類似性を決定するよう構成された類似性決定ユニットとを有する。ここで、前記第一のチャンネルのオーディオ信号からオーディオ・オブジェクトを抽出することは、前記プリ・レンダリング類似性にさらに基づいてオーディオ・オブジェクトを抽出することを含む。

10

【0066】

いくつかの実施形態では、前記オーディオ・コンテンツは、フル帯域オーディオ信号のフレームを、周波数領域および時間領域の少なくとも一方において分割することによって得られる一つまたは複数のオーディオ・ブロックを含んでいてもよい。

【0067】

20

明確のため、システム300のいくつかの任意的なコンポーネントは図3には示していない。しかしながら、図1～図2を参照して上記した事項はみなシステム300に適用可能であることは理解されるはずである。さらに、システム300のコンポーネントは、ハードウェア・モジュールまたはソフトウェア・ユニット・モジュールでありうる。たとえば、いくつかの実施形態では、システム300は、部分的にまたは完全に、たとえばコンピュータ可読媒体において具現されたコンピュータ・プログラム・プロダクトとして実装されるソフトウェアおよび/またはファームウェアとして実装されてもよい。代替的または追加的に、システム300は部分的または完全に、たとえば集積回路(IC)、特定用途向け集積回路(ASIC)、システムオンチップ(SOC)、フィールド・プログラマブル・ゲート・アレイ(FPGA)などのようなハードウェアに基づいて実装されてもよい。例示的実施形態の範囲はこれに関して限定されない。

30

【0068】

図4は、例示的実施形態を実装するために好適な例示的なコンピュータ・システム400のブロック図を示している。図のように、コンピュータ・システム400は、読み出し専用メモリ(ROM)402に記憶されたプログラムまたは記憶ユニット408からランダム・アクセス・メモリ(RAM)403にロードされたプログラムに従ってさまざまなプロセスを実行することのできる中央処理ユニット(CPU)401を有する。RAM 403では、CPU 401がさまざまなプロセスを実行するときに必要なとされるデータなども必要に応じて記憶される。CPU 401、ROM 402およびRAM 403はバス404を介して互いに接続されている。入出力(I/O)インターフェース405もバス404に接続されている。

40

【0069】

以下のコンポーネントがI/Oインターフェース405に接続される：キーボード、マウスなどを含む入力部406；陰極線管(CRT)、液晶ディスプレイ(LCD)などのようなディスプレイまたはスピーカーなどを含む出力部407；ハードディスクなどを含む記憶部408；およびLANカード、モデムなどのようなネットワーク・インターフェース・カードを含む通信部409である。通信部409は、インターネットのようなネットワークを介して通信プロセスを実行する。ドライブ410も必要に応じてI/Oインターフェース405に接続される。磁気ディスク、光ディスク、光磁気ディスク、半導体メモリなどのような着脱可能な媒体411が必要に応じてドライブ410にマウントされ、それにより必

50

要に応じて、そこから読まれたコンピュータ・プログラムが記憶部 408 にインストールされる。

【0070】

特に、例示的实施形態によれば、図2を参照して上記したプロセスがコンピュータ・ソフトウェア・プログラムとして実装されてもよい。たとえば、例示的实施形態の実施形態は、方法200を実行するためのプログラム・コードを含む、機械可読媒体上に有体に具現されたコンピュータ・プログラムを含むコンピュータ・プログラム・プロダクトを含む。そのような実施形態では、コンピュータ・プログラムは、通信ユニット409を介してネットワークからダウンロードおよびマウントされ、および/または着脱可能な媒体411からインストールされてもよい。

10

【0071】

一般に、さまざまな例示的实施形態はハードウェアまたは特殊目的回路、ソフトウェア、論理またはそれらの任意の組み合わせにおいて実装される。いくつかの側面はハードウェアにおいて実装され、一方で他の側面がコントローラ、マイクロプロセッサまたは他のコンピューティング装置によって実行されるファームウェアまたはソフトウェアにおいて実装されてもよい。例示的实施形態のさまざまな側面がブロック図、フローチャートとしてまたは他のいくつかの絵的表現を使って図示され、記述されているが、本稿に記載されるブロック、装置、システム、技法または方法は、限定しない例として、ハードウェア、ソフトウェア、ファームウェア、特殊目的回路または論理、汎用ハードウェアまたはコントローラまたは他のコンピューティング装置またはそれらの何らかの組み合わせにおいて実装されてもよいことは理解されるであろう。

20

【0072】

さらに、フローチャートに示されるさまざまなブロックを方法ステップとしておよび/またはコンピュータ・プログラム・コードの動作から帰結する動作としておよび/または関連する機能(単数または複数)を実行するよう構築された複数の結合された論理回路要素として見ることができる。たとえば、実施形態は、機械可読媒体上に有体に具現されたコンピュータ・プログラムを有するコンピュータ・プログラム・プロダクトを含み、該コンピュータ・プログラムは、上記で述べた諸方法を実行するために構成されたプログラム・コードを含む。

【0073】

本開示のコンテキストにおいて、機械可読媒体は、命令実行システム、装置またはデバイスによってまたはそれとの関連で使うためのプログラムを含むまたは記憶することができるいかなる有体の媒体であってもよい。機械可読媒体は機械可読信号媒体または機械可読記憶媒体でありうる。機械可読媒体は、電子式、磁気式、光学式、電磁式、赤外線または半導体のシステム、装置またはデバイスまたは上記の任意の好適な組み合わせを含みうるが、それに限られなくてもよい。機械可読記憶媒体のより具体的な例は、一つまたは複数のワイヤを有する電気接続、ポータブルなコンピュータ・ディスクレット、ハードディスク、ランダム・アクセス・メモリ(RAM)、読み出し専用メモリ(ROM)、消去可能なプログラム可能型読み出し専用メモリ(EPROMまたはフラッシュ・メモリ)、光ファイバー、ポータブルなコンパクト・ディスク読み出し専用メモリ(CD-ROM)、光記憶デバイス、磁気記憶デバイスまたは上記の任意の好適な組み合わせを含む。

30

40

【0074】

本稿に開示される例示的实施形態の方法を実行するためのコンピュータ・プログラム・コードは、一つまたは複数のプログラミング言語の任意の組み合わせにおいて書かれる。これらのコンピュータ・プログラム・コードは、汎用コンピュータ、特殊目的コンピュータまたは他のプログラム可能なデータ処理装置のプロセッサに提供されてもよく、それにより該プログラム・コードは、該コンピュータまたは他のプログラム可能なデータ処理装置のプロセッサによって実行されたとき、フローチャートおよび/またはブロック図において規定された機能/動作を実装させる。プログラム・コードは完全にコンピュータ上で、部分的にコンピュータ上で、スタンドアローンのソフトウェア・パッケージとして、

50

部分的にはコンピュータ上で部分的にはリモート・コンピュータ上で、あるいは完全にリモート・コンピュータまたはサーバー上で実行されてもよい。

【0075】

さらに、動作は特定の順序で描かれているが、これは、そのような動作が示される特定の順序で、あるいは逐次順に実行されること、あるいは所望される結果を達成するために示されているすべての動作が実行されることを要求するものと理解されるべきではない。ある種の状況では、マルチタスクおよび並列処理が有利であることがある。同様に、いくつかの個別的な実装詳細が上記の議論に含まれるものの、これらはいずれかの実施形態のまたは特許請求されうるものの範囲に対する限定として解釈されるべきではなく、むしろ特定の実施形態に固有でありうる事項の記述と解釈されるべきである。別個の実施形態のコンテキストにおいて本明細書に記載されるある種の特徴は、単一の実施形態において組み合わせて実装されてもよい。逆に、単一の実施形態のコンテキストにおいて記述されているさまざまな特徴が、複数の実施形態において別個にまたは任意の好適なサブコンビネーションにおいて実装されることもできる。

10

【0076】

付属の図面との関連で読まれるときの上記の記述に鑑み、本稿に開示される上記の例示の実施形態へのさまざまな修正および適応が当業者に明白となることがありうる。任意の、あらゆる修正がそれでも、本稿に開示される、限定しない、例示的な実施形態の範囲内にはいる。さらに、本稿に記載される他の実施形態が、上記の記述および図面に呈示される教示の恩恵をもつ当業者には思いつくであろう。

20

【0077】

例示の実施形態は、本稿に記載される形の任意のもので具現されうる。たとえば、以下の付番実施例（EEE: enumerated example embodiment）は、例示の実施形態のいくつかの側面のいくつかの構造、特徴および機能を記述するものである。

〔EEE1〕

複数のチャンネルに基づくフォーマットのオーディオ・コンテンツからのオーディオ・オブジェクト抽出のための方法であって：

諸投影空間から導出されたフィルタ行列を通じたオブジェクト抽出と；

任意的に、

芸術的な意図を保存するよう、前記抽出されたオブジェクトまたは前記フィルタ行列に対して追加的な利得を提供するソフト・ゲーティングとを含む、方法。

30

〔EEE2〕

オブジェクト抽出を実行するために各オーディオ・ブロックについて、

各チャンネル入力についての投影ベクトル集合を生成し、各対のチャンネルの間の最大相関（類似性スコア）がそれらを諸投影空間に投影することによって計算され；

それらのチャンネルを対応する相関（類似性スコア）に基づいてグループ化し；

グループ内の各チャンネルについて、各オーディオ・ブロックについてのフィルタ行列を導出し；

各チャンネルの入力オーディオ信号にそれ自身のフィルタ行列を乗算することによってオブジェクトを復元し、

40

前記オーディオ・ブロックは、フル帯域オーディオ信号のフレームを周波数領域および時間領域の少なくとも一方において分割することによって得られる、

EEE1記載の方法。

〔EEE3〕

前記投影ベクトル集合が、現在のオーディオ・ブロックおよび近隣のオーディオ・ブロックを使うことによってブロックごとに形成される、EEE2記載の方法。

〔EEE4〕

フィルタ行列Fの生成がWおよびHの選択に関わり、

Hの選択は式(9)を介してなされることができ、

50

Wの選択はグループ内での第二のチャンネルの特定に関わる、
E E E 3 記載の方法。

〔 E E E 5 〕

第二のチャンネルの特定はチャンネル・エネルギーに基づき、たとえばグループのうちで最大のエネルギーをもつチャンネルが選択される、E E E 4 記載の方法。

〔 E E E 6 〕

第一のチャンネルについてのWの選択は、前記第二のチャンネルに関して投影ベクトル集合を選択する、E E E 4 記載の方法。

〔 E E E 7 〕

前記第二のチャンネルについてのWの選択は、前記第二のチャンネルへの最大の類似性を示す、前記グループのうちでの前記チャンネルに関する投影ベクトル集合を選択する、E E E 4 記載の方法。

10

〔 E E E 8 〕

前記ソフト・ゲーティング段階は、各オーディオ・ブロックについての利得ベクトルの生成に関わり、出力を生成するために、前記利得ベクトルはブロックごとに前記オーディオ信号入力を乗算される、E E E 1 記載の方法。

〔 E E E 9 〕

前記利得ベクトルは、それぞれプリ・ダウンミックス処理およびプリ・レンダリング処理から生成される二つのサブ利得ベクトルの積として計算される、つまり式(22)である、E E E 8 記載の方法。

20

〔 E E E 1 0 〕

プリ・ダウンミックス処理からの前記サブ利得ベクトルは式(10)～(13)によって計算できる、E E E 9 記載の方法。

〔 E E E 1 1 〕

プリ・レンダリング処理からの前記サブ利得ベクトルは式(17)～(20)によって計算できる、E E E 9 記載の方法。

〔 E E E 1 2 〕

複数のチャンネルに基づくフォーマットになっているオーディオ・コンテンツからのオーディオ・オブジェクト抽出のためのシステムであって、E E E 1 ないし 1 1 のうちいずれか一項記載の方法を実行するよう構成された諸ユニットを有する、システム。

30

〔 E E E 1 3 〕

オーディオ・コンテンツからのオーディオ・オブジェクト抽出のためのコンピュータ・プログラム・プロダクトであって、当該コンピュータ・プログラム・プロダクトは、非一時的なコンピュータ可読媒体上に有体に記憶されており、実行されたときにE E E 1 ないし 1 1 のうちいずれか一項記載の方法の段階を機械に実行させる機械実行可能命令を有する、コンピュータ・プログラム・プロダクト。

【 0 0 7 8 〕

本稿に開示される例示的实施形態は開示される特定の实施形態に限定されず、他の实施形態も付属の請求項の範囲内に含まれることが意図されていることは理解されるであろう。本稿では個別的な用語が使われているが、それらは一般的で記述的な意味において使われているだけであり、限定のためではない。

40

【図 1】

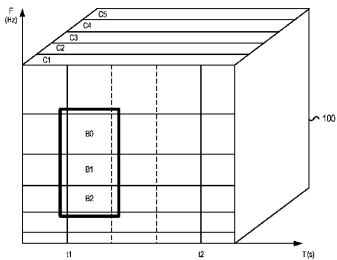
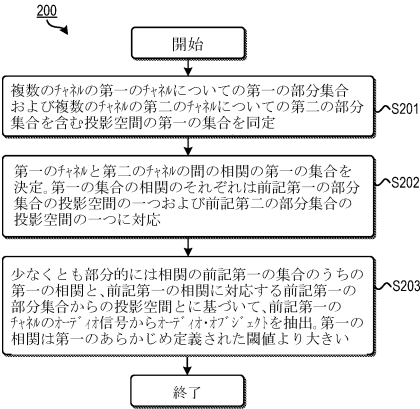
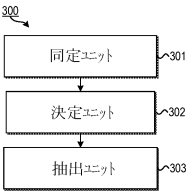


Figure 1

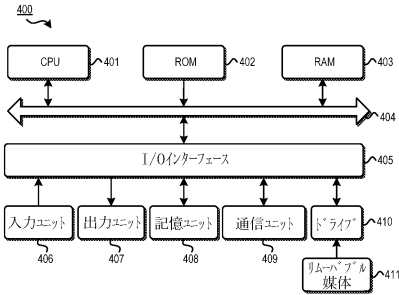
【図 2】



【図 3】



【図 4】



フロントページの続き

(72)発明者 フー, ミンチン

中華人民共和国, ベイジン 10020, シャオヤン・ディストリクト, イースト・サード・リング・ミドル・ロード, ワールド・フィナンシャル・センター・ナンバー 1, ドルビー ラボラトリーズ インターナショナル サービスズ内

(72)発明者 ルー, リエ

中華人民共和国, ベイジン 10020, シャオヤン・ディストリクト, イースト・サード・リング・ミドル・ロード, ワールド・フィナンシャル・センター・ナンバー 1, ドルビー ラボラトリーズ インターナショナル サービスズ内

(72)発明者 チェン, リアンウー

中華人民共和国, ベイジン 10020, シャオヤン・ディストリクト, イースト・サード・リング・ミドル・ロード, ワールド・フィナンシャル・センター・ナンバー 1, ドルビー ラボラトリーズ インターナショナル サービスズ内

審査官 大野 弘

(56)参考文献 再公表特許第2010/092913(JP, A1)

(58)調査した分野(Int.Cl., DB名)

G10L 21/028

G10L 19/008

G10L 21/0308