

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2016-224449

(P2016-224449A)

(43) 公開日 平成28年12月28日(2016.12.28)

| (51) Int.Cl. | F I | テーマコード (参考) |
|---------------------------|--------------------|-------------|
| G 1 O L 21/034 (2013.01) | G 1 O L 21/034 | 5 J O 3 O |
| G 1 O L 25/51 (2013.01) | G 1 O L 25/51 | |
| G 1 O L 21/0364 (2013.01) | G 1 O L 21/0364 | |
| H O 3 G 5/16 (2006.01) | H O 3 G 5/16 1 6 5 | |

審査請求 未請求 請求項の数 18 O L (全 85 頁)

(21) 出願番号 特願2016-145567 (P2016-145567)
 (22) 出願日 平成28年7月25日 (2016.7.25)
 (62) 分割の表示 特願2016-505487 (P2016-505487) の分割
 原出願日 平成26年3月17日 (2014.3.17)
 (31) 優先権主張番号 201310100422.1
 (32) 優先日 平成25年3月26日 (2013.3.26)
 (33) 優先権主張国 中国 (CN)
 (31) 優先権主張番号 61/811,072
 (32) 優先日 平成25年4月11日 (2013.4.11)
 (33) 優先権主張国 米国 (US)

(71) 出願人 507236292
 ドルビー ラボラトリーズ ライセンシング
 グ コーポレイション
 アメリカ合衆国 94103 カリフォル
 ニア州 サンフランシスコ マーケット
 ストリート 1275
 (74) 代理人 100107766
 弁理士 伊東 忠重
 (74) 代理人 100070150
 弁理士 伊東 忠彦
 (74) 代理人 100091214
 弁理士 大貫 進介

(特許庁注：以下のものは登録商標)

1. ウィンドウズ

最終頁に続く

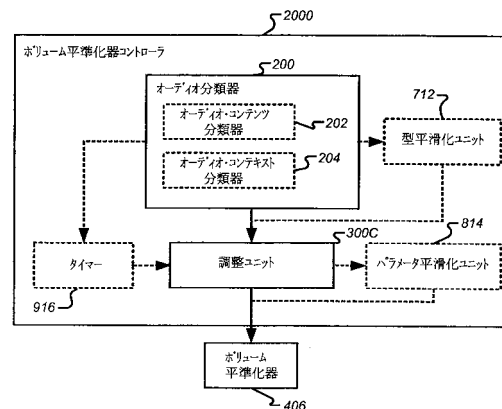
(54) 【発明の名称】 ボリューム平準化器コントローラおよび制御方法

(57) 【要約】

【課題】 ボリューム平準化器コントローラおよび制御方法が提供される。

【解決手段】 ある実施形態では、ボリューム平準化器コントローラは、リアルタイムでオーディオ信号のコンテンツ型を識別するためのオーディオ・コンテンツ分類器と；識別されたコンテンツ型に基づいて連続的な仕方でボリューム平準化器を調整する調整ユニットとを有する。調整ユニットは、ボリューム平準化器の動的な利得を、前記オーディオ信号の情報性のコンテンツ型と正に相関させ、ボリューム平準化器の動的な利得を、前記オーディオ信号の干渉性のコンテンツ型と負に相関させるよう構成されていてもよい。

【選択図】 図 2 0



【特許請求の範囲】**【請求項 1】**

リアルタイムでオーディオ信号のコンテンツ型を識別するためのオーディオ・コンテンツ分類器と；

識別されたコンテンツ型に基づいて連続的な仕方でボリューム平準化器を調整する調整ユニットとを有する、
ボリューム平準化器コントローラ。

【請求項 2】

請求項 1 記載のボリューム平準化器コントローラを有するオーディオ処理装置。

【請求項 3】

オーディオ信号の短期的セグメントのコンテンツ型を識別する段階と；

少なくとも部分的には識別されたコンテンツ型に基づく前記短期的セグメントのコンテキスト型を識別する段階とを含む、
オーディオ分類方法。

【請求項 4】

前記コンテンツ型进行分类する動作が、前記短期的セグメントをコンテンツ型VoIP発話またはコンテンツ型非VoIP発話に分类することを含み、

前記コンテキスト型を識別する動作が、VoIP発話および非VoIP発話の信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分类するよう構成されている、

請求項 3 記載のオーディオ分類方法。

【請求項 5】

前記コンテンツ型进行分类する動作がさらに、

短期的セグメントをコンテンツ型VoIPノイズおよびコンテンツ型非VoIPノイズに分类することを含み、

前記コンテキスト型を識別する動作が、VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズの信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分类するよう構成されている、

請求項 4 記載のオーディオ分類方法。

【請求項 6】

前記コンテキスト型を識別する動作が：

VoIP発話の信頼値が第一の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；

VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；

それ以外の場合には、前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、

請求項 4 記載のオーディオ分類方法。

【請求項 7】

前記コンテキスト型を識別する動作が：

VoIP発話の信頼値が第一の閾値より大きい場合またはVoIPノイズの信頼値が第三の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；

VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合またはVoIPノイズの信頼値が前記第三の閾値より大きくない第四の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；

それ以外の場合には前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、

請求項 2 記載のオーディオ分類方法。

【請求項 8】

前記コンテンツ型の過去の信頼値に基づいて現在の時点での前記コンテンツ型の信頼値

10

20

30

40

50

を平滑化することをさらに含む、請求項 1 記載のオーディオ分類方法。

【請求項 9】

前記の平滑化する動作は、現在の短期的セグメントの信頼値と最後の短期的セグメントの平滑化された信頼値との重み付けされた和を計算することによって、現在の短期的セグメントの平滑化された信頼値を決定するよう構成されている、請求項 8 記載のオーディオ分類方法。

【請求項 10】

前記短期的セグメントからコンテンツ型発話を識別する段階をさらに含み、平滑化前の現在の短期的セグメントについてのVoIP発話の信頼値が、所定の信頼値として、あるいはコンテンツ型発話についての信頼値が第五の閾値より低い最後の短期的セグメントの平滑化された信頼値として、設定される、請求項 9 記載のオーディオ分類方法。

10

【請求項 11】

前記コンテキスト型を識別する動作が、特徴として、前記短期的セグメントのコンテンツ型の信頼値および前記短期的セグメントから抽出された他の特徴を使って、機械学習モデルに基づいて前記短期的セグメントを分類するよう構成されている、請求項 2 記載のオーディオ分類方法。

【請求項 12】

前記コンテキスト型を識別する動作が同じコンテキスト型を連続的に出力する継続時間を測定する段階をさらに含み、当該オーディオ分類方法は、新しいコンテキスト型の継続時間の長さが第六の閾値に達するまで、現在のコンテキスト型を使い続けるよう構成される、請求項 6 記載のオーディオ分類方法。

20

【請求項 13】

あるコンテキスト型から別のコンテキスト型への異なる遷移対について、異なる第六の閾値が設定される、請求項 12 記載のオーディオ分類方法。

【請求項 14】

前記第六の閾値が、前記新しいコンテキスト型の信頼値と負に相関している、請求項 12 記載のオーディオ分類方法。

【請求項 15】

前記第一および/または第二の閾値が、最後の短期的セグメントのコンテキスト型によって異なる、請求項 6 記載のオーディオ分類方法。

30

【請求項 16】

オーディオ信号の短期的セグメントのコンテンツ型を識別するオーディオ・コンテンツ分類器と；

少なくとも部分的には前記オーディオ・コンテンツ分類器によって識別されたコンテンツ型に基づいて前記短期的セグメントのコンテキスト型を識別するオーディオ・コンテキスト分類器とを有しており、

請求項 3 記載の方法を実行するよう構成されている、

オーディオ分類器。

【請求項 17】

オーディオ信号の短期的セグメントのコンテンツ型を識別するオーディオ・コンテンツ分類器と；

40

請求項 3 に従って、少なくとも部分的には前記オーディオ・コンテンツ分類器によって識別されたコンテンツ型に基づいて前記短期的セグメントのコンテキスト型を識別するオーディオ・コンテキスト分類器とを有する、

オーディオ分類器を有するオーディオ処理装置。

【請求項 18】

一つまたは複数のプロセッサによって実行されたときに請求項 3 記載のオーディオ分類方法を実行する命令が記憶されている非一時的なコンピュータ可読媒体。

【発明の詳細な説明】

【技術分野】

50

【 0 0 0 1 】

関連出願への相互参照

本願は2013年3月26日に出願された中国特許出願第201310100422.1号および2013年4月11日に出願された米国仮特許出願第61/811,072号の優先権を主張するものである。これら各出願はここに参照によってその全体において組み込まれる。

【 0 0 0 2 】

技術分野

本願は概括的にはオーディオ信号処理に関する。詳細には、本願の実施形態はオーディオ分類および処理、特にダイアログ向上器、サラウンド仮想化器、ボリューム平準化器および等化器の制御のための装置および方法に関する。

10

【背景技術】

【 0 0 0 3 】

いくつかのオーディオ改善装置は、オーディオの全体的な品質を改善し、相応してユーザーの経験を向上させるために、時間領域またはスペクトル領域のいずれかにおいてオーディオ信号を修正する傾向がある。さまざまなオーディオ改善装置がさまざまな目的のために開発されている。オーディオ改善装置のいくつかの典型的な例は次のものを含む。

【 0 0 0 4 】

ダイアログ向上器〔エンハンサー〕：ダイアログは、映画およびラジオまたはテレビ・プログラムにおいてストーリーを理解するための最も重要な構成要素である。特に聴力が衰えつつある高齢者のために、ダイアログの明瞭性および了解性を高めるためにダイアログを向上させる諸方法が開発された。

20

【 0 0 0 5 】

サラウンド仮想化器：サラウンド仮想化器は、PCの内部スピーカーを通じてまたはヘッドフォンを通じてサラウンド（マルチチャンネル）サウンド信号がレンダリングされるようにする。すなわち、（スピーカーおよびヘッドフォンのような）ステレオ装置を用いて、仮想的にサラウンドの効果を生成し、消費者のために映画館の体験を提供するのである。

【 0 0 0 6 】

ボリューム平準化器：ボリューム平準化器は、再生時にオーディオ・コンテンツのボリュームを調整し、目標ラウドネス値に基づいて時間軸を通じてボリュームがほぼ一貫しているようにすることをねらいとする。

30

【 0 0 0 7 】

等化器〔イコライザー〕：等化器は、「トーン」または「音色」として知られるスペクトル・バランスの一貫性を提供し、ユーザーが、ある種の音を強調したり望ましくない音を除去したりするために、個々の周波数帯域での周波数応答（利得）の全体的なプロファイル（曲線または形状）を構成設定できるようにする。伝統的な等化器では、異なる等化器プリセットが、種々の音楽ジャンルのような種々の音のために提供されてもよい。ひとたびプリセットが選択されたらまたは等化プロファイルが設定されたら、手動で等化プロファイルが修正されるまで、同じ等化利得が信号に対して適用される。対照的に、動的等化器は、オーディオのスペクトル・バランスを連続的にモニタリングし、それを所望されるトーンと比較し、オーディオの元のトーンを所望されるトーンに変換するための等化フィルタを動的に調整することによってスペクトル・バランス一貫性を達成する。

40

【 0 0 0 8 】

一般に、オーディオ改善装置はその独自の応用シナリオ／コンテキストをもつ。すなわち、オーディオ改善装置は、あらゆる可能なオーディオ信号についてではなく、ある種のコンテンツの集合についてのみ好適であることがある。異なるコンテンツは異なる仕方で処理される必要があることがあるからである。たとえば、ダイアログ向上方法は、通例、映画コンテンツに適用される。それがダイアログのない音楽に適用されたとしたら、いくつかの周波数サブバンドを誤ってブーストし、重度の音色変化および知覚的な非一貫性を導入することがありうる。同様に、音楽信号に対してノイズ抑制方法が適用されたとした

50

ら、強いアーチファクトが可聴となるであろう。

【0009】

しかしながら、通例はオーディオ改善装置の集合を含むオーディオ処理システムについて、その入力は、必然的に、あらゆる可能な型のオーディオ信号でありうる。たとえば、PCに統合されたオーディオ処理システムは映画、音楽、VoIPおよびゲームを含む多様な源からオーディオ・コンテンツを受領することになる。よって、対応するコンテンツに対してよりよいアルゴリズムまたは各アルゴリズムのよりよいパラメータを適用するために、処理されるコンテンツを識別または区別することが重要になる。

【0010】

オーディオ・コンテンツを区別して、相応してよりよいパラメータまたはよりよいオーディオ改善アルゴリズムを適用するために、伝統的なシステムは、通例、プリセットの集合を事前設計し、ユーザーは再生されるコンテンツについてのプリセットを選ぶことを求められる。プリセットは通例、映画または音楽再生のために特に設計された「映画」プリセットおよび「音楽」プリセットのように、適用されるオーディオ改善アルゴリズムおよび/またはその最良のパラメータの集合をエンコードする。

10

【先行技術文献】

【特許文献】

【0011】

【特許文献1】国際公開第2008/106036号 (H. Muesch, "Speech Enhancement in Entertainment Audio")

20

【特許文献2】米国特許出願公開第2009/0097676A1号 (A. J. Seefeldt et al., "Calculating and Adjusting the Perceived Loudness and/or the Perceived Spectral Balance of an Audio Signal")

【特許文献3】国際公開第2007/127023号 (B.G. Grockett et al., "Audio Gain Control Using Specific-Loudness-Based Auditory Event Detection")

【特許文献4】国際公開第2009/011827号、(A. Seefeldt et al., "Audio Processing Using Auditory Scene Analysis and Spectral Skewness")

【非特許文献】

【0012】

【非特許文献1】L. Lu, H.-J. Zhang, and S. Li, "Content-based Audio Classification and Segmentation by Using Support Vector Machines", ACM Multimedia Systems Journal 8 (6), pp. 482-492, March, 2003

30

【非特許文献2】L. Lu, D. Liu, and H.-J. Zhang, "Automatic mood detection and tracking of music audio signals", IEEE Transactions on Audio, Speech, and Language Processing, 14(1):5 - 18, 2006

【非特許文献3】M.F McKinney and J. Breebaart, "Features for audio and music classification", Proc. ISMIR, 2003

【非特許文献4】G. H. Wakefield, "Mathematical representation of joint time Chroma distributions", SPIE, 1999

【非特許文献5】Ludovic Malfait, Jens Berger, and Martin Kastner, "P.563-The ITU-T Standard for Single-Ended Speech Quality Assessment", IEEE Transaction on Audio, Speech, and Language Processing, VOL. 14, NO. 6, November 2006

40

【発明の概要】

【発明が解決しようとする課題】

【0013】

しかしながら、手動選択はユーザーにとって不便である。ユーザーは通例、あらかじめ定義されたプリセットの間で頻繁に切り換えることはせず、単にすべてのコンテンツについて一つのプリセットを使い続ける。さらに、いくつかの自動ソリューションにおいてさえ、プリセットにおけるパラメータまたはアルゴリズム・セットアップは通例離散的で

50

あり（たとえば特定のコンテンツに関して特定のアルゴリズムについてオンまたはオフにするなど）、コンテンツ・ベースの連続的な仕方でパラメータを調整することはできない。

【課題を解決するための手段】

【0014】

本願の第一の側面は、再生時にオーディオ・コンテンツに基づいて連続的な仕方でオーディオ改善装置を自動的に構成設定することである。この「自動」モードでは、ユーザーは、わざわざ異なるプリセットを選ぶことなく、簡単にコンテンツを享受することができる。他方、連続的に調整することは、遷移点における耳に聞こえるアーチファクトを避けるために、より重要である。

10

【0015】

第一の側面のある実施形態によれば、オーディオ処理装置は、リアルタイムでオーディオ信号を少なくとも一つのオーディオ型に分類するオーディオ分類器と；聴衆の経験を改善するためのオーディオ改善装置と；前記少なくとも一つのオーディオ型の信頼値に基づいて連続的な仕方で前記オーディオ改善装置の少なくとも一つのパラメータを調整するための調整ユニットとを含む。

【0016】

オーディオ改善装置は、ダイアログ向上器、サラウンド仮想化器、ボリューム平準化器および等化器のうちの任意のものであってもよい。

【0017】

20

対応して、オーディオ処理方法は、リアルタイムでオーディオ信号を少なくとも一つのオーディオ型に分類し；前記少なくとも一つのオーディオ型の信頼値に基づいて連続的な仕方でオーディオ改善のための少なくとも一つのパラメータを調整することを含む。

【0018】

第一の側面のもう一つの実施形態によれば、ボリューム平準化器コントローラが、リアルタイムでオーディオ信号のコンテンツ型を識別するためのオーディオ・コンテンツ分類器と；識別されたコンテンツ型に基づいて連続的な仕方でボリューム平準化器を調整する調整ユニットとを含む。調整ユニットは、ボリューム平準化器の動的な利得を、オーディオ信号の情報性のコンテンツ型と正に相関させ、ボリューム平準化器の動的な利得を、オーディオ信号の干渉性のコンテンツ型と負に相関させるよう構成されていてもよい。

30

【0019】

上記のようなボリューム平準化器コントローラを有するオーディオ処理装置も開示される。

【0020】

対応して、ボリューム平準化器制御方法は、リアルタイムでオーディオ信号のコンテンツ型を識別し；識別されたコンテンツ型に基づいて連続的な仕方でボリューム平準化器を調整することを含む。該調整は、ボリューム平準化器の動的な利得を、オーディオ信号の情報性のコンテンツ型と正に相関させ、ボリューム平準化器の動的な利得を、オーディオ信号の干渉性のコンテンツ型と負に相関させることによる。

【0021】

40

第一の側面のさらにもう一つの実施形態によれば、等化器コントローラが、リアルタイムでオーディオ信号のオーディオ型を識別するためのオーディオ分類器と；識別されたオーディオ型の信頼値に基づいて連続的な仕方で等化器を調整する調整ユニットとを含む。

【0022】

上記のような等化器コントローラを有するオーディオ処理装置も開示される。

【0023】

対応して、等化器制御方法は、リアルタイムでオーディオ信号のオーディオ型を識別し；識別されたオーディオ型の信頼値に基づいて連続的な仕方で等化器を調整することを含む。

【0024】

50

本願は、コンピュータ・プログラム命令が記録されたコンピュータ可読媒体であって、前記命令は、プロセッサによって実行されると、前記プロセッサが上述したオーディオ処理方法または前記ボリューム平準化器制御方法または前記等化器制御方法を実行できるようにするものをも提供する。

【0025】

第一の側面の実施形態によれば、ダイアログ向上器、サラウンド仮想化器、ボリューム平準化器および等化器のうちの一つであってもよいオーディオ改善装置は、オーディオ信号の型および/または該型の信頼値に従って連続的に調整されてもよい。

【0026】

本願の第二の側面は、複数のオーディオ型を識別するコンテンツ識別コンポーネントを開発することである。検出結果は、連続的な仕方によりよいパラメータを見出すことにおいてさまざまなオーディオ改善装置の挙動を操縦/案内するために使われてもよい。

10

【0027】

第二の側面のある実施形態によれば、オーディオ分類器は：それぞれオーディオ・フレームのシーケンスを含む短期的オーディオ・セグメントから短期的特徴を抽出する短期的特徴抽出器と；長期的オーディオ・セグメント内の短期的セグメントのシーケンスをそれぞれの短期的特徴を使って諸短期的オーディオ型に分類する短期的分類器と；前記長期的オーディオ・セグメント内の短期的セグメントのシーケンスに関して短期的分類器の結果の統計量を長期的特徴として計算する統計抽出器と；前記長期的特徴を使って、前記長期的オーディオ・セグメントを長期的オーディオ型に分類する長期的分類器とを含む。

20

【0028】

上記のようなオーディオ分類器を有するオーディオ処理装置も開示される。

【0029】

対応して、オーディオ分類方法は：それぞれオーディオ・フレームのシーケンスを含む短期的オーディオ・セグメントから短期的特徴を抽出し；長期的オーディオ・セグメント内の短期的セグメントのシーケンスをそれぞれの短期的特徴を使って諸短期的オーディオ型に分類し；前記長期的オーディオ・セグメント内の短期的セグメントのシーケンスに関して分類処理の結果の統計量を長期的特徴として計算し；前記長期的特徴を使って、前記長期的オーディオ・セグメントを長期的オーディオ型に分類することを含む。

【0030】

第二の側面のもう一つの実施形態によれば、オーディオ分類器は、オーディオ信号の短期的セグメントのコンテンツ型を識別するオーディオ・コンテンツ分類器と；少なくとも部分的には前記オーディオ・コンテンツ分類器によって識別されたコンテンツ型に基づいて前記短期的セグメントのコンテキスト型を識別するオーディオ・コンテキスト分類器とを含む。

30

【0031】

上記のようなオーディオ分類器を有するオーディオ処理装置も開示される。

【0032】

対応して、オーディオ分類方法は、オーディオ信号の短期的セグメントのコンテンツ型を識別し；少なくとも部分的には識別されたコンテンツ型に基づいて前記短期的セグメントのコンテキスト型を識別することを含む。

40

【0033】

本願は、コンピュータ・プログラム命令が記録されたコンピュータ可読媒体であって、前記命令は、プロセッサによって実行されると、前記プロセッサが上述したオーディオ分類方法を実行できるようにするものをも提供する。

【0034】

第二の側面の実施形態によれば、オーディオ信号は、短期的な型またはコンテンツ型とは異なる、種々の長期的な型またはコンテキスト型に分類されてもよい。オーディオ信号の型および/または該型の信頼値は、ダイアログ向上器、サラウンド仮想化器、ボリューム平準化器または等化器のようなオーディオ改善装置を調整するためにさらに使われても

50

よい。

【図面の簡単な説明】

【0035】

本願は、限定ではなく例として、付属の図面の図において例解される。図面において同様の参照符号は同様の要素を指す。

【図1】本願のある実施形態に基づくオーディオ処理装置を示す図である。

【図2】図1に示した実施形態の変形を示す図である。

【図3】図1に示した実施形態の変形を示す図である。

【図4】複数のオーディオ型を識別し、信頼値を計算する分類器の可能な構成を示す図である。

10

【図5】複数のオーディオ型を識別し、信頼値を計算する分類器の可能な構成を示す図である。

【図6】複数のオーディオ型を識別し、信頼値を計算する分類器の可能な構成を示す図である。

【図7】本願のオーディオ処理装置のさらなる実施形態を示す図である。

【図8】本願のオーディオ処理装置のさらなる実施形態を示す図である。

【図9】本願のオーディオ処理装置のさらなる実施形態を示す図である。

【図10】種々のオーディオ型の間の変換の遅延を示す図である。

【図11】本願の実施形態に基づくオーディオ処理方法を示すフローチャートである。

【図12】本願の実施形態に基づくオーディオ処理方法を示すフローチャートである。

20

【図13】本願の実施形態に基づくオーディオ処理方法を示すフローチャートである。

【図14】本願の実施形態に基づくオーディオ処理方法を示すフローチャートである。

【図15】本願のある実施形態に基づくダイアログ向上器コントローラを示す図である。

【図16】ダイアログ向上器の制御において本願に基づくオーディオ処理方法の使用を示すフローチャートである。

【図17】ダイアログ向上器の制御において本願に基づくオーディオ処理方法の使用を示すフローチャートである。

【図18】本願のある実施形態に基づくサラウンド仮想化器コントローラを示す図である。

【図19】サラウンド仮想化器の制御において本願に基づくオーディオ処理方法の使用を示すフローチャートである。

30

【図20】本願のある実施形態に基づくボリューム平準化器コントローラを示す図である。

【図21】本願に基づくボリューム平準化器コントローラの効果を示す図である。

【図22】本願のある実施形態に基づく等化器コントローラを示す図である。

【図23】所望されるスペクトル・バランス・プリセットのいくつかの例を示す図である。

【図24】本願のある実施形態に基づくオーディオ分類器を示す図である。

【図25】本願のオーディオ分類器によって使用されるいくつかの特徴を示す図である。

【図26】本願のオーディオ分類器によって使用されるいくつかの特徴を示す図である。

40

【図27】本願に基づくオーディオ分類器のさらなる実施形態を示す図である。

【図28】本願に基づくオーディオ分類器のさらなる実施形態を示す図である。

【図29】本願に基づくオーディオ分類器のさらなる実施形態を示す図である。

【図30】本願の実施形態に基づくオーディオ分類方法を示すフローチャートである。

【図31】本願の実施形態に基づくオーディオ分類方法を示すフローチャートである。

【図32】本願の実施形態に基づくオーディオ分類方法を示すフローチャートである。

【図33】本願の実施形態に基づくオーディオ分類方法を示すフローチャートである。

【図34】本願のもう一つの実施形態に基づくオーディオ分類器を示す図である。

【図35】本願にさらにもう一つの実施形態に基づくオーディオ分類器を示す図である。

【図36】本願のオーディオ分類器において使われるヒューリスティック規則を示す図で

50

ある。

【図37】本願に基づくオーディオ分類器のさらなる実施形態を示す図である。

【図38】本願に基づくオーディオ分類器のさらなる実施形態を示す図である。

【図39】本願の実施形態に基づくオーディオ分類方法を示すフローチャートである。

【図40】本願の実施形態に基づくオーディオ分類方法を示すフローチャートである。

【図41】本願の実施形態を実装する例示的なシステムを示すブロック図である。

【発明を実施するための形態】

【0036】

以下では、本願の実施形態が図面を参照しつつ記述される。明確のために、当業者に知られているが本願を理解するために必要ではないコンポーネントおよびプロセスについての表現および記述は、図面および説明において省略されることを注意しておく。

10

【0037】

当業者は理解するであろうが、本願の諸側面は、システム、装置（たとえば携帯電話、ポータブル・メディア・プレーヤー、パーソナル・コンピュータ、パーサー、テレビジョン・セットトップボックスまたはデジタル・ビデオ・レコーダまたは他の任意のメディア・プレーヤー）、方法またはコンピュータ・プログラム・プロダクトとして具現されうる。よって、本願の諸側面は、ハードウェア実施形態、ソフトウェア実施形態（ファームウェア、常駐ソフトウェア、マイクロコードなどを含む）またはソフトウェアおよびハードウェア側面の両方を組み合わせる実施形態の形と取りうる。これらはみな本稿では「回路」、「モジュール」または「システム」と称されることがある。さらに、本願の諸側面は、コンピュータ可読プログラム・コードが具現された一つまたは複数の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・プロダクトの形を取ることがある。

20

【0038】

一つまたは複数のコンピュータ可読媒体のいかなる組み合わせが利用されてもよい。コンピュータ可読媒体は、コンピュータ可読信号媒体またはコンピュータ可読記憶媒体でありうる。コンピュータ可読記憶媒体は、たとえば、電子的、磁氣的、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以上のものの任意の好適な組み合わせでありうるがそれに限られない。コンピュータ可読記憶媒体のさらなる個別的な例（網羅的なりリストではない）は、次のものを含む：一つまたは複数のワイヤをもつ電気的接続、ポータブル・コンピュータ・ディスク、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・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以上のものの任意の好適な組み合わせ。本稿のコンテキストでは、コンピュータ可読記憶媒体は、命令実行システム、装置またはデバイスによって使うためまたはそれら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。

30

【0039】

コンピュータ可読信号媒体は、たとえばベースバンドにおいてまたは搬送波の一部としてコンピュータ可読プログラム・コードが具現されている伝搬するデータ信号を含みうる。そのような伝搬する信号は、電磁的または光学的な信号またはそれらの任意の好適な組み合わせを含むがそれに限られない多様な形の任意のものを取りうる。

40

【0040】

コンピュータ可読信号媒体は、命令実行システム、装置またはデバイスによって使うためまたはそれらとの関連で使うためのプログラムを通信する、伝搬させるまたは搬送することができる、コンピュータ可読記憶媒体ではないいかなるコンピュータ可読媒体であってもよい。

【0041】

コンピュータ可読媒体上に具現されるプログラム・コードは、無線、有線、光ファイバーケーブル、RFなどまたは以上のものの任意の好適な組み合わせを含むがそれに限られな

50

いいかなる適切な媒体を使って伝送されてもよい。

【 0 0 4 2 】

本願の諸側面の動作を実行するためのコンピュータ・プログラム・コードは、ジャバ、スモルトーク、C++などといったオブジェクト指向プログラミング言語および「C」プログラミング言語といった従来型の手続き型プログラミング言語または同様のプログラミング言語を含む、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で書かれてもよい。プログラム・コードは、完全にユーザーの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、部分的にユーザーの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る。この最後のシナリオでは、リモート・コンピュータはユーザーのコンピュータに、ローカル・エリア・ネットワーク（LAN）または広域ネットワーク（WAN）を含む任意の型のネットワークを通じて接続されてもよく、あるいは（たとえばインターネット・サービス・プロバイダーを使ってインターネットを通じて）外部コンピュータに接続がされてもよい。

10

【 0 0 4 3 】

本発明の諸側面は、本発明の実施形態に基づく方法、装置（システム）およびコンピュータ・プログラム・プロダクトのフローチャート図および/またはブロック図を参照して記述される。フローチャート図および/またはブロック図の各ブロックならびにフローチャート図および/またはブロック図のブロックの組み合わせは、コンピュータ・プログラム命令によって実装されることができるとは理解されるであろう。これらのコンピュータ・プログラム命令は、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与えられて、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る該命令が前記フローチャートおよび/またはブロック図の単数または複数のブロックにおいて特定されている機能/工程を実装する手段を作り出すよう、機械を生成してもよい。

20

【 0 0 4 4 】

これらのコンピュータ・プログラム命令は、コンピュータ、他のプログラム可能なデータ処理装置または他のデバイスが特定の仕方で機能するよう指令することができるコンピュータ可読媒体に記憶され、それにより、該コンピュータ可読媒体に記憶される命令は、前記フローチャートおよび/またはブロック図の単数または複数のブロックにおいて特定されている機能/工程を実装する命令を含む製造物を作り出してもよい。

30

【 0 0 4 5 】

コンピュータ・プログラム命令はコンピュータ、他のプログラム可能なデータ処理装置または他のデバイスにロードされて、該コンピュータ、他のプログラム可能な装置または他のデバイス上で一連の動作処理を実行させて、前記コンピュータまたは他のプログラム可能な装置上で実行される前記命令が前記フローチャートおよび/またはブロック図の単数または複数のブロックにおいて特定されている機能/工程を実装するためのプロセスを提供するようなコンピュータ実装されたプロセスを作り出してもよい。

【 0 0 4 6 】

下記では、本願の実施形態が詳細に記述される。明確のため、記述は次の構成に編成される：

40

第一部：オーディオ処理装置および方法

- 1 . 1 節 オーディオ型
- 1 . 2 節 オーディオ型の信頼値および分類器の構成
- 1 . 3 節 オーディオ型の信頼値の平滑化
- 1 . 4 節 パラメータ調整
- 1 . 5 節 パラメータ平滑化
- 1 . 6 節 オーディオ型の遷移
- 1 . 7 節 実施形態の組み合わせおよび応用シナリオ
- 1 . 8 節 オーディオ処理方法

50

| | |
|-----------------------------|----|
| 第二部：ダイアログ向上器コントローラおよび制御方法 | |
| 2.1節 ダイアログ向上のレベル | |
| 2.2節 向上させるべき周波数帯域の決定のための閾値 | |
| 2.3節 背景レベルへの調整 | |
| 2.4節 実施形態の組み合わせおよび応用シナリオ | |
| 2.5節 ダイアログ向上器制御方法 | |
| 第三部：サラウンド仮想化器コントローラおよび制御方法 | |
| 3.1節 サラウンド・ブースト量 | |
| 3.2節 開始周波数 | |
| 3.3節 実施形態の組み合わせおよび応用シナリオ | 10 |
| 3.4節 サラウンド仮想化器制御方法 | |
| 第四部：ボリューム平準化器コントローラおよび制御方法 | |
| 4.1節 情報性および干渉性のコンテンツ型 | |
| 4.2節 種々のコンテキストにおけるコンテンツ型 | |
| 4.3節 コンテキスト型 | |
| 4.4節 実施形態の組み合わせおよび応用シナリオ | |
| 4.5節 ボリューム平準化器制御方法 | |
| 第五部：等化器コントローラおよび制御方法 | |
| 5.1節 コンテンツ型に基づく制御 | |
| 5.2節 音楽における優勢な源の確からしさ | 20 |
| 5.3節 等化器プリセット | |
| 5.4節 コンテキスト型に基づく制御 | |
| 5.5節 実施形態の組み合わせおよび応用シナリオ | |
| 5.6節 等化器制御方法 | |
| 第六部：オーディオ分類器および分類方法 | |
| 6.1節 コンテンツ型分類に基づくコンテキスト分類器 | |
| 6.2節 長期的特徴の抽出 | |
| 6.3節 短期的特徴の抽出 | |
| 6.4節 実施形態の組み合わせおよび応用シナリオ | |
| 6.5節 オーディオ分類方法 | 30 |
| 第七部：VoIP分類器および分類方法 | |
| 7.1節 短期的セグメントに基づくコンテキスト分類 | |
| 7.2節 VoIP発話およびVoIPノイズを使った分類 | |
| 7.3節 平滑化ゆらぎ | |
| 7.4節 実施形態の組み合わせおよび応用シナリオ | |
| 7.5節 VoIP分類方法。 | |

【0047】

第一部：オーディオ処理装置および方法

図1は、再生時にオーディオ・コンテンツに基づく改善されたパラメータでの少なくとも一つのオーディオ改善装置の自動的な構成設定をサポートするコンテンツ適応的なオーディオ処理装置100の概括的なフレームワークを示している。これは三つの主要なコンポーネントを有する：オーディオ分類器200、調整ユニット300、オーディオ改善装置400である。

【0048】

オーディオ分類器200は、リアルタイムでオーディオ信号を少なくとも一つのオーディオ型に分類するものである。これは再生時にコンテンツのオーディオ型を自動的に識別する。オーディオ・コンテンツを識別するためには、信号処理、機械学習およびパターン認識を通じてなど、いかなるオーディオ分類技術が適用されることもできる。あらかじめ定義された目標オーディオ型の集合に関するオーディオ・コンテンツの確率を表わす信頼値がほぼ同時に推定される。

【 0 0 4 9 】

オーディオ改善装置 4 0 0 は、オーディオ信号に対して処理を実行することによって聴衆の経験を改善するものであり、のちに詳細に論じる。

【 0 0 5 0 】

調整ユニット 3 0 0 は、前記少なくとも一つのオーディオ型の信頼値に基づいて連続的な仕方で前記オーディオ改善装置の少なくとも一つのパラメータを調整するものである。これは、オーディオ改善装置 4 0 0 の挙動を操縦するよう設計される。これは、オーディオ分類器 2 0 0 から得られた結果に基づいて対応するオーディオ改善装置の最も好適なパラメータを推定する。

【 0 0 5 1 】

さまざまなオーディオ改善装置がこの装置において適用できる。図 2 は、ダイアログ向上器 (DE: Dialog Enhancer) 4 0 2、サラウンド仮想化器 (SV: Surround Virtualizer) 4 0 4、ポリウム準化器 (VL: Volume Leveler) 4 0 6 および等化器 (EQ: Equalizer) 4 0 8 を含む四つのオーディオ改善装置を含む例示的なシステムを示している。各オーディオ改善装置は、オーディオ分類器 2 0 0 において得られる結果 (オーディオ型および/または信頼値) に基づいて連続的な仕方で自動的に調整されることができる。

【 0 0 5 2 】

むしろ、オーディオ処理装置は、必ずしもすべての種類のオーディオ改善装置を含まなくてもよく、そのうち一つまたは複数を含むだけでもよい。他方、オーディオ改善装置は本開示において与えられている装置に限定されず、さらなる種類のオーディオ改善装置を含んでいてもよく、それらも本願の範囲内である。さらに、ダイアログ向上器 (DE) 4 0 2、サラウンド仮想化器 (SV) 4 0 4、ポリウム準化器 (VL) 4 0 6 および等化器 (EQ) 4 0 8 を含む本開示において論じられるオーディオ改善装置の名称は限定をなすものではなく、そのそれぞれは同じまたは同様の機能を実現する他の任意の装置をカバーすると解釈される。

【 0 0 5 3 】

1 . 1 節 オーディオ型

さまざまな種類のオーディオ改善装置を適正に制御するために、本願はさらに、オーディオ型の新たな構成を提供する。ただし、従来技術におけるオーディオ型も本願で適用可能である。

【 0 0 5 4 】

具体的には、オーディオ信号中の基本成分を表わす低レベルのオーディオ要素と、現実のユーザー娯楽アプリケーションにおけるたいいていの一般的なオーディオ・コンテンツを表わす高レベルのオーディオ・ジャンルとを含め、異なる意味的レベルからのオーディオ型がモデル化される。前者は「コンテンツ型」と称されてもよい。基本的オーディオ・コンテンツ型は、発話、音楽 (歌を含む)、背景音 (または効果音) およびノイズを含んでいてもよい。

【 0 0 5 5 】

発話および音楽の意味は明らかである。本願におけるノイズは、意味的なノイズではなく、物理的なノイズを意味する。本願における物理的なノイズは、たとえばエアコンからのノイズや、信号伝送経路に起因するピンク・ノイズのような技術的理由により生じるノイズを含みうる。対照的に、本願における「背景音」は、聴取者の注意のコア・ターゲットの周辺で生起する聴覚イベントであってもよい効果音である。たとえば、電話の通話におけるオーディオ信号では、話者の声のほかに、通話に関係ない何らかの他の人物の声、キーボードの音、足音などのような、意図されない何らかの他の音があることがある。これらの望まれない音は、ノイズではなく、「背景音」と称される。つまり、「背景音」は、ターゲット (または聴取者の注意のコア・ターゲット) ではない、あるいはさらに望まれないものであるが、それでも何らかの内容的な意味をもつ音と定義してもよい。一方、「ノイズ」は、ターゲットオンおよび背景音を除く望まれない音と定義されてもよい。

10

20

30

40

50

【0056】

時に、背景音は本当に「望まれない」のではなく、意図的に生成され、何らかの有用な情報を担う。たとえば、映画、テレビ番組またはラジオ放送番組における背景音がそうである。よって、時に、「効果音」と称されることがある。本開示では以下では、簡潔のため「背景音」のみが使用され、さらに「背景」と短縮されることもある。

【0057】

さらに、音楽はさらに、優勢な源のない音楽と優勢な源のある音楽に分類されてもよい。音楽片において他の源よりずっと強い源（声または楽器）がある場合には、「優勢な源のある音楽」と称される。そうでない場合には、「優勢な源のない音楽」と称される。たとえば、歌声およびさまざまな楽器を伴う多声音楽では、和声的にバランスが取れているまたはいくつかの最も顕著な源のエネルギーが互いに匹敵する場合には、優勢な源のない音楽と考えられる。対照的に、ある源（たとえば声）がずっと大きく、一方他の源がずっと静かである場合には、優勢な源を含んでいると考えられる。もう一つの例として、飛び抜けた、あるいは目立つ楽器トーンは「優勢な源をもつ音楽」である。

10

【0058】

音楽はさらに、種々の標準に基づく種々の型に分類されうる。音楽は、これに限られないがロック、ジャズ、ラップおよびフォークのような音楽のジャンルに基づいて分類されることができる。音楽は、声楽および器楽のように、楽器に基づいて分類されることもできる。器楽は、ピアノの音楽およびギター音楽など、種々の楽器を用いて演奏されるさまざまな音楽を含みうる。他の例示的な標準は、リズム、テンポ、音楽の音色および/または属性の類似性に基づいて音楽がグループ化されることのできる他の任意の音楽的な属性を含む。たとえば、音色に基づいて、声楽はテノール、バリトン、バス、ソプラノ、メゾソプラノおよびアルトに分類されうる。

20

【0059】

オーディオ信号のコンテンツ型は、複数のフレームから構成されるような短期的オーディオ・セグメントに関して分類されてもよい。一般に、オーディオ・フレームは20msのような複数ミリ秒の長さであり、オーディオ分類器によって分類されるべき短期的オーディオ・セグメントの長さは、数百ミリ秒から数秒、たとえば1秒の長さをもちうる。

【0060】

コンテンツ適応的な仕方でオーディオ改善装置を制御するために、オーディオ信号はリアルタイムで分類されてもよい。上記のコンテンツ型については、現在の短期的オーディオ・セグメントのコンテンツ型は現在のオーディオ信号のコンテンツ型を表わす。短期的オーディオ・セグメントの長さはそれほど長くないので、オーディオ信号は、順次の、重なり合わない短期的オーディオ・セグメントとして分割されてもよい。しかしながら、短期的オーディオ・セグメントは、オーディオ信号の時間軸に沿って連続的/半連続的にサンプリングされてもよい。すなわち、短期的オーディオ・セグメントは、オーディオ信号の時間軸に沿って一つまたは複数のフレームのステップ・サイズで動く所定の長さ（短期的オーディオ・セグメントの意図される長さ）をもつ窓を用いてサンプリングされてもよい。

30

【0061】

高レベルのオーディオ・ジャンルは、オーディオ信号の長期的な型を示すので「コンテキスト型」と称されることもあり、その時のサウンド・イベントの環境またはコンテキストと見なされてもよく、それは上記のようなコンテンツ型に分類されてもよい。本願によれば、コンテキスト型は、映画メディア、音楽（歌を含む）、ゲームおよびVoIP（インターネット・プロトコル上での音声）のようなたいていの一般的なオーディオ・アプリケーションを含みうる。

40

【0062】

音楽、ゲームおよびVoIPの意味は自明である。映画メディアは映画、テレビ番組、ラジオ放送番組または上記のものと同様の他の任意のオーディオ・メディアを含んでいてもよい。映画メディアの主要な特徴は、可能な発話、音楽およびさまざまな種類の背景音

50

(効果音)の混合である。

【0063】

コンテンツ型およびコンテキスト型はいずれも音楽(歌を含む)を含むことを注意してもよいであろう。以下、本願では、それらを区別するためにそれぞれ「短期的音楽」および「長期的音楽」という言い方を使う。

【0064】

本願のいくつかの実施形態については、他のいくつかのコンテキスト型構成も提案される。

【0065】

たとえば、オーディオ信号は高品質オーディオ(映画的メディアおよび音楽CDなど)または低品質オーディオ(VoIP、低ビットレート・オンライン・ストリーミング・オーディオおよびユーザー生成コンテンツなど)と分類されてもよい。これらはまとめて「オーディオ品質型」と称されてもよい。

10

【0066】

もう一つの例として、オーディオ信号はVoIPまたは非VoIPとして分類されてもよい。これらは上述した4コンテキスト型構成(VoIP、映画的メディア、(長期的)音楽およびゲーム)の変換と見なされてもよい。VoIPまたは非VoIPのコンテキストとの関連で、オーディオ信号は、VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズのようなVoIPに関係したオーディオ・コンテンツ型として分類されてもよい。VoIPオーディオ・コンテンツ型の構成は、VoIPおよび非VoIPコンテキストを区別するために特に有用である。通例、VoIPコンテキストは、ボリューム平準化器(オーディオ改善装置の一つの種類)の最も困難な応用シナリオだからである。

20

【0067】

一般に、オーディオ信号のコンテキスト型は、短期的オーディオ・セグメントより長い長期的オーディオ・セグメントに関して分類されてもよい。長期的オーディオ・セグメントは、短期的オーディオ・セグメントにおけるフレーム数より多くの数の複数のフレームから構成される。長期的オーディオ・セグメントは、複数の短期的オーディオ・セグメントから構成されてもよい。一般に、長期的オーディオ・セグメントは、数秒から数十秒、たとえば10秒などの秒のオーダーの長さをもちうる。

【0068】

同様に、適応的な仕方でオーディオ改善装置を制御するために、オーディオ信号は、リアルタイムでコンテキスト型に分類されてもよい。同様に、現在の長期的オーディオ・セグメントのコンテキスト型は、現在のオーディオ信号のコンテキスト型を表わす。長期的オーディオ・セグメントの長さは比較的長いので、オーディオ信号は、そのコンテキスト型の急激な変化を、よってオーディオ改善装置(単数または複数)の作動パラメータの急激な変化を避けるために、オーディオ信号の時間軸に沿って連続的/半連続的にサンプリングされてもよい。すなわち、長期的オーディオ・セグメントは、オーディオ信号の時間軸に沿って一つまたは複数のフレームまたは一つまたは複数の短期的セグメントのステップ・サイズで動く所定の長さ(長期的オーディオ・セグメントの意図される長さ)をもつ窓を用いてサンプリングされてもよい。

30

40

【0069】

上記では、コンテンツ型およびコンテキスト型の両方が記述された。本願の実施形態では、調整ユニット300は、さまざまなコンテンツ型の少なくとも一つおよび/またはさまざまなコンテキスト型の少なくとも一つに基づいて、オーディオ改善装置(単数または複数)の少なくとも一つのパラメータを調整してもよい。したがって、図3に示されるように、図1に示した実施形態のある変形では、オーディオ分類器200は、オーディオ・コンテンツ分類器202またはオーディオ・コンテキスト分類器204またはその両方を有してもよい。

【0070】

上記では、(たとえばコンテキスト型について)種々の標準に基づく種々のオーディオ

50

型および（たとえばコンテンツ型について）種々の階層レベルでの種々のオーディオ型が言及された。しかしながら、これらの標準および階層レベルは単にここでの記述の便宜のためであって、全く限定するものではない。つまり、本願では、上述したオーディオ型の任意の二つ以上は、同時にオーディオ分類器 200 によって識別され、同時に調整ユニット 300 によって考慮されることができる。これについては後述する。つまり、種々の階層レベルにおけるすべてのオーディオ型は並列、あるいは同じレベルであってもよい。

【0071】

1.2節 オーディオ型の信頼値および分類器の構成

オーディオ分類器 200 は、硬判定結果を出力してもよく、あるいは調整ユニット 300 はオーディオ分類器 200 の結果を硬判定結果と見なしてもよい。硬判定についてでも、複数のオーディオ型がオーディオ・セグメントに割り当てられることができる。たとえば、オーディオ・セグメントは、発話および短期的音楽の混合信号でありうるので、「発話」および「短期的音楽」の両方によってラベル付けされることができる。得られたラベルは、オーディオ改善装置（単数または複数）400 を操縦するために直接使われることができる。簡単な例は、発話が存在するときにダイアログ向上器 402 を有効にし、発話が存在しないときにオフにするというものである。しかしながら、この硬判定方法は、注意深い平滑化方式（後述）なしの場合には、あるオーディオ型から別のオーディオ型への遷移点においていくらかの不自然さを導入することができる。

10

【0072】

より柔軟性をもち、連続的な仕方でオーディオ改善装置のパラメータを調整するために、各ターゲット・オーディオ型の信頼値が推定されることができる（軟判定）。信頼値は、識別されるべきオーディオ・コンテンツとターゲット・オーディオ型の間の一致レベルを 0 から 1 の値で表わす。

20

【0073】

先述したように、多くの分類技法は直接、信頼値を出力してもよい。信頼値は、分類器の一部と見なされてもよいさまざまな方法から計算されることもできる。たとえば、オーディオ・モデルがガウシアン混合モデル（GMM: Gaussian Mixture Models）のようないくつかの確率的モデル化技術によってトレーニングされる場合、信頼値を表わすために、次のように事後確率が使われることができる。

30

【0074】

【数 1】

$$p(c_i | x) = \frac{p(x | c_i)}{\sum_{i=1}^N p(x | c_i)} \quad (1)$$

ここで、 x はオーディオ・セグメントの一片であり、 c_i はターゲット・オーディオ型であり、 N はターゲット・オーディオ型の数であり、 $p(c_i)$ はオーディオ・セグメント x がオーディオ型 c_i である確からしさであり、 $p(c_i | x)$ は対応する事後確率である。

40

【0075】

他方、オーディオ・モデルがサポートベクターマシン（SVM: Support Vector Machine）およびアダブースト（AdaBoost）のようないくつかの弁別的方法からトレーニングされる場合には、モデル比較からはスコア（実数値）だけが得られる。これらの場合、得られたスコア（理論的には - から ）を期待される信頼度（0 から 1）にマッピングするために、通例、次のシグモイド関数が使われる。

【0076】

【数 2】

$$conf = \frac{1}{1 + e^{Ay+B}} \quad (2)$$

ここで、 y はSVMまたはアダブーストからの出力スコアであり、 A および B は何らかのよく知られた技術を使ってトレーニング・データ・セットから推定される必要のある二つのパラメータである。

【0077】

10

本願のいくつかの実施形態については、調整ユニット300は、三つ以上のコンテンツ型および/または三つ以上のコンテキスト型を使ってもよい。その場合、オーディオ・コンテンツ分類器202は三つ以上のコンテンツ型を識別する必要がある、および/またはオーディオ・コンテキスト分類器204は三つ以上のコンテキスト型を識別する必要がある。そのような状況では、オーディオ・コンテンツ分類器202またはオーディオ・コンテキスト分類器204は、ある構成で編成された分類器の群であってもよい。

【0078】

たとえば、調整ユニット300が四種類のコンテキスト型、映画のメディア、長期的音楽、ゲームおよびVoIPの全部を必要とする場合には、オーディオ・コンテキスト分類器204は以下の種々の構成をもちうる。

20

【0079】

第一に、オーディオ・コンテキスト分類器204は、図4に示されるように編成された6個の対一二項分類器（各分類器は一つのターゲット・オーディオ型を別のあるターゲット・オーディオ型から弁別する）と、図5に示されるように編成された3個の対他の二項分類器（各分類器はターゲット・オーディオ型を他のオーディオ型から弁別する）と、図6に示されるように編成された4個の対他の分類器を有していてもよい。判定有向非環状グラフ（DDAG: Decision Directed Acyclic Graph）構成のような他の構成もある。図4～図6および対応する以下の記述において、簡潔のため、「映画のメディア」の代わりに「映画」が使われていることを注意しておく。

【0080】

30

各二項分類器は、その出力について信頼スコア $H(x)$ を与える（ x はオーディオ・セグメントを表わす）。各二項分類器の出力が得られたのち、それらを、識別された諸コンテキスト型の最終的な諸信頼値にマッピングする必要がある。

【0081】

一般に、オーディオ信号は M 個のコンテキスト型に分類されるとする（ M は正の整数）。通常の一対一構成は $M(M-1)/2$ 個の分類器を構築する。ここで、各分類器は二つのクラスからのデータでトレーニングされる。次いで、それぞれの対一分類器はその好ましいクラスについて一票を投じ、最終結果は、 $M(M-1)/2$ 個の分類器の分類のうちの最多票をもつクラスである。通常の一対一構成と比べ、図4における階層的な構成も $M(M-1)/2$ 個の分類器を構築することを必要とする。しかしながら、セグメント x は各階層レベルにおいて対応するクラスにある/ないと判定され、全体的なレベル・カウントは $M-1$ なので、試験反復工程は $M-1$ に短縮されることができ。さまざまなコンテキスト型についての最終的な信頼値は、たとえば二項分類信頼値 $H_k(x)$ から計算されてもよい（ $k=1,2,\dots,6$ は種々のコンテキスト型を表わす）。

40

【0082】

【数 3】

$$\begin{aligned}
C_{MOVIE} &= (1 - H_1(x)) \cdot (1 - H_3(x)) \cdot (1 - H_6(x)) \\
C_{VOIP} &= H_1(x) \cdot H_2(x) \cdot H_4(x) \\
C_{MUSIC} &= H_1(x) \cdot (1 - H_2(x)) \cdot (1 - H_5(x)) + H_3(x) \cdot (1 - H_1(x)) \cdot (1 - H_5(x)) \\
&\quad + H_6(x) \cdot (1 - H_1(x)) \cdot (1 - H_3(x)) \\
C_{GAME} &= H_1(x) \cdot H_2(x) \cdot (1 - H_4(x)) + H_1(x) \cdot H_5(x) \cdot (1 - H_2(x)) + H_3(x) \cdot H_5(x) \\
&\quad \cdot (1 - H_1(x))
\end{aligned}
\tag{10}$$

図 5 に示した構成では、二項分類結果 $H_k(x)$ から最終的な信頼値へのマッピング関数は次の例のように定義できる。

【0083】

【数 4】

$$\begin{aligned}
C_{MOVIE} &= H_1(x) \\
C_{MUSIC} &= H_2(x) \cdot (1 - H_1(x)) \\
C_{VOIP} &= H_3(x) \cdot (1 - H_2(x)) \cdot (1 - H_1(x)) \\
C_{GAME} &= (1 - H_3(x)) \cdot (1 - H_2(x)) \cdot (1 - H_1(x))
\end{aligned}
\tag{20}$$

図 6 に示した構成では、最終的な信頼値は、対応する二項分類結果 $H_k(x)$ に等しくてもよく、あるいはすべてのクラスについての信頼値の和が1であることが要求されるならば、最終的な信頼値は単に推定された $H_k(x)$ に基づいて規格化されることができる。

【0084】

【数 5】

$$\begin{aligned}
C_{MOVIE} &= H_1(x) / (H_1(x) + H_2(x) + H_3(x) + H_4(x)) \\
C_{MUSIC} &= H_2(x) / (H_1(x) + H_2(x) + H_3(x) + H_4(x)) \\
C_{VOIP} &= H_3(x) / (H_1(x) + H_2(x) + H_3(x) + H_4(x)) \\
C_{GAME} &= H_4(x) / (H_1(x) + H_2(x) + H_3(x) + H_4(x))
\end{aligned}$$

最大の信頼値をもつ一つまたは複数が最終的な識別されたクラスとして決定されることができる。

【0085】

図 4 ~ 図 6 に示される構成では、種々の二項分類器のシーケンスは必ずしも図示したとおりではなく、他のシーケンスであってもよいことを注意しておくべきである。そうしたシーケンスは、さまざまな応用の種々の要求に応じて、手動の割り当てまたは自動学習によって選択されてもよい。

【0086】

上記の記述は、オーディオ・コンテキスト分類器 204 に向けられている。オーディオ・コンテンツ分類器 202 については、状況は同様である。

【0087】

あるいはまた、オーディオ・コンテンツ分類器 202 またはオーディオ・コンテキスト分類器 204 は、同時にすべてのコンテンツ型 / コンテキスト型を識別する一つの単独の

40

50

分類器として実装され、同時に対応する諸信頼値を与えてもよい。これを行なうために多くの既存の技法がある。

【0088】

信頼値を使って、オーディオ分類器200の出力はベクトルとして表現できる。ベクトルの各次元が各ターゲット・オーディオ型の信頼値を表わす。たとえば、ターゲット・オーディオ型が順次（発話、短期的音楽、ノイズ、背景）であれば、例示的な出力結果は(0.9,0.5,0.0,0.0)であることができる。これは、そのオーディオ・コンテンツが発話であることは90%確実であり、そのオーディオが音楽であることは50%確実であることを示す。出力ベクトル中のすべての次元の和が1である必要はないことを注意しておく（たとえば、図6からの結果は必ずしも規格化されない）。つまり、オーディオ信号は発話および短期的音楽の混合信号であってもよい。

10

【0089】

のちに、第六部および第七部において、オーディオ・コンテキスト分類およびオーディオ・コンテンツ分類の新規な実装を詳細に論じる。

【0090】

1.3節 オーディオ型の信頼値の平滑化

任意的に、各オーディオ・セグメントがあらかじめ定義されたオーディオ型に分類された後、追加的なステップは、ある型から別の型への急激なジャンプを避け、オーディオ改善装置におけるパラメータのよりなめらかな推定をするために、時間軸に沿って分類結果を平滑化するというものである。たとえば、長い抜粋が、VoIPとして分類された一つだけのセグメントを除いて映画のメディアと分類されると、急激なVoIP判定は、平滑化によって映画のメディアに修正されることができる。

20

【0091】

したがって、図7に示した実施形態のある変形では、各オーディオ型について、現在の時間でのオーディオ信号の信頼値を平滑化するために、型平滑化ユニット712がさらに設けられる。

【0092】

一般的な平滑化方法は、現在の実際の信頼値と最後の時点の平滑化された信頼値の重み付けされた和を計算するなど、重み付けされた平均に基づく。

【0093】

$$\text{smoothConf}(t) = \alpha \cdot \text{smoothConf}(t-1) + (1-\alpha) \cdot \text{conf}(t) \quad (3)$$

30

ここで、 t は現在の時間（現在のオーディオ・セグメント）、 $t-1$ は最後の時間（最後のオーディオ・セグメント）、 α は重み、 conf および smoothConf はそれぞれ平滑化の前および後の信頼値（confidence value）である。

【0094】

信頼値の観点からは、分類器の硬決定からの結果は、値が0か1のいずれかである信頼値を用いて表わされることもできる。すなわち、ターゲット・オーディオ型が選ばれ、オーディオ・セグメントに割り当てられる場合、対応する信頼値は1であり、そうでなければ信頼値は0である。したがって、たとえオーディオ分類器200が信頼値を与えず、オーディオ型に関する硬決定を与えるだけであっても、調整ユニット300の連続的な調整は、型平滑化ユニット712の平滑化動作を通じて、相変わらず可能である。

40

【0095】

平滑化アルゴリズムは、異なる場合について異なる平滑化重みを使うことによって「非対称」であることができる。たとえば、重み付けされた和を計算するための重みは、オーディオ信号のオーディオ型の信頼値に基づいて適応的に変えられてもよい。現在のセグメントの信頼値がより大きければ、その重みもより大きくなる。

【0096】

別の観点からは、重み付けされた和を計算するための重みは、特に、オーディオ改善装置（単数または複数）が、一つの単独のコンテンツ型の有無に基づくのではなく、オーディオ分類器200にいおって識別される複数のコンテンツ型に基づいて調整されるときは

50

、あるオーディオ型から別のオーディオ型への異なる遷移対に基づいて適応的に変更されてもよい。たとえば、あるコンテキストにおいてより頻繁に現われるオーディオ型からそのコンテキストにおいてそれほど頻繁に現われない別のオーディオ型への遷移については、後者の信頼値は、あまり速く増大しないよう、平滑化されてもよい。たまたまの突発であることもありうるからである。

【0097】

もう一つの要因は、変化レートを含む変化（増大または減少）トレンドである。あるオーディオ型が存在するようになるときに（すなわちその信頼値が増すときに）我々が遅延についてより気にするとすると、次のようにして平滑化アルゴリズムを設計できる。

【0098】

【数6】

$$\text{smoothConf}(t) = \begin{cases} \text{conf}(t) & \text{conf}(t) \geq \text{smoothConf}(t-1) \\ \beta \cdot \text{smoothConf}(t-1) + (1-\beta) \cdot \text{conf}(t) & \text{それ以外の場合} \end{cases} \quad (4)$$

上記の公式は、平滑化された信頼値が、信頼値が増加するときには迅速に現在状態に回答し、信頼値が減少するときにはゆっくりと平滑化されていくことを許容する。平滑化関数の諸変形が同様にして簡単に設計できる。たとえば、公式(4)は、 $\text{conf}(t) \geq \text{smoothConf}(t-1)$ であるときに $\text{conf}(t)$ の重みがより大きくなるよう修正されてもよい。実際、公式(4)では、 $\beta = 0$ であり、 $\text{conf}(t)$ の重みは最大、すなわち1であると見なすことができる。

【0099】

異なる観点からは、あるオーディオ型の変化トレンドを考えることは、オーディオ型の異なる遷移対を考えることの単なる個別的な例である。たとえば、型Aの信頼値を増すことは、非AからAへの遷移と見なされてもよく、型Aの信頼値を減らすことは、Aから非Aへの遷移と見なされてもよい。

【0100】

1.4節 パラメータ調整

調整ユニット300は、オーディオ分類器200からの得られた結果に基づいてオーディオ改善装置（単数または複数）400についての適正なパラメータを推定または調整するよう設計される。コンテンツ型またはコンテキスト型のいずれかをまたは合同判定のために両方を使って、異なるオーディオ改善装置について異なる調整アルゴリズムが設計されてもよい。たとえば、映画のメディアおよび長期的音楽のようなコンテキスト型情報では、上述したようなプリセットが自動的に選択され、対応するコンテンツに適用されることができる。利用可能なコンテンツ型情報を用いて、各オーディオ改善装置のパラメータは、後述する部において示されるように、より細かい仕方で調整されることができる。コンテンツ型情報およびコンテキスト情報はさらに、長期的および短期的情報をバランスさせるために、調整ユニット300において合同で使われることができる。特定のオーディオ改善装置についての特定の調整アルゴリズムは、別個の調整ユニットと見なされてもよい。あるいは、異なる調整アルゴリズムがまとめて連合した調整ユニットと見なされてもよい。

【0101】

すなわち、調整ユニット300は、少なくとも一つのコンテンツ型の信頼値および/または少なくとも一つのコンテキスト型の信頼値に基づいて、オーディオ改善装置の前記少なくとも一つのパラメータを調整するよう構成されていてもよい。特定のオーディオ改善装置について、オーディオ型のいくつかは情報性であり、オーディオ型のいくつかは干渉性である。よって、特定のオーディオ改善装置のパラメータは、情報性のオーディオ型（単数または複数）または干渉性のオーディオ型（単数または複数）の信頼値（単数または複数）と正または負に相関しうる。ここで「正に相関する」とは、オーディオ型の信頼値の増大または減少とともに、パラメータが線形または非線型に増大または減少することを

10

20

30

40

50

意味する。「負に相関する」とは、オーディオ型の信頼値の減少または増大とともに、パラメータが線形または非線型にそれぞれ増大または減少することを意味する。

【0102】

ここで、信頼値の減少および増大は、正または負の相関によって、調整されるべきパラメータに直接「伝達」される。数学では、そのような相関または「伝達」は、正比例または反比例、プラスまたはマイナス（加算または減算）演算、乗算または除算または非線型関数として具現されうる。これらすべての形の相関は「伝達関数」と称されてもよい。信頼値の増大または減少を決定するために、現在の信頼値またはその数学的変換を最後の信頼値もしくは複数の履歴信頼値またはその数学的変換と比較することもできる。本願のコンテキストにおいて、用語「比較」は、減算を通じた比較または除算を通じた比較を意味する。差が0より大きいか否かまたは比が1より大きいか否かを判定することによって増大または減少を判定できる。

10

【0103】

個別的な実装では、適正なアルゴリズム（伝達関数など）を通じてパラメータを信頼値またはその比もしくは差と直接関係させることができ、「外部観察者」が特定の信頼値および/または特定のパラメータが増大したまたは減少したことを明示的に知ることは必要でない。いくつかの個別的な例が、個別的なオーディオ改善装置についての後述する第二～五部において与えられる。

【0104】

前節で述べたように、同じオーディオ・セグメントに関し、分類器200は、それぞれの信頼値をもつ複数のオーディオ型を識別することがある。オーディオ・セグメントは音楽と発話と背景音など、同時に複数の成分を含むことがあるので、それらの信頼値は必ずしも総計1にならないことがある。そのような状況では、オーディオ改善装置のパラメータは、異なるオーディオ型の間でバランスを取る必要がある。たとえば、調整ユニット300は、前記少なくとも一つのオーディオ型の重要性に基づいて前記少なくとも一つのオーディオ型の信頼値を重み付けすることを通じて、複数のオーディオ型の少なくともいくつかを考慮するよう構成されていてもよい。特定のオーディオ型がより重要であるほど、パラメータはそれによってより大きく影響される。

20

【0105】

重みは、オーディオ型の情報性および干渉性の効果を反映することもできる。たとえば、干渉性のオーディオ型については、マイナスの重みが与えられてもよい。いくつかの個別的な例が、個別的なオーディオ改善装置についての後述する第二～五部において与えられる。

30

【0106】

本願のコンテキストにおいて、「重み」は、多項式における係数よりも広い意味をもつことを注意しておく。多項式における係数のほか、「重み」は指数または冪の形を取ることでもできる。多項式における係数であるとき、重み付け係数は規格化されてもされなくてもよい。一言で言うと、重みは単に重み付けされた対象が調整されるべきパラメータに対してどのくらいの影響をもつかを表わすものである。

【0107】

いくつかの他の実施形態では、同じオーディオ・セグメントに含まれる複数のオーディオ型について、その信頼値が、規格化されることを通じて重みに変換されてもよい。次いで、最終的なパラメータが、各オーディオ型についてあらかじめ定義され、信頼値に基づいて重みによって重み付けされたパラメータ・プリセット値の和を計算することを通じて決定されてもよい。すなわち、調整ユニット300は、信頼値に基づいて複数のオーディオ型の効果を重み付けすることを通じて複数のオーディオ型を考慮するよう構成されていてもよい。

40

【0108】

重み付けの個別的な例として、調整ユニットは、信頼値に基づいて少なくとも一つの優勢なオーディオ型を考慮するよう構成される。あまりに低い（閾値より低い）信頼値をも

50

つオーディオ型については、考慮されなくてもよい。これは、信頼値が閾値より小さい他のオーディオ型の重みが0に設定されることと等価である。いくつかの個別的な例が、個別的なオーディオ改善装置についての後述する第二～五部において与えられる。

【0109】

コンテンツ型およびコンテキスト型と一緒に考慮されることができる。ある実施形態では、それらは同じレベルにあると見なされることができ、それらの信頼値はそれぞれの重みをもちうる。もう一つの実施形態では、まさに名称が示すように、「コンテキスト型」は、「コンテキスト型」が位置しているコンテキストまたは環境であり、よって調整ユニット200は、異なるコンテキスト型のオーディオ信号におけるコンテンツ型が、オーディオ信号のコンテキスト型に依存して異なる重みを割り当てられるよう構成されてもよい。一般に、いかなるオーディオ型も、別のオーディオ型のコンテキストを構成することができ、調整ユニット200はあるオーディオ型の重みを、別のオーディオ型の信頼値を用いて修正するよう構成されてもよい。いくつかの個別的な例が、個別的なオーディオ改善装置についての後述する第二～五部において与えられる。

10

【0110】

本願のコンテキストでは、「パラメータ」は、その文字通りの意味より広い意味をもつ。単一の値をもつパラメータのほかに、パラメータは、種々のパラメータの集合、種々のパラメータからなるベクトルまたはプロファイルを含め、上述したようなプリセットを意味することもある。特に、後述する第二～五部においては、次のパラメータが論じられるが、本願はそれに限定されるものではない：ダイアログ向上のレベル、ダイアログ向上されるべき周波数帯域を決定するための閾値、背景レベル、サラウンド・ブースト量、サラウンド仮想化器についての開始周波数、ボリューム平準化器の動的利得または動的利得の範囲〔ダイナミック・ゲインのレンジ〕、オーディオ信号が新しい知覚可能なオーディオ・イベントである度合いを示すパラメータ、等化レベル、等化プロファイルおよびスペクトル・バランス・プリセット。

20

【0111】

1.5節 パラメータ平滑化

1.3節では、急激な変化を避け、よってオーディオ改善装置のパラメータの急激な変化を避けるためにオーディオ型の信頼値を平滑化することを論じた。他の措置も可能である。一つは、オーディオ型に基づいて調整されるパラメータを平滑化することであり、本節で論じる。他方は、オーディオ分類器および/または調整ユニットを、オーディオ分類器の結果の変化を遅らせるよう構成することであり、これについては1.6節で論じる。

30

【0112】

ある実施形態では、パラメータは、遷移点における可聴アーチファクトを導入しうる素早い変化を避けるために、次のように、さらに平滑化されることができる。

【0113】

【数7】

$$\tilde{L}(t) = \tau \tilde{L}(t-1) + (1-\tau)L(t) \quad (3')$$

40

ここで、チルダ付きの $L(t)$ は平滑化されたパラメータ、 $L(t)$ は平滑化されていないパラメータ、 τ は時定数を表わす係数、 t は現在の時間、 $t-1$ は最後の時間である。

【0114】

すなわち、図8に示されるように、オーディオ処理装置は、パラメータ平滑化ユニット814を有していてもよい。これは、調整ユニット300によって調整される（ダイアログ向上器402、サラウンド仮想化器404、ボリューム平準化器406および等化器408のうち少なくとも一つのような）オーディオ改善装置のパラメータについて、現在の時間において調整ユニットによって決定されたパラメータ値および最後の時間の平滑化されたパラメータ値の重み付けされた和を計算することによって、現在の時間における調

50

整ユニット 300 によって決定されるパラメータ値を平滑化する。

【0115】

時定数は、応用の個別的な要求および/またはオーディオ改善装置 400 の実装に基づく固定値であることができる。時定数は、オーディオ型に基づいて、特に、音楽から発話、発話から音楽など、あるオーディオ型から別のオーディオ型への種々の遷移型に基づいて適応的に変更されてもよい。

【0116】

等化器を例に取る（さらなる詳細は第五部で言及されうる）。等化は、音楽コンテンツに適用するには良好だが、発話コンテンツに適用するのはよくない。よって、等化のレベルを平滑化するためには、オーディオ信号が音楽から発話に遷移するときは時定数は比較的小さくてもよく、それにより発話コンテンツに対してより小さな等化レベルがより迅速に適用されることができる。他方、発話から音楽への遷移についての時定数は、遷移点における可聴アーチファクトを避けるために、比較的大きいことができる。

10

【0117】

遷移型（たとえば、発話から音楽または音楽から発話）を推定するために、コンテンツ分類結果は直接使われることができる。すなわち、オーディオ・コンテンツを音楽または発話に分類するれば、遷移型を得ることは単純になる。より連続的な仕方で遷移を推定できるよう、オーディオ型の硬決定を直接比較する代わりに、推定された平滑化されない等化レベルにも頼ることができる。一般的な発想は、平滑化されない等化レベルが増大する場合には、それは発話から音楽への（またはより音楽的への）遷移を示し、そうでない場合にはそれは音楽から発話への（またはより発話的への）遷移により近い。異なる遷移型を区別することにより、時定数は対応して設定されることができる。一例は次のようなものである。

20

【0118】

【数 8】

$$\tau(t) = \begin{cases} \tau_1 & L(t) \geq L(t-1) \\ \tau_2 & L(t) < L(t-1) \end{cases} \quad (4')$$

30

ここで、 $\tau(t)$ は、コンテンツに依存する時間変化する時定数であり、 τ_1 および τ_2 は二つのプリセット時定数値であり、通例 $\tau_1 > \tau_2$ を満たす。直観的には、上記の関数は、等化レベルが増加するときには比較的遅い遷移を示し、等化レベルが減少するときには比較的速い遷移を示す。だが本願はこれに限定されるものではない。さらに、パラメータは等化レベルに限定されず、他のパラメータであってもよい。すなわち、パラメータ平滑化ユニット 814 は、重み付けされた和を計算するための重みが、調整ユニット 300 によって決定されるパラメータ値の増加トレンドまたは減少トレンドに基づいて適応的に変えられるよう構成されてもよい。

【0119】

40

1.6 節 オーディオ型の遷移

図 9 および図 10 を参照して、オーディオ型の急激な変化を避け、よってオーディオ改善装置のパラメータの急激な変化を避けるためのもう一つの方式が記述される。

【0120】

図 9 に示されるように、オーディオ処理装置 100 はさらに、オーディオ分類器 200 が連続的に同じ新しいオーディオ型を出力する持続時間を測定するためのタイマー 916 を有していてもよい。調整ユニット 300 は、新しいオーディオ型の持続時間の長さが閾値に達するまで、現在のオーディオ型を使い続けるよう構成されてもよい。

【0121】

換言すれば、図 10 に示されるように、観察（または維持）フェーズが導入される。調

50

整ユニット300が実際に新しいオーディオ型を使う前に、(持続時間の長さの閾値に対応する)観察フェーズにおいて、オーディオ型が本当に変化したのかどうかを確認するために、ある連続量の時間にわたってオーディオ型の変化がさらにモニタリングされる。

【0122】

図10に示されるように、矢印(1)は、現在状態が型Aであり、オーディオ分類器200の結果が変わらない状況を示す。

【0123】

現在状態が型Aであり、オーディオ分類器200の結果が型Bになる場合、タイマー916は計時を開始する、あるいは図10に示されるように、プロセスは観察フェーズにはいり(矢印(2))、残存(hangover)カウントcntの初期値が設定される。これは観察継続時間の長さ(前記閾値に等しい)を示す。

10

【0124】

次いで、オーディオ分類器200が連続的に型Bを出力し、cntは連続的に減少して(矢印(3))、しまいにはcntは0に等しくなる(すなわち、新しい型Bの持続時間の長さが閾値に達する)場合、調整ユニット300は新しいオーディオ型Bを使用しうる(矢印(4))。あるいは、換言すれば、この時点になってはじめて、オーディオ型は本当に型Bに変わったと見なされうる。

【0125】

そうでなく、cntが0になる前に(持続時間の長さが閾値に達する前に)オーディオ分類器200の出力がもとの型Aに戻る場合には、観察フェーズは打ち切られ、調整ユニット300は相変わらずもとの型Aを使う(矢印(5))。

20

【0126】

型Bから型Aへの変化は、上記のプロセスと同様であってもよい。

【0127】

上記のプロセスでは、閾値(または残存カウント)は用途の要件に基づいて設定される。これはあらかじめ定義された固定値であってもよい。これは適応的に設定されてもよい。ある変形では、閾値は、あるオーディオ型から別のオーディオ型への異なる遷移対については異なる。たとえば、型Aから型Bが変わるとき、閾値は第一の値であってもよく、型Bから型Aが変わるとき、閾値は第二の値であってもよい。

【0128】

もう一つの変形では、残存カウント(閾値)は、新しいオーディオ型の信頼値と負に相関していてもよい。一般的な発想は、信頼値が二つの型の間の混乱を示す場合には(たとえば、信頼値が約0.5しかないときは)、観察継続時間は長い必要がある。そうでない場合には、継続時間は比較的短くてもよい。このガイドラインに従い、例示的な残存カウント(hangover count)は、次の公式によって設定されることができる。

30

【0129】

$$\text{HangCnt} = C \cdot |0.5 - \text{Conf}| + D$$

ここで、HangCntは残存継続時間または閾値であり、CおよびDは用途の要求に基づいて設定されることのできる二つのパラメータであり、通例、Cは負、Dは正の値である。

【0130】

なお、タイマー916(よって上記の遷移プロセス)はオーディオ処理装置の一部だがオーディオ分類器200の外部として記述した。他のいくつかの実施形態では、まさに7.3節で述べるように、オーディオ分類器200の一部と見なされてもよい。

40

【0131】

1.7節 実施形態の組み合わせおよび応用シナリオ

上記で論じたすべての実施形態およびその変形は、そのいかなる組み合わせにおいて実装されてもよく、異なる部/実施形態において言及されるが同じまたは同様の機能をもついかなる構成要素も同じまたは別個の構成要素として実装されてもよい。

【0132】

具体的には、以上において実施形態およびその変形を記述するとき、前の実施形態また

50

は変形においてすでに記述されたものと同様の参照符号をもつ構成要素は省略され、異なる構成よそが記述されるだけである。実際、これらの異なる構成要素は、他の実施形態または変形の構成要素と組み合わせられたり、あるいは単独で別個の解決策を構成したりすることができる。たとえば、図 1 ないし図 10 を参照して述べた解決策の任意の二つ以上が互いと組み合わせられてもよい。最も完備な解決策として、オーディオ処理装置はオーディオ・コンテンツ分類器 202 およびオーディオ・コンテキスト分類器 204 の両方ならびに平滑化ユニット 712、パラメータ平滑化ユニット 814 およびタイマー 916 を有していてもよい。

【0133】

先述したように、オーディオ改善装置 400 は、ダイアログ向上器 402、サラウンド仮想化器 404、ボリューム準化器 406 および等化器 408 を含んでいてもよい。オーディオ処理装置 100 は、それらの任意の一つまたは複数を含んでいてもよく、調整ユニット 300 がそれに適応されてもよい。複数のオーディオ改善装置 400 に関わるとき、調整ユニット 300 は、それぞれのオーディオ改善装置 400 に固有の複数のサブユニット 300A ~ 300D (図 15、図 18、図 20 および図 22) を含むものと見なされてもよし、あるいは相変わらず一つの連合した調整ユニットと見なされてもよい。あるオーディオ改善装置に固有であるとき、調整ユニット 300 はオーディオ分類器 200 および他の可能なコンポーネントと一緒に、その特定のオーディオ改善装置のコントローラと見なされてもよい。これについては、後述する第二部 ~ 第五部において詳細に論じる。

【0134】

さらに、オーディオ改善装置 400 は、上述した例に限定されず、他のいかなるオーディオ改善装置を含んでいてもよい。

【0135】

さらに、すでに論じた任意の解決策またはそれらの任意の組み合わせは、本開示の他の部分において記述または含意される任意の実施形態とさらに組み合わせられてもよい。特に、第六部および第七部において論じられるオーディオ分類器の実施形態は、オーディオ処理装置において使用されてもよい。

【0136】

1.8 節 オーディオ処理方法

上記の実施形態におけるオーディオ処理装置を記述する過程で、いくつかのプロセスまたは方法も開示されていることは明らかである。以下では、これらの方法の概要が与えられるが、上記ですでに論じた詳細の一部は繰り返さない。ただし、これらの方法はオーディオ処理装置を記述する過程において開示されているものの、これらの方法は必ずしも記載されるコンポーネントを採用するものではなく、必ずしもそうしたコンポーネントによって実行されるのではない。たとえば、オーディオ処理装置の実施形態は、部分的または完全にハードウェアおよび/またはファームウェアを用いて実現されてもよく、一方、以下で論じるオーディオ処理方法は、オーディオ処理装置のハードウェアおよび/またはファームウェアを採用してもよいが、完全にコンピュータ実行可能プログラムによって実現されてもよい。

【0137】

図 11 ~ 図 14 を参照して以下でこれらの方法について述べる。本方法がリアルタイムで実装されるときは、オーディオ信号のストリーミング属性に対応して、さまざまな動作が繰り返され、異なる動作は必ずしも同じオーディオ・セグメントに関してではないことに注意されたい。

【0138】

図 11 に示される実施形態では、オーディオ処理方法が提供される。まず、処理されるべきオーディオ信号がリアルタイムで少なくとも一つのオーディオ型に分類される (動作 1102)。前記少なくとも一つのオーディオ型の信頼値に基づいて、オーディオ改善のための少なくとも一つのパラメータが連続的に調整されることことができる (動作 1104)

。オーディオ改善は、ダイアログ向上（動作 1 1 0 6）、サラウンド仮想化（動作 1 1 0 8）、ボリューム平準化（動作 1 1 1 0）および/または等化（動作 1 1 1 2）であってもよい。対応して、前記少なくとも一つのパラメータは、ダイアログ向上処理、サラウンド仮想化処理、ボリューム平準化処理および等化処理のうちの少なくとも一つについての少なくとも一つのパラメータを含んでいてもよい。

【0 1 3 9】

ここで、「リアルタイムで」および「連続的に」はオーディオ型が、よってパラメータも、オーディオ信号の特定の内容とともにリアルタイムで変化することを意味する。「連続的に」は、調整が、急激または離散的な調整ではなく、信頼値に基づく連続的な調整であることをも意味する。

10

【0 1 4 0】

オーディオ型はコンテンツ型および/またはコンテキスト型を含んでいてもよい。対応して、調整の動作 1 1 0 4 は、少なくとも一つのコンテンツ型の信頼値および少なくとも一つのコンテキスト型の信頼値に基づいて前記少なくとも一つのパラメータを調整するよう構成されていてもよい。コンテンツ型はさらに、短期的音楽、発話、背景音およびノイズのコンテンツ型のうちの少なくとも一つを含んでいてもよい。コンテキスト型はさらに、長期的音楽、映画のメディア、ゲームおよびVoIPのコンテキスト型のうちの少なくとも一つを含んでいてもよい。

【0 1 4 1】

他のいくつかのコンテキスト型スキームも提案される。たとえば、VoIPおよび非VoIPを含むVoIP関係コンテキスト型および高品質オーディオまたは低品質オーディオを含むオーディオ品質型などである。

20

【0 1 4 2】

短期的音楽は、種々の標準に従ってサブ型にさらに分類されてもよい。優勢な源の存在に依存して、優勢な源のない音楽および優勢な源のある音楽を含んでいてもよい。さらに、短期的音楽は、少なくとも一つのジャンル・ベースのクラスターまたは少なくとも一つの楽器ベースのクラスターまたは音楽のリズム、テンポ、音色および/または他の任意の音楽的属性に基づいて分類された少なくとも一つの音楽クラスターを含んでいてもよい。

【0 1 4 3】

コンテンツ型およびコンテキスト型の両方が識別されたとき、コンテンツ型の重要性は、そのコンテンツ型が位置しているところのコンテキスト型によって決定されてもよい。すなわち、異なるコンテキスト型のオーディオ信号におけるコンテンツ型は、オーディオ信号のコンテキスト型に依存して異なる重みを割り当てられる。より一般には、あるオーディオ型が別のオーディオ型に影響してもよく、別のオーディオ型の前提であってもよい。したがって、調整 1 1 0 4 の動作は、あるオーディオ型の重みを別のオーディオ型の信頼値を用いて修正するよう構成されていてもよい。

30

【0 1 4 4】

オーディオ信号が同時に（すなわち、同じオーディオ・セグメントに関して）複数のオーディオ型に分類されるとき、調整 1 1 0 4 の動作は、そのオーディオ・セグメントを改善するためのパラメータ（単数または複数）を調整するための識別されたオーディオ型の一部または全部を考慮してもよい。たとえば、調整 1 1 0 4 の動作は、前記少なくとも一つのオーディオ型の重要性に基づいて前記少なくとも一つのオーディオ型の信頼値に重み付けするよう構成されていてもよい。あるいは、調整 1 1 0 4 の動作は、前記オーディオ型の少なくともいくつかを、その信頼値に基づいてそれらに重み付けすることを通じて考慮するよう構成されていてもよい。ある特別な場合には、調整 1 1 0 4 の動作は、信頼値に基づいて少なくとも一つの優勢なオーディオ型を考慮するよう構成されていてもよい。

40

【0 1 4 5】

結果の急激な変化を避けるために、平滑化方式が導入されてもよい。

【0 1 4 6】

調整されたパラメータ値は平滑化されてもよい（図 1 2 の動作 1 2 1 4）。たとえば、

50

現在の時間における調整 1 1 0 4 の動作によって決定されたパラメータ値が、現在の時間における調整の動作によって決定されたパラメータ値と、最後の時間の平滑化されたパラメータ値の重み付けされた和で置き換えられてもよい。このように、逐次反復される平滑化動作を通じて、パラメータ値は時間軸上で平滑化される。

【 0 1 4 7 】

重み付けされた和を計算するための重みは、オーディオ信号のオーディオ型に基づいてまたはあるオーディオ型から別のオーディオ型への異なる遷移対に基づいて適応的に変更されてもよい。あるいはまた、重み付けされた和を計算するための重みは、調整の動作によって決定されたパラメータ値の増加または減少トレンドに基づいて適応的に変更される。

10

【 0 1 4 8 】

もう一つの平滑化方式が図 1 3 に示されている。すなわち、本方法はさらに、各オーディオ型について、現在における実際の信頼値と最後の時間の平滑化された信頼値の重み付けされた和を計算することによって、現在の時間におけるオーディオ信号の信頼値を平滑化することを含んでいてもよい（動作 1 3 0 3）。パラメータ平滑化動作 1 2 1 4 と同様に、重み付けされた和を計算するための重みは、オーディオ信号のオーディオ型の信頼値に基づいて、またはあるオーディオ型から別のオーディオ型への異なる遷移対に基づいて適応的に変更されてもよい。

【 0 1 4 9 】

もう一つの平滑化方式は、オーディオ分類動作 1 1 0 2 の出力が変わったとしてもあるオーディオ型から別のオーディオ型への遷移を遅らせるバッファ機構である。すなわち、調整 1 1 0 4 の動作は、すぐに新しいオーディオ型を使うのではなく、オーディオ分類動作 1 1 0 2 の出力の安定化を待つ。

20

【 0 1 5 0 】

具体的には、本方法は、分類動作が同じ新しいオーディオ型を連続的に出力する持続時間を測定することを含んでいてもよい（図 1 4 の動作 1 4 0 3）。ここで、調整 1 1 0 4 の動作は、新しいオーディオ型の持続時間の長さがある閾値に達する（動作 1 4 0 3 5 における「Y」および動作 1 1 0 4 2）まで、現在のオーディオ型を使い続ける（動作 1 4 0 3 5 における「N」および動作 1 1 0 4 1）よう構成される。具体的には、オーディオ分類動作 1 1 0 2 から出力されるオーディオ型が、オーディオ・パラメータ調整動作 1 1 0 4 において使われている現在のオーディオ型に関して変化するとき（動作 1 4 0 3 において「Y」）、計時が始まる（動作 1 4 0 3 2）。オーディオ分類動作 1 1 0 2 が新しいオーディオ型を出力することを続ければ、すなわち、動作 1 4 0 3 1 における判断が「Y」であり続ければ、計時は続く（動作 1 4 0 3 2）。最終的には、新しいオーディオ型の持続時間が閾値に達すると（動作 1 4 0 3 5 における「Y」）、調整動作 1 1 0 4 は新しいオーディオ型を使い（動作 1 1 0 4 2）、計時は、オーディオ型の次の切り換えに備えてリセットされる（動作 1 4 0 3 4）。閾値に達するまでは（動作 1 4 0 3 5 における「N」）、調整動作 1 1 0 4 は現在のオーディオ型を使い続ける（動作 1 1 0 4 1）。

30

【 0 1 5 1 】

ここで、計時はタイマーの機構（カウントアップまたはカウントダウン）を用いて実装されてもよい。計時が始まった後、ただし閾値に達する前に、オーディオ分類動作 1 1 0 2 の出力が調整動作 1 1 0 4 において使われている現在のオーディオ型に戻る場合には、調整動作 1 1 0 4 によって使用される現在のオーディオ型に関して変化がないと見なされるべきである（動作 1 4 0 3 1 における「N」）。だが（オーディオ信号における分類されるべき現在のオーディオ・セグメントに対応する）現在の分類結果は、（オーディオ信号における分類されるべき前のオーディオ・セグメントに対応する）オーディオ分類動作 1 1 0 2 の前の出力に関して変化し（動作 1 4 0 3 3 における「Y」）、よって計時は、次の変化（動作 1 4 0 3 1 における「Y」）が計時を開始するまで、リセットされる（動作 1 4 0 3 4）。むしろ、オーディオ分類動作 1 1 0 2 の分類結果がオーディオ・パラメータ調整動作 1 1 0 4 によって使用される現在のオーディオ型に関して変化せず（動作 1

40

50

4031における「N」)、前の分類に関しても変化しない(動作14033における「N」)場合には、そのことは、オーディオ分類が安定した状態にあることを示し、現在のオーディオ型が使用され続ける。

【0152】

ここで使われる閾値は、あるオーディオ型から別のオーディオ型への異なる遷移対については異なってもよい。というのも、状態があまり安定でないときは、一般に、オーディオ改善装置が他の状態よりもそのデフォルト状態にあることが好ましいことがあるからである。他方、新しいオーディオ型の信頼値が比較的高い場合には、新しいオーディオ型に遷移するほうが安全である。したがって、閾値は、新しいオーディオ型の信頼値と負に相関していてもよい。信頼値が高いほど、閾値は低く、つまりオーディオ型は新しいオーディオ型により速く遷移しうる。

10

【0153】

オーディオ処理装置の実施形態と同様に、オーディオ処理方法の実施形態およびその変形の任意の組み合わせが現実的である。他方、オーディオ処理方法の実施形態およびその変形のあらゆる側面は別個の解決策であってもよい。特に、オーディオ処理方法のすべてにおいて、第六部および第七部において論じるようなオーディオ分類方法が使われてもよい。

【0154】

第二部：ダイアログ向上器コントローラおよび制御方法

オーディオ改善装置の一例はダイアログ向上器(DE)である。これは、特に聴力が低下しつつある高齢者のために、再生時にオーディオを断続的にモニタリングし、ダイアログの存在を検出し、ダイアログの明瞭性および了解性を高める(ダイアログを聞いて理解しやすくする)ためにダイアログを向上させることをねらいとする。ダイアログが存在するかどうかを検出するほか、ダイアログが存在し、よって(動的スペクトル再均衡化(dynamic spectral rebalancing)を用いて)相応して向上される場合、了解性に最も重要な周波数も検出される。例示的なダイアログ向上方法が特許文献1に呈示されている。その全体はここに参照によって組み込まれる。

20

【0155】

ダイアログ向上器における一般的な手動の構成設定は、通例、映画のメディア・コンテンツについては有効にされるが、音楽コンテンツについては無効にされるというものである。ダイアログ向上は、音楽信号に対しては誤ってトリガーしすぎることがあるからである。

30

【0156】

利用可能なオーディオ型情報を用いて、ダイアログ向上のレベルおよび他のパラメータが、識別されたオーディオ型の信頼値に基づいて調整されることができる。先に論じたオーディオ処理装置および方法の個別的な例として、ダイアログ向上器は、第一部で論じたすべての実施形態およびそれらの実施形態の任意の組み合わせを使用してもよい。特に、ダイアログ向上器を制御する場合、図1~図10に示されるようなオーディオ処理装置100におけるオーディオ分類器200および調整ユニット300は、図15に示されるようなダイアログ向上器コントローラ1500を構成してもよい。この実施形態では、調整ユニットはダイアログ向上器に固有なので、300Aと称されてもよい。先述した部において論じたように、オーディオ分類器200は、オーディオ・コンテンツ分類器202およびオーディオ・コンテキスト分類器204のうち少なくとも一つを含んでいてもよく、ダイアログ向上器コントローラ1500はさらに、型平滑化ユニット712、パラメータ平滑化ユニット814およびタイマー916のうち少なくとも一つを含んでいてもよい。

40

【0157】

したがって、この部においては、先の部ですでに記述した内容を繰り返すことはせず、単にこの部のいくつかの固有の例を与える。

【0158】

50

ダイアログ向上器については、調整可能なパラメータは、ダイアログ向上のレベル、背景レベルおよび向上されるべき周波数帯域を決定するための閾値を含むがそれに限定されない。特許文献1参照。その全体はここに参照によって組み込まれる。

【0159】

2.1節 ダイアログ向上のレベル

ダイアログ向上のレベルに関わるとき、調整ユニット300Aは、ダイアログ向上器のダイアログ向上のレベルを、発話の信頼値と正に相関させるよう構成されていてもよい。追加的または代替的に、レベルは、他のコンテンツ型の信頼値に負に相関させられてもよい。こうして、ダイアログ向上のレベルは発話信頼度に（線形または非線形に）比例するように設定されることができる。よって、ダイアログ向上は、音楽および背景音（効果音）のような非発話信号についてはそれほど効果的ではない。

10

【0160】

コンテキスト型については、調整ユニット300Aは、ダイアログ向上器のダイアログ向上のレベルを、映画的メディアおよび/またはVoIPの信頼値と正に相関させ、ダイアログ向上器のダイアログ向上のレベルを、長期的音楽および/またはゲームの信頼値と負に相関させるよう構成されていてもよい。たとえば、ダイアログ向上のレベルは映画的メディアの信頼値に（線形または非線形に）比例するように設定されることができる。映画的メディア信頼値が0のとき（たとえば音楽コンテンツにおいて）は、ダイアログ向上のレベルも0であり、これはダイアログ向上器を無効にすることと等価である。

【0161】

先の部で述べたように、コンテンツ型およびコンテキスト型は合同して考慮されてもよい。

20

【0162】

2.2節 向上させるべき周波数帯域の決定のための閾値

ダイアログ向上器の作動の間、各周波数帯域について、向上される必要があるかどうかを判定するための閾値（通例、エネルギーまたはラウドネス閾値）がある。すなわち、それぞれのエネルギー/ラウドネス閾値より上の周波数帯域が向上される。それらの閾値を調整するために、調整ユニット300Aは、閾値を、短期的音楽および/またはノイズおよび/または背景音の信頼値と正に相関させるおよび/または閾値を発話の信頼値と負に相関させるよう構成されていてもよい。たとえば、発話信頼度が高ければ、より信頼できる発話検出を想定して閾値を下げることができ、より多くの周波数帯域が向上されることを許容する。他方、音楽信頼度が高ければ、閾値を高くすることができ、より少数の周波数帯域が向上されるようにする（よってアーチファクトをより少なくする）。

30

【0163】

2.3節 背景レベルへの調整

ダイアログ向上器におけるもう一つのコンポーネントは、図15に示されるような最小追跡ユニット4022である。これは、（SNR推定および2.2節で述べた周波数帯域閾値推定のために）オーディオ信号における背景レベルを推定するために使われる。これは、オーディオ・コンテンツ型の信頼値に基づいて調整されることもできる。たとえば、発話信頼値が高い場合、最小追跡ユニットは、背景レベルを現在の最小に設定するのにより自信を持つことができる。音楽信頼度が高い場合には、背景レベルはその現在の最小よりはやや高く設定されることができ、あるいは別の仕方では現在の最小と現在フレームのエネルギーとの、現在の最小に大きな重みをかけた重み付き平均に設定されることができる。ノイズおよび背景信頼度が高い場合には、背景レベルは現在の最小値よりずっと高く設定されることができ、あるいは別の仕方では、現在の最小と現在フレームのエネルギーとの、現在の最小に小さな重みをかけた重み付き平均に設定されることができる。

40

【0164】

こうして、調整ユニット300Aは、最小追跡ユニットによって推定された背景レベルに調整を割り当てるよう構成されてもよい。ここで、調整ユニットはさらに、調整を、短期的音楽および/またはノイズおよび/または背景音の信頼値と正に相関させるおよび/

50

または調整を発話の信頼値と負に相関させるよう構成されている。ある変形では、調整ユニット300Aは、調整を、短期的音楽より、ノイズおよび/または背景音の信頼値と、より正に相関させるよう構成されていてもよい。

【0165】

2.4節 実施形態の組み合わせおよび応用シナリオ

第一部と同様に、上記で論じたすべての実施形態およびその変形は、そのいかなる組み合わせにおいて実装されてもよく、異なる部/実施形態において言及されるが同じまたは同様の機能をもついかなる構成要素も同じまたは別個の構成要素として実装されてもよい。

【0166】

10

たとえば、2.1節ないし2.3節において述べた解決策の任意の二つ以上が互いと組み合わせられてもよい。そして、これらの組み合わせは、第一部および後述する他の部において記載または含意されている任意の実施形態とさらに組み合わせられてもよい。特に、それぞれの種類のオーディオ改善装置または方法に対して多くの公式が実際に適用可能であるが、それらは必ずしも本開示の各部において記載または議論されていない。そのような場合、ある部において論じられる特定の公式を、特定の用途の特定の要件に応じて関連するパラメータ、係数、冪（指数）および重みが適正に調整されるだけで他の部に適用するために、本開示の各部の間で相互参照がなされることがある。

【0167】

2.5節 ダイアログ向上器制御方法

20

第一部と同様に、上記の実施形態におけるダイアログ向上器コントローラを記述する過程で、いくつかのプロセスまたは方法も開示されていることは明らかである。以下では、これらの方法の概要が与えられるが、上記ですでに論じた詳細の一部は繰り返さない。

【0168】

まず、第一部で論じたオーディオ処理方法の実施形態がダイアログ向上器のために使われてもよい。ダイアログ向上器のパラメータ（単数または複数）が、オーディオ処理方法によって調整されるべきターゲットの一つである。この観点から、オーディオ処理方法はダイアログ向上器制御方法でもある。

【0169】

この節では、ダイアログ向上器の制御に固有の側面のみが論じられる。制御方法の一般的な側面については、第一部が参照されうる。

30

【0170】

ある実施形態によれば、オーディオ処理方法はさらに、ダイアログ向上処理を含んでいてもよく、調整1104の動作は、ダイアログ向上のレベルを映画的メディアおよび/またはVoIPの信頼値と正に相関させるおよびまたはダイアログ向上のレベルを長期的音楽および/またはゲームの信頼値と負に相関させることを含む。すなわち、ダイアログ向上は主として、映画的メディアまたはVoIPのコンテキストにおけるオーディオ信号に向けられる。

【0171】

より具体的には、調整1104の動作は、ダイアログ向上器のダイアログ向上のレベルを発話の信頼値と正に相関させることを含んでいてもよい。

40

【0172】

本願は、ダイアログ向上処理において向上されるべき周波数帯域を調整してもよい。図16に示されるように、それぞれの周波数帯域が向上されるべきかどうかを決定するための閾値（通例エネルギーまたはラウドネス）が、本願に従って、識別されたオーディオ型の信頼値（単数または複数）に基づいて調整されてもよい（動作1602）。次いで、ダイアログ向上器内で、調整された閾値に基づいて、それぞれの閾値より上の周波数帯域が選択され（動作1604）、向上される（動作1606）。

【0173】

具体的には、調整1104の動作は、それらの閾値を、短期的音楽および/またはノイ

50

ズおよび/または背景音の信頼値と正に相関させるおよび/またはそれらの閾値を発話の信頼値と負に相関させることを含んでいてもよい。

【0174】

オーディオ処理方法(特にダイアログ向上処理)は一般にさらに、オーディオ信号における背景レベルを推定することを含む。これは一般に、ダイアログ向上器402において実現される最小追跡ユニット4022によって実装され、SNR推定または周波数帯域閾値推定において使われる。本願は、背景レベルを調整するために使われてもよい。そのような状況では、背景レベルが推定された後(動作1702)、背景レベルはまず、オーディオ型(単数または複数)の信頼値(単数または複数)に基づいて調整され(動作1704)、次いでSNR推定および/または周波数帯域閾値推定において使われる(動作1706)。特に、調整1104の動作は、推定された背景レベルに調整を割り当てるよう構成されていてもよい。ここで、調整1104の動作は、調整を、短期的音楽および/またはノイズおよび/または背景音の信頼値と正に相関させるおよび/または調整を発話の信頼値と負に相関させるよう構成されていてもよい。

10

【0175】

より具体的には、調整1104の動作は、調整を、短期的音楽よりも、ノイズおよび/または背景の信頼値と、より正に相関させるよう構成されていてもよい。

【0176】

オーディオ処理装置の実施形態と同様に、オーディオ処理方法の実施形態およびその変形の任意の組み合わせが現実的である。他方、オーディオ処理方法の実施形態およびその変形のあらゆる側面は別個の解決策であってもよい。さらに、この節に記載される解決策の任意の二つ以上が互いと組み合わせられてもよく、これらの組み合わせがさらに、第一部または後述する他の部において記載または含意される任意の実施形態と組み合わせられてもよい。

20

【0177】

第三部：サラウンド仮想化器コントローラおよび制御方法

サラウンド仮想化器は、サラウンドサウンド信号(マルチチャンネル5.1および7.1など)がPCの内部スピーカーを通じてまたはヘッドフォンを通じてレンダリングされることを可能にする。すなわち、内蔵ラップトップ・スピーカーまたはヘッドフォンのようなステレオ装置を用いて、仮想的にサラウンド効果を生成し、消費者のために映画館の体験を提供するのである。サラウンド仮想化器では、マルチチャンネル・オーディオ信号に関連付けられたさまざまなスピーカー位置からくる音の耳への到来をシミュレートするために、通例、頭部伝達関数(HRTF: Head Related Transfer Function)が利用される。

30

【0178】

現在のサラウンド仮想化器はヘッドフォン上でよく機能するが、組み込みスピーカーと異なるコンテンツに対して異なる仕方で機能する。一般に、映画のメディア・コンテンツはスピーカーのためにサラウンド仮想化器を有効にし、一方、音楽はそうしない。あまりに薄っぺらに聞こえることがありうるからである。

【0179】

サラウンド仮想化器における同じパラメータが映画のメディアおよび音楽コンテンツの両方について同時に良好な音像を作り出すことはできないので、パラメータはコンテンツに基づいてより精密に調整される必要がある。利用可能なオーディオ型情報、特に音楽信頼値および発話信頼値ならびに他の何らかのコンテンツ型情報およびコンテキスト情報を用いて、機能は本願とともに行なうことができる。

40

【0180】

第二部と同様に、第一部で論じたオーディオ処理装置および方法の個別的な例として、サラウンド仮想化器404は、第一部で論じたすべての実施形態およびそこで開示されたそれらの実施形態の任意の組み合わせを使用してもよい。特に、サラウンド仮想化器を制御する場合、図1~図10に示されるようなオーディオ処理装置100におけるオーディオ分類器200および調整ユニット300は、図18に示されるようなサラウンド仮想化

50

器コントローラ 1800 を構成してもよい。この実施形態では、調整ユニットはサラウンド仮想化器に固有なので、300B と称されてもよい。第二部と同様に、オーディオ分類器 200 は、オーディオ・コンテンツ分類器 202 およびオーディオ・コンテキスト分類器 204 のうちの少なくとも一つを含んでいてもよく、サラウンド仮想化器コントローラ 1800 はさらに、型平滑化ユニット 712、パラメータ平滑化ユニット 814 およびタイマー 916 のうちの少なくとも一つを含んでいてもよい。

【0181】

したがって、この部においては、第一部ですでに記述した内容を繰り返すことはせず、単にこの部のいくつかの固有の例を与える。

【0182】

サラウンド仮想化器については、調整可能なパラメータは、サラウンド・ブースト量およびサラウンド仮想化器 404 の開始周波数を含むがそれに限定されない。

【0183】

3.1 節 サラウンド・ブースト量

サラウンド・ブースト量に関わるとき、調整ユニット 300B は、サラウンド仮想化器 404 のサラウンド・ブースト量を、ノイズおよび/または背景および/または発話の信頼値と正に相関させるおよび/またはサラウンド・ブースト量を短期的音楽の信頼値と負に相関させるよう構成されていてもよい。

【0184】

特に、音楽（コンテンツ型）が受け入れ可能に聞こえるようにサラウンド仮想化器 404 を修正するために、調整ユニット 300B の例示的な実装は、短期的音楽信頼値に基づいてサラウンド・ブーストの量を調整することができる。たとえば、

$$SB (1 - \text{Conf}_{\text{music}}) \quad (5)$$

ここで、SB はサラウンド・ブースト量、 $\text{Conf}_{\text{music}}$ は短期的音楽の信頼値である。

【0185】

それは、音楽についてサラウンド・ブーストを減少させ、音楽が生氣ないように聞こえることを防ぐ。

【0186】

同様に、発話信頼値も利用でき、たとえば、

$$SB (1 - \text{Conf}_{\text{music}}) * \text{Conf}_{\text{speech}} \quad (6)$$

ここで、 $\text{Conf}_{\text{speech}}$ は発話の信頼値であり、 $\text{Conf}_{\text{speech}}$ は指数の形の重み付け係数であり、1~2 の範囲にあってもよい。この公式は、サラウンド・ブースト量は純粋な発話（高い発話信頼度および低い音楽信頼度）についてのみ高くなることを示す。

【0187】

あるいは、発話の信頼値のみを考慮することもできる

$$SB \text{Conf}_{\text{speech}} \quad (7)$$

さまざまな変形が同様にして設計できる。特に、ノイズまたは背景音について、公式(5)ないし(7)と同様の公式が構築されてもよい。さらに、それら四つのコンテンツ型の効果は任意の組み合わせにおいて一緒に考慮されてもよい。そのような状況において、ノイズおよび背景は周囲音であり、大きなブースト量をもってもより安全である。発話は、話者が通例はスクリーンの前方に座ると想定して、中程度のブースト量をもつことができる。したがって、調整ユニット 300B は、サラウンド・ブースト量を、コンテンツが型発話よりも、ノイズおよび/または背景の信頼値と、より正に相関させるよう構成されていてもよい。

【0188】

各コンテンツ型について期待されるブースト量（これは重みと等化である）をあらかじめ定義していたとして、もう一つの代替を適用することもできる。

【0189】

10

20

30

40

【数 9】

$$\hat{a} = \frac{a_{speech} \cdot Conf_{speech} + a_{music} \cdot Conf_{music} + a_{noise} \cdot Conf_{noise} + a_{bkg} \cdot Conf_{bkg}}{Conf_{speech} + Conf_{music} + Conf_{noise} + Conf_{bkg}} \quad (8)$$

ここで、^付きのaは推定されたブースト量、コンテンツ型の添え字をもつはそのコンテンツ型の期待される / あらかじめ定義されたブースト量（重み）、コンテンツ型の添え字をもつConfはそのコンテンツ型の信頼値である（ここで、bkgは「background sound」（背景音）を表わす）。状況に依存して、 a_{music} は（必須ではないが）0に設定されてもよい。これはサラウンド仮想化器404が純粋な音楽（コンテンツ型）については無効にされることを示す。

10

【0190】

別の観点からは、公式(8)におけるコンテンツ型の添え字をもつはそのコンテンツ型の期待される / あらかじめ定義されたブースト量であり、対応するコンテンツ型の信頼値を、すべての識別されたコンテンツ型の信頼値の和で割った商は、対応するコンテンツ型のあらかじめ定義された / 期待されるブースト量の規格化された重みと見なされてもよい。すなわち、調整ユニット300Bは、それらの信頼値に基づいて複数のコンテンツ型のあらかじめ定義された諸ブースト量に重み付けすることを通じて、複数のコンテンツ型のうち少なくともいくつかを考慮するよう構成されていてもよい。

20

【0191】

コンテキスト型については、調整ユニット300Bは、サラウンド仮想化器404のサラウンド・ブースト量を、映画のメディアおよび / またはゲームの信頼値と正に相関させ、サラウンド・ブースト量を、長期的音楽および / またはVoIPの信頼値と負に相関させるよう構成されていてもよい。次いで、(5)ないし(8)と同様の公式が構築される。

【0192】

特殊な例として、サラウンド仮想化器404は、純粋な映画のメディアおよび / またはゲームについては有効にされ、音楽および / またはVoIPについては無効にされることができ。一方、サラウンド仮想化器404のブースト量は映画のメディアおよびゲームについて異なるように設定されることができ。映画のメディアはより高いブースト量を使い、ゲームはより低い。したがって、調整ユニット300Bは、サラウンド・ブースト量を、ゲームよりも、映画のメディアの信頼値と、より正に相関させるよう構成されてもよい。

30

【0193】

コンテンツ型と同様に、オーディオ信号のブースト量は、コンテキスト型の信頼値の重み付き平均に設定されることもできる。

【0194】

【数 10】

$$\hat{a} = \frac{a_{MOVIE} \cdot Conf_{MOVIE} + a_{MUSIC} \cdot Conf_{MUSIC} + a_{GAME} \cdot Conf_{GAME} + a_{VOIP} \cdot Conf_{VOIP}}{Conf_{MOVIE} + Conf_{MUSIC} + Conf_{GAME} + Conf_{VOIP}} \quad (9)$$

40

ここで、^付きのaは推定されたブースト量、コンテキスト型の添え字をもつはそのコンテキスト型の期待される / あらかじめ定義されたブースト量（重み）、コンテキストの添え字をもつConfはそのコンテキスト型の信頼値である。状況に依存して、 a_{MUSIC} および a_{VOIP} は（必須ではないが）0に設定されてもよい。これはサラウンド仮想化器404が純粋な音楽（コンテキスト型）およびまたは純粋なVoIPについては無効にされることを示す。

【0195】

50

やはりコンテンツ型と同様に、公式(9)におけるコンテキスト型の添え字をもつはそのコンテキスト型の期待される/あらかじめ定義されたブースト量であり、対応するコンテキスト型の信頼値を、すべての識別されたコンテキスト型の信頼値の和で割った商は、対応するコンテキスト型のあらかじめ定義された/期待されるブースト量の規格化された重みと見なされてもよい。すなわち、調整ユニット300Bは、それらの信頼値に基づいて複数のコンテキスト型のあらかじめ定義された諸ブースト量に重み付けすることを通じて、複数のコンテキスト型のうち少なくともいくつかを考慮するよう構成されていてもよい。

【0196】

3.2節 開始周波数

他のパラメータも開始周波数のようなサラウンド仮想化器において修正されることができ、一般に、オーディオ信号中の高周波数成分は空間的にレンダリングされるのに、より好適である。たとえば、音楽では、ベースがより多くのサラウンド効果をもつようにレンダリングされると、おかしく聞こえる。よって、特定のオーディオ信号について、サラウンド仮想化器は、それより上の成分が空間的にレンダリングされ、それより下の成分が保持される周波数閾値を決定する必要がある。周波数閾値は開始周波数である。

【0197】

本願のある実施形態によれば、サラウンド仮想化器についての開始周波数は音楽コンテンツに対しては高められることができ、それにより音楽信号についてはより多くのベースが保持されることができ、すると、調整ユニット300Bは、サラウンド仮想化器の開始周波数を短期的音楽の信頼値と正に相関させるよう構成されうる。

【0198】

3.3節 実施形態の組み合わせおよび応用シナリオ

第一部と同様に、上記で論じたすべての実施形態およびその変形は、そのいかなる組み合わせにおいて実装されてもよく、異なる部/実施形態において言及されるが同じまたは同様の機能をもついかなる構成要素も同じまたは別個の構成要素として実装されてもよい。

【0199】

たとえば、3.1節および3.2節において述べた解決策の任意の二つ以上が互いと組み合わせられてもよい。そして、これらの組み合わせの任意のものが、第一部、第二部および後述する他の部において記載または含意されている任意の実施形態とさらに組み合わせられてもよい。

【0200】

3.4節 サラウンド仮想化器制御方法

第一部と同様に、上記の実施形態におけるサラウンド仮想化器コントローラを記述する過程で、いくつかのプロセスまたは方法も開示されていることは明らかである。以下では、これらの方法の概要が与えられるが、上記ですでに論じた詳細の一部は繰り返さない。

【0201】

まず、第一部で論じたオーディオ処理方法の実施形態がサラウンド仮想化器について使用されてもよい。サラウンド仮想化器のパラメータ(単数または複数)が、オーディオ処理方法によって調整されるべきターゲットの一つである。この観点から、オーディオ処理方法はサラウンド仮想化器制御方法でもある。

【0202】

この節では、サラウンド仮想化器の制御に固有の側面のみが論じられる。制御方法の一般的な側面については、第一部が参照されうる。

【0203】

ある実施形態によれば、オーディオ処理方法はさらに、サラウンド仮想化処理を含んでもよく、調整する動作1104はサラウンド仮想化処理のサラウンド・ブースト量をノイズおよび/または背景および/または発話の信頼値と正に相関させるおよび/またはサラウンド・ブースト量を短期的音楽の信頼値と負に相関させるよう構成されていてもよ

10

20

30

40

50

い。

【0204】

特に、調整する動作1104はサラウンド・ブースト量をノイズおよび/または背景および/または発話の信頼値と、コンテンツ型発話よりも、より正に相関させるよう構成されていてもよい。

【0205】

代替的または追加的に、サラウンド・ブースト量は、コンテキスト型(単数または複数)の信頼値(単数または複数)に基づいて調整されてもよい。特に、調整する動作1104は、サラウンド仮想化処理のサラウンド・ブースト量を映画のメディアおよび/またはゲームの信頼値と正に相関させるおよび/またはサラウンド・ブースト量を長期的音楽および/またはVoIPの信頼値と負に相関させるよう構成されていてもよい。

10

【0206】

特に、調整する動作1104はサラウンド・ブースト量を映画のメディアと、ゲームよりも、より正に相関させるよう構成されていてもよい。

【0207】

調整されるべきもう一つのパラメータは、サラウンド仮想化処理のための開始周波数である。図19に示されるように、開始周波数はオーディオ型(単数または複数)の信頼値(単数または複数)に基づいてまず調整され(動作1902)、次いで、サラウンド仮想化器は、開始周波数より上のオーディオ・コンポーネントを処理する(動作1904)。特に、調整する動作1104は、サラウンド仮想化処理の開始周波数を短期的音楽の信頼値と正に相関させるよう構成されていてもよい。

20

【0208】

オーディオ処理装置の実施形態と同様に、オーディオ処理方法の実施形態およびその変形の任意の組み合わせが現実的である。他方、オーディオ処理方法の実施形態およびその変形のあらゆる側面は別個の解決策であってもよい。さらに、この節に記載される解決策の任意の二つ以上が互いと組み合わせられてもよく、これらの組み合わせがさらに、本開示の他の部において記載または含意される任意の実施形態と組み合わせられてもよい。

【0209】

第四部：ボリューム平準化器コントローラおよび制御方法

異なるオーディオ源または同じオーディオ源の異なるピースのボリュームは時に大きく変化する。ユーザーがボリュームを頻繁に調整しなければならないので、これはわずらわしい。ボリューム平準化器(VL: Volume Leveler)は再生時のオーディオ・コンテンツのボリュームを調整し、ターゲット・ラウドネス値に基づいて時間軸上でほとんど一貫しているようにすることをねらいとする。例示的なボリューム平準化器は特許文献2、特許文献3、特許文献4に記載されている。これら三つの文書はここに全体において参照によって組み込まれる。

30

【0210】

ボリューム平準化器は、何らかの仕方でオーディオ信号のラウドネスを連続的に測定し、次いで利得(gain)の量だけ信号を修正する。gainはオーディオ信号のラウドネスを修正するためのスケール因子であり、通例、測定されたラウドネス、所望されるターゲット・ラウドネスおよび他のいくつかの因子の関数である。適正な利得を推定するためにいくつかの因子が、ターゲット・ラウドネスに接近しかつダイナミックレンジを維持するための基礎になる基準とともに、考慮される必要がある。それは通例、自動利得制御(AGC)、聴覚イベント検出、ダイナミックレンジ制御(DRC)のようないくつかのサブ要素を含む。

40

【0211】

制御信号は、オーディオ信号の「利得」を制御するためにボリューム平準化器において一般に適用される。たとえば、制御信号は、純粋な信号解析によって導出される、オーディオ信号の大きさの変化のインジケータであることができる。制御信号はまた、聴覚シーン解析または特定ラウドネス・ベースの(specific-loudness-based)聴覚イベント・

50

検出といった音響心理学的解析を通じた、新しいオーディオ・イベントが現われるかどうかを表わすオーディオ・イベント・インジケータであることもできる。そのような制御信号は、たとえばオーディオ信号における利得の急激な変化に起因する可能な可聴アーチファクトを低減するために利得が聴覚イベント内でほぼ一定であることを保証することによって、および利得変化の多くをイベント境界の近傍に制約することによって、利得制御のためにボリューム平準化器において適用される。

【0212】

しかしながら、制御信号を導出する従来の方法は、情報性の聴覚イベントを非情報性（干渉性）の聴覚イベントから区別することができない。ここで、情報性の聴覚イベントは、ダイアログおよび音楽のような、意味のある情報を含むオーディオ・イベントを表わし、ユーザーがより注意を払うことがありうる。一方、非情報性の信号は、VoIPにおけるノイズのように、ユーザーにとって意味のある情報を含まない。結果として、非情報性の信号も、大きな利得を適用され、ターゲット・ラウドネス近くにブーストされうる。それはいくつかの応用では快くない。たとえば、VoIP通話では、会話の休止に現われるノイズ信号が、ボリューム平準化器により処理された後に、しばしば大きなボリュームにブーストされる。これはユーザーによって望まれない。

10

【0213】

この問題に少なくとも部分的に対処するために、本願は、第一部で論じた実施形態に基づいてボリューム平準化器を制御することを提案する。

20

【0214】

第二部および第三部と同様に、第一部で論じたオーディオ処理装置および方法の個別的な例として、ボリューム平準化器406は、第一部で論じたすべての実施形態およびそこで開示されたそれらの実施形態の任意の組み合わせを使用してもよい。特に、ボリューム平準化器406を制御する場合、図1～図10に示されるようなオーディオ処理装置100におけるオーディオ分類器200および調整ユニット300は、図20に示されるようなボリューム平準化器406コントローラ2000を構成してもよい。この実施形態では、調整ユニットはボリューム平準化器406に固有なので、300Cと称されてもよい。

【0215】

すなわち、第一部の開示に基づき、ボリューム平準化器コントローラ2000は、オーディオ信号のオーディオ型（コンテンツ型および/またはコンテキスト型など）を連続的に識別するオーディオ分類器200と、識別されたオーディオ型の信頼値に基づいて連続的な仕方でボリューム平準化器を調整する調整ユニット300Cとを有していてもよい。同様に、オーディオ分類器200は、オーディオ・コンテンツ分類器202およびオーディオ・コンテキスト分類器204のうちの少なくとも一つを含んでいてもよく、ボリューム平準化器コントローラ2000はさらに、型平滑化ユニット712、パラメータ平滑化ユニット814およびタイマー916のうちの少なくとも一つを含んでいてもよい。

30

【0216】

したがって、この部においては、第一部ですでに記述した内容を繰り返すことはせず、単にこの部のいくつかの固有の例を与える。

【0217】

分類結果に基づいて、ボリューム平準化器406の種々のパラメータが適応的に調整されることができる。たとえば非情報性信号についての利得を低減することにより、動的利得または動的利得の範囲に直接関係したパラメータを調整することができる。信号が新しい知覚可能なオーディオ・イベントである度合いを示すパラメータを調整し、そして動的利得を間接的に制御することもできる（利得は、オーディオ・イベント内でゆっくり変化するが、二つのオーディオ・イベントの境界では急激に変化することがある）。本願では、パラメータ調整またはボリューム平準化器制御機構のいくつかの実施形態が呈示される。

40

【0218】

4.1節 情報性および干渉性のコンテンツ型

50

上述したように、ボリューム平準化器の制御との関連で、オーディオ・コンテンツ型は情報性のコンテンツ型および干渉性のコンテンツ型として分類される。

調整ユニット300Cは、ボリューム平準化器の動的利得をオーディオ信号の情報性コンテンツ型と正に相関させ、ボリューム平準化器の動的利得をオーディオ信号の干渉性コンテンツ型と負に相関させるよう構成されてもよい。

【0219】

例として、ノイズが干渉性（非情報性）であり、それは大きなボリュームにブーストされるとわずらわしいとする。動的利得を直接制御するパラメータまたは新しいオーディオ・イベントを示すパラメータは、

$$\text{GainControl} = 1 - \text{Conf}_{\text{noise}} \quad (10)$$

のように、ノイズ信頼値（ $\text{Conf}_{\text{noise}}$ ）の減少関数に比例するよう設定されることができ

【0220】

ここで、簡単のため、ボリューム平準化器における利得制御に関係するすべてのパラメータ（またはその効果）を表わすために記号GainControlを用いる。ボリューム平準化器の異なる実装は異なる基礎的な意味をもつパラメータの異なる名前をいうからである。単一の用語GainControlを使うことは、一般性を失うことなく、短い表現をもつことができる。本質的には、これらのパラメータを調整することは、もとの利得に線形または非線形の重みを適用することと等価である。一例として、GainControlは、GainControlが小さければ利得が小さくなるよう、利得をスケールングするために直接使われることができる。もう一つの個別的な例として、利得は、特許文献3に記載されるイベント制御信号をGainControlを用いてスケールングすることによって間接的に制御される。同出願はここにその全体において参照によって組み込まれる。この場合、GainControlが小さいときは、ボリューム平準化器の利得の制御は、利得が時間とともに著しく変化することを防ぐよう修正される。GainControlが大きいときは、制御は、平準化器の利得がより自由に変化することを許容されるように修正される。

【0221】

公式(10)において記述される利得制御（もとの利得またはイベント制御信号を直接スケールングすること）を用いて、オーディオ信号の動的利得はそのノイズ信頼値に（線形または非線形に）相関される。信号が高い信頼値でノイズであれば、採取的な利得は、因子（ $1 - \text{Conf}_{\text{noise}}$ ）のため、小さくなる。このように、ノイズ信号を快くない大きなボリュームにブーストすることを避ける。

【0222】

公式(10)の変形例として、(VoIPなどの)用途において背景音にも関心がない場合には、背景音も同様に扱うことができ、小さな利得によって適用される。制御関数は、ノイズ信頼値（ $\text{Conf}_{\text{noise}}$ ）および背景信頼値（ Conf_{bkg} ）の両方を考慮に入れることができる。たとえば、

$$\text{GainControl} = (1 - \text{Conf}_{\text{noise}}) \cdot (1 - \text{Conf}_{\text{bkg}}) \quad (11)$$

上記の公式においては、ノイズおよび背景音の両方が望まれないので、GainControlはノイズの信頼値および背景の信頼値によって等しく影響される。これは、ノイズおよび背景音が同じ重みをもつと見なしうる。状況に依存して、両者は異なる重みをもってもよい。たとえば、ノイズおよび背景音の信頼値（またはそれらの1との差）に異なる係数または異なる指数（および）を与えてもよい。すなわち、公式(11)は

$$\text{GainControl} = (1 - \text{Conf}_{\text{noise}})^{\alpha} \cdot (1 - \text{Conf}_{\text{bkg}})^{\beta} \quad (12)$$

または

$$\text{GainControl} = (1 - \text{Conf}_{\text{noise}})^{\alpha} \cdot (1 - \text{Conf}_{\text{bkg}})^{\beta} \quad (13)$$

と書き直されてもよい。

【0223】

あるいはまた、調整ユニット300Cは、信頼値に基づいて少なくとも一つの優勢なコンテンツ型を考慮するよう構成されていてもよい。たとえば、

10

20

30

40

50

$$\text{GainControl} = 1 - \max(\text{Conf}_{\text{noise}}, \text{Conf}_{\text{bkg}}) \quad (14)$$

公式(11)（およびその諸変形）および公式(14)の両方は、ノイズ信号および背景音信号についての小さな利得を示し、ボリューム平準化器のもとの挙動は、（発話および音楽信号におけるように）ノイズ信頼値および背景信頼値の両方が小さくGainControlが1に近いときにのみ保持される。

【0224】

上記の例は、優勢な干渉コンテンツ型を考慮する。状況に依存して、調整ユニット300Cは、信頼値に基づいて優勢な情報性コンテンツ型を考慮するよう構成されていてもよい。より一般には、調整ユニット300Cは、識別されたオーディオ型が情報性および/または干渉性オーディオ型であるノを含むか否かに関わりなく、信頼値に基づいて少なくとも一つの優勢なコンテンツ型を考慮するよう構成されていてもよい。

10

【0225】

公式(10)のもう一つの例示的な変形として、発話信号が最も情報性のコンテンツであり、ボリューム平準化器のデフォルトの挙動に対して必要な修正がより少ないとすると、制御関数はノイズ信頼値（ $\text{Conf}_{\text{noise}}$ ）および発話信頼値（ $\text{Conf}_{\text{speech}}$ ）の両方を

$$\text{GainControl} = 1 - \text{Conf}_{\text{noise}} \cdot (1 - \text{Conf}_{\text{speech}}) \quad (15)$$

として考慮することができる。この関数を用いると、小さなGainControlが得られるのは、高いノイズ信頼度および低い発話信頼度をもつ信号（たとえば純粋なノイズ）についてのみであり、発話信頼度が高い場合にはGainControlは1に近くなる（よってボリューム平準化器のもとの挙動を保持する）。より一般には、あるコンテンツ型（ $\text{Conf}_{\text{noise}}$ など）の重みが少なくとも一つの他のコンテンツ型（ $\text{Conf}_{\text{speech}}$ など）をもって修正されうると見なされることができる。上記の公式(15)において、発話の信頼度はノイズの信頼度の重み係数を変化させると見なされることができる（公式(12)および(13)における重みに比べると別の種類の重み）。換言すれば、公式(10)では $\text{Conf}_{\text{noise}}$ の係数が1と見なされることができ、一方、公式(15)では、いくつかの他のオーディオ型（発話などだがそれに限られない）がノイズの信頼値の重要性に影響する。よって、 $\text{Conf}_{\text{noise}}$ の重みが発話の信頼値によって修正されると言うことができる。本開示のコンテキストにおいて、用語「重み」はこれを含むように解釈される。すなわち、値の重要性を示すが、必ずしも規格化されていない。1.4節が参照されてもよい。

20

【0226】

別の観点からは、公式(12)および(13)と同様に、指数の形の重みが上記の関数における信頼値に適用されて、異なるオーディオ信号の優先度（または重要性）を示すことができる。たとえば、公式(15)は次のように変更できる。

30

【0227】

$$\text{GainControl} = 1 - \text{Conf}_{\text{noise}} \cdot (1 - \text{Conf}_{\text{speech}}) \quad (16)$$

ここで、 $\text{Conf}_{\text{noise}}$ および $\text{Conf}_{\text{speech}}$ は二つの重みである。これらは、平準器パラメータを修正するためにより大きく反応することが期待される場合にはより小さく設定されることができる。

【0228】

公式(10)～(16)は、自由に組み合わせられて、異なる応用において好適でありうるさまざまな制御関数を形成することができる。音楽信頼値のような他のオーディオ・コンテンツ型の信頼値も同様の仕方で制御関数に簡単に組み込まれることができる。

40

【0229】

GainControlが信号が新しい知覚可能なオーディオ・イベントである度合いを示すパラメータを調整し、そして動的利得を間接的に制御するために使われる場合には（利得はオーディオ・イベント内ではゆっくり変化するが、二つのオーディオ・イベントの境界では急激に変化する）、コンテンツ型の信頼値と最終的な動的利得との間のもう一つの伝達関数があると見なされてもよい。

【0230】

4.2節 種々のコンテキストにおけるコンテンツ型

公式(10)～(16)における上記の制御関数は、ノイズ、背景音、短期的音楽および発話の

50

ようなオーディオ・コンテンツ型の信頼値を考慮に入れるが、映画のメディアおよびVoIPなど、音がどこからくるかのオーディオ・コンテキストは考慮しない。同じオーディオ・コンテンツ型が、たとえば背景音について、異なるオーディオ・コンテキストでは異なる仕方で処理される必要があることがありうる。背景音は、自動車エンジン、爆発および拍手など、さまざまな音を含む。VoIPでは意味がないかもしれないが、映画のメディアでは重要であることがある。これは、関心のあるオーディオ・コンテキストが識別され、異なるオーディオ・コンテキストについて異なる制御関数が設計される必要があることを示している。

【0231】

したがって、調整ユニット300Cはオーディオ信号のコンテンツ型を、オーディオ信号のコンテキスト型に基づいて情報性または干渉性を見なすよう構成されていてもよい。たとえば、ノイズ信頼値および背景信頼値を考慮し、VoIPおよび非VoIPコンテキストを区別することによって、オーディオ・コンテキスト依存制御関数は次のようなものであることができる。

【0232】

```
if オーディオ・コンテキストがVoIP
    GainControl 1 - max(Confnoise, Confbkg)
else
    GainControl 1 - Confnoise                (17)
```

すなわち、VoIPコンテキストでは、ノイズおよび背景音は干渉性コンテンツ型と見なされ、一方、非VoIPコンテキストでは、背景音は情報性コンテンツ型と見なされる。

【0233】

もう一つの例として、発話、ノイズおよび背景の信頼値を考え、VoIPおよび非VoIPコンテキストを区別するオーディオ・コンテキスト依存制御関数は次のようなものであることができる。

【0234】

```
if オーディオ・コンテキストがVoIP
    GainControl 1 - max(Confnoise, Confbkg)
else
    GainControl 1 - Confnoise · (1 - Confspeech)    (18)
```

ここで、発話は情報性コンテンツ型として強調される。

【0235】

音楽も非VoIPコンテキストにおいて重要な情報性の情報であるとする、公式(18)の後半を

$$\text{GainControl } 1 - \text{Conf}_{\text{noise}} \cdot (1 - \max(\text{Conf}_{\text{speech}}, \text{Conf}_{\text{music}})) \quad (19)$$

と拡張できる。

【0236】

実のところ、(10)~(16)における制御関数のそれぞれまたはその変形は、異なるノ対応するオーディオ・コンテキストにおいて適用されることができる。よって、オーディオ・コンテキスト依存制御関数を形成する多数の組み合わせを生成することができる。

【0237】

公式(17)および(18)において区別され、利用されるVoIPおよび非VoIPコンテキストのほか、映画のメディア、長期的音楽およびゲームまたは低品質オーディオおよび高品質オーディオのような他のオーディオ・コンテキストが同様の仕方で利用されることができる。

【0238】

4.3節 コンテキスト型

コンテキスト型は、ノイズのようなわずらわしい音がブーストされすぎるのを避けるようボリューム平準化器を制御するために直接使われることもできる。たとえば、VoIP信頼値が、ボリューム平準化器を、その信頼値が高いときに感度を低くするよう操縦するために使われることができる。

10

20

30

40

50

【0239】

具体的には、VoIP信頼値 $Conf_{VoIP}$ を用いて、ボリューム平準化器のレベルは $(1 - Conf_{VoIP})$ に比例するよう設定されることができ。すなわち、ボリューム平準化器はVoIPコンテンツでは（VoIP信頼値が高いときは）ほとんど非作動にされる。これは、VoIPコンテキストについてボリューム平準化器を無効にする伝統的な手動のセットアップ（プリセット）と整合する。

【0240】

あるいはまた、オーディオ信号の種々のコンテキストについて異なる動的利得範囲を設定することができる。一般に、VL（ボリューム平準化器）量は、オーディオ信号に適用される利得の量をさらに調整し、利得に対するもう一つの（非線形な）重みと見ることができる。ある実施形態では、セットアップは次のようなものであることができる。

【0241】

【表1】

| | 映画のメディア | 長期的音楽 | VOIP | ゲーム |
|-----|---------|-------|-----------|-----|
| VL量 | 高 | 中 | オフ(または最低) | 低 |

さらに、期待されるVL量が各コンテキスト型についてあらかじめ定義されているとする。たとえば、VL量は映画のメディアについては1、VoIPについては0、音楽については0.6、ゲームについては0.3と設定されるが、本願はそれに限定されない。この例によれば、映画のメディアの動的利得の範囲が100%であれば、VoIPの動的利得の範囲は60%である、などとなる。オーディオ分類器200の分類が硬判定に基づく場合には、動的利得の範囲は上記の例のように直接設定されてもよい。オーディオ分類器200の分類が軟判定に基づく場合には、該範囲はコンテキスト型の信頼値に基づいて調整されてもよい。

【0242】

同様に、オーディオ分類器200は、オーディオ信号から複数のコンテキスト型を識別することがあり、調整ユニット300Cは、前記複数のコンテンツ型の重要性に基づいて前記複数のコンテンツ型の信頼値に重み付けすることによって動的利得の範囲を調整するよう構成されていてもよい。

【0243】

一般に、コンテキスト型についても、適切なVL量を適応的に設定するために、(10)~(16)と同様の関数が、その中のコンテンツ型をコンテキスト型で置き換えて、ここで使用されることができる。実際、表1は異なるコンテキスト型の重要性を反映する。

【0244】

別の観点からは、信頼値は、1.4節で論じた規格化された重みを導出するために使われてもよい。表1の各コンテキスト型について特定の量があらかじめ定義されているとすると、公式(9)と同様の公式が適用されることもできる。なお、同様の解決策が、複数のコンテンツ型および他の任意のオーディオ型に適用されてもよい。

【0245】

4.4節 実施形態の組み合わせおよび応用シナリオ

第一部と同様に、上記で論じたすべての実施形態およびその変形は、そのいかなる組み合わせにおいて実装されてもよく、異なる部/実施形態において言及されるが同じまたは同様の機能をもついかなる構成要素も同じまたは別個の構成要素として実装されてもよい。たとえば、4.1節ないし4.3節において述べた解決策の任意の二つ以上が互いと組み合わせられてもよい。そして、これらの組み合わせの任意のものが、第一部~第三部および後述する他の部において記載または含意されている任意の実施形態とさらに組み合わせられてもよい。

【0246】

図21は、もとの短期的セグメント（図21(A)）、パラメータ修正なしの通常のボ

10

20

30

40

50

リ्यूーム平準化器によって処理された短期的セグメント（図 2 1（B））および本願において呈示されるポリリューーム平準化器によって処理された短期的セグメント（図 2 1（C））を比べることによって、本願で提案されるポリリューーム平準化器コントローラの効果を示している。見て取れるように、図 2 1（B）に示される通常のポリリューーム平準化器では、ノイズ（オーディオ信号の後半）のポリリューームもブーストとされ、わずらわしい。対照的に、図 2 1（C）に示される新しいポリリューーム平準化器では、オーディオ信号の実効部分のポリリューームが、ノイズのポリリューームを一見してブーストすることなく、ブーストとされ、聞き手に対して良好な経験を与える。

【0247】

4.5節 ポリリューーム平準化器制御方法

第一部と同様に、上記の実施形態におけるポリリューーム平準化器コントローラを記述する過程で、いくつかのプロセスまたは方法も開示されていることは明らかである。以下では、これらの方法の概要が与えられるが、上記ですでに論じた詳細の一部は繰り返さない。

【0248】

まず、第一部で論じたオーディオ処理方法の実施形態がポリリューーム平準化器について使用されてもよい。ポリリューーム平準化器のパラメータ（単数または複数）が、オーディオ処理方法によって調整されるべきターゲットの一つである。この観点から、オーディオ処理方法はポリリューーム平準化器制御方法でもある。

【0249】

この節では、ポリリューーム平準化器の制御に固有の側面のみが論じられる。制御方法の一般的な側面については、第一部が参照されうる。

【0250】

本願によれば、ポリリューーム平準化器制御方法であって、リアルタイムでオーディオ信号のコンテンツ型を識別し、ポリリューーム平準化器の動的利得をオーディオ信号の情報性コンテンツ型と正に相関させ、ポリリューーム平準化器の動的利得をオーディオ信号の干渉性コンテンツ型と負に相関させることによって、識別されたコンテンツ型に基づいて連続的な仕方でポリリューーム平準化器を調整する方法が提供される。

【0251】

コンテンツ型は、発話、短期的音楽、ノイズおよび背景音を含んでいてもよい。一般に、ノイズは干渉性コンテンツ型と見なされる。

【0252】

ポリリューーム平準化器の動的利得を調整するとき、コンテンツ型の信頼値に基づいて直接調整されてもよいし、あるいはコンテンツ型の信頼値の伝達関数を介して調整されてもよい。

【0253】

すでに述べたように、オーディオ信号は同時に複数のオーディオ型に分類されうる。複数のコンテンツ型に関わるとき、調整動作 1104 は、複数のコンテンツ型の重要性に基づいて複数のコンテンツ型の信頼値に重み付けすることを通じて、または信頼値に基づいて複数のコンテンツ型の効果に重み付けすることを通じて、複数のオーディオ・コンテンツ型の少なくともいくつかを考慮するよう構成されてもよい。特に、調整動作 1104 は、信頼値に基づいて少なくとも一つの優勢なコンテンツ型を考慮するよう構成されていてもよい。オーディオ信号が干渉性のコンテンツ型（単数または複数）および情報性のコンテンツ型（単数または複数）の両方を g ふくむと器、調整動作は、信頼値に基づいて少なくとも一つの優勢な干渉性コンテンツ型を考慮するおよび / または信頼値に基づいて少なくとも一つの優勢な情報性コンテンツ型を考慮するよう構成されてもよい。

【0254】

異なるオーディオ型が互いに影響してもよい。したがって、調整動作 1104 は、あるコンテンツ型の重みを、少なくとも一つの他のコンテンツ型の信頼値を用いて修正するよう構成されていてもよい。

【0255】

10

20

30

40

50

第一部で述べたように、オーディオ信号のオーディオ型の信頼値が平滑化されてもよい。平滑化動作の詳細については、第一部を参照されたい。

【0256】

本方法はさらに、オーディオ信号のコンテキスト型を識別することを含んでいてもよい。ここで、調整動作1104は、コンテキスト型の信頼値に基づいて動的利得の範囲を調整するよう構成されていてもよい。

【0257】

コンテンツ型の役割は、それが位置しているコンテキスト型によって制限される。したがって、オーディオ信号について同時に（すなわち、同じオーディオ・セグメントについて）コンテンツ型情報およびコンテキスト型情報の両方が得られるとき、オーディオ信号のコンテンツ型は、オーディオ信号のコンテキスト型に基づいて情報性または干渉性と判定されてもよい。さらに、異なるコンテキスト型のオーディオ信号におけるコンテンツ型は、オーディオ信号のコンテキスト型に依存して異なる重みを割り当てられてもよい。別の観点からは、コンテンツ型の情報性の性質または干渉性の性質を反映するよう、異なる重み（より大きいまたはより小さい、正の値または負の値）を使うことができる。

10

【0258】

オーディオ信号のコンテキスト型は、VoIP、映画的メディア、長期的音楽およびゲームを含んでいてもよい。コンテキスト型VoIPのオーディオ信号においては、背景音が干渉性のコンテンツ型と見なされてもよい。一方、コンテキスト型非VoIPのオーディオ信号においては、背景および/または発話および/または音楽；が情報性のコンテンツ型と見なされる。他のコンテキスト型は高品質オーディオまたは低品質オーディオを含んでいてもよい。

20

【0259】

複数のコンテンツ型と同様に、オーディオ信号が同時に（すなわち、同じオーディオ・セグメントについて）複数のコンテキスト型情報に分類されるとき、調整動作1104は、複数のコンテキスト型の重要性に基づいて複数のコンテキスト型の信頼値に重み付けすることを通じて、または信頼値に基づいて複数のコンテキスト型の効果に重み付けすることを通じて、複数のオーディオ・コンテンツ型の少なくともいくつかを考慮するよう構成されてもよい。特に、調整動作は、信頼値に基づいて少なくとも一つの優勢なコンテンツ型を考慮するよう構成されていてもよい。

30

【0260】

最後に、本節で論じた方法の実施形態は、第六部および第七部で論じるオーディオ分類方法を使ってもよい。詳細な記述はここでは割愛する。

【0261】

オーディオ処理装置の実施形態と同様に、オーディオ処理方法の実施形態およびその変形の任意の組み合わせが現実的である。他方、オーディオ処理方法の実施形態およびその変形のあらゆる側面は別個の解決策であってもよい。さらに、この節に記載される任意の二つ以上の解決策が互いと組み合わせられてもよく、これらの組み合わせがさらに、本開示の他の部において記載または含意される任意の実施形態と組み合わせられてもよい。

40

【0262】

第五部：等化器コントローラおよび制御方法

等化は、通例、音楽信号に適用されて、「トーン」または「音色」として知られるそのスペクトル・バランスを調整または修正する。伝統的な等化器は、ある種の楽器を強調したりまたは望まれない音を除去したりするために、ユーザーが個々の周波数帯域における周波数応答（利得）の全体的なプロファイル（曲線または形状）を構成設定できるようにする。ウィンドウズ・メディア・プレーヤーのような一般的な音楽プレーヤーは、種々のジャンルの音楽の最良の聴取経験を得るために、各周波数帯域における利得を調整するためのグラフィック・イコライザーを提供し、ロック、ラップ、ジャズおよびフォークのような種々の音楽ジャンルについての等化器プリセットの集合をも提供する。ひとたびプリセットが選択され、プロファイルが設定されたら、プロファイルが手動で修正されるまで

50

、同じ等化利得が信号に対して適用される。

【0263】

対照的に、動的等化器は、所望される音色またはトーンに関してスペクトル・バランスの全体的な一貫性を保持するために、各周波数帯域における等化利得を自動的に調整するすべを提供する。この一貫性は、オーディオのスペクトル・バランスを連続的にモニタリングし、それを所望されるプリセット・スペクトル・バランスと比較し、オーディオの元のスペクトル・バランスを所望されるスペクトル・バランスに変換するための適用される等化利得を動的に調整することによって達成される。所望されるスペクトル・バランスは、手動で選択されるまたは処理前に事前設定される。

【0264】

両方の種類の等化器は、次の不都合な点を共有する。最良の等化プロファイル、所望されるスペクトル・バランスまたは関係したパラメータは手動で選択される必要があり、再生時にオーディオ・コンテンツに基づいて自動的に修正されることができない。オーディオ・コンテンツ型を決定することは、種々のオーディオ信号について全体的な良好な品質を提供するために非常に重要である。たとえば、異なるジャンルのものなど、異なる音楽片は異なる等化プロファイルを必要とする。

【0265】

(音楽だけでなく)任意の種類のオーディオ信号が入力されることが可能である等化器システムでは、等化器パラメータはコンテンツ型に基づいて調整される必要がある。たとえば、等化器は通例、音楽信号に対しては有効にされるが、発話信号に対しては、発話の音色を大きく変えすぎて信号が不自然に聞こえるようにしてしまうことがあるので、無効にされる。

【0266】

この問題に少なくとも部分的に対処するために、本願は、第一部で論じた実施形態に基づいて等化器を制御することを提案する。

【0267】

第二部～第四部と同様に、第一部で論じたオーディオ処理装置および方法の個別的な例として、等化器408は、第一部で論じたすべての実施形態およびそこで開示されたそれらの実施形態の任意の組み合わせを使用してもよい。特に、等化器408を制御する場合、図1～図10に示されるようなオーディオ処理装置100におけるオーディオ分類器200および調整ユニット300は、図22に示されるような等化器408コントローラ2000を構成してもよい。この実施形態では、調整ユニットは等化器408に固有なので、300Dと称されてもよい。

【0268】

すなわち、第一部の開示に基づき、等化器コントローラ2200は、オーディオ信号のオーディオ型を連続的に識別するオーディオ分類器200と、識別されたオーディオ型の信頼値に基づいて連続的な仕方で等化器を調整する調整ユニット300Dとを有していてもよい。同様に、オーディオ分類器200は、オーディオ・コンテンツ分類器202およびオーディオ・コンテキスト分類器204のうちの少なくとも一つを含んでいてもよく、ボリューム等化器コントローラ2200はさらに、型平滑化ユニット712、パラメータ平滑化ユニット814およびタイマー916のうちの少なくとも一つを含んでいてもよい。

【0269】

したがって、この部においては、第一部ですでに記述した内容を繰り返すことはせず、単にこの部のいくつかの固有の例を与える。

【0270】

5.1節 コンテンツ型に基づく制御

一般に、音楽、発話、背景音およびノイズのような一般的なオーディオ・コンテンツ型について、等化器は異なるコンテンツ型に対して異なるように設定されるべきである。伝統的なセットアップと同様に、等化器は、自動的に音楽信号に対して有効にされるが、発

10

20

30

40

50

話に対しては無効にされることができる。あるいはより連続的な仕方で、音楽信号に対しては高い等化レベルを、発話信号に対しては低い等化レベルを設定することができる。このようにして、等化器の等化レベルは異なるオーディオ・コンテンツについて自動的に設定されることができる。

【0271】

特に音楽について、優勢な源をもつ音楽片に対しては等化器はあまりうまく機能しないことが観察される。不適切な等化が適用されると、優勢な源の音色が著しく変化して不自然に聞こえることがあるからである。これを考えると、優勢の源がある音楽片に対しては強く伊藤かレベルを設定するほうがよい。一方、優勢な源のない音楽片に対しては等化レベルは高く保つことができる。この情報を用いて、等化器は種々の音楽コンテンツについて自動的に等化レベルを設定することができる。

10

【0272】

音楽は、ジャンル、楽器およびリズム、テンポおよび音色を含む一般的な音楽特性といった種々の属性に基づいてグループ化することもできる。異なる音楽ジャンルについて異なる等化器プリセットが使用されるのと同じように、これらの音楽グループ/クラスターもそれぞれ自身の最適等化プロファイルまたは等化器曲線（伝統的な等化器の場合）または最適な所望されるスペクトル・バランス（動的等化器の場合）をもっていてよい。

【0273】

上述したように、等化器は一般には音楽コンテンツに対しては有効にされるが、発話については無効にされる。等化器は、音色変化のため、ダイアログをそれほどよく聞こえさせないことがありうるからである。それを自動的に達成する一つの方法は、等化器レベルをコンテンツに、特にオーディオ・コンテンツ分類モジュールから得られる音楽信頼値および/または発話信頼値に関係付けることである。ここで、等化レベルは、適用される等化器利得の重みとして説明できる。レベルが高いほど、適用される等化は強い。たとえば、等化レベルが1であれば、フル等化プロファイルが適用される。等化レベルが0であれば、すべての利得が対応して0dBとなり、よって非等化が適用される。等化レベルは、等化器アルゴリズムの種々の実装において種々のパラメータによって表わされることがある。このパラメータの例は、特許文献2において実装されるような等化器重みである。同文献はここにその全体において参照によって組み込まれる。

20

【0274】

等化レベルを調整するために、さまざまな制御方式が設計されることができる。たとえば、オーディオ・コンテンツ型情報では、発話信頼値または音楽信頼値が等化レベルを設定するために

30

$$L_{eq} = Conf_{music} \quad (20)$$

または

$$L_{eq} = 1 - Conf_{speech} \quad (21)$$

として使用されることができる。ここで、 L_{eq} は等化レベルであり、 $Conf_{music}$ および $Conf_{speech}$ は音楽および発話の信頼値を表わす。

【0275】

すなわち、調整ユニットは300Dは、等化レベルを短期的音楽の信頼値と正に相関させるまたは等化レベルを発話の信頼値と負に相関させるよう構成されていてもよい。

40

【0276】

発話信頼値および音楽信頼値はさらに、等化レベルを設定するために統合して利用されることができる。一般的な発想は、等化レベルが高いのは、音楽信頼値が高くかつ発話信頼値が低いときにのみであり、他の場合には等化レベルは低いということである。たとえば、

$$L_{eq} = Conf_{music}(1 - Conf_{speech}) \quad (22)$$

ここで、発話信頼値は、頻繁に起こりうる、音楽信号における0でない発話信頼値を扱うために、乗される。上記の公式を用いれば、等化は、発話成分のない純粋な音楽信号に対してはフルに適用される（1に等しいレベルで）。第一部で述べたように、はコンテ

50

ンツ型の重要性に基づく重み付け係数と見なされてもよく、典型的には1ないし2に設定されることができる。

【0277】

発話の信頼値により大きな重みをおくならば、調整ユニット300Dは、コンテンツ型発話についての信頼値がある閾値より大きいときに等化器408を無効にするよう構成されていてもよい。

【0278】

上記の記述では、音楽および発話のコンテンツ型が例に取られている。代替的または追加的に、背景音および/またはノイズの信頼値も考慮されてもよい。特に、調整ユニット300Dは、等化レベルを背景の信頼値と正に相関させるおよび/または等化レベルをノイズの信頼値と負に相関させるよう構成されていてもよい。

10

【0279】

もう一つの例として、信頼値は1.4節で論じた規格化された重みを導出するために使われてもよい。期待される等化レベルが各コンテンツ型についてあらかじめ定義されるとすると(たとえば、音楽については1、発話については0、ノイズおよび背景については0.5)、公式(8)と同様の公式が厳密に適用されることができる。

【0280】

等化レベルはさらに、遷移点において可聴アーチファクトを導入しうる急激な変化を避けるために、平滑化されてもよい。これは、1.5節で述べたパラメータ平滑化ユニット814を用いてできる。

20

【0281】

5.2節 音楽における優勢な源の確からしさ

優勢な源をもつ音楽が高い等化レベルを適用されることを避けるために、等化レベルはさらに、音楽片が優勢な源を含むかどうかを示す信頼値 $Conf_{dom}$ に相関させられてもよい。たとえば、

$$L_{eq} = 1 - Conf_{dom} \quad (23)$$

【0282】

このようにして、等化レベルは優勢な源をもつ音楽片に対しては低く、優勢な源のない音楽片については高い。

【0283】

ここで、優勢な源をもつ音楽の信頼値が記述されているが、優勢な源をもたない音楽の信頼値を使うこともできる。すなわち、調整ユニット300Dは、等化レベルを優勢な源をもたない短期的音楽の信頼値と正に相関させるおよび/または等化レベルを優勢な源をもつ短期的音楽の信頼値と負に相関させるよう構成されてもよい。

30

【0284】

1.1節で述べたように、音楽および発話と、優勢な源をもつまたはもたない音楽とは、異なる階層レベルでのコンテンツ型であるが、並行して考えられることができる。上記のような優勢な源の信頼値および発話および音楽の信頼値を合同して考えることにより、等化レベルは公式(20)~(21)のうちの少なくとも一つを(23)と組み合わせることによって設定されることができる。一例は、これら三つの公式すべてを組み合わせ

40

$$L_{eq} = Conf_{music}(1 - Conf_{speech})(1 - Conf_{dom}) \quad (24)$$

とすることである。

【0285】

一般性のために、公式(22)のようにして、コンテンツ型の重要性に基づく異なる重みがさらに異なる信頼値に適用されることができる。

【0286】

もう一つの例として、 $Conf_{dom}$ はオーディオ信号が音楽であるときにのみ計算されるとして、階段関数が次のように設計されることができる。

【0287】

【数 1 1】

$$L_{eq} = \begin{cases} (1 - Conf_{dom}) & Conf_{music} > \text{閾値} \\ Conf_{music} (1 - conf_{speech}^\alpha) & \text{それ以外} \end{cases} \quad (25)$$

この関数は、分類システムがオーディオが音楽であることをかなり確証する（音楽信頼値が閾値より大きい）場合には優勢なスコアの信頼値に基づいて等化レベルを設定し、そうでない場合には、音楽および発話信頼値に基づいて設定される。すなわち、調整ユニット 300D は、短期的音楽についての信頼値が閾値より大きいときに、優勢な源がない／ある短期的音楽を考慮するよう構成されていてもよい。もちろん、公式(25)における前半または後半が公式(20)ないし(24)のように修正されてもよい。

10

【0288】

1.5節で論じたのと同じ平滑化方式が適用されることもでき、時定数がさらに、優勢な源をもつ音楽から優勢な源のない音楽への遷移または優勢な源のない音楽から優勢な源をもつ音楽への遷移といった、遷移型に基づいて設定されることができる。この目的のために、公式(4')と同様の公式が適用されることもできる。

【0289】

5.3節 等化器プリセット

20

オーディオ・コンテンツ型の信頼値に基づいて等化レベルを適応的に調整することのほか、種々のオーディオ・コンテンツについて、そのジャンル、楽器または他の特性に依存して、適切な等化プロファイルまたは所望されるスペクトル・バランス・プリセットが自動的に選ばれることもできる。同じジャンルをもつ、同じ楽器を含むまたは同じ音楽特性をもつ音楽は同じ等化プロファイルまたは所望されるスペクトル・バランス・プリセットを共有することができる。

【0290】

一般性のために、同じジャンル、同じ楽器または同様の音楽アトリビュートをもつ音楽グループを表わすために用語「音楽クラスター」を使う。これは、1.1節で述べたオーディオ・コンテンツ型のもう一つの階層レベルと見なすことができる。適切な等化プロファイル、等化レベルおよび／または所望されるスペクトル・バランス・プリセットがそれぞれの音楽クラスターに関連付けられてもよい。等化プロファイルは、音楽信号に対して適用される利得曲線であり、異なる音楽ジャンル（クラシック、ロック、ジャズおよびフォークなど）について使用される等化器プリセットの任意のものであり、所望されるスペクトル・バランス・プリセットは各クラスターについての所望される音色を表わす。図23は、ドルビー・ホーム・シアター技術において実装される所望されるスペクトル・バランス・プリセットのいくつかの例を示している。それぞれは、可聴周波数範囲にわたる所望されるスペクトル形状を記述する。この形状は連続的に、はいてくるオーディオのスペクトル形状と比較され、等化利得は、はいてくるオーディオのスペクトル形状の、プリセットのスペクトル形状への変換との比較から、計算される。

30

40

【0291】

新しい音楽片については、最も近いクラスターが決定されることができる（硬判定）または各音楽クラスターに関する信頼値が計算されることができる（軟判定）。この情報に基づいて、適正な等化プロファイルまたは所望されるスペクトル・バランス・プリセットが所与の音楽片について決定されることができる。最も簡単な方法は、

$$P_{eq} = P_c \quad (26)$$

として、最良のマッチしたクラスターの対応するプロファイルを割り当てることである。ここで、 P_{eq} は推定される等化プロファイルまたは所望されるスペクトル・バランス・プリセットであり、 c^* は最良のマッチした音楽クラスター（優勢なオーディオ型）のインデックスであり、これは最も高い信頼値をもつクラスターを拾うことによって得られる。

50

【0292】

さらに、0より大きい信頼値をもつ音楽クラスターが二つ以上あることがある。つまり、音楽片は、それらのクラスターと多少なりとも同様なアトリビュートをもつ。たとえば、音楽片は複数の楽器をもつことがあり、あるいは複数のジャンルのアトリビュートをもつことがある。このことは、最も近いクラスターのみを使うのではなく、すべてのクラスターを考慮することによって適正な等化プロファイルを推定するもう一つの方法を着想させる。たとえば、重み付けされた和が使用されることができ。

【0293】

【数12】

$$P_{eq} = \sum_{c=1}^N w_c P_c \quad (27)$$

10

ここで、Nはあらかじめ定義されたクラスターの数であり、 w_c は各あらかじめ定義された音楽クラスター（インデックスcをもつ）に関する設計されたプロファイル P_c の重みである。この重みは、その信頼値に基づいて1に規格化されるべきである。このようにして、推定されるプロファイルは、音楽クラスターのプロファイルの混合となる。たとえば、ジャズおよびロックのアトリビュートの両方をもつ音楽片について、推定されたプロファイルは中間の何かであろう。

20

【0294】

いくつかの用途では、公式(27)に示されるすべてのクラスターに関与させることを望まないこともある。該クラスターの部分集合 現在の音楽片に最も関係している諸クラスター のみが考慮される必要があり、公式(27)は次のように微修正できる。

【0295】

【数13】

$$P_{eq} = \sum_{c'=1}^{N'} w_{c'} P_{c'} \quad (28)$$

30

ここで、 N' は考慮されるべきクラスターの数であり、 c' はそれらのクラスターを信頼値に基づいて降順にソートしたあとのクラスター・インデックスである。部分集合を使うことにより、最も関係した諸クラスターにより焦点を当てることができ、それほど重要でないものを除外することができる。換言すれば、調整ユニット300Dは、信頼値に基づいて、少なくとも一つの優勢なオーディオ型を考慮するよう構成されてもよい。

【0296】

上記の記述では、音楽クラスターが例として取られている。実のところ、これらの解決策は、1.1節で論じた任意の階層レベルのオーディオ型に適用可能である。よって、一般に、調整ユニット300Dは、等化レベルおよび/または等化プロファイルおよび/またはスペクトル・バランス・プリセットを各オーディオ型に割り当てるよう構成されてもよい。

40

【0297】

5.4節 コンテキスト型に基づく制御

これまでの節では、さまざまなコンテンツ型に焦点を当てている。本節で論じるさらなる実施形態では、代替的または追加的にコンテキスト型が考慮されてもよい。

【0298】

一般に、等化器は、明らかな音色変化のため映画のメディアにおけるダイアログをあまり良好に聞こえなくしてしまうことがあるので、音楽については有効にされるが、映画のメディアについては無効にされる。このことは、等化レベルは、長期的音楽の信頼値およ

50

び / または映画のメディアの信頼値に関係付けられうることを示す :

$$L_{eq} = \text{Conf}_{\text{MUSIC}} \quad (29)$$

または

$$L_{eq} = 1 - \text{Conf}_{\text{MOVIE}} \quad (30)$$

ここで、 L_{eq} は等化レベル、 $\text{Conf}_{\text{MUSIC}}$ および $\text{Conf}_{\text{MOVIE}}$ は長期的音楽および映画のメディアの信頼値を表わす。

【0299】

すなわち、調整ユニット300Dは、等化レベルを長期的音楽の信頼値と正に相関させるまたは等化レベルを映画のメディアの信頼値と負に相関させるよう構成されていてもよい。

10

【0300】

すなわち、映画のメディア信号については、映画のメディア信頼値は高く（または音楽信頼度は低い）、よって等化レベルは低い。他方、音楽信号については、映画のメディア信頼値は低く（または音楽信頼度は高い）、よって等化レベルは高い。

【0301】

公式(29)および(30)に示される解決策は、公式(22)ないし(25)と同様に修正されてもよく、および / または公式(22)ないし(25)に示される解決策の任意のものと組み合わせられてもよい。

【0302】

追加的または代替的に、調整ユニット300Dは、等化レベルをゲームの信頼値と負に相関させるよう構成されていてもよい。

20

【0303】

もう一つの実施形態では、信頼値は、1.4節で論じた規格化された重みを導出するために使われてもよい。期待される等化レベル / プロファイルが各コンテキスト型についてあらかじめ定義されているとすると（等化プロファイルは下記の表2に示されている）、公式(9)と同様の公式適用されることもできる。

【0304】

【表2】

| | 映画のメディア | 長期的音楽 | VOIP | ゲーム |
|----------|---------|---------|---------|---------|
| 等化プロファイル | プロファイル1 | プロファイル2 | プロファイル3 | プロファイル4 |

30

ここで、いくつかのプロファイルでは、映画のメディアおよびゲームのようなある種のコンテキスト型について等化器を無効にする方法として、すべての利得が0に設定されることができ。

【0305】

5.5節 実施形態の組み合わせおよび応用シナリオ

第一部と同様に、上記で論じたすべての実施形態およびその変形は、そのいかなる組み合わせにおいて実装されてもよく、異なる部 / 実施形態において言及されるが同じまたは同様の機能をもついかなる構成要素も同じまたは別個の構成要素として実装されてもよい。

40

【0306】

たとえば、5.1節ないし5.4節において述べた解決策の任意の二つ以上が互いと組み合わせられてもよい。そして、これらの組み合わせの任意のものが、第一部～第四部および後述する他の部において記載または含意されている任意の実施形態とさらに組み合わせられてもよい。

【0307】

5.6節 等化器制御方法

50

第一部と同様に、上記の実施形態における等化器コントローラを記述する過程で、いくつかのプロセスまたは方法も開示されていることは明らかである。以下では、これらの方法の概要が与えられるが、上記ですでに論じた詳細の一部は繰り返さない。

【0308】

まず、第一部で論じたオーディオ処理方法の実施形態が等化器について使用されてもよい。等化器のパラメータ（単数または複数）が、オーディオ処理方法によって調整されるべきターゲットの一つである。この観点から、オーディオ処理方法は等化器制御方法でもある。

【0309】

この節では、等化器の制御に固有の側面のみが論じられる。制御方法の一般的な側面については、第一部が参照されうる。

10

【0310】

諸実施形態によれば、等化器制御方法が、リアルタイムでオーディオ信号のオーディオ型を識別し、識別されたオーディオ型の信頼値に基づいて連続的な仕方で等化器を調整することを含むうる。

【0311】

本願の他の部分と同様に、対応する信頼値をもつ複数のオーディオ型が関与するとき、調整の動作1104は、複数のオーディオ型の重要性に基づいて複数のオーディオ型の信頼値に重み付けすることを通じて、または信頼値に基づいて複数のオーディオ型の効果に重み付けすることを通じて、複数のオーディオ型の少なくともいくつかを考慮するよう構成されてもよい。特に、調整動作1104は、信頼値に基づいて少なくとも一つの優勢なコンテンツ型を考慮するよう構成されていてもよい。

20

【0312】

第一部で述べたように、調整されたパラメータ値が平滑化されてもよい。1.5節および1.8節が参照され、ここでは詳細は記述は割愛する。

【0313】

オーディオ型はコンテンツ型またはコンテキスト型でありうる。コンテンツ型に関わる時、調整動作1104は、等化レベルを短期的音楽の信頼値と正に相関させるおよび/または等化レベルを発話の信頼値と負に相関させるよう構成されていてもよい。追加的または代替的に、調整動作は、等化レベルを背景の信頼値と正に相関させるおよび/または等化レベルをノイズの信頼値と負に相関させるよう構成されていてもよい。

30

【0314】

コンテキスト型に関わる時、調整動作1104は、等化レベルを長期的音楽の信頼値と正に相関させるおよび/または等化レベルを映画のメディアおよび/またはゲームの信頼値と負に相関させるよう構成されていてもよい。

【0315】

短期的音楽のコンテンツ型については、調整動作1104は、等化レベルを、優勢な源のない短期的音楽の信頼値と正に相関させるおよび/または等化レベルを、優勢な源をもつ短期的音楽の信頼値と負に相関させるよう構成されていてもよい。これは、短期的音楽についての信頼値がある閾値より大きいときにのみ行なわれることができる。

40

【0316】

等化レベルを調整することのほか、等化器の他の側面が、オーディオ信号のオーディオ型（単数または複数）の信頼値（単数または複数）に基づいて調整されてもよい。たとえば、調整動作1104は、等化レベルおよび/または等化プロファイルおよび/またはスペクトル・バランス・プリセットをそれぞれのオーディオ型に割り当てるよう構成されていてもよい。

【0317】

オーディオ型の具体例については第一部が参照されうる。

【0318】

オーディオ処理装置の実施形態と同様に、オーディオ処理方法の実施形態およびその変

50

形の任意の組み合わせが現実的である。他方、オーディオ処理方法の実施形態およびその変形のあらゆる側面は別個の解決策であってもよい。さらに、この節に記載される任意の二つ以上の解決策が互いと組み合わせられてもよく、これらの組み合わせがさらに、本開示の他の部において記載または含意される任意の実施形態と組み合わせられてもよい。

【0319】

第六部：オーディオ分類器および分類方法

1.1節および1.2節で述べたように、さまざまな階層レベルのコンテンツ型およびコンテキスト型を含む本願で論じられるオーディオ型は、機械学習ベースの方法を含め何らかの既存の分類方式を用いて分類または識別されることができる。この部および次の部では、本願は、これまでの部で言及されたコンテンツ型を分類するための分類器および方法のいくつかの新たな側面を提案する。

10

【0320】

6.1節 コンテンツ型分類に基づくコンテキスト分類器

これまでの部で述べたように、オーディオ分類器200は、オーディオ信号のコンテンツ型を識別するおよび/またはオーディオ信号のコンテキスト型を識別するために使われる。したがって、オーディオ分類器200は、オーディオ・コンテンツ分類器202および/またはオーディオ・コンテキスト分類器204を有していてもよい。オーディオ・コンテンツ分類器202および/またはオーディオ・コンテキスト分類器204を実装するための既存の技法を採用するとき、両分類器は互いから独立でありうるが、いくつかの特徴を共有していてもよく、よって該特徴を抽出するためのいくつかの方式を共有していてもよい。

20

【0321】

この部および次の第七部では、本願で提案される新たな側面に従って、オーディオ・コンテキスト分類器204は、オーディオ・コンテンツ分類器202の結果を利用してもよい。すなわち、オーディオ分類器200は、オーディオ信号のコンテンツ型を識別するオーディオ・コンテンツ分類器202と；オーディオ・コンテンツ分類器202の結果に基づいてオーディオ信号のコンテキスト型を識別するオーディオ・コンテキスト分類器204とを有する。よって、オーディオ・コンテンツ分類器202の分類結果は、オーディオ・コンテキスト分類器204およびこれまでの部で論じた調整ユニット300（または調整ユニット300Aないし300D）の両方によって使用されてもよい。しかしながら、図面には示されていないものの、オーディオ分類器200は、調整ユニット300およびオーディオ・コンテキスト分類器204によってそれぞれ使用される二つのオーディオ・コンテンツ分類器202を含んでいてもよい。

30

【0322】

さらに、1.2節において論じられるように、特に複数のオーディオ型を分類するとき、オーディオ・コンテンツ分類器202またはオーディオ・コンテキスト分類器204は、互いと協働する分類器のグループをなしてもよい。ただし、一つの単一の分類器として実装されることも可能である。

【0323】

1.1節で論じたように、コンテンツ型は、一般に数フレームないし数十フレームのオーダーの長さ（たとえば1秒）をもつ短期的オーディオ・セグメントに関するオーディオ型の種類であり、コンテキスト型は、一般に数秒ないし数十秒のオーダーの長さ（たとえば10秒）をもつ長期的オーディオ・セグメントに関するオーディオ型の種類である。よって、「コンテンツ型」および「コンテキスト型」に対応して、必要なときはそれぞれ「短期」「長期」を使う。しかしながら、次の第七部で論じるように、コンテキスト型は比較的長い時間スケールでのオーディオ信号の属性を示すものであるものの、短期的オーディオ・セグメントから抽出される特徴に基づいて識別されることもできる。

40

【0324】

ここで、図24を参照して、オーディオ・コンテンツ分類器202およびオーディオ・コンテキスト分類器204の構造に目を向ける。

50

【 0 3 2 5 】

図 2 4 に示されるように、オーディオ・コンテンツ分類器 2 0 2 は、それぞれオーディオ・フレームのシーケンスを含む短期的オーディオ・セグメントから短期的特徴を抽出する短期的特徴抽出器 2 0 2 2 と；長期的オーディオ・セグメント中の短期的セグメントのシーケンスをそれぞれの短期的特徴を使って短期的オーディオ型に分類する短期的分類器 2 0 2 4 とを有していてもよい。短期的特徴抽出器 2 0 2 2 および短期的分類器 2 0 2 4 の両方は、既存の技法を用いて実装されてもよいが、後述する 6 . 3 節において、短期的特徴抽出器 2 0 2 2 についていくつかの修正も提案される。

【 0 3 2 6 】

短期的分類器 2 0 2 4 は、短期的セグメントのシーケンスの各セグメントを次の短期的オーディオ型（コンテンツ型）のうち少なくとも一つに分類するよう構成されていてもよい：発話、短期的音楽、背景音およびノイズ。これについては 1 . 1 節で説明してある。各コンテンツ型はさらに、1 . 1 節で論じたような、ただしそれに限定されないより低い階層レベルでのコンテンツ型にさらに分類されてもよい。

10

【 0 3 2 7 】

当技術分野において知られているように、分類されたオーディオ型の信頼値は、短期的分類器 2 0 2 4 によって得られてもよい。本願では、何らかの分類器の動作に言及するとき、明示的に記録されるか否かによらず、必要であれば同時に信頼値が得られることは理解される。オーディオ型分類の例は、非特許文献 1 に見出されうる。同文献はここにその全体において参照によって組み込まれる。

20

【 0 3 2 8 】

他方、図 2 4 に示されるように、オーディオ・コンテキスト分類器 2 0 4 は、長期的オーディオ・セグメント内の短期的セグメントのシーケンスに関して短期的分類器の結果の統計量を、長期的特徴として計算するための統計量抽出器 2 0 2 4 と；長期的特徴を使って、長期的オーディオ・セグメントを長期的オーディオ型に分類する長期的分類器 2 0 4 4 とを有していてもよい。同様に、統計量抽出器 2 0 4 2 および長期的分類器 2 0 4 4 の両方は、既存の技法を用いて実装されてもよいが、次の 6 . 2 節において統計量抽出器 2 0 4 2 についていくつかの修正が提案される。

【 0 3 2 9 】

長期的分類器 2 0 4 4 は長期的オーディオ・セグメントを次の長期的オーディオ型（コンテキスト型）の少なくとも一つに分類するよう構成されていてもよい：映画的メディア、長期的音楽、ゲームおよびVoIP。これについては 1 . 1 節で説明してある。代替的または追加的に、長期的分類器 2 0 4 4 は長期的オーディオ・セグメントを、1 . 1 節で説明したVoIPまたは非VoIPに分類するよう構成されていてもよい。代替的または追加的に、長期的分類器 2 0 4 4 は長期的オーディオ・セグメントを、1 . 1 節で説明した高品質オーディオまたは低品質オーディオに分類するよう構成されていてもよい。實際上、用途/システムの必要性に基づいて、さまざまなターゲット・オーディオ型が選ばれることができ、トレーニングされることができる。

30

【 0 3 3 0 】

短期的セグメントおよび長期的セグメント（ならびに 6 . 3 節で論じるフレーム）の意味および選択については、1 . 1 節が参照されうる。

40

【 0 3 3 1 】

6 . 2 節 長期的特徴の抽出

図 2 4 に示されるように、ある実施形態では、統計量抽出器 2 0 4 2 のみが、短期的分類器 2 0 2 4 の結果から長期的特徴を抽出するために使用される。長期的特徴として、次のうちの少なくとも一つが統計量抽出器 2 0 4 2 によって計算されてもよい：分類されるべき長期的セグメント内の短期的セグメントの短期的オーディオ型の信頼値の平均および分散、短期的セグメントの重要度によって重み付けされた前記平均および分散、各短期的オーディオ型の出現頻度および分類されるべき長期的セグメント内の種々の短期的オーディオ型の間の遷移の頻度。

50

【0332】

図25において、各短期的セグメント（長さ1s）における発話および短期的音楽信頼値の平均を示す。比較のために、セグメントは三つの異なるオーディオ・コンテキストから抽出されている：映画のメディア（図25（A））、長期的音楽（図25（B））およびVoIP（図25（C））。映画のメディア・コンテキストについては、発話型についてまたは音楽型について高い信頼値が得られ、これら二つのオーディオ型の間で頻りに交替することが観察できる。対照的に、長期的音楽のセグメントは安定した高い短期的音楽の信頼値および比較的安定した低い発話信頼値を与える。一方、VoIPのセグメントは安定して低い短期的音楽の信頼値を与えるが、VoIP会話の間の休止のため揺動する発話信頼値を与える。

10

【0333】

各オーディオ型についての信頼値の分散も種々のオーディオ・コンテキストを分類するための重要な特徴である。図26は、遺影が的メディア、長期的音楽およびVoIPオーディオ・コンテキストにおける発話、短期的音楽、背景およびノイズの信頼値の分散のヒストグラムを与えている（横軸はデータセット中の信頼値の分散であり、縦軸はデータセット中の分散値sの各ビンにおける生起数であり、これは分散値の各ビンの正規確率を示すよう規格化されることができる）。映画のメディアについては、発話、短期的音楽および背景の信頼値の分散のすべては比較的高く、幅広く分布している。これは、これらのオーディオ型の信頼値が強く変化していることを示す。長期的な音楽については、発話、短期的音楽、背景およびノイズの信頼値の分散のすべては比較的低く、狭く分布している。これは、これらのオーディオ型の信頼値が安定を保っていることを示す。発話信頼値は一定して低く保たれ、音楽信頼値は一定した高く保たれる。VoIPについては、短期的音楽の信頼値の分散は低く、狭く分布している一方、発話の信頼値の分散は比較的幅広く分布している。これは、VoIP会話の間の頻りに休止のためである。

20

【0334】

重み付けされた平均および分散を計算する際に使われる重みについて、これらは各短期的セグメントの重要度に基づいて決定される。短期的セグメントの重要度は、そのエネルギーまたはラウドネスによって測定されてもよい。エネルギーやラウドネスは多くの既存の技法を用いて推定されることができる。

【0335】

分類されるべき長期的セグメントにおける各短期的オーディオ型の出現頻度は、該長期的セグメント内の短期的セグメントが分類された各オーディオ型のカウントを、長期的セグメントの長さで規格化したものである。

30

【0336】

分類されるべき長期的セグメント内の種々の短期的オーディオ型の間での遷移の頻度は、分類されるべき長期的セグメント内の隣り合う短期的セグメント間のオーディオ型変化のカウントを、長期的セグメントの長さで規格化したものである。

【0337】

図25を参照して信頼値の平均および分散を論じるとき、各短期的オーディオ型の出現頻度およびそれら種々の短期型オーディオ型の間での遷移頻度も実際には触れられる。これらの特徴は、オーディオ・コンテキスト分類にも大きく関わってくる。たとえば、長期的音楽はほとんど短期的音楽オーディオ型を含み、よって短期的音楽の高い出現頻度をもつ。一方、VoIPはほとんど発話および休止を含み、よって発話またはノイズの高い出現頻度をもつ。もう一つの例として、映画のメディアは、長期的音楽またはVoIPよりも頻りに異なる短期的オーディオ型の間で遷移し、よって一般に短期的音楽、発話および背景の間でより高い遷移頻度をもつ。VoIPは通例、他よりも発話とノイズの間でより頻りに遷移し、よって発話とノイズの間でのより高い遷移頻度をもつ。

40

【0338】

一般に、長期的セグメントは同じアプリケーション/システムでは同じ長さであると想定する。そうであれば、各短期的オーディオ型の出現カウントおよび長期的セグメント内

50

の異なる短期的オーディオ型の間の変換カウントは規格化なしで直接使用されてもよい。長期的セグメントの長さが可変であれば、上述した出現頻度および変換頻度が使われるべきである。本願の請求項は、両方の状況をカバーするものと解釈されるべきである。

【0339】

追加的または代替的に、オーディオ分類器200（またはオーディオ・コンテキスト分類器204）はさらに、長期的オーディオ・セグメント内の短期的セグメントのシーケンスの短期的特徴に基づいて長期的オーディオ・セグメントからさらなる長期的特徴を抽出するための長期的特徴抽出器2046（図27）を含んでいてもよい。換言すれば、長期的特徴抽出器2046は、短期的分類器2024の分類結果を使わず、短期的特徴抽出器2022によって抽出された短期的特徴を直接使って、長期的分類器2044によって使用されるべきいくつかの長期的特徴を導出する。長期的特徴抽出器2046および統計量抽出器2042は独立してまたは合同して使われてもよい。換言すれば、オーディオ分類器200は、長期的特徴抽出器2046または統計量抽出器2042の一方または両方を含んでいてもよい。

10

【0340】

任意の特徴が、長期的特徴抽出器2046によって抽出されることができる。本願では、長期的特徴として、短期的特徴抽出器2022からの短期的特徴の次の統計量の少なくとも一つを計算することが提案される：平均、分散、重み付けされた平均、重み付けされた分散、高平均、低平均および高平均と低平均の間の比（コントラスト）。

20

【0341】

分類されるべき長期的セグメント内の短期的セグメントから抽出された短期的特徴の平均および分散。

【0342】

分類されるべき長期的セグメント内の短期的セグメントから抽出された短期的特徴の重み付けされた平均および分散。短期的特徴は、各短期的セグメントの、たった今述べたそのエネルギーまたはラウドネスを用いて測定される重要度に基づいて重み付けされる。

【0343】

高平均（high average）：分類されるべき長期的セグメント内の短期的セグメントから抽出された、選択された短期的特徴の平均。短期的特徴は、次の条件のうちの少なくとも一つを満たすときに選択される：ある閾値より大きい；または他のすべての短期的特徴より低くない、短期的特徴のあらかじめ決定された割合以内、たとえば短期的特徴の上位10%以内。

30

【0344】

低平均（low average）：分類されるべき長期的セグメント内の短期的セグメントから抽出された、選択された短期的特徴の平均。短期的特徴は、次の条件のうちの少なくとも一つを満たすときに選択される：ある閾値より小さい；または他のすべての短期的特徴より高くない、短期的特徴のあらかじめ決定された割合以内、たとえば短期的特徴の下位10%以内。

【0345】

コントラスト：高平均と低平均の間の比。長期的セグメント内の諸短期的特徴のダイナミックを表わす。

40

【0346】

短期的特徴抽出器2022は、既存の技法を用いて実装されてもよく、それによりいかなる特徴が抽出されることもできる。にもかかわらず、短期的特徴抽出器2022のためにいくつかの修正が次の6.3節で提案される。

【0347】

6.3節 短期的特徴の抽出

図24および図27に示されるように、短期的特徴抽出器2022は、短期的特徴として、次の特徴のうちの少なくとも一つを、各短期的オーディオ・セグメントから直接抽出するよう構成されていてもよい：リズム特性、中断/ミュート特性および短期的オーディ

50

オ品質特徴。

【0348】

リズム特性は、リズム強さ、リズム規則性、リズム明確性（ここにその全体において参照によって組み込まれる非特許文献2参照）および2Dサブバンド変調（ここにその全体において参照によって組み込まれる非特許文献3参照）を含みうる。

【0349】

中断/ミュート特性は、発話中断、シャープな減衰、ミュート長さ、不自然な無音、不自然な無音の平均、不自然な無音の全エネルギーを含みうる。

【0350】

短期的オーディオ品質特徴は、短期的セグメントに関するオーディオ品質特徴であり、これは下記で論じる、オーディオ・フレームから抽出されるオーディオ品質特徴と同様である。

【0351】

代替的または追加的に、図28に示されるように、オーディオ分類器200は、短期的セグメントに含まれるオーディオ・フレームのシーケンスの各フレームからフレーム・レベル特徴を抽出するフレーム・レベル特徴抽出器2012を有していてもよい。短期的特徴抽出器2022は、オーディオ・フレームのシーケンスから抽出されるフレーム・レベル特徴に基づいて短期的特徴を計算するよう構成されていてもよい。

【0352】

前処理として、入力オーディオ信号は、モノ・オーディオ信号にダウンミックスされてもよい。この前処理は、オーディオ信号がすでにモノ信号であれば不要である。次いで、あらかじめ定義された長さ（典型的には10ないし25ミリ秒）でフレームに分割される。対応して、各フレームからフレーム・レベル特徴が抽出される。

【0353】

フレーム・レベル特徴抽出器2012は、次の特徴のうちの少なくとも一つを抽出するよう構成されていてもよい：さまざまな短期的オーディオ型の属性を特徴付ける特徴、カットオフ周波数、静的な信号雑音比（SNR）特性、セグメントの信号雑音比（SNR）特性、基本的発話記述子および声道特性。

【0354】

さまざまな短期的オーディオ型（特に、発話、短期的音楽、背景音およびノイズ）の属性を特徴付ける特徴は、次の特徴のうちの少なくとも一つを含んでいてもよい：フレーム・エネルギー、サブバンド・スペクトル分布、スペクトル・フラックス（spectral flux）、メル周波数ケプストラム係数（MFCC: Mle-frequency Cepstral Coefficient）、ベース（bass）、残差情報（residual information）、クロマ（Chroma）特徴および零交差レート（zero-crossing rate）。

【0355】

MFCCの詳細については、ここに参照によってその全体において組み込まれる非特許文献1参照。クロマ特徴の詳細については、ここに参照によってその全体において組み込まれる非特許文献4参照。

【0356】

カットオフ周波数は、それより上ではコンテンツのエネルギーが0に近い、オーディオ信号の最高周波数を表わす。これは、帯域が限定されているコンテンツを検出するよう設計され、これは本願ではオーディオ・コンテキスト分類のために有用である。カットオフ周波数は通例、符号化によって引き起こされる。たいていの符号化器は、低ビットレートまたは中ビットレートでは高周波数を破棄するからである。たとえば、MP3コーデックは128kbpsにおいて16kHzのカットオフ周波数をもつ。もう一つの例として、多くの一般的なVoIPコーデックは8kHzまたは16kHzのカットオフ周波数をもつ。

【0357】

カットオフ周波数のほかに、オーディオ・エンコード・プロセスの間の信号劣化が、VoIP対非VoIPコンテキスト、高品質対低品質オーディオ・コンテキストのようなさまざまな

10

20

30

40

50

オーディオ・コンテキストを区別するためのさらなる特性として考慮される。よりリッチな特性を捕捉するために、客観的発話品質評価（ここに参照によってその全体において組み込まれる非特許文献5参照）のためのもののようなオーディオ品質を表わす特徴が、複数レベルにおいてさらに抽出されてもよい。オーディオ品質特徴の例は次のものを含む：

- a) 推定された背景ノイズ・レベル、スペクトル明瞭性などを含む静的なSNR特性
- b) スペクトル・レベル偏差、スペクトル・レベル範囲、相対ノイズ・フロアなどを含むセグメントSNR特性
- c) ピッチ平均、発話セクション・レベル変動、発話レベルなどを含む基本的発話記述子
- d) ロボット化 (robotization)、ピッチ・クロス・パワー (pitch cross power) などを含む声道特性。

10

【0358】

フレーム・レベル特徴から短期的特徴を導出するために、短期的特徴抽出器2022は、フレーム・レベル特徴の統計量を、短期的特徴として計算するよう構成されていてもよい。

【0359】

フレーム・レベル特徴の統計量の例は、平均および標準偏差を含む。これは、短期的音楽、発話、背景およびノイズのようなさまざまなオーディオ型を区別するためのリズム属性を捕捉する。たとえば、発話は通例、有声音と無声音の間で音節レートで交替し、一方音楽はそのような交替がなく、発話のフレーム・レベル特徴の変動は通例、音楽より大きい。

20

【0360】

統計量のもう一つの例は、フレーム・レベル特徴の重み付けされた平均である。たとえば、カットオフ周波数について、短期的セグメント内の全オーディオ・フレームから導出されたカットオフ周波数の間の、各フレームのエネルギーまたはラウドネスを重みとして重み付けされた平均は、短期的セグメントについてのカットオフ周波数となる。

【0361】

代替的または追加的に、図29に示されるように、オーディオ分類器200は、オーディオ・フレームからフレーム・レベル特徴を抽出するフレーム・レベル特徴抽出器2012と、それぞれのフレーム・レベル特徴を使ってオーディオ・フレームのシーケンスの各フレームをフレーム・レベル・オーディオ型に分類するフレーム・レベル分類器2014とを有していてもよい。ここで、短期的特徴抽出器2022は、前記シーケンスのオーディオ・フレームに関するフレーム・レベル分類器2014の結果に基づいて短期的特徴を計算するよう構成されていてもよい。

30

【0362】

換言すれば、オーディオ・コンテンツ分類器202およびオーディオ・コンテキスト分類器204に加えて、オーディオ分類器200はさらに、フレーム分類器201を有していてもよい。そのような構成では、オーディオ・コンテンツ分類器202は、フレーム分類器201のフレーム・レベルの分類結果に基づいて短期的セグメントを分類し、オーディオ・コンテキスト分類器204は、オーディオ・コンテンツ分類器202の短期的分類結果に基づいて長期的セグメントを分類する。

40

【0363】

フレーム・レベル分類器2014は、オーディオ・フレームのシーケンスの各フレームを何らかのクラスに分類するよう構成されていてもよい。そのクラスは、「フレーム・レベルのオーディオ型」と称されてもよい。

ある実施形態では、フレーム・レベルのオーディオ型は、上記で論じたコンテンツ型の構成と同様の構成を有していてもよく、コンテンツ型と同様の意味を有していてもよい。唯一の相違は、フレーム・レベルのオーディオ型とコンテンツ型はオーディオ信号の異なるレベルで、すなわちフレーム・レベルおよび短期的セグメント・レベルで分類されるということである。たとえば、フレーム・レベル分類器2014は、オーディオ・フレームのシーケンスの各フレームを次のフレーム・レベルのオーディオ型のうちの少なくとも一つ

50

に分類するよう構成されていてもよい：発話、音楽、背景音およびノイズ。他方、フレーム・レベルのオーディオ型は、部分的または完全にコンテンツ型の構成とは異なる、フレーム・レベルの分類により好適であり短期的分類のための短期的特徴として使われるのにより好適な構成を有していてもよい。たとえば、フレーム・レベル分類器 2014 は、オーディオ・フレームのシーケンスの各フレームを、次のフレーム・レベルのオーディオ型のうちの少なくとも一つに分類するよう構成されていてもよい：有声、無声および休止。

【0364】

フレーム・レベル分類の結果からいかにして短期的特徴を導出するかについて、6.2 節の記述を参照することによって、同様の方式が採用されてもよい。

【0365】

代替として、フレーム・レベル分類器 2014 の結果に基づく短期的特徴およびフレーム・レベル特徴抽出器 2012 によって得られたフレーム・レベル特徴に直接基づく短期的特徴が短期的分類器 2024 によって使用されてもよい。したがって、短期的特徴抽出器 2022 は、オーディオ・フレームのシーケンスから抽出されたフレーム・レベル特徴およびオーディオ・フレームのシーケンスに関するフレーム・レベル分類器の結果の両方に基づく短期的特徴を計算するよう構成されていてもよい。

【0366】

換言すれば、フレーム・レベル特徴抽出器 2012 は、6.2 節で論じたのと同様の統計量および図 28 との関連で述べた、下記の特徴のうちの一つを含む短期的特徴の両方を計算するよう構成されていてもよい：さまざまな短期的オーディオ型の属性を特徴付ける特徴、カットオフ周波数、静的な信号雑音比特性、セグメントの信号雑音比特性、基本的発話記述子および声道特性。

【0367】

リアルタイムで作業するために、すべての実施形態において、短期的特徴抽出器 2022 は、所定のステップ長さで長期的オーディオ・セグメントの時間次元内をスライドする移動窓を用いて形成される短期的オーディオ・セグメントに対して作用するよう構成されてもよい。短期的オーディオ・セグメントについての移動窓および長期的オーディオ・セグメントについての移動窓について、詳細は 1.1 節が参照されうる。

【0368】

6.4 節 実施形態の組み合わせおよび応用シナリオ

第一部と同様に、上記で論じたすべての実施形態およびその変形は、そのいかなる組み合わせにおいて実装されてもよく、異なる部 / 実施形態において言及されるが同じまたは同様の機能をもついかなる構成要素も同じまたは別個の構成要素として実装されてもよい。

【0369】

たとえば、6.1 節ないし 6.3 節において述べた解決策の任意の二つ以上が互いと組み合わせられてもよい。そして、これらの組み合わせの任意のものが、第一部～第五部および後述する他の部において記載または含意されている任意の実施形態とさらに組み合わせられてもよい。特に、第一部で論じた型平滑化ユニット 712 がオーディオ分類器 200 のコンポーネントとしてこの部において、フレーム分類器 2014 またはオーディオ・コンテンツ分類器 202 またはオーディオ・コンテキスト分類器 204 の結果を平滑化するために使用されてもよい。さらに、オーディオ分類器 200 の出力の急激な変化を避けるために、オーディオ分類器 200 のコンポーネントとしてタイマー 916 も役割を果たしてもよい。

【0370】

6.5 節 オーディオ分類方法

第一部と同様に、上記の実施形態におけるオーディオ分類器を記述する過程で、いくつかのプロセスまたは方法も開示されていることは明らかである。以下では、これらの方法の概要が与えられるが、上記ですでに論じた詳細の一部は繰り返さない。

【0371】

10

20

30

40

50

図30に示されるある実施形態では、オーディオ分類方法が提供される。(互いに重なり合うまたは重なり合わない)短期的オーディオ・セグメントのシーケンスからなる長期的オーディオ・セグメントの長期的オーディオ型(つまりコンテンツ型)を識別するために、短期的オーディオ・セグメントがまず短期的オーディオ型に、つまりコンテンツ型に分類され(動作3004)、長期的オーディオ・セグメント内の短期的セグメントのシーケンスに関する分類動作の結果の統計量を計算する(動作3006)ことによって長期的特徴が得られる。次いで、朝敵特徴を使って長期的分類(動作3008)が実行されてもよい。短期的オーディオ・セグメントは、オーディオ・フレームのシーケンスを含んでいてもよい。むしろ、短期的セグメントの短期的オーディオ型を識別するために、それらのセグメントから短期的特徴が抽出される必要がある。

10

【0372】

短期的オーディオ型(コンテンツ型)は、発話、短期的音楽、背景音およびノイズを含んでいてもよいが、それに限られない。

【0373】

長期的特徴は、それらの短期的セグメントの信頼値の平均および分散、それらの短期的セグメントの重要度によって重み付けされた前記平均および分散、各短期的オーディオ型の出現頻度および異なる短期的オーディオ型の間の遷移の頻度。

【0374】

ある平均では、図31に示されるように、さらなる長期的特徴が、長期的オーディオ・セグメント内の短期的セグメントのシーケンスの短期的特徴に直接基づいて得られてもよい(動作3107)。そのような長期的特徴は、短期的特徴の次の統計量を含んでいてもよいが、それに限られない:平均、分散、重み付けされた平均、重み付けされた分散、高平均、低平均および高平均と低平均の間の比。

20

【0375】

短期的特徴を抽出するためには種々の方法がある。一つは、分類されるべき短期的オーディオ・セグメントから短期的特徴を直接抽出することである。そのような特徴は、リズム特性、中断/ミュート特性および短期的オーディオ品質特徴を含むがそれに限られない。

【0376】

第二の方法は、各短期的セグメントに含まれるオーディオ・フレームからフレーム・レベルの特徴を抽出して(図32の動作3201)、次いでフレーム・レベルの特徴に基づいて短期的特徴を計算する、たとえば短期的特徴としてフレーム・レベルの特徴の統計量を計算することである。フレーム・レベル特徴は:さまざまな短期的オーディオ型の属性を特徴付ける特徴、カットオフ周波数、静的な信号雑音比特性、セグメントの信号雑音比特性、基本的発話記述子および声道特性を含んでいてもよいがそれに限られない。さまざまな短期的オーディオ型の属性を特徴付ける特徴はさらに、フレーム・エネルギー、サブバンド・スペクトル分布、スペクトル・フラックス(spectral flux)、メル周波数ケプストラム係数(MFCC: Mle-frequency Cepstral Coefficient)、ベース(bass)、残差情報(residual information)、クロマ(Chroma)特徴および零交差レート(zero-crossing rate)を含んでいてもよい。

30

40

【0377】

第三の方法は、長期的特徴の抽出と同様の仕方で短期的特徴を抽出することである:分類されるべき短期的セグメント内のオーディオ・フレームからフレーム・レベル特徴を抽出した(動作3201)後、それぞれのフレーム・レベル特徴を使って各オーディオ・フレームをフレーム・レベルのオーディオ型に分類する(図33における動作32011);短期的特徴は、フレーム・レベル・オーディオ型(任意的に信頼値を含む)に基づいて短期的特徴を計算することによって抽出されうる(動作3002)。フレーム・レベル・オーディオ型は、短期的オーディオ型(コンテンツ型)と同様の属性および構成を有していてもよく、発話、音楽、背景音およびノイズをも含んでいてもよい。

【0378】

50

上記第二の方法および第三の方法は、図33において破線矢印で示されるように一緒に組み合わせられてもよい。

【0379】

第一部で論じたように、短期的オーディオ・セグメントおよび長期的オーディオ・セグメントの両方が移動窓を用いてサンプリングされてもよい。すなわち、短期的特徴を抽出する動作(動作3002)は、所定のステップ長さで長期的オーディオ・セグメントの時間次元においてスライドする移動窓を用いて形成される短期的オーディオ・セグメントに対して実行されてもよく、長期的特徴を抽出する動作(動作3107)および短期的オーディオ型の統計量を計算する動作(動作3006)は、所定のステップ長さでオーディオ信号の時間次元においてスライドする移動窓を用いて形成される長期的オーディオ・セグメントに対して実行されてもよい。

10

【0380】

オーディオ処理装置の実施形態と同様に、オーディオ処理方法の実施形態およびその変形の任意の組み合わせが現実的である。他方、オーディオ処理方法の実施形態およびその変形のあらゆる側面は別個の解決策であってもよい。さらに、この節に記載される任意の二つ以上の解決策が互いと組み合わせられてもよく、これらの組み合わせがさらに、本開示の他の部において記載または含意される任意の実施形態と組み合わせられてもよい。特に、6.4節においてすでに論じたように、オーディオ型の平滑化方式および遷移方式がここで論じたオーディオ分類方法の一部であってもよい。

【0381】

第七部：VoIP分類器および分類方法

第六部では、少なくとも部分的にはコンテンツ型分類の結果に基づいてオーディオ信号をオーディオ・コンテキスト型に分類するための新規なオーディオ分類器が提案されている。第六部で論じた実施形態において、長期的特徴は、数秒ないし数十秒の長さの長期的セグメントから抽出される。よって、オーディオ・コンテキスト分類は長いレイテンシーを引き起こしうる。オーディオ・コンテキストがリアルタイムでまたはほぼリアルタイムでたとえば短期的セグメント・レベルにおいて分類されることが望まれる。

20

【0382】

7.1節 短期的セグメントに基づくコンテキスト分類

したがって、図34に示されるように、オーディオ信号の短期的セグメントのコンテンツ型を同定するためのオーディオ・コンテンツ分類器202Aと、少なくとも部分的には前記オーディオ・コンテンツ分類器によって識別されたコンテンツ型に基づいて短期的セグメントのコンテキスト型を識別するオーディオ・コンテキスト分類器204Aとを有するオーディオ分類器200Aが提供される。

30

【0383】

ここで、オーディオ・コンテンツ分類器202Aは、第六部ですでに述べた技法を採用してもよいが、下記の7.2節で論じる種々の技法を採用してもよい。また、オーディオ・コンテキスト分類器204Aは、第六部ですでに述べた技法を採用してもよいが、違いとして、コンテキスト分類器204Aは、オーディオ・コンテンツ分類器202Aからの結果の統計量を使うのではなく、オーディオ・コンテンツ分類器202Aの結果を直接使ってもよい。オーディオ・コンテキスト分類器204およびオーディオ・コンテンツ分類器202Aはいずれも同じ短期的セグメントを分類しているからである。さらに、第六部と同様に、オーディオ・コンテンツ分類器202Aからの結果に加えて、オーディオ・コンテキスト分類器204は、短期的セグメントから直接抽出された他の特徴を使ってもよい。すなわち、オーディオ・コンテキスト分類器204Aは、特徴として、短期的セグメントのコンテンツ型の信頼値および短期的セグメントから抽出された他の特徴を使って、機械学習モデルに基づいて短期的セグメントを分類するよう構成されてもよい。短期的セグメントから抽出される特徴については、第六部が参照されうる。

40

【0384】

オーディオ・コンテンツ分類器200Aは、同時に、短期的セグメントを、VoIP発話 /

50

ノイズおよび/または非VoIP発話/ノイズよりも多くのオーディオ型としてラベル付けしてもよく(VoIP発話/ノイズおよび/または非VoIP発話/ノイズについては下記の7.2節で論じる)、複数のオーディオ型のそれぞれは、1.2節において論じたその独自の信頼値をもちうる。これは、よりリッチな情報が補足できるので、よりよい分類精度を達成できる。たとえば、発話および短期的音楽の信頼値の合同情報は、どの程度そのオーディオ・コンテンツが発話および背景音楽の混合でありそうかを明らかにし、それにより純粋なVoIPコンテンツから弁別されることができる。

【0385】

7.2節 VoIP発話およびVoIPノイズを使った分類

本願のこの側面は、短い判断レイテンシーのために現在の短期的セグメントを分類することを要求されるVoIP/非VoIP分類システムにおいて特に有用である。

10

【0386】

この目的のために、図34に示されるように、オーディオ分類器200Aは、VoIP/非VoIP分類のために特別に設計される。VoIP/非VoIPを分類するために、オーディオ・コンテキスト分類器204による最終的な堅牢なVoIP/非VoIP分類のための中間結果を生成するようVoIP発話分類器2026および/またはVoIPノイズ分類器が開発される。

【0387】

VoIP短期的セグメントはVoIP発話およびVoIPノイズを交互に含むであろう。発話の短期的セグメントをVoIP発話または非VoIP発話に分類するためには高い精度が達成できるが、ノイズの短期的セグメントをVoIPノイズまたは非VoIPノイズに分類するためにはそうではないことが観察される。このように、発話とノイズの間の差を考慮せず、よってこれら二つのコンテンツ型(発話およびノイズ)の特徴が一緒に混ざったままで短期的セグメントをVoIP(VoIP発話およびVoIPノイズを含むが、VoIP発話およびVoIPノイズは個々に同定されていない)および非VoIPに直接分類することによって、弁別性がぼかされていると結論できる。

20

【0388】

分類器にとって、VoIPノイズ/非VoIPノイズ分類よりもVoIP発話/非VoIP発話分類についてより高い精度を達成することは理にかなっている。発話がノイズよりも多くの情報を含み、カットオフ周波数のような特徴は発話を分類するためにより効果的だからである。アダプスト・トレーニング/プロセスから得られる重みランク付けに従って、VoIP/非VoIP分類についての上位の重みを付けられた諸短期的特徴は：対数エネルギーの標準偏差、カットオフ周波数、リズム強さの標準偏差およびスペクトル・フラックスの標準偏差である。対数エネルギーの標準偏差、リズム強さの標準偏差およびスペクトル・フラックスの標準偏差は一般に、非VoIP発話についてよりVoIP発話について高くなる。一つのありそうな理由は、映画のメディアまたはゲームのような非VoIPコンテンツにおける多くの短期的発話セグメントは通例、背景音楽または効果音のような他の音と混ざっており、該背景音楽または効果音については上記の特徴の値はより低いということである。一方、カットオフ特徴は一般に、非VoIP発話についてよりVoIP発話について低くなる。このことは、多くの一般的なVoIPコーデックによって導入される低いカットオフ周波数を示す。

30

【0389】

したがって、ある実施形態では、オーディオ・コンテンツ分類器202Aは、短期的セグメントをコンテンツ型VoIP発話またはコンテンツ型非VoIP発話に分類するVoIP発話分類器2026を有していてもよく、オーディオ・コンテキスト分類器204は、VoIP発話および非VoIP発話の信頼値に基づいて短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されていてもよい。

40

【0390】

もう一つの実施形態では、オーディオ・コンテンツ分類器202Aはさらに、短期的セグメントをコンテンツ型VoIPノイズまたはコンテンツ型非VoIPノイズに分類するVoIPノイズ分類器2028を有していてもよく、オーディオ・コンテキスト分類器204は、VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズの信頼値に基づいて短期的セグメント

50

をコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されていてもよい。

【0391】

VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズのコンテンツ型は、第六部、1.2節および7.1節で論じた既存の技法を用いて同定されてもよい。

【0392】

あるいはまた、オーディオ・コンテンツ分類器202Aはさらに図35に示される階層構造を有していてもよい。すなわち、発話/ノイズ分類器2025からの結果を活用してまず短期的セグメントを発話またはノイズ/背景に分類するのである。

【0393】

単にVoIP発話分類器2026を使う実施形態に基づいて、短期的セグメントが発話/ノイズ分類器2025（このような状況ではこれは単に発話分類器である）によって発話と判定される場合、VoIP発話分類器2026は、それがVoIP発話または非VoIP発話のいずれであるかを分類することに進み、二値分類結果を計算する。それ以外の場合には、VoIP発話の信頼値が低い、あるいはVoIP発話についての決定が不確かであると見なされてもよい。

【0394】

単にVoIPノイズ分類器2028を使う実施形態に基づいて、短期的セグメントが発話/ノイズ分類器2025（このような状況ではこれは単にノイズ（背景）分類器である）によってノイズと判定される場合、VoIPノイズ分類器2028は、それをVoIPノイズまたは非VoIPノイズに分類し、二値分類結果を計算することに進む。それ以外の場合には、VoIPノイズの信頼値が低い、あるいはVoIPノイズについての決定が不確かであると見なされてもよい。

【0395】

ここで、一般に発話は情報性のコンテンツ型であり、ノイズ/背景は干渉性のコンテンツ型であるので、たとえ短期的セグメントがノイズでなくても、前段の実施形態において、短期的セグメントがコンテキスト型VoIPではないと確定的に判定することはできない。短期的セグメントが発話でなければ、単にVoIP発話分類器2026を使う実施形態では、それはおそらくコンテキスト型VoIPではないであろう。したがって、一般に、担任VoIP発話分類器2026を使う実施形態は独立して実現されうる。一方、単位VoIPノイズ分類器2028を使う他方の実施形態は、たとえばVoIP発話分類器2026を使う実施形態と協働する補足的実施形態として使用されうる。

【0396】

すなわち、VoIP発話分類器2026およびVoIPノイズ分類器2028の両方が使用されてもよい。短期的セグメントが発話/ノイズ分類器2025によって発話と判定される場合、VoIP発話分類器2026はそれがVoIP発話または非VoIP発話のいずれであるかを分類することに進み、二値分類結果を計算する。短期的セグメントが発話/ノイズ分類器2025によってノイズと判定される場合、VoIPノイズ分類器2028はそれをVoIPノイズまたは非VoIPノイズに分類し、二値分類結果を計算することに進む。それ以外の場合には、短期的セグメントが非VoIPと分類されうると見なされてもよい。

【0397】

発話/ノイズ分類器2025、VoIP発話分類器2026およびVoIPノイズ分類器2028の実装は、いかなる既存の技法を採用してもよく、第一部ないし第六部で論じたオーディオ・コンテンツ分類器202であってもよい。

【0398】

上記に従って実装されたオーディオ・コンテンツ分類器202Aが最終的に短期的セグメントを発話、ノイズおよび背景のいずれにも、あるいはVoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズのいずれにも分類しない場合、つまり、すべての関連する信頼値が低い場合、オーディオ・コンテンツ分類器202A（およびオーディオ・コンテキスト分類器204）はその短期的セグメントを非VoIPと分類してもよい。

10

20

30

40

50

【0399】

短期的セグメントを、VoIP発話分類器2026およびVoIPノイズ分類器2028の結果に基づいてVoIPまたは非VoIPのコンテキスト型に分類するために、オーディオ・コンテキスト分類器204は、7.1節で論じた機械学習ベースの技法を採用してもよく、修正として、7.1節ですでに論じたような、短期的セグメントから直接抽出された短期的特徴および/またはVoIP関係のコンテンツ型以外のコンテンツ型に向けられた他のオーディオ・コンテンツ分類器(単数または複数)の結果を含むより多くの特徴が使用されてもよい。

【0400】

上記の機械学習ベースの技法のほかに、VoIP/非VoIP分類への代替的なアプローチは、ドメイン知識を活用し、分類結果をVoIP発話およびVoIPノイズとの関連で利用するヒューリスティック規則であることができる。

10

【0401】

時刻 t の現在の短期的セグメントがVoIP発話または非VoIP発話として決定される場合、分類結果はVoIP/非VoIP分類結果として直接取られる。VoIP/非VoIP発話分類は先に論じたように堅牢だからである。すなわち、短期的セグメントがVoIP発話であると判定される場合は、それはコンテキスト型VoIPである。短期的セグメントが非VoIP発話であると判定される場合は、それはコンテキスト型非VoIPである。

【0402】

VoIP発話分類器2026が先述した発話/ノイズ分類器2025によって判定される発話に関してVoIP発話/非VoIP発話に関する二値判定をするとき、VoIP発話および非VoIP発話の信頼値は相補的でありうる。すなわち、それらの和は1である(0が100%否を表わし、1が100%肯定を表わす場合)。そしてVoIP発話と非VoIP発話を区別するための信頼値の閾値は実際に同じ点を指しうる。VoIP発話分類器2026が二値分類器でない場合には、VoIP発話および非VoIP発話の信頼値は相補的ではなく、VoIP発話と非VoIP発話を区別するための信頼値の閾値は必ずしも同じ点を指さないことがある。

20

【0403】

しかしながら、VoIP発話または非VoIP発話信頼度が閾値に近く、該閾値のまわりで揺動する場合、VoIP/非VoIP分類結果は、頻繁に切り替わりすぎることがありうる。そのような揺動を避けるために、バッファ方式が提供されてもよい。VoIP発話および非VoIP発話についての両方の閾値がより大きく設定されてもよく、よって現在のコンテンツ型から他方のコンテンツ型に切り替わるのはそれほど容易ではなくなる。記述の簡単のため、非VoIP発話についての信頼値をVoIP発話の信頼値に変換してもよい。すなわち、信頼値が高ければ、短期的セグメントはVoIP発話により近いと見なされ、信頼値が低ければ、短期的セグメントは非VoIP発話により近いと見なされる。上記のような非二値分類器については、非VoIP発話についての高い信頼値は必ずしもVoIP発話の低い信頼値を意味しないが、そのような単純化は、本解決策の本質をよく反映でき、二値分類器の言辞を用いて記述される関連する請求項は、非二値分類器についての等価な解決策をカバーすると解釈される。

30

【0404】

バッファ方式は図36に示されている。二つの閾値 $Th1$ および $Th2$ ($Th1 > Th2$)の間にバッファ領域がある。VoIP発話の信頼値 $v(t)$ がこの領域において低下するとき、コンテキスト分類は変化しない。これは図36における左側および右側の矢印によって示されている。信頼値 $v(t)$ が大きいほうの閾値 $Th1$ より大きいときにのみ、短期的セグメントはVoIPと分類され(図36の下部の矢印によって示されるように)、信頼値が小さいほうの閾値 $Th2$ より大きくないときにのみ、短期的セグメントは非VoIPと分類されることになる(図36の上部の矢印によって示されるように)。

40

【0405】

VoIPノイズ分類器2028が代わりに使われる場合、状況は同様である。解決策をより堅牢にするために、VoIP発話分類器2026およびVoIPノイズ分類器2028が合同して使用されてもよい。次いで、オーディオ・コンテキスト分類器204Aは: VoIP発話の信

50

信頼値が第一の閾値より大きい場合またはVoIPノイズの信頼値が第三の閾値より大きい場合、短期的セグメントをコンテキスト型VoIPとして分類し；VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合またはVoIPノイズの信頼値が前記第三の閾値より大きくない第四の閾値より大きくない場合、短期的セグメントをコンテキスト型非VoIPとして分類し；それ以外の場合には短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されていてもよい。

【0406】

ここで、第一の閾値は、第二の閾値に等しくてもよく、第三の閾値は第四の閾値に等しくてもよい。これは特に二値VoIP発話分類器および二値VoIPノイズ分類器についてそうであるが、それに限られない。しかしながら、一般に、VoIPノイズ分類結果はそれほど堅牢ではないので、第三および第四の閾値は互いに等しくないほうがよいであろう。両者は0.5から遠いべきである（0は非VoIPノイズである高い信頼度を示し、1はVoIPノイズである高い信頼度を示す）。

10

【0407】

7.3節 平滑化ゆらぎ

急速な揺動〔ゆらぎ〕を避けるために、もう一つの解決策は、オーディオ・コンテンツ分類器によって決定される信頼値を平滑化することである。したがって、図37に示されるように、型平滑化ユニット203Aがオーディオ分類器200Aに含まれてもよい。先に論じた四つのVoIP関係のコンテンツ型のそれぞれの信頼値について、1.3節で論じた平滑化方式が採用されてもよい。

20

【0408】

あるいはまた、7.2節と同様に、VoIP発話および非VoIP発話は、相補的な信頼値を有する対と見なされてもよい。VoIPノイズおよび非VoIPノイズも相補的な信頼値をもつ対と見なされてもよい。そのような状況では、各対のうち一つのみが平滑化される必要があり、1.3節で論じた平滑化方式が採用されてもよい。

【0409】

VoIP発話の信頼値を例にとると、公式(3)が次のように書き直されてもよい。

【0410】

$$v(t) = \alpha \cdot v(t-1) + (1 - \alpha) \cdot \text{voipSpeechConf}(t) \quad (3'')$$

ここで、 $v(t)$ は時刻 t における平滑化されたVoIP発話信頼値であり、 $v(t-1)$ は最後の時点における平滑化されたVoIP発話信頼値であり、 voipSpeechConf は平滑化前の現在時刻 t におけるVoIP発話信頼度であり、 α は重み付け係数である。

30

【0411】

ある変形では、上記のように発話ノイズ分類器2025がある場合、短期的セグメントについての発話の信頼値が低ければ、その短期的セグメントはVoIP発話として堅牢に分類されることはできず、VoIP発話分類器2026を実際に機能させることなく、 $\text{voipSpeechConf}(t) = v(t-1)$ と直接設定することができる。

【0412】

あるいはまた、上記の状況において、不確かなケースを示して $\text{voipSpeechConf}(t) = 0.5$ （または0.4~0.5のような0.5より高くない他の値）と設定することができる（ここで、信頼度=1がVoIPである高い信頼性を示し、信頼度=0がVoIPではない高い信頼性を示す）。

40

【0413】

したがって、この変形によれば、図37に示されるように、オーディオ・コンテンツ分類器200Aはさらに、短期的セグメントの発話のコンテンツ型を識別するための発話ノイズ分類器2025を有していてもよく、型平滑化ユニット203Aは、平滑化前の現在の短期的セグメントについてのVoIP発話の信頼値を、（0.5または0.4~0.5などの他の値のような）所定の信頼値として、あるいは発話ノイズ分類器によって分類されるコンテンツ型発話についての信頼値が第五の閾値より低い最後の短期的セグメントの平滑化された信頼値として、設定するよう構成されていてもよい。そのような状況では、VoIP発話

50

分類器 2026 は作用してもしなくてもよい。あるいはまた、信頼値の設定は、VoIP 発話分類器 2026 によって行なわれてもよい。これは、この作用が型平滑化ユニット 203A によってなされる解決策と等価であり、請求項は両方の状況をカバーするものと解釈されるべきである。さらに、ここで、「発話 / ノイズ分類器によって分類されるコンテンツ型発話についての信頼値が第五の閾値より低い」という言辞を使っているが、保護範囲はそれに限定されず、短期的セグメントが発話以外のコンテンツ型に分類される状況に等価である。

【0414】

VoIP ノイズの信頼値について、状況は同様であり、詳細な説明はここでは割愛する。

【0415】

急激なゆらぎを避けるために、さらにもう一つの解決策は、オーディオ・コンテキスト分類器 204A によって決定される信頼値を平滑化するというものであり、1.3 節で論じた平滑化方式が採用されてもよい。

【0416】

急激なゆらぎを避けるために、さらにもう一つの解決策は、VoIP と非 VoIP の間のコンテキスト型の遷移を遅らせるというものであり、1.6 節で述べたのと同じ方式が使われてもよい。1.6 節で述べたように、タイマー 916 はオーディオ分類器の外部であってもよいし、あるいはオーディオ分類器の一部としてその内部であってもよい。したがって、図 38 に示されるように、オーディオ分類器 200A はさらにタイマー 916 を有していてもよい。そして、オーディオ分類器は、新しいコンテキスト型の持続時間の長さが第六の閾値に達するまで現在のコンテキスト型を出力し続けるよう構成される（コンテキスト型はオーディオ型のインスタンスである）。1.6 節を参照することにより、詳細な説明はここでは割愛しうる。

【0417】

追加的または代替的に、VoIP と非 VoIP の間の遷移を遅らせるもう一つの方式として、VoIP / 非 VoIP 分類について先述した第一および / または第二の閾値が、最後の短期的セグメントのコンテキスト型に依存して異なってもよい。すなわち、新しい短期的セグメントのコンテキスト型が最後の短期的セグメントのコンテキスト型と異なるときは、第一および / または第二の閾値はより大きくなり、一方、新しい短期的セグメントのコンテキスト型が最後の短期的セグメントのコンテキスト型と同じときはより小さくなる。このようにして、コンテキスト型は現在のコンテキスト型に維持される傾向があり、よってコンテキスト型の急激なゆらぎがある程度抑制されうる。

【0418】

7.4 節 実施形態の組み合わせおよび応用シナリオ

第一部と同様に、上記で論じたすべての実施形態およびその変形は、そのいかなる組み合わせにおいて実装されてもよく、異なる部 / 実施形態において言及されるが同じまたは同様の機能をもついかなる構成要素も同じまたは別個の構成要素として実装されてもよい。

【0419】

たとえば、7.1 節ないし 7.3 節において述べた解決策の任意の二つ以上が互いと組み合わせられてもよい。そして、これらの組み合わせの任意のものが、第一部～第六部において記載または含意されている任意の実施形態とさらに組み合わせられてもよい。特に、この部で論じた諸実施形態およびそれらの任意の組み合わせは、第四部で論じたオーディオ処理装置 / 方法またはボリューム平滑化器 / 制御方法の実施形態と組み合わせられてもよい。

【0420】

7.5 節 VoIP 分類方法

第一部と同様に、上記の実施形態におけるオーディオ分類器を記述する過程で、いくつかのプロセスまたは方法も開示されていることは明らかである。以下では、これらの方法の概要が与えられるが、上記ですでに論じた詳細の一部は繰り返さない。

10

20

30

40

50

【 0 4 2 1 】

図 3 9 に示されるある実施形態では、オーディオ分類方法が、オーディオ信号の短期的セグメントのコンテンツ型を識別し（動作 4 0 0 4）、次いで、少なくとも識別されたコンテンツ型に基づいて短期的セグメントのコンテキスト型を識別する（動作 4 0 0 8）ことを含む。

【 0 4 2 2 】

オーディオ信号のコンテキスト型を動的にかつ高速に識別するために、この部におけるオーディオ分類方法は、コンテキスト型VoIPおよび非VoIPを識別する際に特に有用である。そのような状況では、短期的セグメントはまずコンテンツ型VoIP発話またはコンテンツ型非VoIP発話に分類されてもよく、コンテキスト型を識別する上記動作は、VoIP発話および非VoIP発話の信頼値に基づいて、短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成される。

10

【 0 4 2 3 】

あるいはまた、短期的セグメントはまずコンテンツ型VoIPノイズまたはコンテンツ型非VoIPノイズに分類されてもよく、コンテキスト型を識別する上記動作は、VoIPノイズおよび非VoIPノイズの信頼値に基づいて、短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されてもよい。

【 0 4 2 4 】

発話およびノイズは合同して考えられてもよい。そのような状況では、コンテキスト型を識別する上記動作は、VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズの信頼値に基づいて、短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されてもよい。

20

【 0 4 2 5 】

短期的セグメントのコンテキスト型を識別するために、短期的セグメントのコンテンツ型の信頼値および短期的セグメントから抽出された他の特徴の両方を特徴として取り入れて、機械学習モデルが使われてもよい。

【 0 4 2 6 】

コンテキスト型を識別する上記動作は、ヒューリスティック規則に基づいて実現されてもよい。VoIP発話および非VoIP発話のみが関わる場合は、ヒューリスティック規則は次のようなものである：VoIP発話の信頼値が第一の閾値より大きければ短期的セグメントをコンテキスト型VoIPとして分類し；VoIP発話の信頼値が第二の閾値より大きくなければ短期的セグメントをコンテキスト型非VoIPとして分類し、ここで、第二の閾値は第一の閾値より小さくなく；それ以外の場合には、短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類する。

30

【 0 4 2 7 】

VoIPノイズおよび非VoIPノイズのみが関わる状況についてのヒューリスティック規則も同様である。

【 0 4 2 8 】

発話およびノイズの両方が関わる場合は、ヒューリスティック規則は次のようなものである：VoIP発話の信頼値が第一の閾値より大きい場合またはVoIPノイズの信頼値が第三の閾値より大きい場合には短期的セグメントをコンテキスト型VoIPとして分類し；VoIP発話の信頼値が、第一の閾値より大きくない第二の閾値より大きくない場合またはVoIPノイズの信頼値が、第三の閾値より大きくない第四の閾値より大きくない場合には、短期的セグメントをコンテキスト型非VoIPとして分類し；それ以外の場合には、短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類する。

40

【 0 4 2 9 】

1 . 3 節および 1 . 8 節で論じた平滑化方式がここで採用されてもよく、詳細な説明は割愛する。1 . 3 節で述べた平滑化方式への修正として、平滑化動作 4 1 0 6 の前に、本方法はさらに、短期的セグメントからのコンテンツ型発話を識別する段階（図 4 0 の動作 4 0 0 4 0）を含んでいてもよい。ここで、コンテンツ型発話についての信頼値は第五の

50

閾値より低い（動作40041における「No」）場合には、平滑化前の現在の短期的セグメントについてのVoIP発話の信頼値は、所定の信頼値または最後の短期的セグメントの平滑化された信頼値として設定される（図40の動作40044）。

【0430】

そうではなく、コンテンツ型発話を識別する動作が短期的セグメントを発話として堅牢に判断する場合には（動作40041における「Yes」）、短期的セグメントは、平滑化動作4106の前に、さらにVoIP発話または非VoIP発話に分類される（動作40042）。

【0431】

実のところ、平滑化方式を使わなくても、本方法は、コンテンツ型発話および/またはノイズを最初に識別してもよく、端的セグメントが発話またはノイズとして分類される時、短期的セグメントをVoIP発話および非VoIP発話の一方またはVoIPノイズまたは非VoIPノイズの一方に分類するようさらなる分類が実装される。次いで、コンテキスト型を識別する動作が行なわれる。

【0432】

1.6節および1.8節で述べたように、そこで論じた遷移方式は、ここで記述されるオーディオ分類方法の一部として取り入れられてもよく、詳細は割愛する。手短かに言うと、本方法はさらに、コンテキスト型を識別する動作が同じコンテキスト型を連続的に出力する継続時間を測定することを含んでいてもよい。オーディオ分類方法は、新しいコンテキスト型の継続時間の長さが第六の閾値に達するまで、現在のコンテキスト型を出力し続けるよう構成される。

【0433】

同様に、あるコンテキスト型から別のコンテキスト型への異なる遷移対について、異なる第六の閾値が設定されてもよい。さらに、第六の閾値は、新しいコンテキスト型の信頼値と負に相関されてもよい。

【0434】

VoIP/非VoIP分類に特に向けられるオーディオ分類方法における遷移方式への修正として、現在の短期的セグメントについての第一ないし第四の閾値の任意の一つまたは複数が、最後の短期的セグメントのコンテキスト型に依存して異なるように設定されてもよい。

【0435】

オーディオ処理装置の実施形態と同様に、オーディオ処理方法の実施形態およびその変形の任意の組み合わせが現実的である。他方、オーディオ処理方法の実施形態およびその変形のあらゆる側面は別個の解決策であってもよい。さらに、この節に記載される任意の二つ以上の解決策が互いと組み合わせられてもよく、これらの組み合わせがさらに、本開示の他の部において記載または含意される任意の実施形態と組み合わせられてもよい。特に、ここに記載したオーディオ分類方法は、先述したオーディオ処理方法、特にボリューム平準化器制御方法において使用されてもよい。

【0436】

本願の「発明を実施するための形態」の冒頭で論じたように、本願の実施形態はハードウェアまたはソフトウェアまたは両方において具現されうる。図41は、本願の諸側面を実装する例示的なシステムを示すブロック図である。

【0437】

図41において、中央処理ユニット（CPU）4201は、読み出し専用メモリ（ROM）4202に記憶されたプログラムまたは記憶セクション4208からランダム・アクセス・メモリ（RAM）4203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。RAM 4203では、CPU 4201が該さまざまなプロセスなどを実行するときに必要なとされるデータも必要に応じて記憶される。

【0438】

CPU 4201、ROM 4202およびRAM 4203はバス4204を介して互いに接続される。入出力インターフェース4205もバス4204に接続される。

10

20

30

40

50

【0439】

次のコンポーネントが入出力インターフェース4205に接続される：キーボード、マウスなどを含む入力部4206；陰極線管（CRT）、液晶ディスプレイ（LCD）などのディスプレイおよびラウドスピーカーなどを含む出力部4207；ハードディスクなどを含む記憶部4208；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4209。通信部4209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。

【0440】

ドライブ4210も必要に応じて入出力インターフェース4205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リムーバブル媒体4211がドライブ4210に必要に応じてマウントされる。それにより、そこから読み込まれるコンピュータ・プログラムが必要に応じて記憶部4208にインストールされる。

【0441】

上記のコンポーネントがソフトウェアによって実装される場合、該ソフトウェアをなすプログラムはインターネットのようなネットワークまたはリムーバブル媒体4211のような記憶媒体からインストールされる。

【0442】

本願で使われる用語は単に具体的な実施形態を記述するためのものであり、本願を限定することは意図されていないことを注意されたい。本稿での用法では、文脈がそうでないことを明確に示すのでない限り、単数形は複数も含むことが意図されている。さらに、本明細書において使われるときの「含む」および/または「有する」という用語は、述べられている特徴、整数、動作、段階、要素および/またはコンポーネントの存在を示すが、一つまたは複数の他の特徴、整数、動作、段階、要素、コンポーネントおよび/またはそれらの群の存在または追加を排除するものではない。

【0443】

請求項におけるあらゆる手段または動作に機能を加えた要素の対応する構造、材料、工程および等価物は、はっきりと請求項に記載されている他の請求項記載の要素との組み合わせにおいて機能を実行するための任意の構造、材料または工程を含むことが意図されている。本願の記述は、例解および説明のために提示されたが、網羅的であることや開示される形の応用に限定されることは意図されていない。本願の範囲および精神から外れることなく、多くの修正および変形が当業者には明白となるであろう。実施形態は、本願の原理および実際の応用を最もよく説明するためおよび当業者が、考えられている具体的な用途に適したさまざまな修正をもつさまざまな実施形態について本願を理解できるようにするために選ばれ、記述された。

【0444】

いくつかの態様を記載しておく。

〔態様1〕

リアルタイムでオーディオ信号のコンテンツ型を識別するためのオーディオ・コンテンツ分類器と；

識別されたコンテンツ型に基づいて連続的な仕方でボリューム平準化器を調整する調整ユニットとを有するボリューム平準化器コントローラであって、

前記調整ユニットは、前記ボリューム平準化器の動的な利得を、前記オーディオ信号の情報性のコンテンツ型と正に相関させ、前記ボリューム平準化器の動的な利得を、前記オーディオ信号の干渉性のコンテンツ型と負に相関させるよう構成されている、ボリューム平準化器コントローラ。

〔態様2〕

前記オーディオ信号の前記コンテンツ型が、発話、短期的音楽、ノイズおよび背景音のうちの一つを含む、態様1記載のボリューム平準化器コントローラ。

〔態様3〕

ノイズが干渉性のコンテンツ型と見なされる、態様1記載のボリューム平準化器コント

10

20

30

40

50

ローラ。

〔態様４〕

前記調整ユニットが、前記コンテンツ型の信頼値に基づいて前記ボリューム平準化器の動的な利得を調整するよう構成されている、態様１記載のボリューム平準化器コントローラ。

〔態様５〕

前記調整ユニットが、前記コンテンツ型の信頼値の伝達関数を介して前記動的な利得を調整するよう構成されている、態様４記載のボリューム平準化器コントローラ。

〔態様６〕

前記オーディオ・コンテンツ分類器が前記オーディオ信号を、対応する信頼値をもつ複数のコンテンツ型に分類するよう構成されており、前記調整ユニットが、前記複数のコンテンツ型の重要性に基づいて前記複数のコンテンツ型の前記信頼値を重み付けすることを通じて前記複数のオーディオ型の少なくともいくつかを考慮するよう構成されている、態様１記載のボリューム平準化器コントローラ。

10

〔態様７〕

前記オーディオ・コンテンツ分類器が前記オーディオ信号を、対応する信頼値をもつ複数のコンテンツ型に分類するよう構成されており、前記調整ユニットが、あるコンテンツ型の重みを少なくとも一つの他のコンテンツ型の信頼値を用いて修正するよう構成されている、態様１記載のボリューム平準化器コントローラ。

〔態様８〕

前記オーディオ・コンテンツ分類器が前記オーディオ信号を、対応する信頼値をもつ複数のコンテンツ型に分類するよう構成されており、前記調整ユニットが、前記信頼値に基づいて前記複数のコンテンツ型の効果を重み付けすることを通じて前記複数のオーディオ型の少なくともいくつかを考慮するよう構成されている、態様１記載のボリューム平準化器コントローラ。

20

〔態様９〕

前記調整ユニットが、前記信頼値に基づいて少なくとも一つの優勢なコンテンツ型を考慮するよう構成されている、態様８記載のボリューム平準化器コントローラ。

〔態様１０〕

前記オーディオ・コンテンツ分類器が前記オーディオ信号を、対応する信頼値をもつ複数の干渉性のコンテンツ型および／または複数の情報性のコンテンツ型に分類するよう構成されており、前記調整ユニットが、前記信頼値に基づいて少なくとも一つの優勢な干渉性のコンテンツ型および／または少なくとも一つの優勢な情報性のコンテンツ型を考慮するよう構成されている、態様１記載のボリューム平準化器コントローラ。

30

〔態様１１〕

各コンテンツ型について、前記オーディオ信号の過去の信頼値に基づいて現在の時点での前記オーディオ信号の信頼値を平滑化するための型平滑化ユニットをさらに有する、態様１ないし１０のうちいずれか一項記載のボリューム平準化器コントローラ。

〔態様１２〕

前記型平滑化ユニットは、現在の実際の信頼値と最後の時点での平滑化された信頼値との重み付けされた和を計算することによって、現時点での前記オーディオ信号の平滑化された信頼値を決定するよう構成されている、態様１１記載のボリューム平準化器コントローラ。

40

〔態様１３〕

前記オーディオ信号のコンテキスト型を識別するオーディオ・コンテキスト分類器をさらに有しており、前記調整ユニットは、前記コンテキスト型の信頼値に基づいて前記動的な利得の範囲を調整するよう構成されている、態様１ないし１０のうちいずれか一項記載のボリューム平準化器コントローラ。

〔態様１４〕

前記オーディオ信号のコンテキスト型を識別するオーディオ・コンテキスト分類器をさ

50

らに有しており、前記調整ユニットは、前記オーディオ信号の前記コンテキスト型に基づいて前記オーディオ信号の前記コンテンツ型を情報性または干渉性で見なすよう構成されている、態様 1 ないし 10 のうちいずれか一項記載のボリューム平準化器コントローラ。

〔態様 15〕

前記オーディオ信号の前記コンテキスト型が、VoIP、映画のメディア、長期的音楽およびゲームのうちの一つを含む、態様 14 記載のボリューム平準化器コントローラ。

〔態様 16〕

コンテキスト型VoIPのオーディオ信号においては、背景音が干渉性コンテンツ型と見なされ、一方、コンテキスト型非VoIPのオーディオ信号においては、背景音および/または発話および/または音楽が情報性コンテンツ型と見なされる、態様 14 記載のボリューム平準化器コントローラ。

10

〔態様 17〕

前記オーディオ信号の前記コンテキスト型が高品質オーディオまたは低品質オーディオを含む、態様 14 記載のボリューム平準化器コントローラ。

〔態様 18〕

異なるコンテキスト型のオーディオ信号におけるコンテンツ型が、前記オーディオ信号のコンテキスト型に依存して異なる重みを割り当てられる、態様 14 記載のボリューム平準化器コントローラ。

〔態様 19〕

前記オーディオ・コンテキスト分類器が前記オーディオ信号を、対応する信頼値をもつ複数のコンテキスト型に分類するよう構成されており、前記調整ユニットが、前記複数のコンテキスト型の重要性に基づいて前記複数のコンテキスト型の前記信頼値を重み付けすることを通じて前記複数のコンテキスト型の少なくともいくつかを考慮するよう構成されている、態様 14 記載のボリューム平準化器コントローラ。

20

〔態様 20〕

前記オーディオ・コンテキスト分類器が前記オーディオ信号を、対応する信頼値をもつ複数のコンテキスト型に分類するよう構成されており、前記調整ユニットが、前記信頼値に基づいて前記複数のコンテキスト型の効果を重み付けすることを通じて前記複数のコンテキスト型の少なくともいくつかを考慮するよう構成されている、態様 14 記載のボリューム平準化器コントローラ。

30

〔態様 21〕

前記オーディオ・コンテンツ分類器が前記オーディオ信号の短期的セグメントに基づいて前記コンテンツ型を識別するよう構成されており、

前記オーディオ・コンテキスト分類器が少なくとも部分的には前記オーディオ・コンテンツ分類器によって識別されたコンテンツ型に基づく前記オーディオ信号の短期的セグメントに基づいて前記コンテキスト型を識別するよう構成されている、態様 14 記載のボリューム平準化器コントローラ。

〔態様 22〕

前記オーディオ・コンテンツ分類器が、短期的セグメントをコンテンツ型VoIP発話またはコンテンツ型非VoIP発話に分類するVoIP発話分類器を有しており、

40

前記オーディオ・コンテキスト分類器は、VoIP発話および非VoIP発話の信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されている、態様 21 記載のボリューム平準化器コントローラ。

〔態様 23〕

前記オーディオ・コンテンツ分類器がさらに、

短期的セグメントをVoIPノイズのコンテンツ型および非VoIPノイズのコンテンツ型に分類するVoIPノイズ分類器を有しており、

前記オーディオ・コンテキスト分類器は、VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズの信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されている、

50

態様 2 2 記載のボリューム平準化器コントローラ。

〔態様 2 4〕

前記オーディオ・コンテキスト分類器が：

VoIP発話の信頼値が第一の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；それ以外の場合には、前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、態様 2 2 記載のボリューム平準化器コントローラ。

〔態様 2 5〕

前記オーディオ・コンテキスト分類器が：

VoIP発話の信頼値が第一の閾値より大きい場合またはVoIPノイズの信頼値が第三の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合またはVoIPノイズの信頼値が前記第三の閾値より大きくない第四の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；それ以外の場合には前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、態様 2 3 記載のボリューム平準化器コントローラ。

〔態様 2 6〕

前記コンテンツ型の過去の信頼値に基づいて現在の時点での前記コンテンツ型の信頼値を平滑化するための型平滑化ユニットをさらに有する、態様 2 1 ないし 2 5 のうちいずれか一項記載のボリューム平準化器コントローラ。

〔態様 2 7〕

前記型平滑化ユニットは、現在の短期的セグメントの信頼値と最後の短期的セグメントの平滑化された信頼値との重み付けされた和を計算することによって、現在の短期的セグメントの平滑化された信頼値を決定するよう構成されている、態様 2 6 記載のボリューム平準化器コントローラ。

〔態様 2 8〕

前記オーディオ・コンテンツ分類器が前記短期的セグメントの発話のコンテンツ型を識別する発話／ノイズ分類器をさらに有しており、前記型平滑化ユニットは、平滑化前の現在の短期的セグメントについてのVoIP発話の信頼値を、所定の信頼値として、あるいは前記発話／ノイズ分類器によって分類されるコンテンツ型発話についての信頼値が第五の閾値より低い最後の短期的セグメントの平滑化された信頼値として、設定するよう構成されている、態様 2 6 記載のボリューム平準化器コントローラ。

〔態様 2 9〕

前記オーディオ・コンテキスト分類器が、特徴として、前記短期的セグメントのコンテンツ型の信頼値および前記短期的セグメントから抽出された他の特徴を使って、機械学習モデルに基づいて前記短期的セグメントを分類するよう構成されている、態様 2 2 または 2 3 記載のボリューム平準化器コントローラ。

〔態様 3 0〕

前記オーディオ・コンテキスト分類器が同じコンテキスト型を連続的に出力する継続時間を測定するタイマーをさらに有しており、前記調整ユニットは、新しいコンテキスト型の継続時間の長さが第六の閾値に達するまで、現在のコンテキスト型を使い続けるよう構成される、態様 1 4 ないし 2 9 のうちいずれか一項記載のボリューム平準化器コントローラ。

〔態様 3 1〕

あるコンテキスト型から別のコンテキスト型への異なる遷移対について、異なる第六の閾値が設定される、態様 3 0 記載のボリューム平準化器コントローラ。

〔態様 3 2〕

前記第六の閾値が、前記新しいコンテキスト型の信頼値と負に相関している、態様 3 0

10

20

30

40

50

記載のボリューム平準化器コントローラ。

〔態様 33〕

前記第一および/または第二の閾値が、最後の短期的セグメントのコンテキスト型によって異なる、態様 24 または 25 記載のボリューム平準化器コントローラ。

〔態様 34〕

態様 1 ないし 33 のうちいずれか一項記載のボリューム平準化器コントローラを有するオーディオ処理装置。

〔態様 35〕

オーディオ信号の短期的セグメントのコンテンツ型を識別するオーディオ・コンテンツ分類器と；

少なくとも部分的には前記オーディオ・コンテンツ分類器によって識別されたコンテンツ型に基づいて前記短期的セグメントのコンテキスト型を識別するオーディオ・コンテキスト分類器とを有する、

オーディオ分類器。

〔態様 36〕

前記オーディオ・コンテンツ分類器が、前記短期的セグメントをコンテンツ型VoIP発話またはコンテンツ型非VoIP発話に分類するVoIP発話分類器を有しており、

前記オーディオ・コンテキスト分類器は、VoIP発話および非VoIP発話の信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されている、

態様 35 記載のオーディオ分類器。

〔態様 37〕

前記オーディオ・コンテンツ分類器がさらに、

前記短期的セグメントをコンテンツ型VoIPノイズおよびコンテンツ型非VoIPノイズに分類するVoIPノイズ分類器を有しており、

前記オーディオ・コンテキスト分類器は、VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズの信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されている、

態様 36 記載のオーディオ分類器。

〔態様 38〕

前記オーディオ・コンテキスト分類器が：

VoIP発話の信頼値が第一の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；それ以外の場合には、前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、態様 37 記載のオーディオ分類器。

〔態様 39〕

前記オーディオ・コンテキスト分類器が：

VoIP発話の信頼値が第一の閾値より大きい場合またはVoIPノイズの信頼値が第三の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合またはVoIPノイズの信頼値が前記第三の閾値より大きくない第四の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；それ以外の場合には前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、

態様 37 記載のオーディオ分類器。

〔態様 40〕

前記コンテンツ型の過去の信頼値に基づいて現在の時点での前記コンテンツ型の信頼値を平滑化するための型平滑化ユニットをさらに有する、態様 35 ないし 39 のうちいずれか一項記載のオーディオ分類器。

10

20

30

40

50

〔態様 4 1〕

前記型平滑化ユニットは、現在の短期的セグメントの信頼値と最後の短期的セグメントの平滑化された信頼値との重み付けされた和を計算することによって、現在の短期的セグメントの平滑化された信頼値を決定するよう構成されている、態様 4 0 記載のオーディオ分類器。

〔態様 4 2〕

前記オーディオ・コンテンツ分類器が前記短期的セグメントからコンテンツ型発話を識別する発話/ノイズ分類器をさらに有しており、前記型平滑化ユニットは、平滑化前の現在の短期的セグメントについてのVoIP発話の信頼値を、所定の信頼値として、あるいは前記発話/ノイズ分類器によって分類されるコンテンツ型発話についての信頼値が第五の閾値より低い最後の短期的セグメントの平滑化された信頼値として、設定するよう構成されている、態様 4 1 記載のオーディオ分類器。

10

〔態様 4 3〕

前記オーディオ・コンテキスト分類器が、特徴として、前記短期的セグメントのコンテンツ型の信頼値および前記短期的セグメントから抽出された他の特徴を使って、機械学習モデルに基づいて前記短期的セグメントを分類するよう構成されている、態様 3 6 または 3 7 記載のオーディオ分類器。

〔態様 4 4〕

前記オーディオ・コンテキスト分類器が同じコンテキスト型を連続的に出力する継続時間を測定するタイマーをさらに有しており、当該オーディオ分類器は、新しいコンテキスト型の継続時間の長さが第六の閾値に達するまで、現在のコンテキスト型を出力し続けるよう構成される、態様 3 8 または 3 9 記載のオーディオ分類器。

20

〔態様 4 5〕

あるコンテキスト型から別のコンテキスト型への異なる遷移対について、異なる第六の閾値が設定される、態様 4 4 記載のオーディオ分類器。

〔態様 4 6〕

前記第六の閾値が、前記新しいコンテキスト型の信頼値と負に相関している、態様 4 4 記載のオーディオ分類器。

〔態様 4 7〕

前記第一および/または第二の閾値が、最後の短期的セグメントのコンテキスト型によって異なる、態様 3 8 または 3 9 記載のオーディオ分類器。

30

〔態様 4 8〕

態様 3 5 ないし 4 7 のうちいずれか一項記載のオーディオ分類器を有するオーディオ処理装置。

〔態様 4 9〕

リアルタイムでオーディオ信号のコンテンツ型を識別する段階と；

識別されたコンテンツ型に基づいて連続的な仕方でボリューム平準化器を調整することを、前記ボリューム平準化器の動的な利得を、前記オーディオ信号の情報性のコンテンツ型と正に相関させ、前記ボリューム平準化器の動的な利得を、前記オーディオ信号の干渉性のコンテンツ型と負に相関させることによって行なう段階とを含む、ボリューム平準化器制御方法。

40

〔態様 5 0〕

前記オーディオ信号の前記コンテンツ型が、発話、短期的音楽、ノイズおよび背景音のうちの一つを含む、態様 4 9 記載のボリューム平準化器制御方法。

〔態様 5 1〕

ノイズが干渉性のコンテンツ型と見なされる、態様 4 9 記載のボリューム平準化器制御方法。

〔態様 5 2〕

前記調整する動作が、前記コンテンツ型の信頼値に基づいて前記ボリューム平準化器の動的な利得を調整するよう構成されている、態様 4 9 記載のボリューム平準化器制御方法

50

。

〔態様 5 3〕

前記調整する動作が、前記コンテンツ型の信頼値の伝達関数を介して前記動的な利得を調整するよう構成されている、態様 5 2 記載のボリューム平準化器制御方法。

〔態様 5 4〕

前記オーディオ信号が、対応する信頼値をもつ複数のコンテンツ型に分類され、前記調整する動作が、前記複数のコンテンツ型の重要性に基づいて前記複数のコンテンツ型の前記信頼値を重み付けすることを通じて前記複数のオーディオ型の少なくともいくつかを考慮するよう構成されている、態様 4 9 記載のボリューム平準化器制御方法。

〔態様 5 5〕

前記オーディオ信号が、対応する信頼値をもつ複数のコンテンツ型に分類され、前記調整する動作が、あるコンテンツ型の重みを少なくとも一つの他のコンテンツ型の信頼値を用いて修正するよう構成されている、態様 4 9 記載のボリューム平準化器制御方法。

10

〔態様 5 6〕

前記オーディオ信号が、対応する信頼値をもつ複数のコンテンツ型に分類され、前記調整する動作が、前記信頼値に基づいて前記複数のコンテンツ型の効果を重み付けすることを通じて前記複数のコンテンツ型の少なくともいくつかを考慮するよう構成されている、態様 4 9 記載のボリューム平準化器制御方法。

〔態様 5 7〕

前記調整する動作が、前記信頼値に基づいて少なくとも一つの優勢なコンテンツ型を考慮するよう構成されている、態様 5 6 記載のボリューム平準化器制御方法。

20

〔態様 5 8〕

前記オーディオ信号が、対応する信頼値をもつ複数の干渉性のコンテンツ型および/または複数の情報性のコンテンツ型に分類され、前記調整する動作が、前記信頼値に基づいて少なくとも一つの優勢な干渉性のコンテンツ型および/または少なくとも一つの優勢な情報性のコンテンツ型を考慮するよう構成されている、態様 5 6 記載のボリューム平準化器制御方法。

〔態様 5 9〕

各コンテンツ型について、前記オーディオ信号の過去の信頼値に基づいて現在の時点での前記オーディオ信号の信頼値を平滑化する段階をさらに含む、態様 4 9 ないし 5 8 のうちいずれか一項記載のボリューム平準化器制御方法。

30

〔態様 6 0〕

前記の型平滑化の動作は、現在の実際の信頼値と最後の時点での平滑化された信頼値との重み付けされた和を計算することによって、現時点での前記オーディオ信号の平滑化された信頼値を決定するよう構成されている、態様 5 9 記載のボリューム平準化器制御方法。

〔態様 6 1〕

前記オーディオ信号のコンテキスト型を識別する段階をさらに含み、前記調整する動作は、前記コンテキスト型の信頼値に基づいて前記動的な利得の範囲を調整するよう構成されている、態様 4 9 ないし 5 8 のうちいずれか一項記載のボリューム平準化器制御方法。

40

〔態様 6 2〕

前記オーディオ信号のコンテキスト型を識別する段階をさらに含み、前記調整する動作は、前記オーディオ信号の前記コンテキスト型に基づいて前記オーディオ信号の前記コンテンツ型を情報性または干渉性に見なすよう構成されている、態様 4 9 ないし 5 8 のうちいずれか一項記載のボリューム平準化器制御方法。

〔態様 6 3〕

前記オーディオ信号の前記コンテキスト型が、VoIP、映画のメディア、長期的音楽およびゲームのうちの一つを含む、態様 6 2 記載のボリューム平準化器制御方法。

〔態様 6 4〕

コンテキスト型VoIPのオーディオ信号においては、背景音が干渉性コンテンツ型と見な

50

され、一方、コンテキスト型非VoIPのオーディオ信号においては、背景音および/または発話および/または音楽が情報性コンテンツ型と見なされる、態様62記載のボリューム平準化器制御方法。

〔態様65〕

前記オーディオ信号の前記コンテキスト型が高品質オーディオまたは低品質オーディオを含む、態様62記載のボリューム平準化器制御方法。

〔態様66〕

異なるコンテキスト型のオーディオ信号におけるコンテンツ型が、前記オーディオ信号のコンテキスト型に依存して異なる重みを割り当てられる、態様62記載のボリューム平準化器制御方法。

〔態様67〕

前記オーディオ信号が、対応する信頼値をもつ複数のコンテキスト型に分類され、前記調整する動作が、前記複数のコンテキスト型の重要性に基づいて前記複数のコンテキスト型の前記信頼値を重み付けすることを通じて前記複数のコンテキスト型の少なくともいくつかを考慮するよう構成されている、態様62記載のボリューム平準化器制御方法。

〔態様68〕

前記オーディオ信号が、対応する信頼値をもつ複数のコンテキスト型に分類され、前記調整する動作が、前記信頼値に基づいて前記複数のコンテキスト型の効果を重み付けすることを通じて前記複数のコンテキスト型の少なくともいくつかを考慮するよう構成されている、態様62記載のボリューム平準化器制御方法。

〔態様69〕

前記コンテンツ型を識別する動作が、前記オーディオ信号の短期的セグメントに基づいて前記コンテンツ型を識別するよう構成されており、

前記コンテキスト型を識別する動作が、少なくとも部分的には識別されたコンテンツ型に基づく前記オーディオ信号の短期的セグメントに基づいて前記コンテキスト型を識別するよう構成されている、

態様62記載のボリューム平準化器制御方法。

〔態様70〕

コンテンツ型を識別する動作が、短期的セグメントをコンテンツ型VoIP発話またはコンテンツ型非VoIP発話に分類することを含み、

コンテキスト型を識別する動作が、VoIP発話および非VoIP発話の信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されている、態様69記載のボリューム平準化器制御方法。

〔態様71〕

コンテンツ型を識別する動作がさらに、

短期的セグメントをコンテンツ型VoIPノイズおよびコンテンツ型非VoIPノイズに分類することを含み、

コンテキスト型を識別する動作は、VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズの信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されている、

態様70記載のボリューム平準化器制御方法。

〔態様72〕

コンテキスト型を識別する動作が：

VoIP発話の信頼値が第一の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；

VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；

それ以外の場合には、前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、

態様70記載のボリューム平準化器制御方法。

10

20

30

40

50

〔態様 7 3〕

コンテキストを識別する動作が：

VoIP発話の信頼値が第一の閾値より大きい場合またはVoIPノイズの信頼値が第三の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；

VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合またはVoIPノイズの信頼値が前記第三の閾値より大きくない第四の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；

それ以外の場合には前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、態様 7 1 記載のボリューム平準化器制御方法。

10

〔態様 7 4〕

前記コンテンツ型の過去の信頼値に基づいて現在の時点での前記コンテンツ型の信頼値を平滑化する段階をさらに含む、態様 6 9 ないし 7 3 のうちいずれか一項記載のボリューム平準化器制御方法。

〔態様 7 5〕

前記の型平滑化の動作は、現在の短期的セグメントの信頼値と最後の短期的セグメントの平滑化された信頼値との重み付けされた和を計算することによって、現在の短期的セグメントの平滑化された信頼値を決定するよう構成されている、態様 7 4 記載のボリューム平準化器制御方法。

〔態様 7 6〕

前記短期的セグメントの発話のコンテンツ型を識別する段階をさらに含み、平滑化前の現在の短期的セグメントについてのVoIP発話の信頼値が、所定の信頼値として、あるいはコンテンツ型発話についての信頼値が第五の閾値より低い最後の短期的セグメントの平滑化された信頼値として、設定される、態様 7 5 記載のボリューム平準化器制御方法。

20

〔態様 7 7〕

特徴として、前記短期的セグメントのコンテンツ型の信頼値および前記短期的セグメントから抽出された他の特徴を使って、機械学習モデルに基づいて、前記短期的セグメントが分類される、態様 7 0 または 7 1 記載のボリューム平準化器制御方法。

〔態様 7 8〕

コンテキスト型を識別する動作が同じコンテキスト型を連続的に出力する継続時間を測定することをさらに含み、前記調整する動作は、新しいコンテキスト型の継続時間の長さが第六の閾値に達するまで、現在のコンテキスト型を使い続けるよう構成される、態様 6 2 ないし 7 7 のうちいずれか一項記載のボリューム平準化器制御方法。

30

〔態様 7 9〕

あるコンテキスト型から別のコンテキスト型への異なる遷移対について、異なる第六の閾値が設定される、態様 7 8 記載のボリューム平準化器制御方法。

〔態様 8 0〕

前記第六の閾値が、前記新しいコンテキスト型の信頼値と負に相関している、態様 7 8 記載のボリューム平準化器制御方法。

〔態様 8 1〕

前記第一および/または第二の閾値が、最後の短期的セグメントのコンテキスト型によって異なる、態様 7 2 または 7 3 記載のボリューム平準化器制御方法。

40

〔態様 8 2〕

オーディオ信号の短期的セグメントのコンテンツ型を識別する段階と；

少なくとも部分的には識別されたコンテンツ型に基づいて前記短期的セグメントのコンテキスト型を識別する段階とを含む、

オーディオ分類方法。

〔態様 8 3〕

コンテンツ型を分類する動作が、前記短期的セグメントをコンテンツ型VoIP発話またはコンテンツ型非VoIP発話に分類することを含み、

50

コンテキスト型を識別する動作が、VoIP発話および非VoIP発話の信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されている、

態様 8 2 記載のオーディオ分類方法。

〔態様 8 4〕

コンテンツ型を分類する動作がさらに、

前記短期的セグメントをコンテンツ型VoIPノイズまたはコンテンツ型非VoIPノイズに分類することを含み、

コンテキスト型を識別する動作が、VoIP発話、非VoIP発話、VoIPノイズおよび非VoIPノイズの信頼値に基づいて、前記短期的セグメントをコンテキスト型VoIPまたはコンテキスト型非VoIPに分類するよう構成されている、

態様 8 3 記載のオーディオ分類方法。

〔態様 8 5〕

コンテキスト型を識別する動作が：

VoIP発話の信頼値が第一の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；

VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；

それ以外の場合には、前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、

態様 8 3 記載のオーディオ分類方法。

〔態様 8 6〕

コンテキスト型を識別する動作が：

VoIP発話の信頼値が第一の閾値より大きい場合またはVoIPノイズの信頼値が第三の閾値より大きい場合、前記短期的セグメントをコンテキスト型VoIPとして分類し；

VoIP発話の信頼値が、前記第一の閾値より大きくない第二の閾値より大きくない場合またはVoIPノイズの信頼値が前記第三の閾値より大きくない第四の閾値より大きくない場合、前記短期的セグメントをコンテキスト型非VoIPとして分類し；

それ以外の場合には前記短期的セグメントを最後の短期的セグメントについてのコンテキスト型として分類するよう構成されている、

態様 8 4 記載のオーディオ分類方法。

〔態様 8 7〕

前記コンテンツ型の過去の信頼値に基づいて現在の時点での前記コンテンツ型の信頼値を平滑化する段階をさらに含む、態様 8 2 ないし 8 6 のうちいずれか一項記載のオーディオ分類方法。

〔態様 8 8〕

前記の型平滑化の動作は、現在の短期的セグメントの信頼値と最後の短期的セグメントの平滑化された信頼値との重み付けされた和を計算することによって、現在の短期的セグメントの平滑化された信頼値を決定するよう構成されている、態様 8 7 記載のオーディオ分類方法。

〔態様 8 9〕

前記短期的セグメントからコンテンツ型発話を識別する段階をさらに含み、平滑化前の現在の短期的セグメントについてのVoIP発話の信頼値が、所定の信頼値として、あるいはコンテンツ型発話についての信頼値が第五の閾値より低い最後の短期的セグメントの平滑化された信頼値として、設定される、態様 8 8 記載のオーディオ分類方法。

〔態様 9 0〕

コンテキスト型を識別する動作が、特徴として、前記短期的セグメントのコンテンツ型の信頼値および前記短期的セグメントから抽出された他の特徴を使って、機械学習モデルに基づいて前記短期的セグメントを分類するよう構成されている、態様 8 3 または 8 4 記載のオーディオ分類方法。

10

20

30

40

50

〔態様 9 1〕

コンテキスト型を識別する動作が同じコンテキスト型を連続的に出力する継続時間を測定する段階をさらに含み、当該オーディオ分類方法は、新しいコンテキスト型の継続時間の長さが第六の閾値に達するまで、現在のコンテキスト型を出力し続けるよう構成される、態様 8 5 または 8 6 記載のオーディオ分類方法。

〔態様 9 2〕

あるコンテキスト型から別のコンテキスト型への異なる遷移対について、異なる第六の閾値が設定される、態様 9 1 記載のオーディオ分類方法。

〔態様 9 3〕

前記第六の閾値が、前記新しいコンテキスト型の信頼値と負に相関している、態様 9 1 記載のオーディオ分類方法。

10

〔態様 9 4〕

前記第一および/または第二の閾値が、最後の短期的セグメントのコンテキスト型によって異なる、態様 8 5 または 8 6 記載のオーディオ分類方法。

〔態様 9 5〕

プロセッサによって実行されると該プロセッサがボリューム平準化器制御方法を実行できるようにするコンピュータ・プログラム命令が記録されたコンピュータ可読媒体であって、前記ボリューム平準化器制御方法は、

リアルタイムでオーディオ信号のコンテンツ型を識別する段階と；

識別されたコンテンツ型に基づいて連続的な仕方でボリューム平準化器を調整することを、前記ボリューム平準化器の動的な利得を、前記オーディオ信号の情報性のコンテンツ型と正に相関させ、前記ボリューム平準化器の動的な利得を、前記オーディオ信号の干渉性のコンテンツ型と負に相関させることによって行なう段階とを含む、コンピュータ可読媒体。

20

〔態様 9 6〕

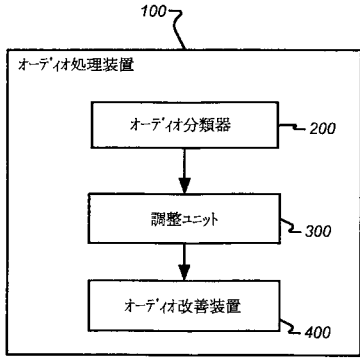
プロセッサによって実行されると該プロセッサがオーディオ分類方法を実行できるようにするコンピュータ・プログラム命令が記録されたコンピュータ可読媒体であって、前記オーディオ分類方法は、

オーディオ信号の短期的セグメントのコンテンツ型を識別する段階と；

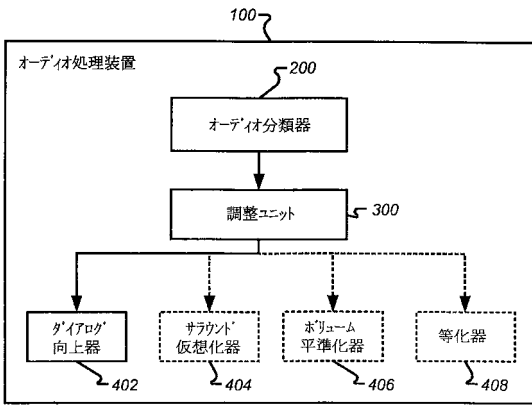
少なくとも部分的には識別されたコンテンツ型に基づいて前記短期的セグメントのコンテキスト型を識別する段階とを含む、コンピュータ可読媒体。

30

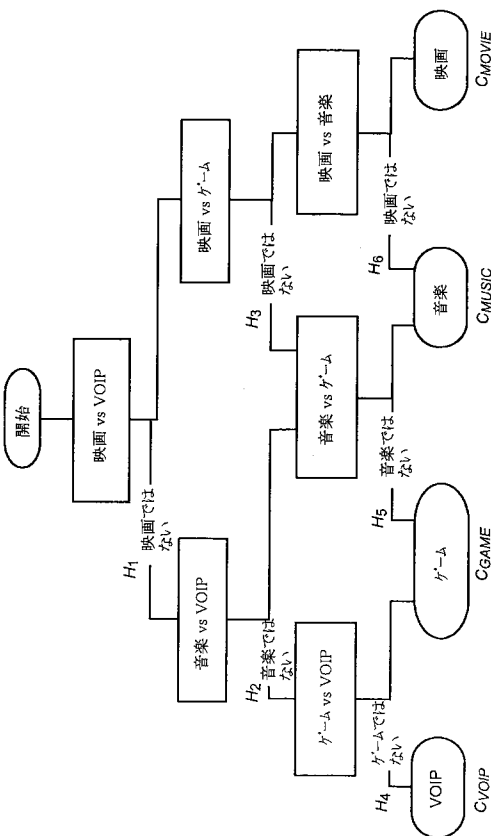
【 図 1 】



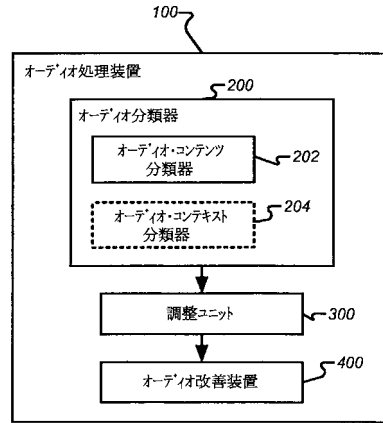
【 図 2 】



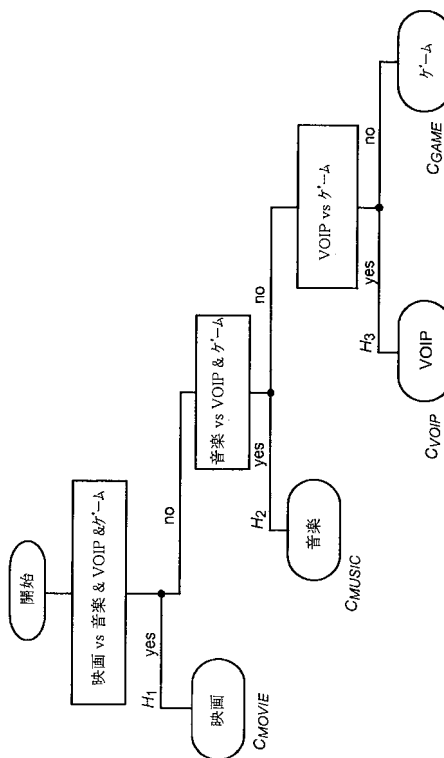
【 図 4 】



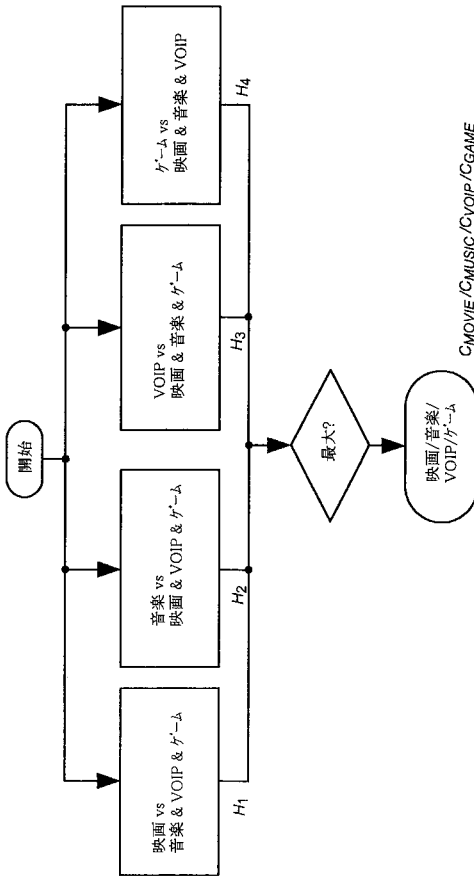
【 図 3 】



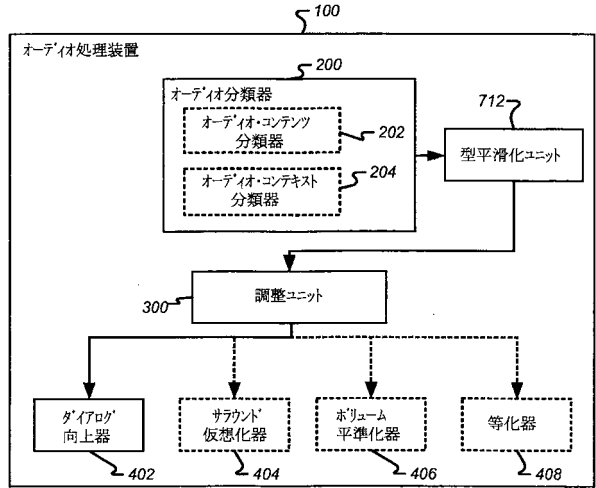
【 図 5 】



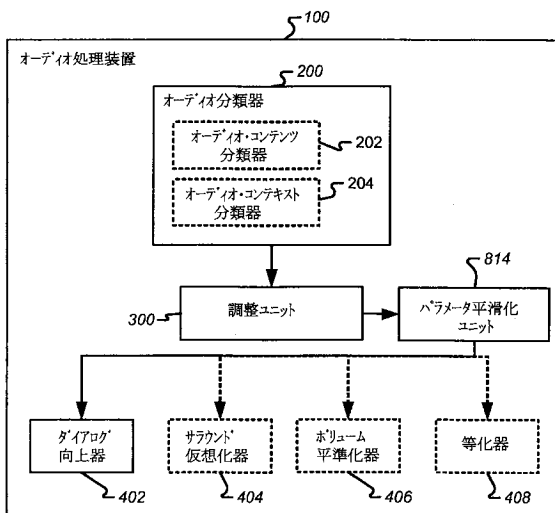
【 図 6 】



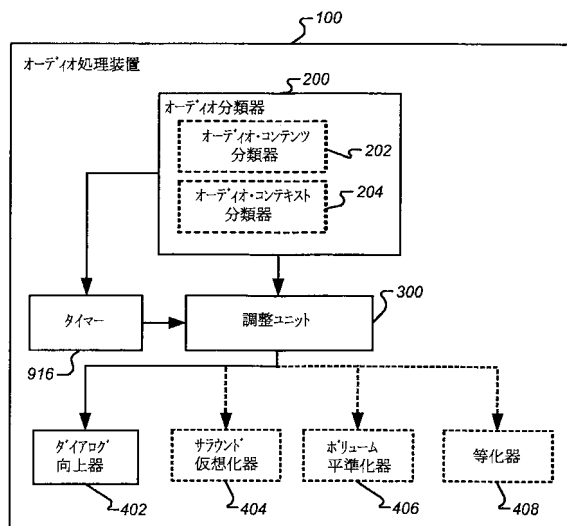
【 図 7 】



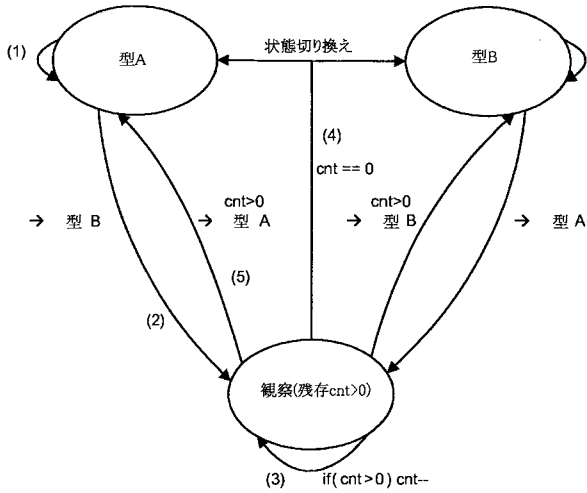
【 図 8 】



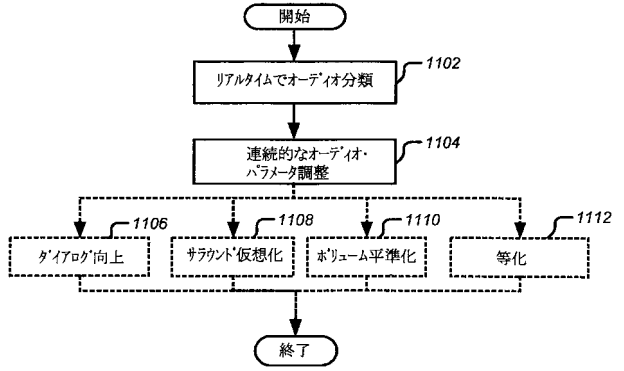
【 図 9 】



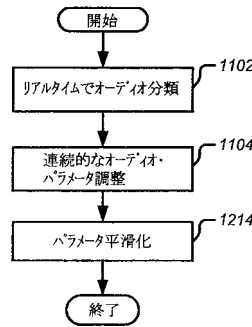
【図10】



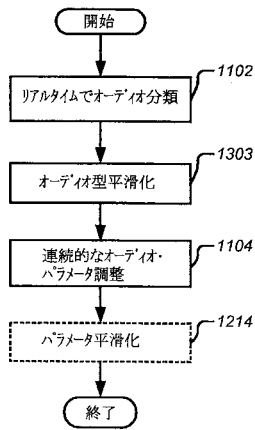
【図11】



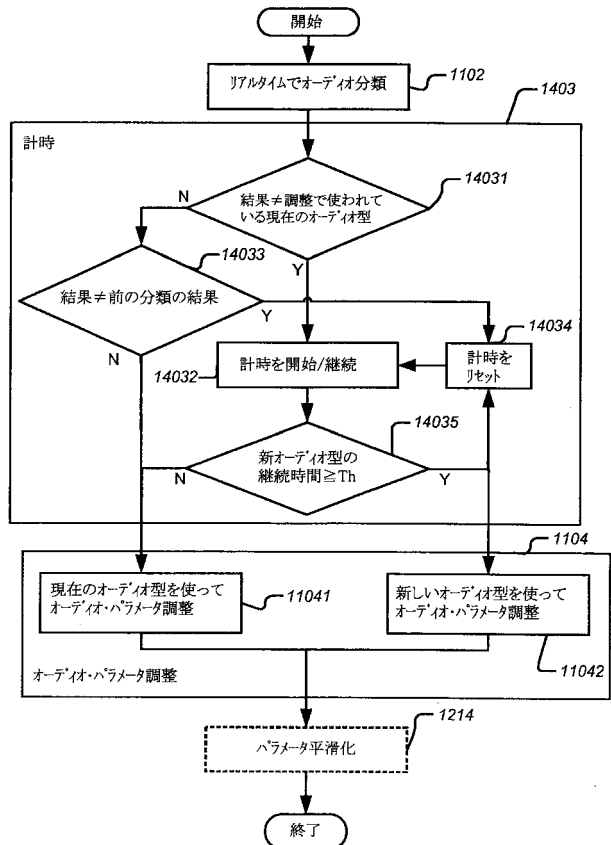
【図12】



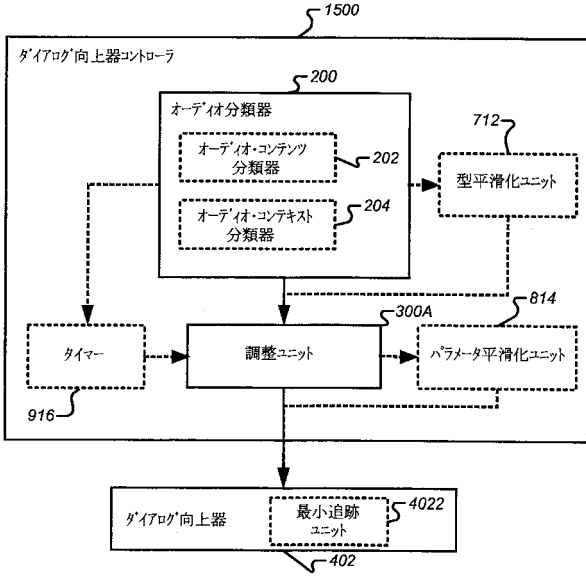
【図13】



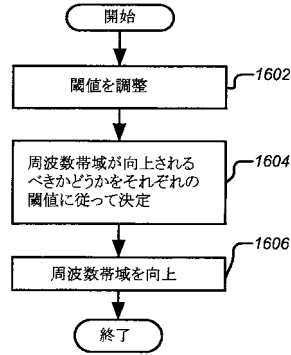
【図14】



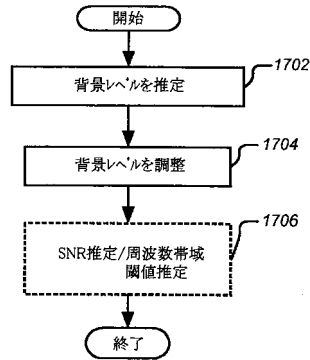
【図15】



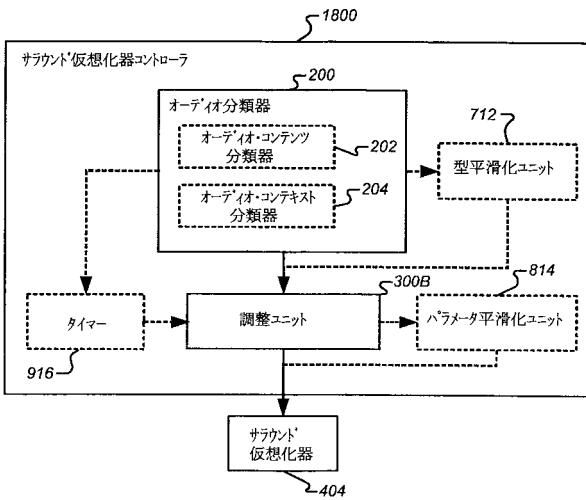
【図16】



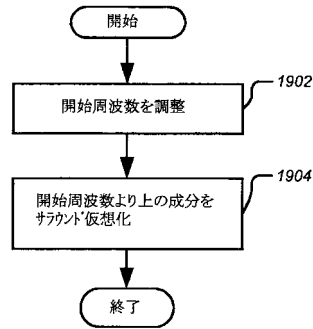
【図17】



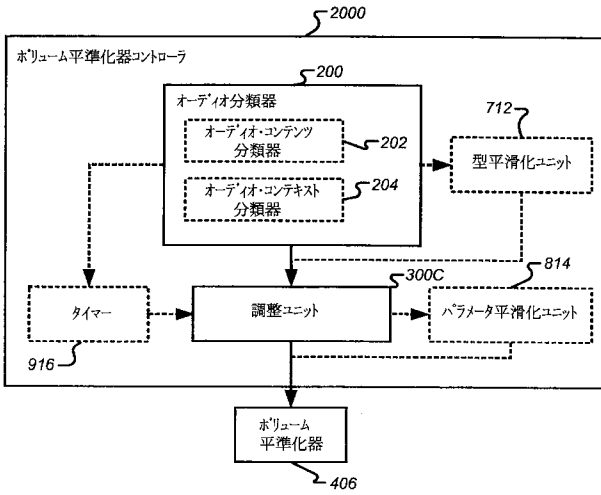
【図18】



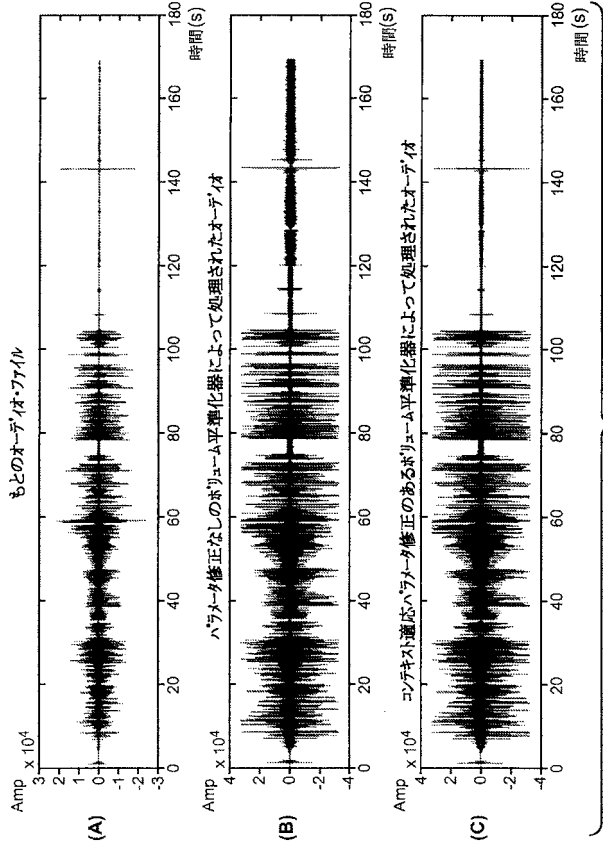
【図19】



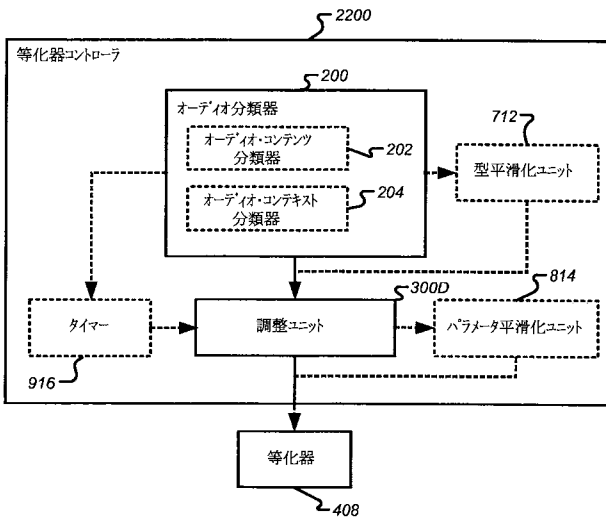
【図20】



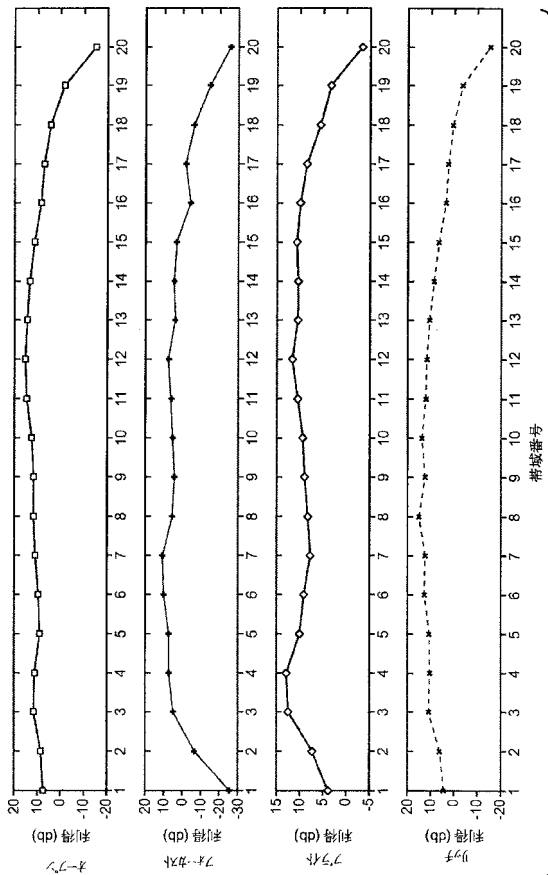
【図21】



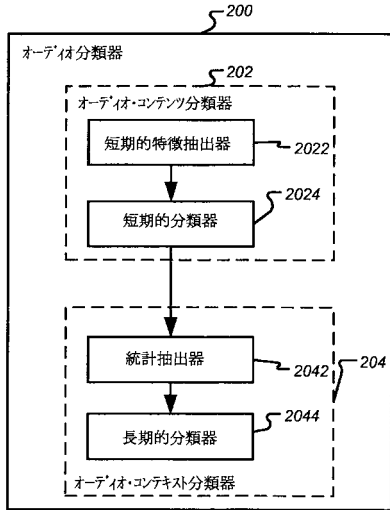
【図22】



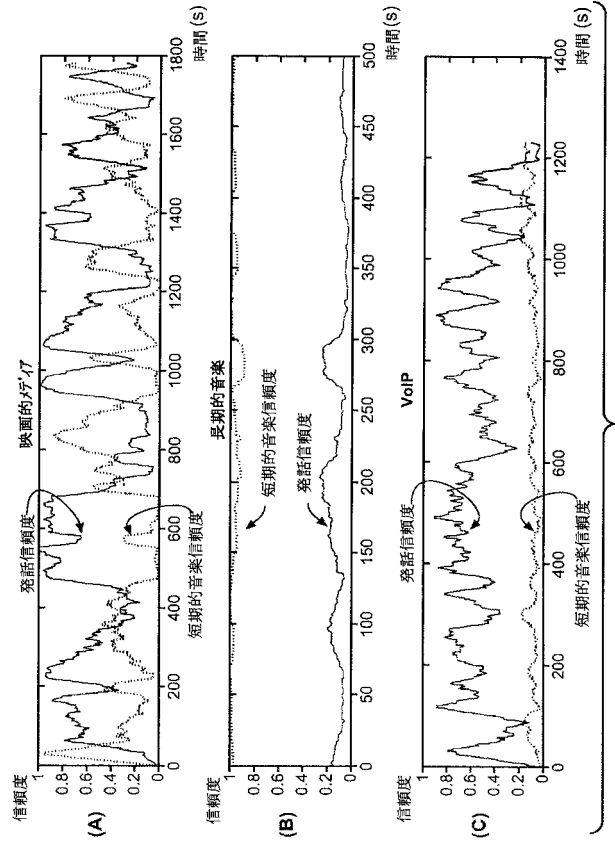
【図23】



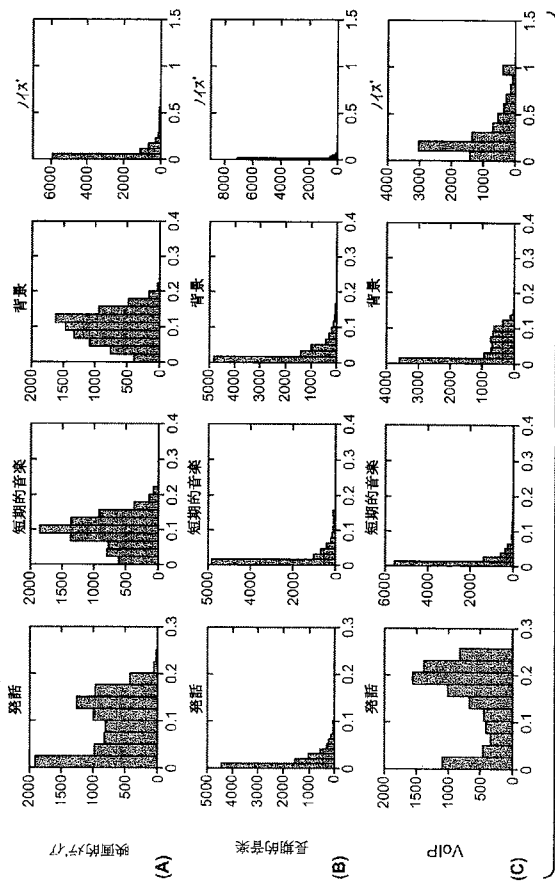
【図 2 4】



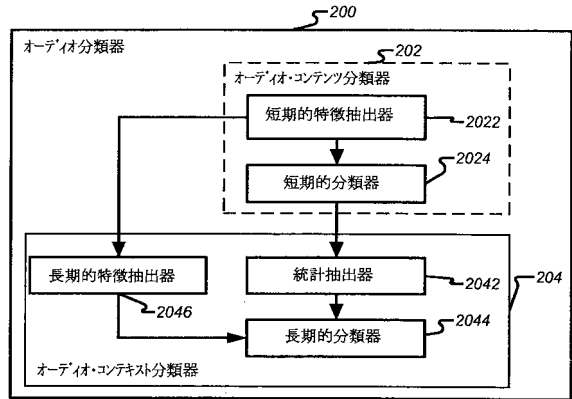
【図 2 5】



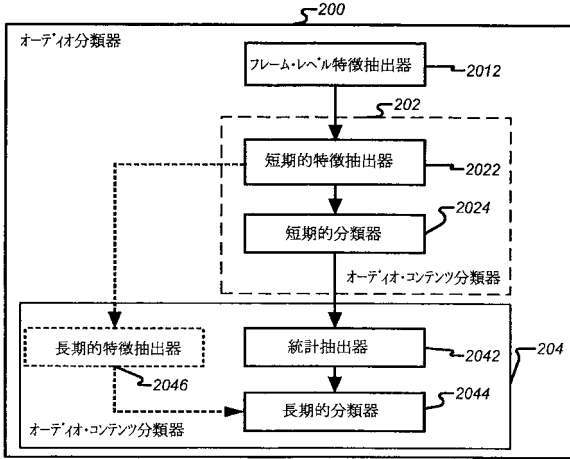
【図 2 6】



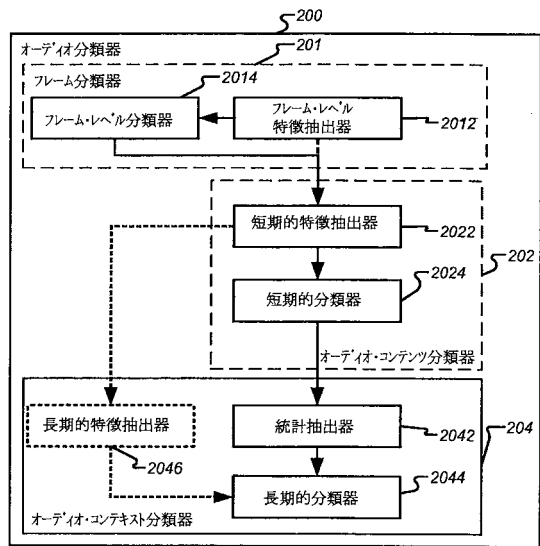
【図 2 7】



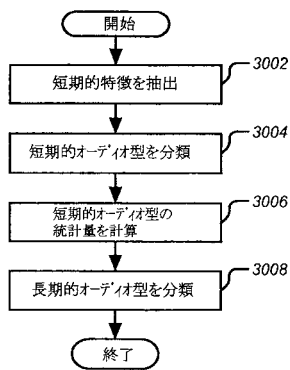
【図 28】



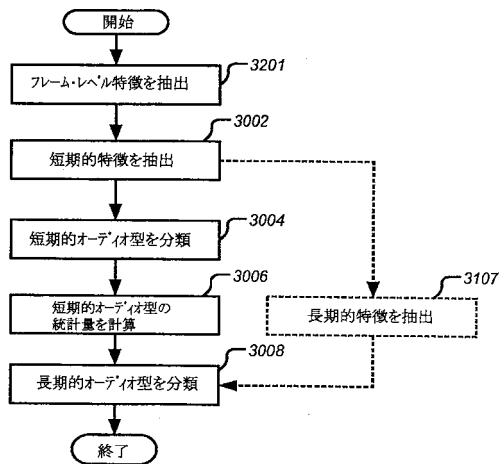
【図 29】



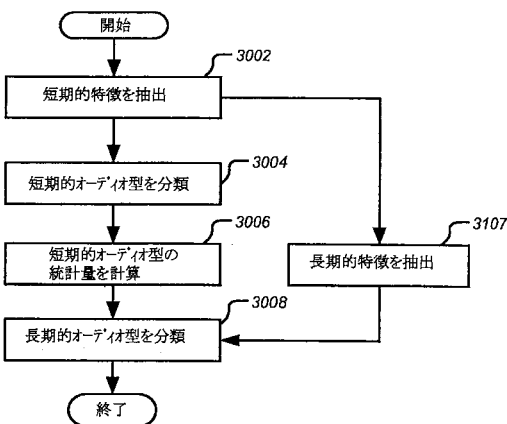
【図 30】



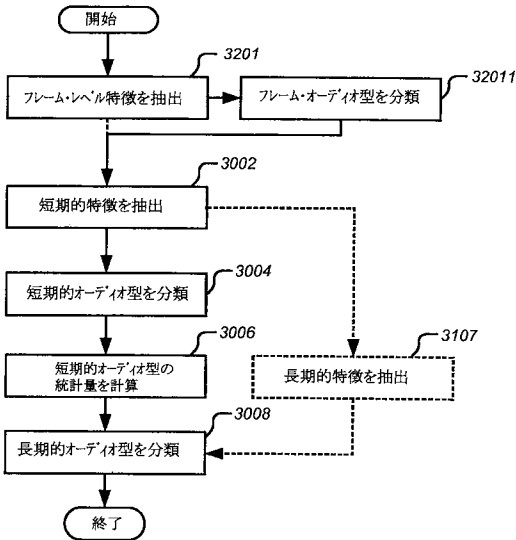
【図 32】



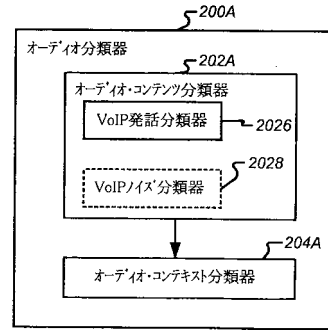
【図 31】



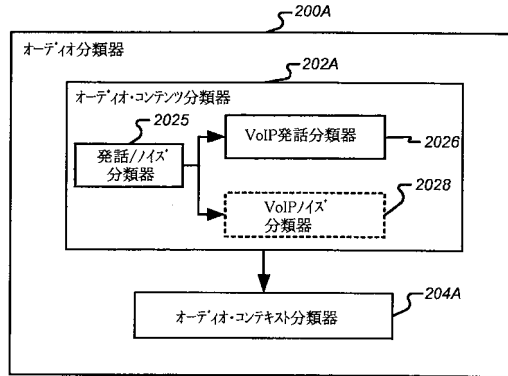
【図33】



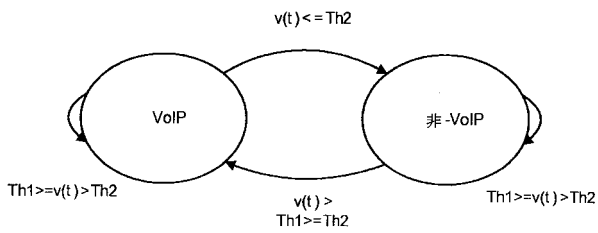
【図34】



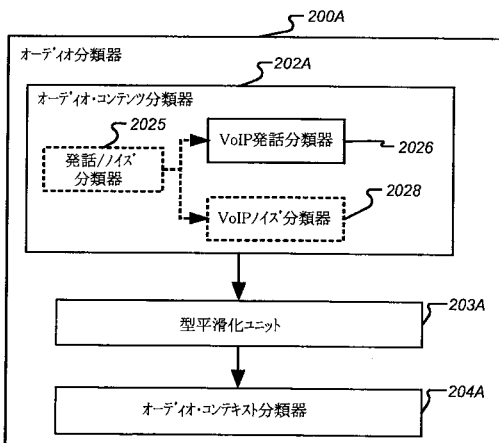
【図35】



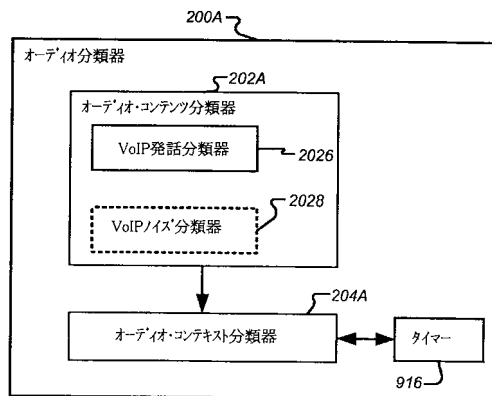
【図36】



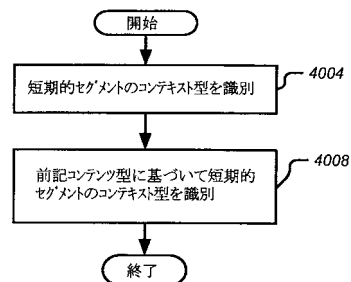
【図37】



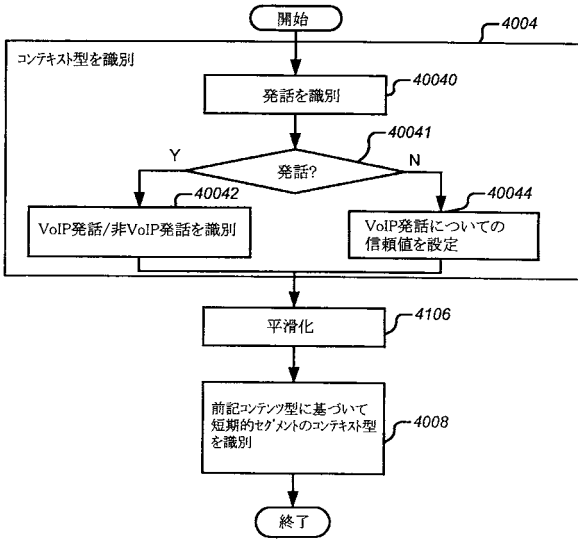
【図38】



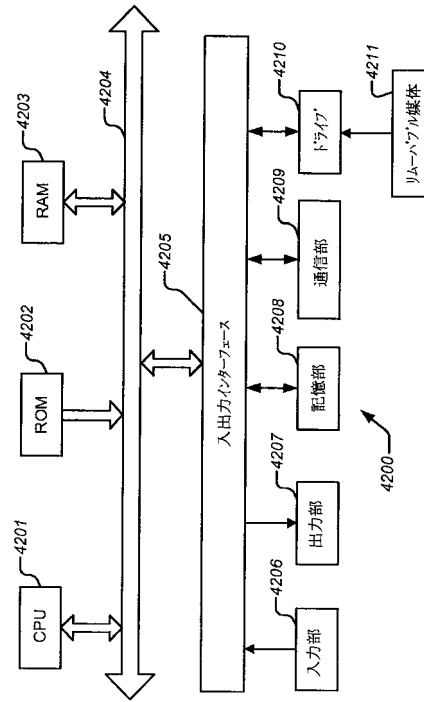
【図39】



【図40】



【図41】



フロントページの続き

(72)発明者 ワン, ジュン

中華人民共和国, ベイジン 100020, 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, ナンバー 1, ワールド・フィナンシャル・センター 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(ベイジン) カンパニー リミテッド内

(72)発明者 ルー, リエ

中華人民共和国, ベイジン 100020, 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, ナンバー 1, ワールド・フィナンシャル・センター 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(ベイジン) カンパニー リミテッド内

(72)発明者 シーフエルドット, アラン

アメリカ合衆国, カリフォルニア州 94103-4813, サンフランシスコ, ポットレロ Avenue 100, ドルビー ラボラトリーズ, インコーポレイテッド内

Fターム(参考) 5J030 AA01 AA15 AC10 AC20