

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7443823号
(P7443823)

(45)発行日 令和6年3月6日(2024.3.6)

(24)登録日 令和6年2月27日(2024.2.27)

(51)国際特許分類	F I
G 1 0 L 21/0272(2013.01)	G 1 0 L 21/0272 1 0 0 Z
G 1 0 L 25/30 (2013.01)	G 1 0 L 25/30

請求項の数 1 (全24頁)

(21)出願番号	特願2020-33347(P2020-33347)	(73)特許権者	000004075 ヤマハ株式会社 静岡県浜松市中央区中沢町10番1号
(22)出願日	令和2年2月28日(2020.2.28)	(74)代理人	110003177 弁理士法人旺知国際特許事務所
(65)公開番号	特開2021-135446(P2021-135446 A)	(72)発明者	北村 大地 香川県三豊市詫間町香田551 香川高 等専門学校内
(43)公開日	令和3年9月13日(2021.9.13)	(72)発明者	渡辺 瑠伊 香川県三豊市詫間町香田551 香川高 等専門学校内
審査請求日	令和4年12月20日(2022.12.20)	審査官	大野 弘

最終頁に続く

(54)【発明の名称】 音響処理方法

(57)【特許請求の範囲】

【請求項1】

第1音源に対応する第1音のうち第1周波数帯域の成分を表す第1入力データと、前記第1音源とは異なる第2音源に対応する第2音のうち前記第1周波数帯域の成分を表す第2入力データと、前記第1音と前記第2音との混合音のうち前記第1周波数帯域とは異なる第2周波数帯域を含む周波数帯域の成分を含む音を表す混合音データと、を含む入力データを取得し、

学習済の推定モデルに前記入力データを入力することで、前記第1音のうち前記第2周波数帯域を含む周波数帯域の成分を表す第1出力データと、前記第2音のうち前記第2周波数帯域を含む周波数帯域の成分を表す第2出力データとの少なくとも一方を生成する

コンピュータにより実現される音響処理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、音響処理に関する。

【背景技術】

【0002】

相異なる音源が発生した複数の音の混合音を音源毎に分離する音源分離技術が従来から提案されている。例えば非特許文献1には、信号の独立性と音源の低ランク性とを同時に考慮することで高精度な音源分離を実現する独立低ランク行列分析(ILRMA: Independ

10

20

ent Low-Rank Matrix Analysis) が開示されている。また、非特許文献 2 には、振幅スペクトログラムをニューラルネットワークに入力することで、音源分離のための時間-周波数領域マスクを生成する技術が開示されている。

【先行技術文献】

【非特許文献】

【0003】

【文献】 Daichi Kitamura, Nobutaka Ono, Hiroshi Sawada, Hirokazu Kameoka, and Hiroshi Saruwatari, "Determined blind source separation unifying independent vector analysis and nonnegative matrix factorization," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 24, no. 9, pp. 1626-1641, September 2016

10

【0004】

【文献】 Andreas Jansson, Eric J. Humphrey, Nicola Montecchio, Rachel Bittner, Aparna Kumar, Tillman Weyde, "Singing Voice Separation with Deep U-Net Convolutional Networks," Proceedings of the 18th International Society for Music Information Retrieval Conference (ISMIR), 2017

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかし、非特許文献 1 および非特許文献 2 に開示された技術においては、音源分離のための処理負荷が過大であるという問題がある。以上の事情を考慮して、本開示のひとつの態様は、音源分離のための処理負荷を軽減することを目的とする。

20

【課題を解決するための手段】

【0006】

以上の課題を解決するために、本開示のひとつの態様に係る音響処理方法は、第 1 音源に対応する第 1 音のうち第 1 周波数帯域の成分を表す第 1 入力データと、前記第 1 音源とは異なる第 2 音源に対応する第 2 音のうち前記第 1 周波数帯域の成分を表す第 2 入力データと、前記第 1 音と前記第 2 音との混合音のうち前記第 1 周波数帯域とは異なる第 2 周波数帯域を含む周波数帯域の成分を含む音を表す混合音データと、を含む入力データを取得し、学習済の推定モデルに前記入力データを入力することで、前記第 1 音のうち前記第 2 周波数帯域を含む周波数帯域の成分を表す第 1 出力データと、前記第 2 音のうち前記第 2 周波数帯域を含む周波数帯域の成分を表す第 2 出力データとの少なくとも一方を生成する。

30

【図面の簡単な説明】

【0007】

【図 1】音響処理システムの構成を例示するブロック図である。

【図 2】音響処理システムの機能的な構成を例示するブロック図である。

【図 3】入力データおよび出力データの説明図である。

【図 4】推定モデルの構成を例示するブロック図である。

【図 5】音響処理の具体的な手順を例示するフローチャートである。

【図 6】訓練データの説明図である。

40

【図 7】学習処理の具体的な手順を例示するフローチャートである。

【図 8】第 2 実施形態における入力データおよび出力データの説明図である。

【図 9】第 3 実施形態における入力データの模式図である。

【図 10】第 3 実施形態における音響処理システムの機能的な構成を例示するブロック図である。

【図 11】第 1 実施形態および第 3 実施形態による効果の説明図である。

【図 12】第 1 実施形態から第 3 実施形態に関する観測結果の図表である。

【図 13】第 5 実施形態における入力データおよび出力データの説明図である。

【図 14】第 5 実施形態における訓練データの説明図である。

【図 15】第 5 実施形態に係る音響処理システムの機能的な構成を例示するブロック図で

50

ある。

【発明を実施するための形態】

【0008】

A：第1実施形態

図1は、本開示の第1実施形態に係る音響処理システム100の構成を例示するブロック図である。音響処理システム100は、制御装置11と記憶装置12と放音装置13とを具備するコンピュータシステムである。音響処理システム100は、例えばスマートフォン、タブレット端末またはパーソナルコンピュータ等の情報端末により実現される。なお、音響処理システム100は、単体の装置で実現されるほか、相互に別体で構成された複数の装置（例えばクライアントサーバシステム）でも実現される。

10

【0009】

記憶装置12は、制御装置11が実行するプログラムと制御装置11が使用する各種のデータとを記憶する単数または複数のメモリである。記憶装置12は、例えば磁気記録媒体もしくは半導体記録媒体等の公知の記録媒体、または、複数種の記録媒体の組合せで構成される。なお、音響処理システム100とは別体の記憶装置12（例えばクラウドストレージ）を用意し、例えば移動体通信網またはインターネット等の通信網を介して、制御装置11が記憶装置12に対する書込および読出を実行してもよい。すなわち、記憶装置12は音響処理システム100から省略されてもよい。

【0010】

記憶装置12は、音波形を表す時間領域の音響信号 S_x を記憶する。音響信号 S_x は、第1音源から発音される音（以下「第1音」という）と第2音源から発音される音（以下「第2音」という）とが混合された音（以下「混合音」という）を表す。第1音源と第2音源とは別個の音源である。第1音源および第2音源の各々は、歌唱者または楽器等の発音源である。例えば、第1音は、歌唱者（第1音源）が発音する歌唱音声であり、第2音は、打楽器等の楽器（第2音源）が発音する楽器音である。音響信号 S_x は、第1音源と第2音源とが並列に発音する環境において例えばマイクロホンアレイ等の収音装置を利用して収録される。ただし、公知の合成技術により合成された信号が音響信号 S_x として利用されてもよい。すなわち、第1音源および第2音源の各々は仮想的な音源でもよい。

20

【0011】

なお、単体の音源のほか複数の音源の集合を第1音源または第2音源として把握してもよい。また、第1音源と第2音源とは基本的には別種の音源であり、第1音と第2音とは音響特性が相違する。ただし、第1音源と第2音源とが相異なる位置に設置された場合のように、各音源の位置を利用して第1音と第2音とを分離可能であれば、第1音源と第2音源とは同種の音源でもよい。すなわち、第1音の音響特性と第2音の音響特性とは、相互に近似または一致してもよい。

30

【0012】

制御装置11は、音響処理システム100の各要素を制御する単数または複数のプロセッサである。具体的には、例えばCPU（Central Processing Unit）、SPU（Sound Processing Unit）、DSP（Digital Signal Processor）、FPGA（Field Programmable Gate Array）、またはASIC（Application Specific Integrated Circuit）等の1種類以上のプロセッサにより、制御装置11が構成される。制御装置11は、記憶装置12に記憶された音響信号 S_x から音響信号 S_z を生成する。音響信号 S_z は、第1音および第2音の一方が他方に対して強調された音を表す時間領域の信号である。すなわち、音響処理システム100は、音響信号 S_x を音源毎に分離する音源分離を実行する。

40

【0013】

放音装置13は、制御装置11が生成した音響信号 S_z が表す音を放音する。放音装置13は、例えばスピーカまたはヘッドホンである。なお、音響信号 S_z をデジタルからアナログに変換するD/A変換器と、音響信号 S_z を増幅する増幅器とは、便宜的に図示が省略されている。また、図1においては、放音装置13を音響処理システム100に搭載した構成を例示したが、音響処理システム100とは別体の放音装置13が有線または無線によ

50

り音響処理システム 100 に接続されてもよい。

【0014】

[1] 音響処理部 20

図 2 は、音響処理システム 100 の機能的な構成を例示するブロック図である。図 2 に例示される通り、制御装置 11 は、記憶装置 12 に記憶された音響処理プログラム P1 を実行することで音響処理部 20 として機能する。音響処理部 20 は、音響信号 S_x から音響信号 S_z を生成する。音響処理部 20 は、周波数解析部 21 と音源分離部 22 と帯域拡張部 23 と波形合成部 24 と音量調整部 25 とを具備する。

【0015】

周波数解析部 21 は、音響信号 S_x の強度スペクトル $X(m)$ を時間軸上の単位期間（フレーム）毎に順次に生成する。記号 m は、時間軸上の 1 個の単位期間を意味する。強度スペクトル $X(m)$ は、例えば振幅スペクトルまたはパワースペクトルである。強度スペクトル $X(m)$ の生成には、例えば短時間フーリエ変換またはウェーブレット変換等の公知の周波数分析が任意に採用される。なお、音響信号 S_x から算定される複素スペクトルが強度スペクトル $X(m)$ とされてもよい。

10

【0016】

図 3 には、音響信号 S_x から生成される強度スペクトル $X(m)$ の時系列（... , $X(m-1)$, $X(m)$, $X(m+1)$, ...）が例示されている。強度スペクトル $X(m)$ は、周波数軸上の所定の周波数帯域（以下「全帯域」という）BF 内に分布する。全帯域 BF は、例えば 0 kHz から 8 kHz までの範囲である。

20

【0017】

音響信号 S_x が表す混合音は、周波数帯域 BL の成分と周波数帯域 BH の成分とを含む。周波数帯域 BL および周波数帯域 BH は、全帯域 BF 内の相異なる周波数帯域である。周波数帯域 BL は周波数帯域 BH よりも低域側に位置する。具体的には、周波数帯域 BL は、全帯域 BF のうち周波数軸上の所定の周波数を下回る帯域であり、周波数帯域 BH は、全帯域 BF のうち当該周波数を上回る帯域である。したがって、周波数帯域 BL と周波数帯域 BH とは相互に重複しない。例えば、周波数帯域 BL は 0 kHz から 4 kHz までの範囲であり、周波数帯域 BH は 4 kHz から 8 kHz までの範囲である。なお、周波数帯域 BL の帯域幅と周波数帯域 BH の帯域幅との異同は不問である。混合音を構成する第 1 音および第 2 音の各々は、周波数帯域 BL の成分と周波数帯域 BH の成分との双方を含む。なお、周波数帯域 BL は「第 1 周波数帯域」の一例であり、周波数帯域 BH は「第 2 周波数帯域」の一例である。

30

【0018】

図 2 の音源分離部 22 は、強度スペクトル $X(m)$ に対する音源分離を実行する。具体的には、音源分離部 22 は、全帯域 BF にわたる強度スペクトル $X(m)$ のうち周波数帯域 BL の成分を対象として音源分離を実行する。すなわち、強度スペクトル $X(m)$ のうち周波数帯域 BH の成分については音源分離の処理対象から除外される。

【0019】

音源分離部 22 による強度スペクトル $X(m)$ の処理には、公知の音源分離が任意に採用される。例えば、独立成分分析（ICA：Independent Component Analysis）、独立ベクトル分析（IVA：Independent Vector Analysis）、非負行列因子分解（NMF：Non-negative Matrix Factorization）、多チャンネル非負行列因子分解（MNMF：Multichannel NMF）、独立低ランク行列分析（ILRMA：Independent Low-Rank Matrix Analysis）、独立低ランクテンソル分析（ILRTA：Independent Low-Rank Tensor Analysis）、または独立深層学習行列分析（IDLMA：Independent Deeply-Learned Matrix Analysis）等の技術が、音源分離部 22 による音源分離に利用される。なお、以上の説明では周波数領域における音源分離を例示したが、音源分離部 22 は、時間領域における音源分離を音響信号 S_x に対して実行してもよい。

40

【0020】

音源分離部 22 は、強度スペクトル $X(m)$ のうち周波数帯域 BL の成分に対する音源分離

50

により強度スペクトル $Y_1(m)$ と強度スペクトル $Y_2(m)$ とを生成する。図3に例示される通り、強度スペクトル $Y_1(m)$ は、混合音に含まれる第1音のうち周波数帯域 BL 内の成分（以下「第1成分」という）のスペクトルを意味する。すなわち、強度スペクトル $Y_1(m)$ は、混合音のうち周波数帯域 BL 内の成分に含まれる第1音を第2音に対して強調した結果（理想的には第2音を除去した結果）を表すスペクトルである。他方、強度スペクトル $Y_2(m)$ は、混合音に含まれる第2音のうち周波数帯域 BL 内の成分（以下「第2成分」という）のスペクトルを意味する。すなわち、強度スペクトル $Y_2(m)$ は、混合音のうち周波数帯域 BL 内の成分に含まれる第2音を第1音に対して強調した結果（理想的には第1音を除去した結果）を表すスペクトルである。以上の説明から理解される通り、混合音のうち周波数帯域 BH の成分は、強度スペクトル $Y_1(m)$ および強度スペクトル $Y_2(m)$ には含まれない。

10

【0021】

以上の通り、第1実施形態においては、音響信号 S_x が表す混合音のうち周波数帯域 BH の成分が音源分離の対象から除外される。したがって、周波数帯域 BL および周波数帯域 BH の双方を含む全帯域 BF を対象として混合音の音源分離を実行する構成と比較して、音源分離部22による処理負荷が軽減される。

【0022】

図2の帯域拡張部23は、混合音の強度スペクトル $X(m)$ と第1成分の強度スペクトル $Y_1(m)$ と第2成分の強度スペクトル $Y_2(m)$ とを利用して出力データ $O(m)$ を生成する。出力データ $O(m)$ は、第1出力データ $O_1(m)$ と第2出力データ $O_2(m)$ とで構成される。第1出力データ $O_1(m)$ は、強度スペクトル $Z_1(m)$ を表すデータであり、第2出力データ $O_2(m)$ は、強度スペクトル $Z_2(m)$ を表すデータである。

20

【0023】

第1出力データ $O_1(m)$ が表す強度スペクトル $Z_1(m)$ は、図3に例示される通り、周波数帯域 BL と周波数帯域 BH とを含む全帯域 BF にわたる第1音のスペクトルである。すなわち、音源分離において周波数帯域 BL に制限された第1音の強度スペクトル $Y_1(m)$ が、帯域拡張部23の処理により、全帯域 BF にわたる強度スペクトル $Z_1(m)$ に変換される。他方、第2出力データ $O_2(m)$ が表す強度スペクトル $Z_2(m)$ は、全帯域 BF にわたる第2音のスペクトルである。すなわち、音源分離において周波数帯域 BL に制限された第2音の強度スペクトル $Y_2(m)$ が、帯域拡張部23の処理により、全帯域 BF にわたる強度スペクトル $Z_2(m)$ に変換される。以上の説明から理解される通り、帯域拡張部23は、第1音および第2音の各々の周波数帯域を、周波数帯域 BL から全帯域 BF （周波数帯域 BL および周波数帯域 BH ）に拡張する。

30

【0024】

図2に例示される通り、帯域拡張部23は、取得部231と生成部232とを具備する。取得部231は、単位期間毎に入力データ $D(m)$ を生成する。入力データ $D(m)$ は、混合音の強度スペクトル $X(m)$ と第1成分の強度スペクトル $Y_1(m)$ と第2成分の強度スペクトル $Y_2(m)$ とに応じたベクトルを表すデータである。

【0025】

図3に例示される通り、入力データ $D(m)$ は、混合音データ $D_x(m)$ と第1入力データ $D_1(m)$ と第2入力データ $D_2(m)$ とを含む。混合音データ $D_x(m)$ は、混合音の強度スペクトル $X(m)$ を表すデータである。具体的には、任意の1個の単位期間（以下「目標期間」という）について生成される混合音データ $D_x(m)$ は、当該目標期間の強度スペクトル $X(m)$ と、目標期間の周囲に位置する他の単位期間の強度スペクトル $X(X(m-4), X(m-2), X(m+2), X(m+4))$ とを含む。具体的には、混合音データ $D_x(m)$ は、目標期間の強度スペクトル $X(m)$ と、目標期間の2個前の単位期間の強度スペクトル $X(m-2)$ と、目標期間の4個前の単位期間の強度スペクトル $X(m-4)$ と、目標期間の2個後の単位期間の強度スペクトル $X(m+2)$ と、目標期間の4個後の単位期間の強度スペクトル $X(m+4)$ とを含む。

40

【0026】

第1入力データ $D_1(m)$ は、第1音の強度スペクトル $Y_1(m)$ を表すデータである。具体

50

的には、任意の1個の目標期間について生成される第1入力データ $D_1(m)$ は、当該目標期間の強度スペクトル $Y_1(m)$ と、目標期間の周囲に位置する他の単位期間の強度スペクトル $Y_1(Y_1(m-4), Y_1(m-2), Y_1(m+2), Y_1(m+4))$ とを含む。具体的には、第1入力データ $D_1(m)$ は、目標期間の強度スペクトル $Y_1(m)$ と、目標期間の2個前の単位期間の強度スペクトル $Y_1(m-2)$ と、目標期間の4個前の単位期間の強度スペクトル $Y_1(m-4)$ と、目標期間の2個後の単位期間の強度スペクトル $Y_1(m+2)$ と、目標期間の4個後の単位期間の強度スペクトル $Y_1(m+4)$ とを含む。以上の説明から理解される通り、第1入力データ $D_1(m)$ は、第1音のうち周波数帯域 BL 内の第1成分を表すデータである。

【0027】

第2入力データ $D_2(m)$ は、第2音の強度スペクトル $Y_2(m)$ を表すデータである。具体的には、任意の1個の目標期間について生成される第2入力データ $D_2(m)$ は、当該目標期間の強度スペクトル $Y_2(m)$ と、目標期間の周囲に位置する他の単位期間の強度スペクトル $Y_2(Y_2(m-4), Y_2(m-2), Y_2(m+2), Y_2(m+4))$ とを含む。具体的には、第2入力データ $D_2(m)$ は、目標期間の強度スペクトル $Y_2(m)$ と、目標期間の2個前の単位期間の強度スペクトル $Y_2(m-2)$ と、目標期間の4個前の単位期間の強度スペクトル $Y_2(m-4)$ と、目標期間の2個後の単位期間の強度スペクトル $Y_2(m+2)$ と、目標期間の4個後の単位期間の強度スペクトル $Y_2(m+4)$ とを含む。以上の説明から理解される通り、第2入力データ $D_2(m)$ は、第2音のうち周波数帯域 BL 内の第2成分を表すデータである。

【0028】

入力データ $D(m)$ の全体で表現されるベクトル V の各要素は、当該ベクトル V の大きさが1(すなわち単位ベクトル)となるように正規化される。例えば、正規化前の入力データ $D(m)$ において、第1入力データ $D_1(m)$ と第2入力データ $D_2(m)$ と混合音データ $D_x(m)$ とにより、 N 個の要素 $e_1 \sim e_N$ が配列された N 次元のベクトル V が構成されると想定する。正規化後の入力データ $D(m)$ を構成する N 個の要素 $E_1 \sim E_N$ の各々は、以下の数式(1)で表現される($n = 1 \sim N$)。

【数1】

$$E_n = \frac{e_n}{\|V\|_2} \quad (1)$$

【0029】

数式(1)の記号 $\| \cdot \|_2$ は、以下の数式(2)で表現される L_2 ノルムを意味し、ベクトル V の大きさを表す指標(以下「強度指標」という)に相当する。

【数2】

$$\|V\|_2 = \alpha = \left(\sum_{n=1}^N |e_n|^2 \right)^{\frac{1}{2}} \quad (2)$$

【0030】

図2の生成部232は、入力データ $D(m)$ から出力データ $O(m)$ を生成する。出力データ $O(m)$ は、単位期間毎に順次に生成される。具体的には、生成部232は、各単位期間の

入力データ $D(m)$ から当該単位期間の出力データ $O(m)$ を生成する。出力データ $O(m)$ の生成には推定モデル M が利用される。推定モデル M は、入力データ $D(m)$ を入力として出力データ $O(m)$ を出力する統計的モデルである。すなわち、推定モデル M は、入力データ $D(m)$ と出力データ $O(m)$ との関係を学習した学習済モデルである。

【 0 0 3 1 】

推定モデル M は、例えばニューラルネットワークで構成される。図 4 は、推定モデル M の構造を例示するブロック図である。推定モデル M は、例えば、入力層 L_{in} と出力層 L_{out} との間の隠れ層 L_h に 4 層の全結合層 L_a を含む深層ニューラルネットワークである。活性化関数は、例えば ReLU (Rectified Linear Unit) である。入力データ $D(m)$ は、隠れ層 L_h の第 1 層において出力層 L_{out} と同等の次元数に圧縮される。なお、推定モデル M の構造は以上の例示に限定されない。例えば、再帰型ニューラルネットワーク (RNN: Recurrent Neural Network)、または畳込ニューラルネットワーク (CNN: Convolutional Neural Network) 等の任意の形式のニューラルネットワークが推定モデル M として利用される。複数種のニューラルネットワークの組合せが推定モデル M として利用されてもよい。また、長短期記憶 (LSTM: Long Short-Term Memory) 等の付加的な要素が推定モデル M に搭載されてもよい。

10

【 0 0 3 2 】

推定モデル M は、入力データ $D(m)$ から出力データ $O(m)$ を生成する演算を制御装置 1 1 に実行させる推定プログラムと、当該演算に適用される複数の変数 K (具体的には加重値およびバイアス) との組合せで実現される。推定プログラムと複数の変数 K とは記憶装置 1 2 に記憶される。複数の変数 K の各々の数値は、機械学習により事前に設定される。

20

【 0 0 3 3 】

図 2 の波形合成部 2 4 は、帯域拡張部 2 3 が順次に生成する出力データ $O(m)$ の時系列から音響信号 Sz_0 を生成する。具体的には、波形合成部 2 4 は、第 1 出力データ $O_1(m)$ および第 2 出力データ $O_2(m)$ の何れかの時系列から音響信号 Sz_0 を生成する。例えば、第 1 音の強調が利用者から指示された場合、波形合成部 2 4 は、第 1 出力データ $O_1(m)$ (強度スペクトル $Z_1(m)$) の時系列から音響信号 Sz_0 を生成する。すなわち、第 1 音が強調された音響信号 Sz_0 が生成される。他方、第 2 音の強調が利用者から指示された場合、波形合成部 2 4 は、第 2 出力データ $O_2(m)$ (強度スペクトル $Z_2(m)$) の時系列から音響信号 Sz_0 を生成する。すなわち、第 2 音が強調された音響信号 Sz_0 が生成される。音響信号 Sz_0 の生成には、例えば短時間逆フーリエ変換が利用される。

30

【 0 0 3 4 】

前述の通り、入力データ $D(m)$ を構成する各要素 E_n は、強度指標 を利用して正規化された数値である。したがって、音響信号 Sz_0 の音量は、音響信号 S_x とは相違する可能性がある。音量調整部 2 5 は、音響信号 Sz_0 の音量を音響信号 S_x と同等の音量に調整すること (すなわちスケールング) で音響信号 Sz を生成する。音響信号 Sz が放音装置 1 3 に供給されることで音波として放射される。具体的には、音量調整部 2 5 は、音響信号 S_x の音量と音響信号 Sz_0 の音量との相違に応じた調整値 G を音響信号 Sz_0 に乗算することで音響信号 Sz を生成する。調整値 G は、音響信号 S_x と音響信号 Sz との音量差が最小化されるように設定される。

40

【 0 0 3 5 】

図 5 は、制御装置 1 1 が音響信号 S_x から音響信号 Sz を生成する処理 (以下「音響処理 S_a 」という) の具体的な手順を例示するフローチャートである。例えば音響処理システム 1 0 0 に対する利用者からの指示を契機として音響処理 S_a が開始される。

【 0 0 3 6 】

音響処理 S_a が開始されると、制御装置 1 1 (周波数解析部 2 1) は、複数の単位期間の各々について音響信号 S_x の強度スペクトル $X(m)$ を生成する (S_{a1})。制御装置 1 1 (音源分離部 2 2) は、強度スペクトル $X(m)$ のうち周波数帯域 B_L 内の成分に対する音源分離により各単位期間の強度スペクトル $Y_1(m)$ と強度スペクトル $Y_2(m)$ とを生成する (S_{a2})。

50

【 0 0 3 7 】

制御装置 1 1 (取得部 2 3 1) は、強度スペクトル $X(m)$ と強度スペクトル $Y_1(m)$ と強度スペクトル $Y_2(m)$ とから各単位期間の入力データ $D(m)$ を生成する (Sa3)。制御装置 1 1 (生成部 2 3 2) は、入力データ $D(m)$ を推定モデル M に入力することで各単位期間の出力データ $O(m)$ を生成する (Sa4)。制御装置 1 1 (波形合成部 2 4) は、第 1 出力データ $O_1(m)$ または第 2 出力データ $O_2(m)$ の時系列から音響信号 Sz_0 を生成する (Sa5)。制御装置 1 1 (音量調整部 2 5) は、音響信号 Sz_0 に調整値 G を乗算することで音響信号 Sz を生成する (Sa6)。

【 0 0 3 8 】

以上に説明した通り、第 1 実施形態においては、周波数帯域 BL の成分を表す第 1 入力データ $D_1(m)$ および第 2 入力データ $D_2(m)$ を含む入力データ $D(m)$ から、周波数帯域 BL を含む全帯域 BF の音を表す出力データ $O(m)$ が生成される。すなわち、音響信号 Sx が表す混合音のうち周波数帯域 BL についてのみ限定的に音源分離を実行する構成にも関わらず、全帯域 BF の成分を含む出力データ $O(m)$ が生成される。したがって、音源分離のための処理負荷を軽減できる。

【 0 0 3 9 】

[2] 学習処理部 3 0

図 2 に例示される通り、制御装置 1 1 は、記憶装置 1 2 に記憶された機械学習プログラム P_2 を実行することで学習処理部 3 0 として機能する。学習処理部 3 0 は、音響処理 Sa に利用される推定モデル M を機械学習により確立する。学習処理部 3 0 は、取得部 3 1 と訓練部 3 2 とを具備する。

【 0 0 4 0 】

記憶装置 1 2 には、推定モデル M の機械学習に利用される複数の訓練データ T が記憶される。図 6 は、訓練データ T の説明図である。複数の訓練データ T の各々は、訓練用の入力データ $D_t(m)$ と訓練用の出力データ $O_t(m)$ との組合せで構成される。図 3 の入力データ $D(m)$ と同様に、訓練用の入力データ $D_t(m)$ は、混合音データ $D_x(m)$ と第 1 入力データ $D_1(m)$ と第 2 入力データ $D_2(m)$ とを含む。

【 0 0 4 1 】

図 6 には、参照信号 S_r と第 1 信号 S_{r1} と第 2 信号 S_{r2} とが図示されている。参照信号 S_r は、第 1 音源から発音される第 1 音と第 2 音源から発音される第 2 音との混合音を表す時間領域の信号である。参照信号 S_r が表す混合音は、周波数帯域 BL と周波数帯域 BH とを含む全帯域 BF にわたる。参照信号 S_r は、例えば、第 1 音源と第 2 音源とが並列に発音する環境において収録装置を利用して収録される。また、第 1 信号 S_{r1} は、第 1 音を表す時間領域の信号であり、第 2 信号 S_{r2} は、第 2 音を表す時間領域の信号である。第 1 音および第 2 音の各々は、周波数帯域 BL と周波数帯域 BH とを含む全帯域 BF にわたる。第 1 信号 S_{r1} は、第 1 音源のみが発音する環境において収録され、第 2 信号 S_{r2} は、第 2 音源のみが発音する環境において収録される。なお、相互に個別に収録された第 1 信号 S_{r1} と第 2 信号 S_{r2} とを混合することで参照信号 S_r が生成されてもよい。

【 0 0 4 2 】

図 6 には、参照信号 S_r の強度スペクトル $X(m)$ の時系列 (... , $X(m-1)$, $X(m)$, $X(m+1)$, ...) と、第 1 信号 S_{r1} の強度スペクトル $R_1(m)$ の時系列 (... , $R_1(m-1)$, $R_1(m)$, $R_1(m+1)$, ...) と、第 2 信号 S_{r2} の強度スペクトル $R_2(m)$ の時系列 (... , $R_2(m-1)$, $R_2(m)$, $R_2(m+1)$, ...) とが図示されている。訓練用の入力データ $D_t(m)$ のうちの混合音データ $D_x(m)$ は、参照信号 S_r の強度スペクトル $X(m)$ から生成される。具体的には、任意の 1 個の目標期間の混合音データ $D_x(m)$ は、図 3 の例示と同様に、当該目標期間の強度スペクトル $X(m)$ と、目標期間の周囲に位置する他の単位期間の強度スペクトル $X(X(m-4)$, $X(m-2)$, $X(m+2)$, $X(m+4)$) とを含む。

【 0 0 4 3 】

第 1 信号 S_{r1} は、周波数帯域 BL の成分と周波数帯域 BH の成分とを含む。第 1 信号 S_{r1} の強度スペクトル $R_1(m)$ は、周波数帯域 BL 内の強度スペクトル $Y_1(m)$ と周波数帯域 B

10

20

30

40

50

H内の強度スペクトル $H_1(m)$ とで構成される。訓練用の入力データ $D_t(m)$ の第1入力データ $D_1(m)$ は、周波数帯域 BL の強度スペクトル $Y_1(m)$ を表すデータである。具体的には、目標期間の第1入力データ $D_1(m)$ は、当該目標期間の強度スペクトル $Y_1(m)$ と、当該目標期間の周囲に位置する他の単位期間の強度スペクトル $Y_1(Y_1(m-4), Y_1(m-2), Y_1(m+2), Y_1(m+4))$ とを含む。

【0044】

第1信号 S_{r1} と同様に、第2信号 S_{r2} は、周波数帯域 BL の成分と周波数帯域 BH の成分とを含む。第2信号 S_{r2} の強度スペクトル $R_2(m)$ は、周波数帯域 BL 内の強度スペクトル $Y_2(m)$ と周波数帯域 BH 内の強度スペクトル $H_2(m)$ とで構成される。訓練用の入力データ $D_t(m)$ の第2入力データ $D_{t2}(m)$ は、周波数帯域 BL の強度スペクトル $Y_2(m)$ を表すデータである。具体的には、目標期間の第2入力データ $D_{t2}(m)$ は、当該目標期間の強度スペクトル $Y_2(m)$ と、目標期間の周囲に位置する他の単位期間の強度スペクトル $Y_2(Y_2(m-4), Y_2(m-2), Y_2(m+2), Y_2(m+4))$ とを含む。

10

【0045】

他方、各訓練データ T を構成する訓練用の出力データ $O_t(m)$ は、第1出力データ $O_{t1}(m)$ と第2出力データ $O_{t2}(m)$ とで構成される正解データである。第1出力データ $O_{t1}(m)$ は、第1信号 S_{r1} の強度スペクトル $R_1(m)$ を表す。すなわち、第1出力データ $O_{t1}(m)$ は、参照信号 S_r が表す混合音のうち全帯域 BF にわたる第1音のスペクトルである。第2出力データ $O_{t2}(m)$ は、第2信号 S_{r2} の強度スペクトル $R_2(m)$ を表す。すなわち、第2出力データ $O_{t2}(m)$ は、参照信号 S_r が表す混合音のうち全帯域 BF にわたる第2音のスペクトルである。

20

【0046】

訓練用の入力データ $D_t(m)$ の全体で表現されるベクトル V の各要素は、前述の入力データ $D_t(m)$ と同様に、当該ベクトル V の大きさが1となるように正規化される。同様に、訓練用の出力データ $O_t(m)$ の全体で表現されるベクトル V の各要素は、当該ベクトル V の大きさが1となるように正規化される。

【0047】

図2の取得部31は、複数の訓練データ T の各々を記憶装置12から取得する。なお、参照信号 S_r と第1信号 S_{r1} と第2信号 S_{r2} とが記憶装置12に記憶された構成においては、取得部31が参照信号 S_r と第1信号 S_{r1} と第2信号 S_{r2} とから複数の訓練データ T を生成する。すなわち、取得部31による「取得」は、事前に用意された訓練データ T を記憶装置12から読出する処理のほか、当該取得部31自身が訓練データ T を生成する処理も包含する。

30

【0048】

訓練部32は、複数の訓練データ T を利用した処理(以下「学習処理 S_b 」という)により推定モデル M を確立する。学習処理 S_b は、複数の訓練データ T を利用した教師あり機械学習である。具体的には、訓練部32は、各訓練データ T の入力データ $D_t(m)$ を入力した場合に暫定的な推定モデル M が生成する出力データ $O(m)$ と、当該訓練データ T に含まれる出力データ $O_t(m)$ との誤差を表す損失関数 L が低減(理想的には最小化)されるように、推定モデル M を規定する複数の変数 K を反復的に更新する。したがって、推定モデル M は、複数の訓練データ T における入力データ $D_t(m)$ と出力データ $O_t(m)$ との間に潜在する関係を学習する。すなわち、訓練部32による訓練後の推定モデル M は、未知の入力データ $D(m)$ に対して当該関係のもとで統計的に妥当な出力データ $O(m)$ を出力する。

40

【0049】

損失関数 L は、例えば以下の数式(3)で表現される。

【数3】

$$L = \varepsilon[O_1(m), O_{t1}(m)] + \varepsilon[O_2(m), O_{t2}(m)] \dots (3)$$

50

数式(3)の記号 $[a, b]$ は、要素 a と要素 b との誤差（例えば平均二乗誤差またはクロスエントロピー関数）である。

【0050】

図7は、学習処理 S_b の具体的な手順を例示するフローチャートである。例えば音響処理システム100に対する利用者からの指示を契機として学習処理 S_b が開始される。

【0051】

制御装置11（取得部31）は、訓練データ T を記憶装置12から取得する（ S_{b1} ）。制御装置11（訓練部32）は、当該訓練データ T を利用した機械学習を実行する（ S_{b2} ）。すなわち、訓練データ T の入力データ $D_t(m)$ から推定モデル M が生成する出力データ $O(m)$ と、当該訓練データ T の出力データ $O_t(m)$ （すなわち正解値）との間の損失関数 L が低減されるように、推定モデル M の複数の変数 K を反復的に更新する。損失関数 L に応じた複数の変数 K の更新には、例えば誤差逆伝播法が利用される。

10

【0052】

制御装置11は、学習処理 S_b に関する終了条件が成立したか否かを判定する（ S_{b3} ）。終了条件は、例えば、損失関数 L が所定の閾値を下回ること、または、損失関数 L の変化量が所定の閾値を下回ることである。終了条件が成立しない場合（ $S_{b3} : NO$ ）、制御装置11（取得部31）は、未取得の訓練データ T を記憶装置12から取得する（ S_{b1} ）。すなわち、終了条件の成立まで、訓練データ T の取得（ S_{b1} ）と当該訓練データ T を利用した複数の変数 K の更新（ S_{b2} ）とが反復される。終了条件が成立した場合（ $S_{b3} : YES$ ）、制御装置11は学習処理 S_b を終了する。

20

【0053】

以上に説明した通り、第1実施形態においては、周波数帯域 BL の成分を表す第1入力データ $D_1(m)$ および第2入力データ $D_2(m)$ を含む入力データ $D(m)$ から、周波数帯域 BL および周波数帯域 BH の音を表す出力データ $O(m)$ が生成されるように、推定モデル M が確立される。すなわち、音響信号 S_x が表す混合音のうち周波数帯域 BL についてのみ限定的に音源分離を実行する構成でも、推定モデル M を利用することで、周波数帯域 BH の成分を含む出力データ $O(m)$ が生成される。したがって、音源分離のための処理負荷を軽減できる。

【0054】

B：第2実施形態

第2実施形態について以下に説明する。なお、以下に例示する各形態において機能が第1実施形態と同様である要素については、第1実施形態の説明で使用した符号を流用して各々の詳細な説明を適宜に省略する。

30

【0055】

第1実施形態においては、混合音データ $D_x(m)$ が周波数帯域 BL の成分と周波数帯域 BH の成分とを双方を含む構成を例示した。しかし、第1音のうち周波数帯域 BL 内の成分は第1入力データ $D_1(m)$ に含まれ、第2音のうち周波数帯域 BH 内の成分は第2入力データ $D_2(m)$ に含まれるから、混合音データ $D_x(m)$ が周波数帯域 BL の成分を含む構成は必須ではない。以上の事情を考慮して、第2実施形態においては、混合音データ $D_x(m)$ が混合音のうち周波数帯域 BL の成分を含まない。

40

【0056】

図8は、第2実施形態における入力データ $D(m)$ の模式図である。音響信号 S_x の強度スペクトル $X(m)$ は、周波数帯域 BL 内の強度スペクトル $X_L(m)$ と周波数帯域 BH 内の強度スペクトル $X_H(m)$ とに分割される。入力データ $D(m)$ の混合音データ $D_x(m)$ は、周波数帯域 BH の強度スペクトル $X_H(m)$ を表すデータである。具体的には、1個の目標期間について生成される混合音データ $D_x(m)$ は、当該目標期間の強度スペクトル $X_H(m)$ と、当該目標期間の周囲に位置する他の単位期間の強度スペクトル $X_H(X_H(m-4), X_H(m-2), X_H(m+2), X_H(m+4))$ とを含む。すなわち、第2実施形態の混合音データ $D_x(m)$ は、混合音のうち周波数帯域 BL の成分（強度スペクトル $X_L(m)$ ）を含まない。なお、音源分離部22が強度スペクトル $X(m)$ のうち周波数帯域 BL の成分を対象として音源分離を実行する点

50

は第1実施形態と同様である。

【0057】

以上の説明においては、音響処理 S_a に利用される入力データ $D(m)$ を例示したが、学習処理 S_b に利用される訓練用の入力データ $D_t(m)$ についても同様に、参照信号 S_r が表す混合音のうち周波数帯域 B_H の成分を表す混合音データ $D_x(m)$ が含まれる。すなわち、訓練用の混合音データ $D_x(m)$ は、参照信号 S_r の強度スペクトル $X(m)$ のうち周波数帯域 B_H 内の強度スペクトル $X_H(m)$ を表し、周波数帯域 B_L 内の強度スペクトル $X_L(m)$ は混合音データ $D_x(m)$ に反映されない。

【0058】

第2実施形態においても第1実施形態と同様の効果が実現される。また、第2実施形態においては、混合音データ $D_x(m)$ が混合音のうち周波数帯域 B_L の成分を含まない。したがって、混合音データ $D_x(m)$ が全帯域 B_F の成分を含む構成と比較して、学習処理 S_b の処理負荷および推定モデル M の規模が低減されるという利点がある。

10

【0059】

第1実施形態においては、全帯域 B_F にわたる混合音を表す混合音データ $D_x(m)$ を例示した。第2実施形態においては、混合音のうち周波数帯域 B_H の成分を表す混合音データ $D_x(m)$ を例示した。以上の例示から理解される通り、混合音データ $D_x(m)$ は、混合音のうち周波数帯域 B_H を含む周波数帯域の成分を表すデータとして包括的に表現される。

【0060】

C：第3実施形態

20

図9は、第3実施形態における入力データ $D(m)$ の模式図である。第3実施形態の入力データ $D(m)$ は、混合音データ $D_x(m)$ と第1入力データ $D_1(m)$ と第2入力データ $D_2(m)$ とに加えて強度指標を含む。強度指標は、前述の通り、入力データ $D(m)$ の全体で表現されるベクトル V の大きさ（例えば L_2 ノルム）を表す指標であり、前掲の数式(2)で算定される。学習処理 S_b に利用される訓練用の入力データ $D_t(m)$ についても同様に、混合音データ $D_x(m)$ と第1入力データ $D_1(m)$ と第2入力データ $D_2(m)$ とに加えて、当該入力データ $D_t(m)$ で表現されるベクトル V の大きさに応じた強度指標が含まれる。なお、混合音データ $D_x(m)$ と第1入力データ $D_1(m)$ と第2入力データ $D_2(m)$ とは、第1実施形態または第2実施形態と同様である。

【0061】

30

図10は、第3実施形態に係る音響処理システム100の機能的な構成を例示するブロック図である。第3実施形態の入力データ $D(m)$ には強度指標が含まれるから、当該強度指標が反映された出力データ $O(t)$ が推定モデル M から出力される。具体的には、波形合成部24が出力データ $O(t)$ から生成する音響信号 S_z は、音響信号 S_x と同等の音量となる。したがって、第1実施形態において例示した音量調整部25（図5のステップ S_a6 ）が第3実施形態においては省略される。すなわち、波形合成部24による出力信号（第1実施形態における音響信号 S_z0 ）が最終的な音響信号 S_z として出力される。

【0062】

第3実施形態においても第1実施形態と同様の効果が実現される。また、第3実施形態においては、強度指標が入力データ $D(m)$ に含まれるから、混合音に対応する音量の音を表す出力データ $O(m)$ が生成される。したがって、第1出力データ $O_1(m)$ および第2出力データ $O_2(m)$ が表す音の強度を調整する処理（音量調整部25）が不要であるという利点がある。

40

【0063】

図11は、第1実施形態および第3実施形態による効果の説明図である。図11の結果Aは、第1実施形態により生成された音響信号 S_z の振幅スペクトログラムであり、図11の結果Bは、第3実施形態により生成された音響信号 S_z の振幅スペクトログラムである。結果Aおよび結果Bにおいては、打楽器音（第1音）と歌唱音声（第2音）との混合音を表す音響信号 S_x に対して音響処理 S_a を実行することで、打楽器音を表す音響信号 S_z を生成した場合が想定されている。図11の正解Cは、単独で発音された打楽器音の振幅ス

50

ペクトログラムである。

【 0 0 6 4 】

図 1 1 の結果 A からは、第 1 実施形態により、正解 C に近い音響信号 S_z を生成できることが確認できる。また、図 1 1 の結果 B からは、入力データ $D(m)$ が強度指標 を含む第 3 実施形態により、第 1 実施形態と比較しても正解 C に十分に近い音響信号 S_z を生成できることが確認される。

【 0 0 6 5 】

図 1 2 は、第 1 実施形態から第 3 実施形態に関する観測結果の図表である。図 1 2 においては、打楽器音（第 1 音）と歌唱音声（第 2 音）との混合音を表す音響信号 S_x に対して音響処理 S_a を実行することで、打楽器音（Drums）を表す音響信号 S_z と、歌唱音声（Vocals）を表す音響信号 S_z とを生成した場合が想定されている。図 1 2 には、評価指標として有効な SAR （信号対非線形歪比：Sources to Artifacts Ratio）および SAR 改善量が、第 1 実施形態から第 3 実施形態の各々について図示されている。 SAR 改善量は、比較例を基準とした SAR の改善量である。比較例については、音響信号 S_z のうち周波数帯域 BH の成分を一律にゼロとした場合の SAR が基準として例示されている。

【 0 0 6 6 】

第 1 実施形態および第 2 実施形態においても SAR が改善することが図 1 2 から確認できる。また、第 3 実施形態によれば、打楽器音および歌唱音声の何れについても、第 1 実施形態および第 2 実施形態と比較して非常に高精度な音源分離が実現されることが図 1 2 から確認できる。

【 0 0 6 7 】

D：第 4 実施形態

第 4 実施形態の学習処理 S_b においては、前掲の数式(3)で表現される損失関数 L が、以下の数式(4)で表現される損失関数 L に置換される。

【数 4】

$$L = \varepsilon[O_1(m), O_{t1}(m)] + \varepsilon[O_2(m), O_{t2}(m)] \\ + \varepsilon[X_H(m), O_{1H}(m) + O_{2H}(m)] \quad (4)$$

【 0 0 6 8 】

数式(4)における記号 $O_{1H}(m)$ は、第 1 出力データ $O_1(m)$ が表す強度スペクトル $Z_1(m)$ のうち周波数帯域 BH 内の強度スペクトルであり、記号 $O_{2H}(m)$ は、第 2 出力データ $O_2(m)$ が表す強度スペクトル $Z_2(m)$ のうち周波数帯域 BH 内の強度スペクトルである。すなわち、数式(4)の右辺における第 3 項は、参照信号 S_r の強度スペクトル $X(m)$ のうち周波数帯域 BH 内の強度スペクトル $X_H(m)$ と、強度スペクトル $H_1(m)$ および強度スペクトル $H_2(m)$ の合計 ($H_1(m) + H_2(m)$) との誤差を意味する。以上の説明から理解される通り、第 4 実施形態の訓練部 3 2 は、強度スペクトル $Z_1(m)$ のうち周波数帯域 BH 内の成分と、強度スペクトル $Z_2(m)$ のうち周波数帯域 BH 内の成分とを混合した結果が、混合音の強度スペクトル $X(m)$ のうち周波数帯域 BH の成分（強度スペクトル $X_H(m)$ ）に近似または一致するという条件（以下「追加条件」という）のもとで、推定モデル M を訓練する。

【 0 0 6 9 】

第 4 実施形態においても第 1 実施形態と同様の効果を実現される。また、第 4 実施形態によれば、追加条件なしで訓練された推定モデル M を利用する構成と比較して、第 1 音のうち周波数帯域 BH の成分（第 1 出力データ $O_1(m)$ ）と第 2 音のうち周波数帯域 BH の成分（第 2 出力データ $O_2(m)$ ）とを高精度に推定できる。なお、第 4 実施形態の構成は、第 2 実施形態および第 3 実施形態にも同様に適用される。

【 0 0 7 0 】

E : 第 5 実施形態

図 1 3 は、第 5 実施形態における入力データ $D(m)$ および出力データ $O(m)$ の模式図である。第 1 実施形態の出力データ $O(m)$ における第 1 出力データ $O_1(m)$ は、全帯域 BF にわたる強度スペクトル $Z_1(m)$ を表し、第 2 出力データ $O_2(m)$ は、全帯域 BF にわたる強度スペクトル $Z_2(m)$ を表す。第 5 実施形態における第 1 出力データ $O_1(m)$ は、第 1 音のうち周波数帯域 BH の成分を表す。すなわち、第 1 出力データ $O_1(m)$ は、第 1 音の強度スペクトル $Z_1(m)$ のうち周波数帯域 BH 内の強度スペクトル $H_1(m)$ を表し、周波数帯域 BL 内の強度スペクトルを含まない。同様に、第 5 実施形態における第 2 出力データ $O_2(m)$ は、第 2 音のうち周波数帯域 BH の成分を表す。すなわち、第 2 出力データ $O_2(m)$ は、第 2 音の強度スペクトル $Z_2(m)$ のうち周波数帯域 BH 内の強度スペクトル $H_2(m)$ を表し、周波数帯域 BL 内の強度スペクトルを含まない。

10

【 0 0 7 1 】

図 1 4 は、第 5 実施形態における訓練用の入力データ $D_t(m)$ および出力データ $O_t(m)$ の模式図である。第 1 実施形態において、訓練用の出力データ $O_t(m)$ における第 1 出力データ $O_{t1}(m)$ は、全帯域 BF にわたる第 1 音の強度スペクトル $R_1(m)$ を表し、第 2 出力データ $O_{t2}(m)$ は、全帯域 BF にわたる第 2 音の強度スペクトル $R_2(m)$ を表す。第 5 実施形態における第 1 出力データ $O_{t1}(m)$ は、第 1 音のうち周波数帯域 BH の成分を表す。すなわち、第 1 出力データ $O_{t1}(m)$ は、第 1 音の強度スペクトル $R_1(m)$ のうち周波数帯域 BH 内の強度スペクトル $H_1(m)$ を表し、周波数帯域 BL 内の強度スペクトル $Y_1(m)$ を含まない。同様に、第 5 実施形態における第 2 出力データ $O_{t2}(m)$ は、第 2 音のうち周波数帯域 BH の成分を表す。すなわち、第 2 出力データ $O_{t2}(m)$ は、第 2 音の強度スペクトル $R_2(m)$ のうち周波数帯域 BH 内の強度スペクトル $H_2(m)$ を表し、周波数帯域 BL 内の強度スペクトル $Y_2(m)$ を含まない。

20

【 0 0 7 2 】

図 1 5 は、第 5 実施形態における音響処理部 2 0 の部分的な構成を例示するブロック図である。第 5 実施形態の波形合成部 2 4 には、第 1 音のうち周波数帯域 BH 内の強度スペクトル $H_1(m)$ を表す第 1 出力データ $O_1(m)$ が音響処理部 2 0 から供給されるほか、第 1 音のうち周波数帯域 BL 内の強度スペクトル $Y_1(m)$ が音源分離部 2 2 から供給される。第 1 音の強調が利用者から指示された場合、波形合成部 2 4 は、強度スペクトル $H_1(m)$ と強度スペクトル $Y_1(m)$ とを合成することで全帯域 BF にわたる強度スペクトル $Z_1(m)$ を生成し、強度スペクトル $Z_1(m)$ の時系列から音響信号 S_{z0} を生成する。

30

【 0 0 7 3 】

また、第 5 実施形態の波形合成部 2 4 には、第 2 音のうち周波数帯域 BH 内の強度スペクトル $H_2(m)$ を表す第 2 出力データ $O_2(m)$ が音響処理部 2 0 から供給されるほか、第 2 音のうち周波数帯域 BL 内の強度スペクトル $Y_2(m)$ が音源分離部 2 2 から供給される。第 2 音の強調が利用者から指示された場合、波形合成部 2 4 は、強度スペクトル $H_2(m)$ と強度スペクトル $Y_2(m)$ とを合成することで全帯域 BF にわたる強度スペクトル $Z_2(m)$ を生成し、強度スペクトル $Z_2(m)$ の時系列から音響信号 S_{z0} を生成する。

40

【 0 0 7 4 】

第 5 実施形態においても第 1 実施形態と同様の効果が実現される。また、第 5 実施形態においては、出力データ $O(m)$ が周波数帯域 BL の成分を含まない。したがって、出力データ $O(m)$ が全帯域 BF の成分を含む構成（例えば第 1 実施形態）と比較して、学習処理 S_b の処理負荷および推定モデル M の規模が低減されるという利点がある。他方、出力データ $O(m)$ が全帯域 BF の成分を含む第 1 実施形態によれば、第 5 実施形態と比較して、全帯域 BF にわたる音響を簡便に生成できるという利点がある。

【 0 0 7 5 】

第 1 実施形態においては、第 1 音のうち周波数帯域 BL と周波数帯域 BH とを含む全帯域 BF の成分を表す第 1 出力データ $O_1(m)$ を例示した。第 5 実施形態においては、第 1 音のうち周波数帯域 BH の成分を表す第 1 出力データ $O_1(m)$ を例示した。以上の例示から理解

50

される通り、第1出力データ $O_1(m)$ は、第1音のうち周波数帯域 B_H を含む周波数帯域の成分を表すデータとして包括的に表現される。同様に、第2出力データ $O_2(m)$ は、第2音のうち周波数帯域 B_H を含む周波数帯域の成分を表すデータとして包括的に表現される。

【0076】

F：変形例

以上に例示した各態様に付加される具体的な変形の態様を以下に例示する。以下の例示から任意に選択された2以上の態様を、相互に矛盾しない範囲で適宜に併合してもよい。

【0077】

(1) 前述の各形態においては、目標期間の強度スペクトル $X(m)$ と他の単位期間の強度スペクトル X とを含む混合音データ $D_x(m)$ を例示したが、混合音データ $D_x(m)$ の内容は以上の例示に限定されない。例えば、目標期間の混合音データ $D_x(m)$ が当該目標期間の強度スペクトル $X(m)$ のみを含む構成が想定される。目標期間の混合音データ $D_x(m)$ が、当該目標期間に対して過去および未来の一方の単位期間の強度スペクトル X を含んでもよい。また、前述の各形態においては、目標期間の混合音データ $D_x(m)$ が、当該目標期間に間隔をあけて前後する他の単位期間の強度スペクトル $X(X(m-4), X(m-2), X(m+2), X(m+4))$ を含む構成を例示したが、目標期間の直前の単位期間の強度スペクトル $X(m-1)$ または直後の単位期間の強度スペクトル $X(m+1)$ を混合音データ $D_x(m)$ が含んでもよい。

【0078】

以上の説明においては混合音データ $D_x(m)$ に着目したが、第1入力データ $D_1(m)$ および第2入力データ $D_2(m)$ についても同様である。例えば、目標期間の第1入力データ $D_1(m)$ は、当該目標期間の強度スペクトル $Y_1(m)$ のみで構成されてもよいし、当該目標期間の過去および未来の一方の単位期間の強度スペクトル Y_1 を含んでもよい。また、目標期間の第1入力データ $D_1(m)$ が、当該目標期間の直前の単位期間の強度スペクトル $Y_2(m-1)$ 、または直後の単位期間の強度スペクトル $Y_1(m+1)$ を含んでもよい。第2入力データ $D_2(m)$ についても同様である。

【0079】

(2) 前述の各形態においては、所定の周波数を下回る周波数帯域 B_L と当該周波数を上回る周波数帯域 B_H とに着目したが、周波数帯域 B_L と周波数帯域 B_H との関係は以上の例示に限定されない。例えば、周波数帯域 B_L が所定の周波数を上回り、周波数帯域 B_H が当該周波数を下回る構成も想定される。また、周波数帯域 B_L および周波数帯域 B_H の各々は、周波数軸上で連続する周波数帯域に限定されない。例えば、周波数軸を区分した複数の周波数帯域のうち奇数番目および偶数番目の一方に属する2以上の周波数帯域の集合が周波数帯域 B_L とされ、奇数番目および偶数番目の他方に属する2以上の周波数帯域の集合が周波数帯域 B_H とされてもよい。

【0080】

(3) 前述の各形態においては、事前に用意された音響信号 S_x を処理する場合を例示したが、音響処理部20は、音響信号 S_x の収録に並行して実時間的に、音響信号 S_x に対する音響処理 S_a を実行してもよい。なお、前述の各形態における例示のように混合音データ $D_x(m)$ が目標期間の後方の強度スペクトル $X(m+4)$ を含む構成では、単位期間の4個分に相当する時間長の遅延が発生する。

【0081】

(4) 前述の各形態においては、第1音が強調された強度スペクトル $Z_1(m)$ を表す第1出力データ $O_1(m)$ と第2音が強調された強度スペクトル $Z_2(m)$ を表す第2出力データ $O_2(m)$ との双方を帯域拡張部23が生成したが、第1出力データ $O_1(m)$ および第2出力データ $O_2(m)$ の一方のみを出力データ $O(m)$ として帯域拡張部23が生成してもよい。例えば、歌唱音声(第1音)と楽器音(第2音)との混合音に対する音響処理 S_a で歌唱音声を抑制するという用途に使用される音響処理システム100においては、第2音が強調された強度スペクトル $Z_2(m)$ を表す出力データ $O(m)$ (第2出力データ $O_2(m)$)を帯域拡張部23が生成すれば充分である。すなわち、第1音が強調された強度スペクトル $Z_1(m)$ の生成は省略される。以上の説明から理解される通り、生成部232は、第1出力データ $O_1(m)$

10

20

30

40

50

および第 2 出力データ O2(m)の少なくとも一方を生成する要素として表現される。

【 0 0 8 2 】

(5) 前述の各形態においては、第 1 音および第 2 音の一方が強調された音響信号 Sz を生成したが、音響処理部 2 0 による処理の内容は以上の例示に限定されない。例えば、第 1 出力データ O1(m)の時系列から生成される第 1 音響信号と第 2 出力データ O2(m)の時系列から生成される第 2 音響信号との加重和を、音響処理部 2 0 が音響信号 Sz として出力してもよい。第 1 音響信号は第 1 音が強調された信号であり、第 2 音響信号は第 2 音が強調された信号である。また、第 1 音響信号および第 2 音響信号の各々に対して、例えば効果付与等の音響処理を相互に独立に実行し、処理後の第 1 音響信号と第 2 音響信号とを加算することで、音響処理部 2 0 が音響信号 Sz を生成してもよい。

10

【 0 0 8 3 】

(6) 携帯電話機またはスマートフォン等の端末装置との間で通信するサーバ装置により音響処理システム 1 0 0 が実現されてもよい。例えば、音響処理システム 1 0 0 は、端末装置から受信した音響信号 Sx に対する音響処理 Sa により音響信号 Sz を生成し、当該音響信号 Sz を端末装置に送信する。端末装置に搭載された周波数解析部 2 1 が生成した強度スペクトル X(m) を音響処理システム 1 0 0 が受信する構成においては、音響処理システム 1 0 0 から周波数解析部 2 1 が省略される。また、波形合成部 2 4 (および音量調整部 2 5) が端末装置に搭載された構成においては、帯域拡張部 2 3 が生成した出力データ O(m) が音響処理システム 1 0 0 から端末装置に送信される。したがって、波形合成部 2 4 および音量調整部 2 5 は音響処理システム 1 0 0 から省略される。

20

【 0 0 8 4 】

また、周波数解析部 2 1 および音源分離部 2 2 は端末装置に搭載されてもよい。音響処理システム 1 0 0 は、周波数解析部 2 1 が生成した強度スペクトル X(m) と、音源分離部 2 2 が生成した強度スペクトル Y1(m) および強度スペクトル Y2(m) とを、端末装置から受信する。以上の説明から理解される通り、音響処理システム 1 0 0 から音源分離部 2 2 が省略されてもよい。音響処理システム 1 0 0 が音源分離部 2 2 を具備しない構成でも、端末装置等の外部装置において実行される音源分離の処理負荷を軽減できる、という所期の効果は実現される。

【 0 0 8 5 】

(7) 前述の各形態においては、音響処理部 2 0 と学習処理部 3 0 とを具備する音響処理システム 1 0 0 を例示したが、音響処理部 2 0 および学習処理部 3 0 の一方が省略されてもよい。学習処理部 3 0 を具備するコンピュータシステムは、推定モデル訓練システム (機械学習システム) とも換言される。推定モデル訓練システムにおける音響処理部 2 0 の有無は不問である。

30

【 0 0 8 6 】

(8) 以上に例示した音響処理システム 1 0 0 の機能は、前述の通り、制御装置 1 1 を構成する単数または複数のプロセッサと、記憶装置 1 2 に記憶されたプログラム (P1 , P2) との協働により実現される。本開示に係るプログラムは、コンピュータが読取可能な記録媒体に格納された形態で提供されてコンピュータにインストールされ得る。記録媒体は、例えば非一過性 (non-transitory) の記録媒体であり、CD-ROM 等の光学式記録媒体 (光ディスク) が好例であるが、半導体記録媒体または磁気記録媒体等の公知の任意の形式の記録媒体も包含される。なお、非一過性の記録媒体とは、一過性の伝搬信号 (transitory, propagating signal) を除く任意の記録媒体を含み、揮発性の記録媒体も除外されない。また、配信装置が通信網を介してプログラムを配信する構成では、当該配信装置においてプログラムを記憶する記憶装置 1 2 が、前述の非一過性の記録媒体に相当する。

40

【 0 0 8 7 】

G : 付記

以上に例示した形態から、例えば以下の構成が把握される。

【 0 0 8 8 】

本開示のひとつの態様 (態様 1) に係る音響処理方法は、第 1 音源に対応する第 1 音の

50

うち第1周波数帯域の成分を表す第1入力データと、前記第1音源とは異なる第2音源に対応する第2音のうち前記第1周波数帯域の成分を表す第2入力データと、前記第1音と前記第2音との混合音のうち前記第1周波数帯域とは異なる第2周波数帯域を含む周波数帯域の成分を含む音を表す混合音データと、を含む入力データを取得し、学習済の推定モデルに前記入力データを入力することで、前記第1音のうち前記第2周波数帯域を含む周波数帯域の成分を表す第1出力データと、前記第2音のうち前記第2周波数帯域を含む周波数帯域の成分を表す第2出力データとの少なくとも一方を生成する。

【0089】

以上の構成によれば、第1音のうち第1周波数帯域の成分を表す第1入力データと、第2音のうち第1周波数帯域の成分を表す第2入力データとを含む入力データから、第1音のうち第2周波数帯域を含む周波数帯域の成分を表す第1出力データと、第2音のうち第2周波数帯域を含む周波数帯域の成分を表す第2出力データとの少なくとも一方が生成される。すなわち、第1入力データが表す音は第1音のうち第1周波数帯域の成分であれば足り、第2入力データが表す音は第2音のうち第1周波数帯域の成分であれば足りる。以上の構成によれば、第1音源に対応する第1音と第2音源に対応する第2音との混合音を第1音と第2音とに分離する音源分離を、第1周波数帯域についてのみ限定的に実行すれば足りる。したがって、音源分離のための処理負荷が軽減される。

【0090】

「第1音源に対応する第1音」は、第1音源から発音される音を優勢に含む音を意味する。すなわち、第1音源から発音される音単独のほか、例えば第1音源から発音される第1音に加えて第2音源からの第2音（例えば音源分離により完全には除去されなかった第2音）が僅かに含まれる音も、「第1音源に対応する第1音」の概念には包含される。同様に、「第2音源に対応する第2音」は、第2音源から発音される音を優勢に含む音を意味する。すなわち、第2音源から発音される音単独のほか、例えば第2音源から発音される第2音に加えて第1音源からの第1音（例えば音源分離により完全には除去されなかった第1音）が僅かに含まれる音も、「第2音源に対応する第2音」の概念には包含される。

【0091】

混合音データが表す音は、混合音のうち第1周波数帯域および第2周波数帯域の双方の成分を含む音（例えば全帯域にわたる混合音）と、混合音のうち第1周波数帯域の成分を含まない音とを包含する。

【0092】

第1周波数帯域および第2周波数帯域は、周波数軸上の相異なる周波数帯域である。典型的には、第1周波数帯域と第2周波数帯域とは相互に重複しない。ただし、第1周波数帯域と第2周波数帯域とが部分的に重複してもよい。第1周波数帯域の周波数軸上の位置と第2周波数帯域の周波数軸上の位置との関係は任意である。また、第1周波数帯域の帯域幅と第2周波数帯域の帯域幅との異同は不問である。

【0093】

第1出力データは、第1音のうち第2周波数帯域の成分のみを表すデータ、または、第1音のうち第1周波数帯域および第2周波数帯域を含む周波数帯域の成分を表すデータである。同様に、第2出力データは、第2音のうち第2周波数帯域の成分のみを表すデータ、または、第2音のうち第1周波数帯域および第2周波数帯域を含む周波数帯域の成分を表すデータである。

【0094】

推定モデルは、入力データと出力データ（第1出力データおよび第2出力データ）との関係を学習した統計的モデルである。推定モデルの典型例はニューラルネットワークであるが、推定モデルの種類は以上の例示に限定されない。

【0095】

態様1の具体例（態様2）において、前記混合音は、前記第1周波数帯域の成分と前記第2周波数帯域の成分とを含み、前記混合音データは、前記混合音のうち前記第1周波数帯域の成分を含まない音を表す。以上の構成によれば、混合音データが表す音が第1周波

10

20

30

40

50

数帯域の成分を含まないから、混合音データが表す音が第1周波数帯域の成分と第2周波数帯域の成分とを含む構成と比較して、推定モデルの機械学習に必要な処理負荷および当該推定モデルの規模が低減されるという利点がある。

【0096】

態様1または態様2の具体例(態様3)において、前記第1入力データは、前記第1音のうち前記第1周波数帯域の成分の強度スペクトルを表し、前記第2入力データは、前記第2音のうち前記第1周波数帯域の成分の強度スペクトルを表し、前記混合音データは、前記混合音のうち前記第2周波数帯域を含む周波数帯域の成分の強度スペクトルを表し、前記入力データは、前記第1入力データと前記第2入力データと前記混合音データとで構成される正規化されたベクトルと、当該ベクトルの大きさを表す強度指標とを含む。以上の構成によれば、強度指標が入力データに含まれるから、混合音に対応する音量の音を表す第1出力データおよび第2出力データが生成される。したがって、第1出力データおよび第2出力データが表す音の強度を調整する処理(スケーリング)が不要であるという利点がある。

10

【0097】

態様1から態様3の何れかの具体例(態様4)において、前記推定モデルは、前記第1出力データが表す音のうち前記第2周波数帯域の成分と、前記第2出力データが表す音のうち前記第2周波数帯域の成分とを混合した結果が、前記混合音のうち前記第2周波数帯域の成分に近似するように訓練されたモデルである。以上の構成によれば、第1出力データが表す音のうち第2周波数帯域の成分と、第2出力データが表す音のうち第2周波数帯域の成分とを混合した結果が、混合音のうち第2周波数帯域の成分に近似するように、推定モデルが訓練される。したがって、以上の条件を加味せずに訓練された推定モデルを利用する構成と比較して、第1音のうち第2周波数帯域の成分(第1出力データ)と第2音のうち第2周波数帯域の成分(第2出力データ)とを高精度に推定できる。

20

【0098】

態様1から態様4の何れかの具体例(態様5)において、さらに、前記混合音のうち前記第1周波数帯域の成分に対する音源分離により、前記第1音のうち第1周波数帯域の第1成分と、前記第2音のうち前記第1周波数帯域の第2成分とを生成し、前記入力データの取得においては、前記第1成分を表す前記第1入力データと、前記第2成分を表す前記第2入力データとを取得する。以上の構成によれば、混合音のうち第1周波数帯域の成分に対して音源分離が実行されるから、混合音の全帯域を対象として音源分離を実行する構成と比較して、音源分離のための処理負荷が軽減される。

30

【0099】

態様1から態様5の何れかの具体例(態様6)において、前記第1出力データは、前記第1音のうち前記第1周波数帯域の成分と前記第2周波数帯域の成分とを表し、前記第2出力データは、前記第2音のうち前記第1周波数帯域の成分と前記第2周波数帯域の成分とを表す。以上の構成によれば、第1周波数帯域および第2周波数帯域の双方の成分を含む第1出力データおよび第2出力データが生成される。したがって、第1出力データが第1音のうち第2周波数帯域の成分のみを表すデータであり、第2出力データが第2音のうち第2周波数帯域の成分のみを表すデータである構成と比較して、第1周波数帯域および第2周波数帯域の双方にわたる音響を簡便に生成できる。

40

【0100】

本開示のひとつの態様(態様7)に係る推定モデルの訓練方法は、入力データと出力データとを各々が含む複数の訓練データを取得し、前記複数の訓練データを利用した機械学習により、前記入力データと前記出力データとを関係を学習した推定モデルを確立し、前記入力データは、第1音源に対応する第1音のうち第1周波数帯域の成分を表す第1入力データと、前記第1音源とは異なる第2音源に対応する第2音のうち前記第1周波数帯域の成分を表す第2入力データと、前記第1音と前記第2音との混合音のうち前記第1周波数帯域とは異なる第2周波数帯域を含む周波数帯域の成分を含む音を表す混合音データとを含み、前記出力データは、前記第1音のうち前記第2周波数帯域を含む周波数帯域の成

50

分を表す第 1 出力データと、前記第 2 音のうち前記第 2 周波数帯域を含む周波数帯域の成分を表す第 2 出力データとを含む。

【 0 1 0 1 】

以上の構成によれば、第 1 音のうち第 1 周波数帯域の成分を表す第 1 入力データと、第 2 音のうち第 1 周波数帯域の成分を表す第 2 入力データとを含む入力データから、第 1 音のうち第 2 周波数帯域を含む周波数帯域の成分を表す第 1 出力データと、第 2 音のうち第 2 周波数帯域を含む周波数帯域の成分を表す第 2 出力データとの少なくとも一方を生成する推定モデルが確立される。以上の構成によれば、第 1 音源に対応する第 1 音と第 2 音源に対応する第 2 音との混合音を第 1 音と第 2 音とに分離する音源分離を、第 1 周波数帯域についてのみ限定的に実行すれば足りる。したがって、音源分離のための処理負荷が軽減される。

10

【 0 1 0 2 】

なお、本開示は、以上に例示した各態様（態様 1 から態様 6）に係る音響処理方法を実現する音響処理システム、または、当該音響処理方法をコンピュータに実行させるプログラム、としても実現される。また、本開示は、前述の態様 7 に係る訓練方法を実現する推定モデル訓練システム、または、当該訓練方法をコンピュータに実行させるプログラム、としても実現される。

【 符号の説明 】

【 0 1 0 3 】

1 0 0 ... 音響処理システム、 1 1 ... 制御装置、 1 2 ... 記憶装置、 1 3 ... 放音装置、 2 0 ... 音響処理部、 2 1 ... 周波数解析部、 2 2 ... 音源分離部、 2 3 ... 帯域拡張部、 2 3 1 ... 取得部、 2 3 2 ... 生成部、 2 4 ... 波形合成部、 2 5 ... 音量調整部、 3 0 ... 学習処理部、 3 1 ... 取得部、 3 2 ... 訓練部、 M ... 推定モデル。

20

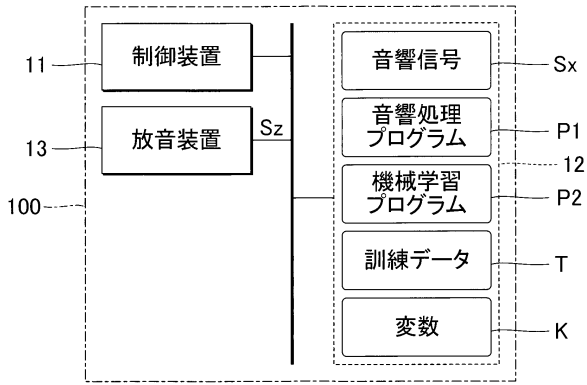
30

40

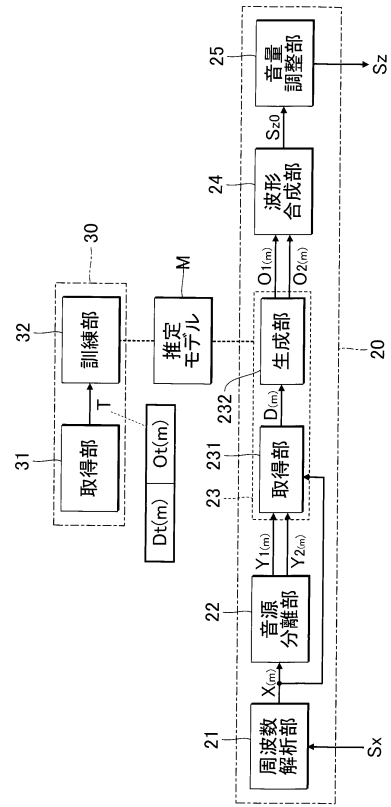
50

【図面】

【図 1】



【図 2】



10

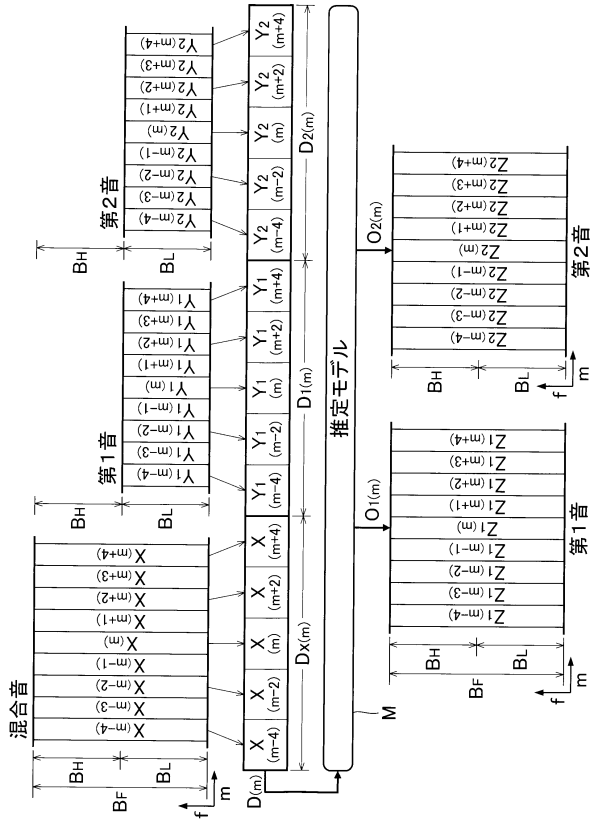
20

30

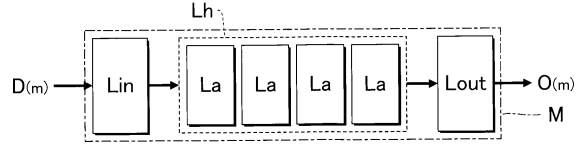
40

50

【図3】



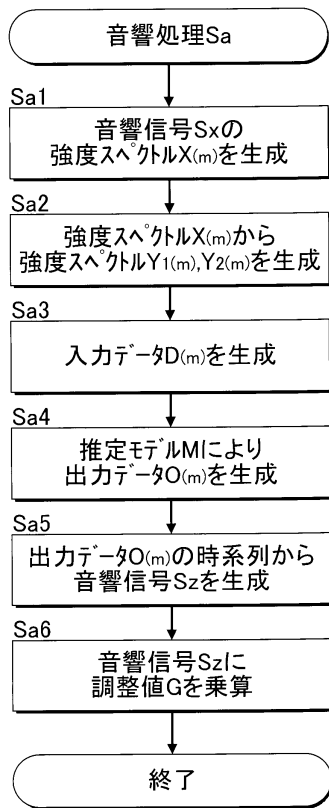
【図4】



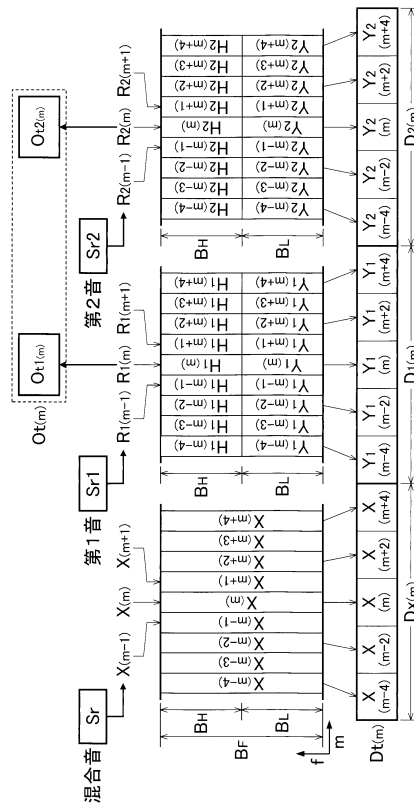
10

20

【図5】



【図6】

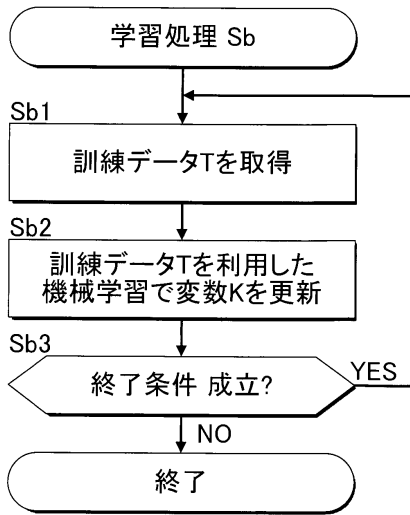


30

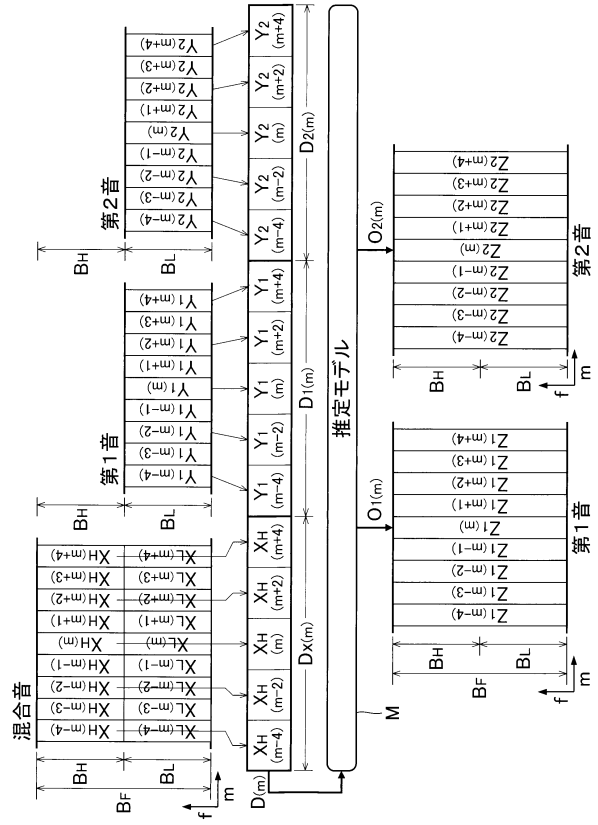
40

50

【 図 7 】



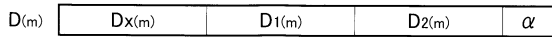
【 図 8 】



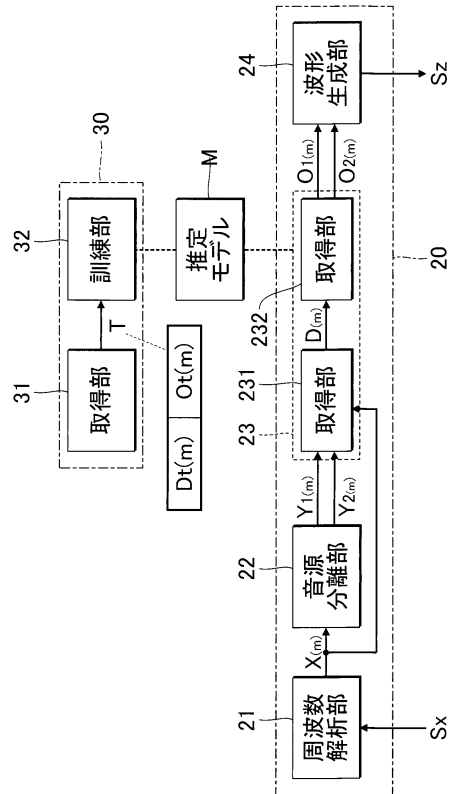
10

20

【 図 9 】



【 図 10 】

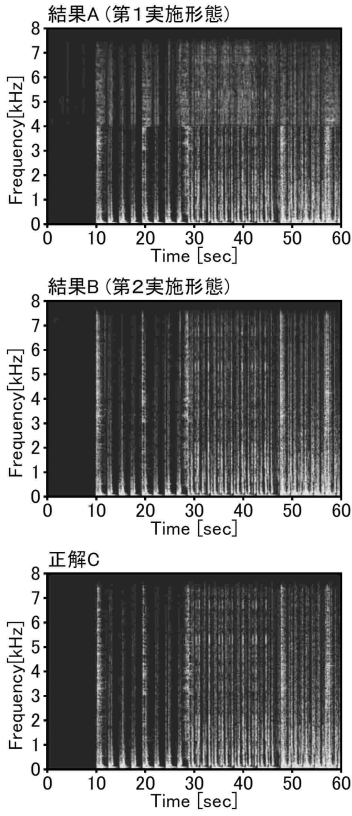


30

40

50

【 図 1 1 】



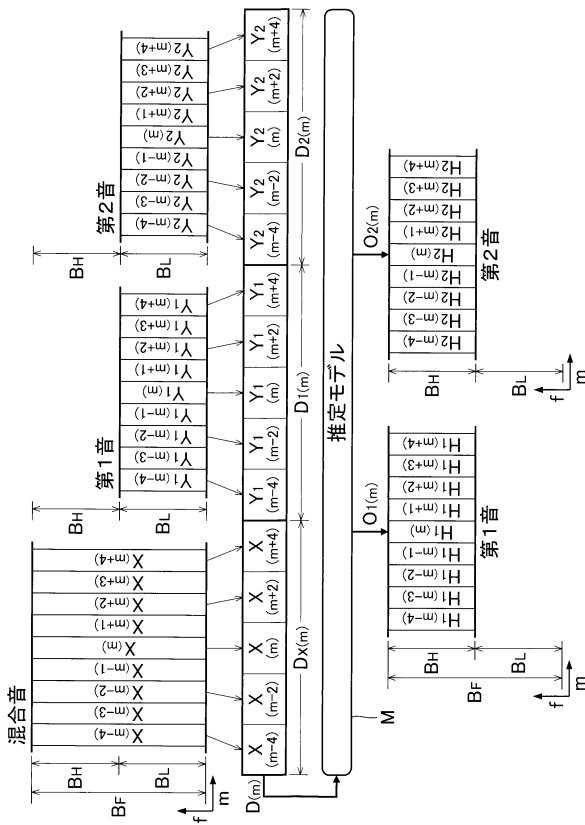
【 図 1 2 】

		SAR [dB]	改善量 [dB]
Drums	比較例	21.96	—
	第1実施形態	22.51	0.55
	第2実施形態	22.50	0.54
	第3実施形態	27.60	5.64
Vocals	比較例	25.00	—
	第1実施形態	25.62	0.62
	第2実施形態	25.62	0.62
	第3実施形態	29.28	4.29

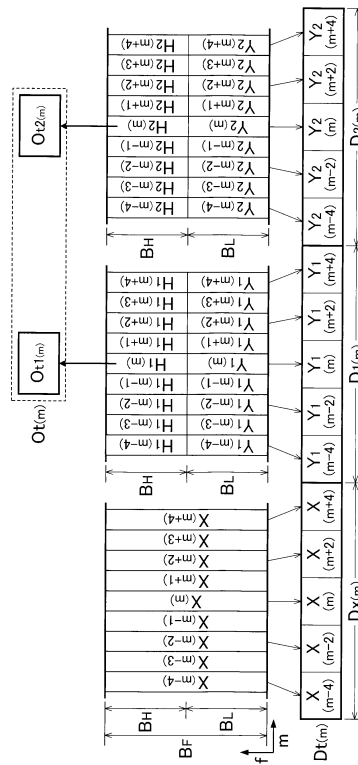
10

20

【 図 1 3 】



【 図 1 4 】

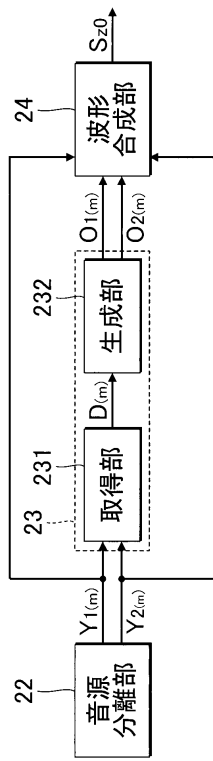


30

40

50

【図 15】



10

20

30

40

50

フロントページの続き

- (56)参考文献 特開 2012 - 22120 (JP, A)
特開 2008 - 278406 (JP, A)
- (58)調査した分野 (Int.Cl., DB名)
- G10L 21/0272
G10L 25/30