

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7129331号

(P7129331)

(45)発行日 令和4年9月1日(2022.9.1)

(24)登録日 令和4年8月24日(2022.8.24)

(51)国際特許分類

F I

G 1 0 L 21/0224(2013.01)

G 1 0 L 21/0224

A 6 3 F 13/54 (2014.01)

A 6 3 F 13/54

A 6 3 F 13/60 (2014.01)

A 6 3 F 13/60

G 1 0 L 21/034(2013.01)

G 1 0 L 21/034

G 1 0 L 21/0364(2013.01)

G 1 0 L 21/0364

請求項の数 8 (全12頁)

(21)出願番号 特願2018-241640(P2018-241640)

(22)出願日 平成30年12月25日(2018.12.25)

(65)公開番号 特開2020-101767(P2020-101767
A)

(43)公開日 令和2年7月2日(2020.7.2)

審査請求日 令和3年10月15日(2021.10.15)

(73)特許権者 595000427

株式会社コーエーテクモゲームス

神奈川県横浜市西区みなとみらい四丁目
3番6号

(74)代理人 100107766

弁理士 伊東 忠重

(74)代理人 100070150

弁理士 伊東 忠彦

(72)発明者 小池 雅人

神奈川県横浜市港北区箕輪町一丁目18

番12号 株式会社コーエーテクモゲー

ムス内

審査官 大石 剛

最終頁に続く

(54)【発明の名称】 情報処理装置、情報処理方法、及びプログラム

(57)【特許請求の範囲】

【請求項1】

録音された音声データから第1閾値以上の帯域の音が抽出された第1データに基づいて、リップノイズが録音されている第1区間を検出する検出部と、

前記音声データから第2閾値以下の帯域の音が抽出された前記第1区間の第2データに基づいて、前記音声データにおける前記第1区間の音のデータを修正する修正部と、
を有する情報処理装置。

【請求項2】

第2閾値は、前記第1閾値よりも小さい、
請求項1に記載の情報処理装置。

【請求項3】

前記検出部は、
前記第1データにおいて、前記第1区間よりも前の所定時間長の第2区間の音量の最大値よりも、前記第1区間の音量の最大値が第3閾値以上増加しており、
前記第1区間よりも後の所定数の所定時間長の各区間の音量の最大値の少なくとも一つが前記第2区間の音量の最大値以下の場合、前記第1区間を検出する、
請求項1または2に記載の情報処理装置。

【請求項4】

前記検出部は、
前記第1データにおいて、前記第1区間よりも前の所定時間長の第2区間の音量の最大

値よりも、前記第 1 区間の音量の最大値が第 3 閾値以上増加しており、

前記第 1 区間よりも後の所定数の所定時間長の各区間のうち音量の最大値が最も大きい第 3 区間の音量の最大値よりも、当該各区間のうち前記第 3 区間よりも後の第 4 区間の音量の最大値が第 4 閾値以上低下している場合、前記第 1 区間を検出する、
請求項 1 から 3 のいずれか一項に記載の情報処理装置。

【請求項 5】

前記修正部は、

前記第 2 データに対して、前記第 1 区間に含まれる最初の時間帯の音量を、前記音声データにおける前記第 1 区間の前の時間帯の音量に基づいて修正する処理、及び前記第 1 区間に含まれる最後の時間帯の音量を、前記音声データにおける前記第 1 区間の後の時間帯の音量に基づいて修正する処理の少なくとも一方を実行する、
請求項 1 から 4 のいずれか一項に記載の情報処理装置。

10

【請求項 6】

前記修正部は、前記音声データに対して、前記音声データにおける前記第 1 区間の前の時間帯の音量を、前記第 2 データの前記第 1 区間に含まれる最初の時間帯の音量に基づいて修正する処理、及び前記音声データにおける前記第 1 区間の後の時間帯の音量を、前記第 1 区間に含まれる最後の時間帯の音量に基づいて修正する処理の少なくとも一方を実行する、
請求項 1 から 5 のいずれか一項に記載の情報処理装置。

【請求項 7】

20

情報処理装置が、

録音された音声データから第 1 閾値以上の帯域の音が抽出された第 1 データに基づいて、リップノイズが録音されている第 1 区間を検出する処理と、
前記音声データから第 2 閾値以下の帯域の音が抽出された前記第 1 区間の第 2 データに基づいて、前記音声データにおける前記第 1 区間の音のデータを修正する処理と、
を実行する情報処理方法。

【請求項 8】

情報処理装置に、

録音された音声データから第 1 閾値以上の帯域の音が抽出された第 1 データに基づいて、リップノイズが録音されている第 1 区間を検出する処理と、
前記音声データから第 2 閾値以下の帯域の音が抽出された前記第 1 区間の第 2 データに基づいて、前記音声データにおける前記第 1 区間の音のデータを修正する処理と、
を実行させるプログラム。

30

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理装置、情報処理方法、及びプログラムに関する。

【背景技術】

【0002】

従来、コンピュータゲーム等において、ゲームの状況に応じて、ゲームのキャラクタのセリフを、予め録音されている音声（ボイス）により出力する技術が知られている（例えば、特許文献 1 を参照）。

40

【0003】

声優や歌手等の発声者が発声した音声を録音する場合、発声者が口を開閉した際に生じる雑音（リップノイズ）が、発声者の口元に設置されたマイクにより集音される場合がある。この場合、録音されている音声を出力すると、ユーザにとって耳障りな雑音となる場合がある。従来、音声を修正する者が、録音されている音声をスピーカに出力させてリップノイズを耳で聞き取り、リップノイズが発生した時間の音声の波形を画面に表示させ、波形を手入力により修正することが知られている。

【先行技術文献】

50

【特許文献】

【 0 0 0 4 】

【文献】特開 2 0 1 7 - 1 8 4 8 4 2 号公報

【発明の概要】

【発明が解決しようとする課題】

【 0 0 0 5 】

しかしながら、従来技術では、職人の経験と勘に基づいて手作業により周波数成分や音量等を修正するため、作業に手間がかかると共に、修正の品質にばらつきがあるという問題がある。

【 0 0 0 6 】

そこで、一側面では、より適切に音声を修正することができる技術を提供することを目的とする。

【課題を解決するための手段】

【 0 0 0 7 】

一つの案では、録音された音声データから第 1 閾値以上の帯域の音が抽出された第 1 データに基づいて、リップノイズが録音されている第 1 区間を検出する検出部と、前記音声データから第 2 閾値以下の帯域の音が抽出された前記第 1 区間の第 2 データに基づいて、前記音声データにおける前記第 1 区間の音のデータを修正する修正部と、を有する。

【発明の効果】

【 0 0 0 8 】

一側面によれば、より適切に音声を修正することができる。

【図面の簡単な説明】

【 0 0 0 9 】

【図 1】実施形態に係る情報処理装置のハードウェア構成例を示す図である。

【図 2】実施形態に係る情報処理装置の機能ブロック図である。

【図 3】実施形態に係る情報処理装置 10 の処理の一例を示すフローチャートである。

【図 4】実施形態に係る高域部分に基づいてリップノイズ区間を特定する処理の一例を示すフローチャートである。

【図 5】実施形態に係る高域部分に基づいてリップノイズ区間を特定する処理について説明する図である。

【図 6】実施形態に係る特定したノイズ区間の音声を修正する処理の一例を示すフローチャートである。

【図 7 A】ノイズ区間を含む所定区間における、録音された音声データの一例を示す図である。

【図 7 B】実施形態に係るリップノイズ区間の中低域の音を抽出する処理について説明する図である。

【図 7 C】実施形態に係るリップノイズ区間の音を修正する処理について説明する図である。

【発明を実施するための形態】

【 0 0 1 0 】

以下、図面に基づいて本発明の実施形態を説明する。

【 0 0 1 1 】

< ハードウェア構成 >

図 1 は、実施形態に係る情報処理装置 10 のハードウェア構成例を示す図である。図 1 に示す情報処理装置 10 は、それぞれバス B で相互に接続されているドライブ装置 100、補助記憶装置 102、メモリ装置 103、CPU 104、インタフェース装置 105、表示装置 106、及び入力装置 107 等を有する。

【 0 0 1 2 】

情報処理装置 10 での処理を実現するゲームプログラムは、記録媒体 101 によって提供される。ゲームプログラムを記録した記録媒体 101 がドライブ装置 100 にセットさ

10

20

30

40

50

れると、ゲームプログラムが記録媒体 101 からドライブ装置 100 を介して補助記憶装置 102 にインストールされる。但し、ゲームプログラムのインストールは必ずしも記録媒体 101 より行う必要はなく、ネットワークを介して他のコンピュータよりダウンロードするようにしてもよい。補助記憶装置 102 は、インストールされたゲームプログラムを格納すると共に、必要なファイルやデータ等を格納する。

【0013】

メモリ装置 103 は、例えば、DRAM (Dynamic Random Access Memory)、または SRAM (Static Random Access Memory) 等のメモリであり、プログラムの起動指示があった場合に、補助記憶装置 102 からプログラムを読み出して格納する。CPU 104 は、メモリ装置 103 に格納されたプログラムに従って情報処理装置 10 に係る機能を実現する。インタフェース装置 105 は、ネットワークに接続するためのインタフェースとして用いられる。表示装置 106 はプログラムによる GUI (Graphical User Interface) 等を表示する。入力装置 107 は、コントローラ等、キーボード及びマウス等、またはタッチパネル及びボタン等で構成され、様々な操作指示を入力させるために用いられる。

10

【0014】

なお、記録媒体 101 の一例としては、CD-ROM、DVD ディスク、ブルーレイディスク、又は USB メモリ等の可搬型の記録媒体が挙げられる。また、補助記憶装置 102 の一例としては、HDD (Hard Disk Drive)、SSD (Solid State Drive)、又はフラッシュメモリ等が挙げられる。記録媒体 101 及び補助記憶装置 102 のいずれについても、コンピュータ読み取り可能な記録媒体に相当する。

20

【0015】

<機能構成>

次に、図 2 を参照し、情報処理装置 10 の機能構成について説明する。図 2 は、実施形態に係る情報処理装置 10 の機能ブロック図である。

【0016】

情報処理装置 10 は、記憶部 11 を有する。記憶部 11 は、例えば、補助記憶装置 102 等を用いて実現される。記憶部 11 は、録音されたセリフの音声データ等を記憶する。

【0017】

また、情報処理装置 10 は、取得部 12、検出部 13、及び修正部 14 を有する。これら各部は、情報処理装置 10 にインストールされた 1 以上のプログラムが、情報処理装置 10 の CPU 104 に実行させる処理により実現される。

30

【0018】

取得部 12 は、録音されたセリフ等の音声データを記憶部 11 から取得する。検出部 13 は、取得部 12 により取得された音声データから第 1 閾値以上の帯域の音を抽出することにより高域の音データ(「第 1 データ」の一例。)を生成する。そして、検出部 13 は、生成した高域の音データから、リップノイズが録音されている各区間を検出する。

【0019】

修正部 14 は、取得部 12 により取得された音声データから、検出部 13 により検出された各区間の第 2 閾値以下の帯域の各音を抽出することにより、中低域の各音データ(「第 2 データ」の一例。)を生成する。そして、修正部 14 は、生成した中低域の各音データに基づいて、取得部 12 により取得された音声データにおける検出部 13 により検出された各区間の音のデータを修正する。これにより、例えば、声優等が発声したセリフの音声データに含まれるリップノイズを低減することができる。

40

【0020】

<処理>

次に、図 3 を参照して、情報処理装置 10 の処理について説明する。図 3 は、実施形態に係る情報処理装置 10 の処理の一例を示すフローチャートである。

【0021】

ステップ S1 において、検出部 13 は、取得部 12 により取得された音声データから高

50

域の音声データを抽出する。ここで、検出部 13 は、例えば、ハイパスフィルター（アンチエイリアスフィルタ）を用いて、録音された音声データから第 1 閾値（例えば、9000 Hz）未満の周波数成分を除去した音声データを生成してもよい。

【0022】

続いて、検出部 13 は、抽出した高域の音声データに基づいて、リップノイズが発生している区間（以下で、「リップノイズ区間」と称する。）をそれぞれ特定する（ステップ S2）。ここで、検出部 13 は、例えば、抽出した高域の音声データの音量が、短時間（例えば、0.1 秒）で所定比（例えば、12 dB）以上変化する各区間（時間帯、時間的な位置）を検出する。なお、この処理については後述する。

【0023】

続いて、検出部 13 は、各リップノイズ区間の音声に周期性があるか否かを判定する（ステップ S3）。ここで、検出部 13 は、例えば、リップノイズ区間より前の所定期間内の波形から、リップノイズ区間の波形と類似度が高い区間である第 1 類似区間を特定する。そして、検出部 13 は、リップノイズ区間と第 1 類似区間との間の第 1 距離（時間差）を算出する。また、検出部 13 は、リップノイズ区間より後の所定期間内の波形から、リップノイズ区間の波形と類似度が高い区間である第 2 類似区間を特定する。そして、リップノイズ区間と第 2 類似区間との間の第 2 距離（時間差）を算出する。

【0024】

そして、検出部 13 は、第 1 距離と第 2 距離との差が所定の閾値未満である場合、リップノイズ区間の音声に周期性があると判定する。また、検出部 13 は、第 1 距離と第 2 距離との一方を整数倍した値と、他方との差が所定の閾値未満である場合も、リップノイズ区間の音声に周期性があると判定する。

【0025】

周期性がある場合（ステップ S3 で YES）、処理を終了する。なお、修正部 14 は、周期性がある場合、リップノイズを除去する量を低減してもよい。この場合、修正部 14 は、周期性があると判定した区間に対しては、後述する第 2 閾値を比較的大きな値（例えば、4000 Hz）に決定してリップノイズ区間の音声の修正処理を行うようにしてもよい。

【0026】

周期性がない場合（ステップ S3 で NO）、修正部 14 は、取得部 12 により取得された音声データにおいて、各リップノイズ区間の音声を修正し（ステップ S4）、処理を終了する。なお、この処理については後述する。

【0027】

リップノイズ区間の特定処理

次に、図 4、及び図 5 を参照し、図 3 のステップ S2 の、抽出した高域の音声データに基づいてリップノイズ区間を特定する処理について説明する。図 4 は、実施形態に係る高域部分に基づいてリップノイズ区間を特定する処理の一例を示すフローチャートである。図 5 は、実施形態に係る高域部分に基づいてリップノイズ区間を特定する処理について説明する図である。なお、以下の説明において、音量の最大値は、音量の絶対値の最大値としてもよい。

【0028】

ステップ S101 において、検出部 13 は、音量が所定の閾値以上（例えば、12 dB）変化し始める各時点を検出する。図 5 の例では、ステップ S1 の処理で録音された音声データから抽出された、高域部分のみの音声データの波形 500 が、横軸を時間、縦軸を音量（dB）として示されている。ここで、検出部 13 は、図 5 のように、例えば、ある時点 501 から所定時間（例えば、0.1 秒）の後の時点 502 までの区間 521 の最大音量 511 と、時点 502 から当該一定時間の後の時点 503 までの区間 522 の最大音量 512 との差が所定の閾値以上変化する場合、時点 502 を、音量が所定の閾値以上変化し始める時点として検出する。

【0029】

10

20

30

40

50

なお、検出部 13 は、以下の処理を、検出した各時点に対してそれぞれ実行する。そのため、以下では、検出した各時点に含まれる一の時点（以下で、「処理対象の時点」と称する。）に対しての処理について説明する。図 5 の時点 502 は、処理対象の時点の一例である。

【0030】

続いて、検出部 13 は、処理対象の時点から当該一定時間の後の時点から、所定期間（例えば、0.3 秒）後までの時間帯（以下で、「リリース区間」とも称する。）に含まれる各区間の音量の最大値の少なくとも一つが、音量が所定の閾値以上変化する前の区間の音量の最大値以下であるか否かを判定する（ステップ S102）。ここで、検出部 13 は、処理対象の時点 502 から当該一定時間の後の時点である時点 503 から当該一定時間後の時点 504 までの区間 523 の最大音量 513、時点 504 から当該一定時間後の時点 505 までの区間 524 の最大音量 514、時点 505 から当該一定時間後の時点 506 までの区間 525 の最大音量 515 を特定する。そして、検出部 13 は、最大音量 513、最大音量 514、及び最大音量 515 の少なくとも一つが、最大音量 511 よりも小さいか否かを判定する。

10

【0031】

リリース区間に含まれる各区間の音量の最大値の少なくとも一つが、音量が所定の閾値以上変化する前の区間の音量の最大値以下でない場合（ステップ S102 で NO）、ステップ S104 の処理に進む。

【0032】

一方、最大値以下である場合（ステップ S102 で YES）、検出部 13 は、処理対象の時点から、変化する前の音量まで下がった区間の終了時点までを、リップノイズ区間であると特定する（ステップ S103）。この場合、検出部 13 は、図 5 の例では、最大音量 511 以下となるのが区間 525 の最大音量 515 であるため、時点 502 から区間 525 の終了時点である時点 506 までの区間を、リップノイズ区間とする。これは、音量が所定の閾値以上に変化し始める時点 502 から所定期間内において、一定時間内の最大音量が、時点 502 よりも前の一定時間内の最大音量 511 以下に下がる場合、リップノイズによる音量の変化であると考えられるためである。

20

【0033】

続いて、検出部 13 は、リリース区間において、音量が所定の閾値以上減衰しているか否かを判定する（ステップ S104）。ここで、検出部 13 は、例えば、最大音量 512、最大音量 513、及び最大音量 514 のうち最も大きいものを特定する。そして、区間 523 から 525 のうちの第 1 区間であって、特定した最大音量の区間よりも後の第 1 区間の最大音量が、特定した最大音量よりも所定の閾値（例えば、12 dB）以上減衰しているか否かを判定してもよい。これは、リップノイズとセリフ等の発声が同時に発生したような場合の音量の変化であると考えられるためである。

30

【0034】

音量が所定の閾値以上減衰していない場合（ステップ S104 で NO）、処理対象の時点に対する処理を終了する。一方、音量が所定の閾値以上減衰している場合（ステップ S104 で YES）、検出部 13 は、処理対象の時点から、音量が所定の閾値以上減衰した区間の終了時点までを、リップノイズ区間であると特定する（ステップ S105）。図 5 において、最大音量 512、最大音量 513、及び最大音量 514 のうち最も大きい値は最大音量 513 であり、最大音量 514 は最大音量 513 よりも 12 dB 以上低いものとする。この場合、検出部 13 は、時点 502 から、最大音量 514 の区間 524 の終了時点である時点 505 までの区間を、リップノイズ区間と判定する。

40

【0035】

リップノイズ区間の音声の修正処理

次に、図 6 から図 7C を参照し、図 3 のステップ S4 の、特定したリップノイズ区間の音声を修正する処理について説明する。図 6 は、実施形態に係る特定したリップノイズ区間の音声を修正する処理の一例を示すフローチャートである。図 7A は、リップノイズ区

50

間を含む所定区間における、録音された音声データの一例を示す図である。図 7 B は、実施形態に係るリップノイズ区間の中低域の音を抽出する処理について説明する図である。図 7 C は、実施形態に係るリップノイズ区間の音を修正する処理について説明する図である。

【 0 0 3 6 】

ステップ S 2 0 1 において、修正部 1 4 は、取得部 1 2 により取得された音声データから、検出部 1 3 により検出された各リップノイズ区間の低中域の音声データを抽出する。ここで、修正部 1 4 は、例えば、ローパスフィルター（アンチエイリアスフィルター）を用いて、録音された音声データから第 1 閾値よりも低い第 2 閾値（例えば、2 0 0 0 H z ）以上の周波数成分を除去した音声データを生成してもよい。

10

【 0 0 3 7 】

ここで、図 7 A の例では、録音された音声データの波形 7 0 0 が、横軸を時間、縦軸を音量（d B ）として示されている。図 7 A において、時点 7 0 1 から時点 7 0 2 の区間 7 1 1 が、リップノイズ区間として特定されているものとする。図 7 B の例では、図 7 A の波形 7 0 0 のうち、区間 7 1 1 の部分の波形 7 0 0 A、区間 7 1 1 より前の部分の波形 7 0 0 B、区間 7 1 1 より後の部分の波形 7 0 0 C、及びステップ S 2 0 1 の処理により抽出された区間 7 1 1 の波形 7 2 1 が示されている。図 7 B の例では、ステップ S 2 0 1 の処理により第 2 閾値以上の周波数成分が除去された波形 7 2 1 の各時点における音量の絶対値は、波形 7 0 0 A の音量の絶対値よりも小さくなっている。

【 0 0 3 8 】

続いて、修正部 1 4 は、抽出した低中域の音声データの音量を、録音された音声データのリップノイズ区間前後の音声データと整合するように修正（調整）する（ステップ S 2 0 2 ）。

20

【 0 0 3 9 】

図 7 C の例で、波形 7 2 1 のうち、時点 7 0 1 から所定時間（例えば、0 . 0 3 秒）後の時点 7 4 1 までの区間 7 5 1 の部分の波形を波形 7 2 1 A とする。また、波形 7 2 1 のうち、時点 7 0 2 から所定時間前の時点 7 4 2 から、時点 7 0 2 までの区間 7 5 2 の部分の波形を波形 7 2 1 B とする。

【 0 0 4 0 】

（リップノイズ区間の中低域データの修正）

30

ステップ S 2 0 2 において、修正部 1 4 は、時点 7 0 1 における波形 7 2 1 A の音量 7 3 2 A と、時点 7 0 1 における波形 7 0 0 A の音量 7 3 2 B との中点 7 3 2 を通り、時点 7 4 1 における波形 7 2 1 A の音量 7 3 3 を通るように波形 7 2 1 A を修正した波形 7 3 1 A のデータを生成する。

【 0 0 4 1 】

そして、修正部 1 4 は、時点 7 4 2 における波形 7 2 1 B の音量 7 3 4 を通り、時点 7 0 2 における波形 7 2 1 B の音量 7 3 5 A と、時点 7 0 2 における波形 7 0 0 A の音量 7 3 5 B との中点 7 3 5 を通るように波形 7 2 1 B を修正した波形 7 3 1 B のデータを生成する。

【 0 0 4 2 】

続いて、修正部 1 4 は、録音された音声データのリップノイズ区間前後の音声データの音量を、抽出した低中域の音声データと整合するように修正する（ステップ S 2 0 3 ）。

40

【 0 0 4 3 】

（リップノイズ区間前後の音声データの修正）

また、ステップ S 2 0 2 において、修正部 1 4 は、区間 7 1 1 より前の部分の波形 7 0 0 B を、時点 7 0 1 から所定時間（例えば、0 . 0 3 秒）前の時点 7 4 3 から時点 7 0 1 までの区間において、時点 7 4 3 の音量 7 6 1 と中点 7 3 2 を通る波形 7 7 1 A に修正する。また、区間 7 1 1 より後の部分の波形 7 0 0 C を、時点 7 0 2 から、所定時間（例えば、0 . 0 3 秒）後の時点 7 4 4 までの区間において、中点 7 3 5 と時点 7 4 4 の音量 7 6 2 とを通る波形 7 7 1 B に修正する。

50

【 0 0 4 4 】

続いて、修正部 1 4 は、ステップ S 2 0 3 の処理で修正した音声データのリップノイズ区間の音声データを、ステップ S 2 0 2 の処理で修正した低中域の音声データに置換し（ステップ S 2 0 4 ）、処理を終了する。

【 0 0 4 5 】

< 変形例 >

情報処理装置 1 0 の各機能部は、例えば 1 以上のコンピュータにより構成されるクラウドコンピューティングにより実現されていてもよい。

【 0 0 4 6 】

以上、本発明の実施例について詳述したが、本発明は斯かる特定の実施形態に限定されるものではなく、特許請求の範囲に記載された本発明の要旨の範囲内において、種々の変形・変更が可能である。

【 符号の説明 】

【 0 0 4 7 】

- 1 0 情報処理装置
- 1 1 記憶部
- 1 2 取得部
- 1 3 検出部
- 1 4 修正部

10

20

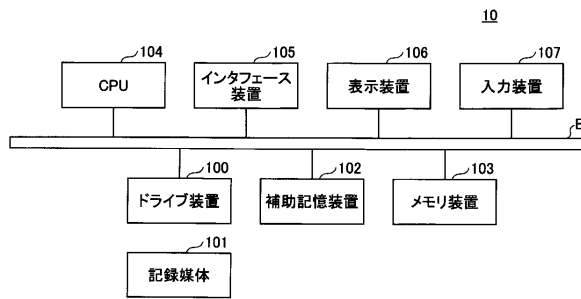
30

40

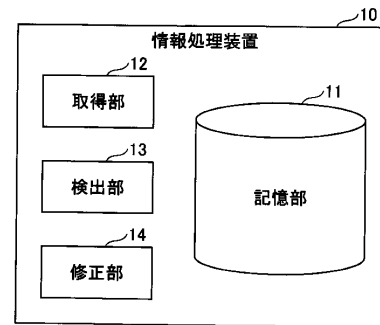
50

【図面】

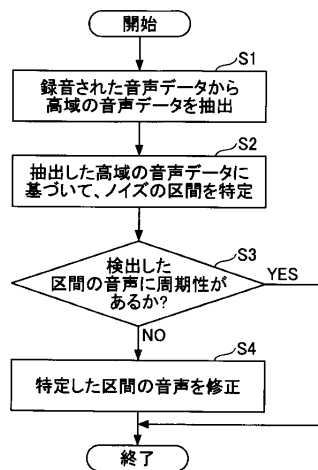
【図 1】



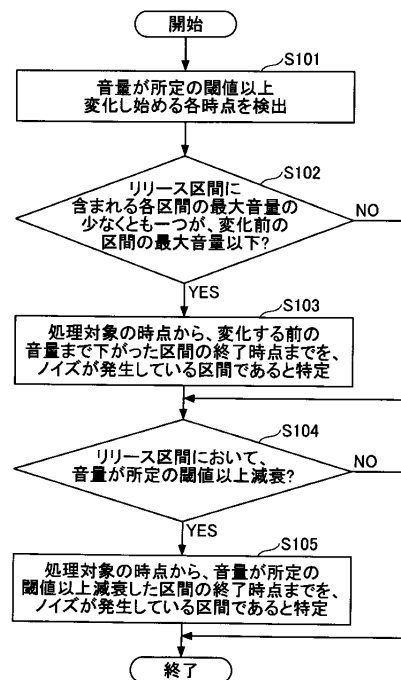
【図 2】



【図 3】



【図 4】



10

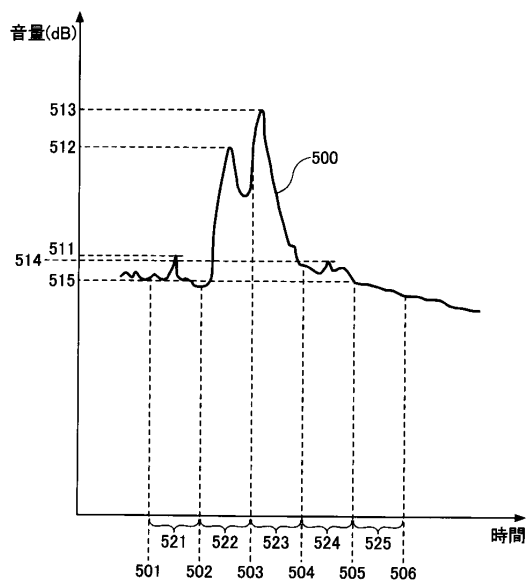
20

30

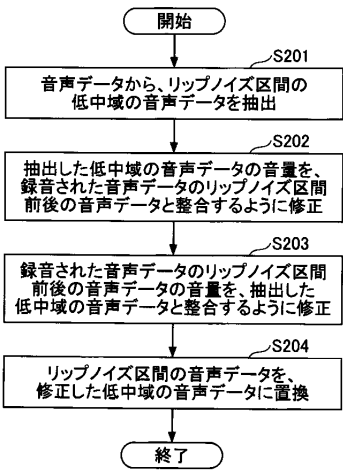
40

50

【図 5】



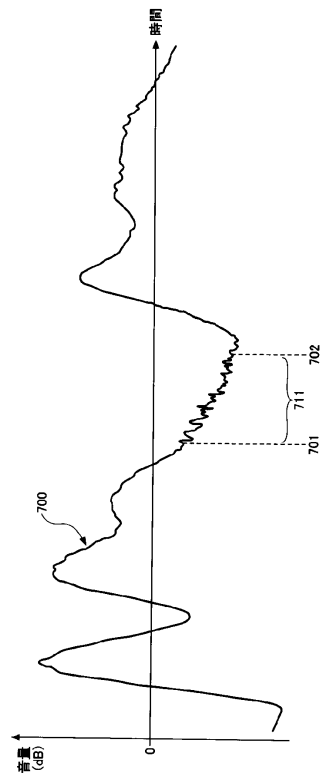
【図 6】



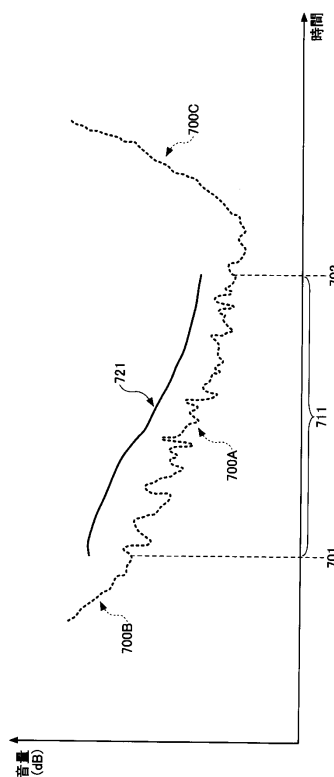
10

20

【図 7 A】



【図 7 B】

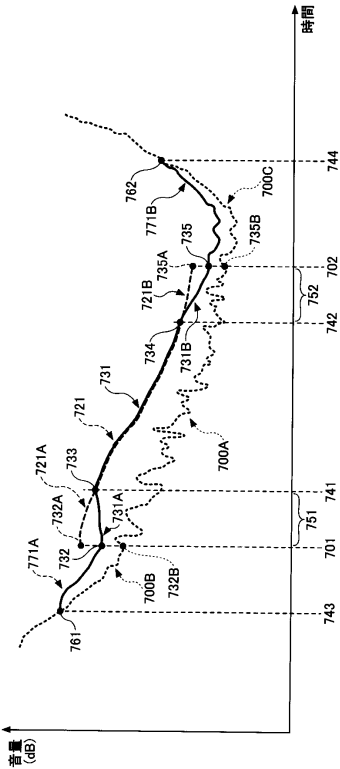


30

40

50

【図 7 C】



10

20

30

40

50

フロントページの続き

- (56)参考文献 特開 2 0 0 3 - 0 6 9 4 3 5 (J P , A)
特開 2 0 1 3 - 0 2 5 2 9 1 (J P , A)
特開 2 0 0 8 - 1 5 8 3 1 6 (J P , A)
特開平 0 7 - 2 6 1 7 7 9 (J P , A)

- (58)調査した分野 (Int.Cl. , D B 名)
G 1 0 L 2 1 / 0 2 2 4
A 6 3 F 1 3 / 5 4
A 6 3 F 1 3 / 6 0
G 1 0 L 2 1 / 0 3 4
G 1 0 L 2 1 / 0 3 6 4