

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6865701号
(P6865701)

(45) 発行日 令和3年4月28日 (2021.4.28)

(24) 登録日 令和3年4月8日 (2021.4.8)

(51) Int. Cl.	F I
G 1 0 L 15/22 (2006.01)	G 1 0 L 15/22 4 7 0 Z
G 1 0 L 19/00 (2013.01)	G 1 0 L 19/00 3 1 2 F
	G 1 0 L 19/00 3 1 2 E
	G 1 0 L 15/22 4 6 0 Z

請求項の数 9 (全 17 頁)

(21) 出願番号	特願2018-23711 (P2018-23711)	(73) 特許権者	000004352
(22) 出願日	平成30年2月14日 (2018.2.14)		日本放送協会
(65) 公開番号	特開2018-180519 (P2018-180519A)		東京都渋谷区神南2丁目2番1号
(43) 公開日	平成30年11月15日 (2018.11.15)	(73) 特許権者	591053926
審査請求日	令和2年10月5日 (2020.10.5)		一般財団法人NHKエンジニアリングシス テム
(31) 優先権主張番号	特願2017-82196 (P2017-82196)		東京都世田谷区砧一丁目10番11号
(32) 優先日	平成29年4月18日 (2017.4.18)	(74) 代理人	110001807
(33) 優先権主張国・地域又は機関	日本国 (JP)		特許業務法人磯野国際特許商標事務所
早期審査対象出願		(72) 発明者	三島 剛
			東京都世田谷区砧一丁目10番11号 日 本放送協会放送技術研究所内
		(72) 発明者	佐藤 庄衛
			東京都世田谷区砧一丁目10番11号 日 本放送協会放送技術研究所内
			最終頁に続く

(54) 【発明の名称】 音声認識誤り修正支援装置およびそのプログラム

(57) 【特許請求の範囲】

【請求項 1】

コンテンツに含まれる音声に対する音声認識の誤りを修正する音声認識誤り修正支援装置であって、

テキストデータである前記音声の認識結果と当該認識結果を構成する単語ごとの時間情報とにより、前記認識結果を予め定めた基準でセグメントに分割する認識結果分割手段と、

項目情報とともに前記セグメントに含まれる単語列を表示するか否かを指定するボタンを表示し、前記ボタンの選択により、編集領域を表示して前記セグメントの単語列を展開するか、前記編集領域を非表示とするかの制御を行う認識結果表示制御手段と、

前記編集領域で前記セグメントの誤りを修正する誤り修正手段と、

前記編集領域の前記セグメントに対応する音声を再生する音声再生手段と、を備え、

前記認識結果分割手段は、前記コンテンツに含まれる位置情報または時間情報の変化点で、前記認識結果を分割し、

前記誤り修正手段は、前記編集領域で指定された単語位置からの前記時間情報に対応する前記コンテンツの音声を前記音声再生手段により再生させることを特徴とする音声認識誤り修正支援装置。

【請求項 2】

コンテンツに含まれる音声に対する音声認識の誤りを修正する音声認識誤り修正支援装置であって、

10

20

テキストデータである前記音声の認識結果と当該認識結果を構成する単語ごとの時間情報とにより、前記認識結果を予め定めた基準でセグメントに分割する認識結果分割手段と

、
項目情報とともに前記セグメントに含まれる単語列を表示するか否かを指定するボタンを表示し、前記ボタンの選択により、編集領域を表示して前記セグメントの単語列を展開するか、前記編集領域を非表示とするかの制御を行う認識結果表示制御手段と、

前記編集領域で前記セグメントの誤りを修正する誤り修正手段と、

前記編集領域の前記セグメントに対応する音声を再生する音声再生手段と、を備え、

前記コンテンツは映像を含み、前記認識結果分割手段は、前記映像のカット点で、前記認識結果を分割し、

10

前記誤り修正手段は、前記編集領域で指定された単語位置からの前記時間情報に対応する前記コンテンツの音声を前記音声再生手段により再生させることを特徴とする音声認識誤り修正支援装置。

【請求項 3】

コンテンツに含まれる音声に対する音声認識の誤りを修正する音声認識誤り修正支援装置であって、

テキストデータである前記音声の認識結果と当該認識結果を構成する単語ごとの時間情報とにより、前記認識結果を予め定めた基準でセグメントに分割する認識結果分割手段と

、
項目情報とともに前記セグメントに含まれる単語列を表示するか否かを指定するボタンを表示し、前記ボタンの選択により、編集領域を表示して前記セグメントの単語列を展開するか、前記編集領域を非表示とするかの制御を行う認識結果表示制御手段と、

20

前記編集領域で前記セグメントの誤りを修正する誤り修正手段と、

前記編集領域の前記セグメントに対応する音声を再生する音声再生手段と、

複数の前記セグメントに含まれる単語から、前記セグメントごとに、TF-IDF法により特徴単語を前記項目情報として抽出する項目情報抽出手段と、を備え、

前記認識結果表示制御手段は、前記セグメントに含まれる単語列を表示するか否かを指定するボタンを含んだ前記項目情報の一覧を表示し、

前記誤り修正手段は、前記編集領域で指定された単語位置からの前記時間情報に対応する前記コンテンツの音声を前記音声再生手段により再生させることを特徴とする音声認識誤り修正支援装置。

30

【請求項 4】

コンテンツに含まれる音声に対する音声認識の誤りを修正する音声認識誤り修正支援装置であって、

テキストデータである前記音声の認識結果と当該認識結果を構成する単語ごとの時間情報とにより、前記認識結果を予め定めた基準でセグメントに分割する認識結果分割手段と

、
項目情報とともに前記セグメントに含まれる単語列を表示するか否かを指定するボタンを表示し、前記ボタンの選択により、編集領域を表示して前記セグメントの単語列を展開するか、前記編集領域を非表示とするかの制御を行う認識結果表示制御手段と、

40

前記編集領域で前記セグメントの誤りを修正する誤り修正手段と、

前記編集領域の前記セグメントに対応する音声を再生する音声再生手段と、を備え、

前記誤り修正手段は、前記編集領域で指定された単語位置からの前記時間情報に対応する前記コンテンツの音声を前記音声再生手段により再生させ、前記コンテンツの音声再生中に前記編集領域の任意の単語位置を指定されることで、前記音声再生手段における音声の再生を停止させ、前記編集領域で指定された単語位置に音声の再生開始を示すポップアップメッセージを表示し、音声が停止した単語位置に音声の再生終了を示すポップアップメッセージを表示することを特徴とする音声認識誤り修正支援装置。

【請求項 5】

前記誤り修正手段は、前記コンテンツの音声再生中に前記編集領域の任意の単語位置を

50

指定されることで、前記音声再生手段における音声の再生を停止することを特徴とする請求項 1 から請求項 3 のいずれか一項に記載の音声認識誤り修正支援装置。

【請求項 6】

前記誤り修正手段は、前記コンテンツの音声再生に連動して、再生される音声に対応する前記編集領域の単語の表示属性を変更することを特徴とする請求項 1 から請求項 5 のいずれか一項に記載の音声認識誤り修正支援装置。

【請求項 7】

前記誤り修正手段は、前記編集領域で指定された単語または指定区間の単語列の前記時間情報に対応する前記コンテンツの音声を、前記音声再生手段により、繰り返して再生することを特徴とする請求項 1 から請求項 6 のいずれか一項に記載の音声認識誤り修正支援装置。

10

【請求項 8】

コンテンツに含まれる音声に対する音声認識の誤りを修正する音声認識誤り修正支援装置であって、

テキストデータである前記音声の認識結果と当該認識結果を構成する単語ごとの時間情報とにより、前記認識結果を発話内容の切り替わりごとのセグメントに分割する認識結果分割手段と、

項目情報とともに前記セグメントに含まれる単語列を表示するか否かを指定するボタンを表示し、前記ボタンの選択により、編集領域を表示して前記セグメントの単語列を展開するか、前記編集領域を非表示とするかの制御を行う認識結果表示制御手段と、

20

前記編集領域で前記セグメントの誤りを修正する誤り修正手段と、

前記編集領域の前記セグメントに対応する音声を再生する音声再生手段と、を備え、

前記誤り修正手段は、前記編集領域で指定された単語位置からの前記時間情報に対応する前記コンテンツの音声を前記音声再生手段により再生させることを特徴とする音声認識誤り修正支援装置。

【請求項 9】

コンピュータを、請求項 1 から請求項 8 のいずれか一項に記載の音声認識誤り修正支援装置として機能させるための音声認識誤り修正支援プログラム。

【発明の詳細な説明】

【技術分野】

30

【0001】

本発明は、音声認識の誤り修正を支援する音声認識誤り修正支援装置およびそのプログラムに関する。

【背景技術】

【0002】

番組取材等で収録した音声素材（映像・音声素材を含む）の音声を文字として利用する場合、音声を文字に書き起こす作業が必須の作業となっている。通常、この作業は、作業者が収録した素材の音声を聴取し、端末のキーボード等で文字を入力することにより行っている。このとき、作業者は、音声の再生と停止を頻繁に繰り返したり、何度も同一箇所の音声を聞き直したりすることになるが、この作業は熟練者であっても素材の収録時間に対して約 6 倍の作業時間がかかるとも言われている。

40

【0003】

従来、音声の書き起こし作業を支援する技術として、入力された音声を任意の単位に区切った文（セル）ごとに音声認識処理を施し、音声認識処理された認識結果と、これに対応する音声とを比較し、音声認識処理の誤りを修正する技術が開示されている（特許文献 1 参照）。

この技術は、音声認識処理においてセル単位で音声を再生し、操作者がセル単位で認識結果を修正し、セルの修正を一般的なテキストエディタの操作で行う。また、この技術では、操作者は、特殊な操作を覚える必要はなく、セルの修正後、セルの先頭から音声を再生して、操作者が認識結果を正しく修正したか否かを確認していた。

50

【 0 0 0 4 】

また、従来の音声の書き起こし作業を支援する技術として、音声の認識結果を、単語ごとに対応付けて、単語単位で修正する技術が開示されている（特許文献 2，3 参照）。

この技術は、字幕放送等のリアルタイム性が要求される誤り修正や、誤りの少ない認識結果を修正する場合には有効である。

【 先行技術文献 】

【 特許文献 】

【 0 0 0 5 】

【 特許文献 1 】 特開 2 0 1 5 - 1 8 4 5 6 4 号公報

【 特許文献 2 】 特開 2 0 0 4 - 2 2 6 9 1 0 号公報

【 特許文献 3 】 特開 2 0 0 5 - 2 2 8 1 7 8 号公報

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 6 】

特許文献 1 で開示されている技術は、セル単位で音声の再生および認識結果の修正を行うため、修正箇所が少なくても、修正箇所の音声と修正結果が合致するか否かを確認するために、セルの先頭から音声を再生する必要がある。

そのため、この技術は、セルの途中にある修正対象箇所の音声再生されるまで、待ち時間が発生してしまうという問題があった。また、この技術は、セル内で、認識結果に対応する音声を操作者が聞き分ける必要があるため、認識結果が悪くなると、音声と修正対象とを対応付けることが困難になってしまうという問題があった。

【 0 0 0 7 】

また、特許文献 2，3 で開示されている技術のように、音声の認識結果を単語単位で修正する技術では、認識結果の修正と音声の確認とを素早く行うことは可能である。しかし、複数の単語に渡って認識誤りがある場合、順番に単語を指定して修正を行わなければならない、手順が複雑となり、その操作に慣れるまでに時間がかかってしまうという問題があった。

【 0 0 0 8 】

そこで、本発明は、音声認識の誤りを修正する際に、修正対象箇所の音声を素早く再生し、簡易な操作で音声認識の誤り修正を行うことが可能な音声認識誤り修正支援装置およびそのプログラムを提供することを課題とする。

【 課題を解決するための手段 】

【 0 0 0 9 】

前記課題を解決するため、本発明に係る音声認識誤り修正支援装置は、コンテンツの音声に対する音声認識の誤りを修正する音声認識誤り修正支援装置であって、認識結果分割手段と、認識結果表示制御手段と、誤り修正手段と、音声再生手段と、を備える構成とした。

【 0 0 1 1 】

かかる構成において、音声認識誤り修正支援装置は、認識結果分割手段によって、テキストデータである音声の認識結果と当該認識結果を構成する単語ごとの時間情報とにより、認識結果を予め定めた基準でセグメントに分割する。

【 0 0 1 3 】

そして、音声認識誤り修正支援装置は、認識結果表示制御手段によって、項目情報とともにセグメントに含まれる単語列を表示するか否かを指定するボタンを表示する。また、音声認識誤り修正支援装置は、認識結果表示制御手段によって、ボタンの選択により、編集領域を表示してセグメントの単語列を展開するか、編集領域を非表示とするかの制御を行う。これによって、認識結果表示制御手段は、音声の認識結果をすべて表示するのではなく、項目一覧によって操作者に編集対象のセグメントを指定させ、対象となったセグメントの単語列を編集領域に展開して操作者に提示する。

【 0 0 1 4 】

そして、音声認識誤り修正支援装置は、誤り修正手段によって、編集領域でセグメントの誤りを修正する。このとき、誤り修正手段は、編集領域で指定された単語位置からの時間情報に対応するコンテンツの音声を生声再生手段により再生させる。これによって、誤り修正手段は、認識結果またはその修正結果に対応する音声を操作者が素早く確認可能なように、指定された位置の単語から音声を再生する。

なお、音声認識誤り修正支援装置は、コンピュータを、前記した各手段として機能させるための音声認識誤り修正支援プログラムで動作させることができる。

【発明の効果】

【0015】

本発明は、以下に示す優れた効果を奏するものである。

10

本発明によれば、素材コンテンツの音声認識結果を分割して、項目の一覧を表示するため、簡易な操作で音声認識の誤りを確認したい認識結果を素早く選択することができる。

また、本発明によれば、編集領域で単語の位置を指定するという簡易な操作で、対応する音声を再生するため、音声認識結果の誤りの発見や、修正確認を素早く行うことができる。

これによって、本発明は、特別なスキルを必要とせずに、音声認識結果の誤りを修正することができる。

【図面の簡単な説明】

【0016】

【図1】本発明の実施形態に係る音声認識誤り修正支援装置の構成を示すブロック構成図である。

20

【図2】素材情報記憶手段が記憶する記憶内容を説明するための説明図である。

【図3】素材コンテンツを選択する素材コンテンツ選択画面の一例を示す画面構成図である。

【図4】素材コンテンツの音声認識結果を分割した項目の一覧を示す項目一覧画面の一例を示す画面構成図である。

【図5】項目一覧画面で編集領域に音声認識結果を展開した例を示す画面構成図である。

【図6】編集領域における編集作業の一例を説明するための説明図である。

【図7】音声再生に連動して編集領域の単語の表示属性を変更する例を説明するための説明図である。

30

【図8】編集領域における編集作業の操作内容を提示する例を説明するための説明図である。

【図9】編集領域における音声の繰り返し再生を指定する例を説明するための説明図である。

【図10】本発明の実施形態に係る音声認識誤り修正支援装置の音声認識結果をセグメント単位で生成するセグメント情報生成動作を示すフローチャートである。

【図11】本発明の実施形態に係る音声認識誤り修正支援装置の音声認識結果をセグメント単位で表示装置に提示するセグメント情報提示動作を示すフローチャートである。

【図12】本発明の実施形態に係る音声認識誤り修正支援装置の音声再生を行いながら認識結果を修正するセグメント修正動作を示すフローチャートである。

40

【図13】本発明の変形例の実施形態に係る音声認識誤り修正支援装置の構成を示すブロック構成図である。

【発明を実施するための形態】

【0017】

以下、本発明の実施形態について図面を参照して説明する。

[音声認識誤り修正支援装置の構成]

最初に、図1を参照して、本発明の実施形態に係る音声認識誤り修正支援装置1の構成について説明する。

音声認識誤り修正支援装置1は、少なくとも音声を含んだ素材コンテンツにおける音声の認識誤りの修正を支援するものである。なお、本実施形態では、素材コンテンツは、映

50

像と音声とからなるコンテンツ、例えば、放送用素材とする。

【0018】

音声認識誤り修正支援装置1は、図1に示すように、素材コンテンツ入力手段10と、音声認識手段11と、認識結果分割手段12と、項目情報抽出手段13と、素材情報記憶手段14と、編集手段15と、書き起こし結果出力手段16と、を備える。

【0019】

素材コンテンツ入力手段10は、素材コンテンツを入力するものである。

素材コンテンツ入力手段10は、例えば、外部の記憶媒体から素材コンテンツを入力するものであってもよいし、通信回線を介して入力するものであってもよい。

この素材コンテンツ入力手段10は、入力した素材コンテンツのうち、音声については、音声認識手段11に出力する。また、素材コンテンツ入力手段10は、入力した素材コンテンツ（映像・音声）を、後記する編集手段15における修正作業に使用するため、素材情報記憶手段14に書き込み記憶する。

【0020】

なお、素材コンテンツ入力手段10は、素材情報記憶手段14に素材コンテンツを書き込んだ後、音声認識手段11に対して、素材コンテンツの書き込み完了を通知し、音声認識手段11が素材情報記憶手段14から音声を読み出すこととしてもよい。

【0021】

音声認識手段11は、素材コンテンツ入力手段10が入力した素材コンテンツの音声を認識し、テキストデータである認識結果と当該認識結果を構成する単語ごとの時間情報とを生成するものである。

この音声認識手段11は、図示を省略した言語モデル、音響モデル、発音辞書により、音声認識を行い、認識した単語と、その単語の音声の先頭からの経過時間を示す時間情報とを生成する。音声認識手段11は、生成した認識結果の単語と時間情報とを認識結果分割手段12に出力する。なお、音声認識手段11における音声認識の手法は、例えば、特開2010-175765等の開示された音声から単語列を認識し、その結果を出力する手法を用いてもよい。

【0022】

認識結果分割手段12は、音声認識手段11で認識された認識結果（単語列）を、予め定めた基準で分割するものである。以下、認識結果分割手段12で生成された分割認識結果のそれぞれのかたまりをセグメントとよぶ。

認識結果分割手段12が用いる分割の基準は、任意の基準を予め定めることができる。

例えば、分割の基準として、音声の無音区間を用いることができる。この場合、認識結果分割手段12は、素材情報記憶手段14に記憶されている音声から音響特徴量であるパワー等によって無音区間を検出し、音声認識手段11による認識結果を、無音区間の前後で分割する。

【0023】

また、例えば、分割の基準として、映像のカット点を用いることができる。この場合、認識結果分割手段12は、素材情報記憶手段14に記憶されている映像から、隣接するフレームの画像特徴が予め定めた基準よりも大きく異なるフレームをカット点として検出し、カット点の時間の前後で認識結果を分割する。

【0024】

また、例えば、分割の基準として、素材コンテンツに予め付加されているメタ情報を用いてもよい。メタ情報としては、GPS（Global Positioning System）の位置情報（ジオタグ）等がある。この場合、認識結果分割手段12は、位置情報によって、素材コンテンツを撮影または集音した場所が異なっている時点で、認識結果を分割する。

【0025】

認識結果分割手段12は、音声認識結果を分割したセグメントを、項目情報抽出手段13に出力する。また、認識結果分割手段12は、セグメントごとに、単語とその時間情報とを素材情報記憶手段14に書き込み記憶する。

【 0 0 2 6 】

項目情報抽出手段 1 3 は、認識結果分割手段 1 2 で分割されたセグメントごとに、当該セグメントに含まれる特徴単語を項目として抽出するものである。

この特徴単語は、セグメント内に含まれる特徴的な単語である。例えば、項目情報抽出手段 1 3 は、T F - I D F 法 (T F : Term Frequency、単語の出現頻度、I D F : Inverse Document Frequency、逆文書頻度) によりセグメントを特徴付ける単語を抽出する。T F - I D F は、文書 (本実施形態では、セグメント) 中の単語に関する重みの一種であり、主に情報検索や文章要約などの分野で利用される。

具体的には、項目情報抽出手段 1 3 は、セグメント s 内の単語 w の出現頻度 $t f (w , s)$ を、以下の式 (1) で算出する。

10

【 0 0 2 7 】

【 数 1 】

$$t f (w , s) = \frac{n_{w,s}}{\sum_{t \in s} n_{t,s}} \quad \cdots \text{式 (1)}$$

【 0 0 2 8 】

この式 (1) で、 $n_{w,s}$ は、ある単語 w のセグメント s 内での出現回数、 $n_{t,s}$ は、セグメント s 内のすべての単語の出現回数の和を示す。

また、項目情報抽出手段 1 3 は、ある単語 w の逆文書頻度 $i d f (w)$ を、以下の式 (2) で算出する。

20

【 0 0 2 9 】

【 数 2 】

$$i d f (w) = \log \frac{N}{d f (w)} + 1 \quad \cdots \text{式 (2)}$$

【 0 0 3 0 】

この式 (2) で、 N は、素材コンテンツ内の全セグメント数、 $d f (w)$ は、ある単語 w が出現する素材コンテンツのセグメントの数 (総セグメント数 [総文書数]) を示す。

そして、項目情報抽出手段 1 3 は、セグメント内の各単語について、以下の式 (3) に示すように、式 (1) の $t f$ 値と式 (2) の $i d f$ 値との積が最も大きい単語、あるいは、予め定めた基準値よりも大きい単語を、当該セグメントの特徴単語とする。

30

【 0 0 3 1 】

【 数 3 】

$$t f (w , s) \times i d f (w) \quad \cdots \text{式 (3)}$$

【 0 0 3 2 】

項目情報抽出手段 1 3 は、抽出した項目を、セグメントに対応付けて素材情報記憶手段 1 4 に書き込み記憶する。

なお、項目情報抽出手段 1 3 は、T F - I D F 法を用いずに、セグメントを形態素解析し、名詞や固有名詞を特徴単語として抽出することとしてもよい。

40

【 0 0 3 3 】

また、項目情報抽出手段 1 3 は、素材コンテンツが映像を含んでいる場合、特徴単語以外に、セグメントに対応する時間区間の映像からサムネイル画像を抽出してもよい。例えば、項目情報抽出手段 1 3 は、セグメントに対応する時間区間の映像の先頭フレームをサムネイル画像として抽出する。項目情報抽出手段 1 3 は、抽出したサムネイル画像を、セグメントに対応付けて素材情報記憶手段 1 4 に書き込み記憶する。

【 0 0 3 4 】

素材情報記憶手段 (記憶手段) 1 4 は、音声認識の誤りを修正する対象となる素材コンテンツと、素材コンテンツをセグメントに分割した各種情報とを記憶するものである。こ

50

の素材情報記憶手段 1 4 は、ハードディスク、半導体メモリ等の一般的な記憶媒体で構成することができる。

【 0 0 3 5 】

ここで、図 2 を参照（適宜図 1 参照）して、素材情報記憶手段 1 4 が記憶する素材情報について具体的に説明する。

図 2 に示すように、素材情報記憶手段 1 4 は、音声認識誤りを修正する対象となる素材コンテンツ（映像・音声）A , B ... を記憶する。この素材コンテンツ（映像・音声）A , B ... は、素材コンテンツ入力手段 1 0 によって、記憶されたものである。

【 0 0 3 6 】

また、図 2 に示すように、素材情報記憶手段 1 4 は、素材コンテンツごとに、音声認識結果をセグメントに分割した情報を記憶する。

図 2 の例では、素材コンテンツの識別情報（ここでは、ファイル名 A , B , ... ）ごとに、セグメント（識別情報 a 1 , a 2 , ... , b 1 , ... ）を対応付けている。

各セグメントは、単語 w と時間情報 t とを複数含み、それぞれは対応付けられている。

このセグメントごとの単語 w および時間情報 t は、音声認識手段 1 1 で対応付けられた単語および時間情報を、認識結果分割手段 1 2 が分割した情報である。

【 0 0 3 7 】

また、各セグメントは、項目 k とサムネイル画像 g とを含む。項目 k は、項目情報抽出手段 1 3 が抽出した特徴単語である。サムネイル画像 g は、項目情報抽出手段 1 3 が当該セグメントの先頭の時間情報に対応した、素材コンテンツの映像から抽出したフレーム画像である。

なお、ここでは、素材コンテンツと、素材コンテンツの音声認識結果を分割したセグメントとを、同一の記憶手段に記憶しているが、別々の記憶手段に記憶することとしてもよい。

図 1 に戻って、音声認識誤り修正支援装置 1 の構成について説明を続ける。

【 0 0 3 8 】

編集手段 1 5 は、外部に接続された修正端末（入力装置 2、表示装置 3、スピーカ 4）を用いて、操作者が、素材情報記憶手段 1 4 に記憶されている音声認識結果を修正するものである。なお、修正端末の表示装置 3 は、タッチパネルを備える構成としてもよい。

編集手段 1 5 は、図 1 に示すように、素材コンテンツ選択手段 1 5 0 と、認識結果表示制御手段 1 5 1 と、誤り修正手段 1 5 2 と、映像 / 音声再生手段 1 5 3 と、を備える。

【 0 0 3 9 】

素材コンテンツ選択手段 1 5 0 は、修正対象となる素材コンテンツを選択するものである。例えば、素材コンテンツ選択手段 1 5 0 は、図 3 に示すように、素材情報記憶手段 1 4 に記憶されている素材コンテンツ A , B , C のいずれかを選択するための選択ボタン 3 0 1 を含んだ素材コンテンツ選択画面 3 0 を表示装置 3 に表示する。そして、素材コンテンツ選択手段 1 5 0 は、素材コンテンツ選択画面 3 0 上の選択ボタン 3 0 1 の押下により、修正対象となる素材コンテンツを選択する。素材コンテンツ選択手段 1 5 0 は、選択された素材コンテンツのファイル名等の識別情報を、認識結果表示制御手段 1 5 1 に出力する。

【 0 0 4 0 】

認識結果表示制御手段 1 5 1 は、セグメントごとに、項目と当該セグメントに含まれる単語列を表示するか否かを指定する選択ボタンとを表示し、選択ボタンの押下により、セグメントの単語列を表示するか否かを制御するものである。

【 0 0 4 1 】

ここで、図 4 および図 5 を参照（適宜図 1 参照）して、認識結果表示制御手段 1 5 1 が表示する画面例について、その制御内容とともに説明する。

図 4 に示すように、認識結果表示制御手段 1 5 1 は、項目一覧画面 3 1 を表示装置 3 の画面上に表示する。

項目一覧画面 3 1 は、選択ボタン 3 1 1 と、項目表示欄 3 1 2 と、サムネイル画像表示

10

20

30

40

50

領域 3 1 3 と、タイムテーブル表示欄 3 1 4 と、スクロールバー表示欄 3 1 5 と、で構成される。

【 0 0 4 2 】

選択ボタン 3 1 1 は、セグメントごとに単語列を表示するか否かの選択を行うボタンである。

項目表示欄 3 1 2 は、セグメント内で抽出された項目を表示する領域である。認識結果表示制御手段 1 5 1 は、素材情報記憶手段 1 4 から、当該セグメントに対応する項目（図 2 の項目 k ）を読み出して、項目表示欄 3 1 2 に表示する。

サムネイル画像表示領域 3 1 3 は、セグメント内で抽出されたサムネイル画像を表示する領域である。認識結果表示制御手段 1 5 1 は、素材情報記憶手段 1 4 から、当該セグメントに対応するサムネイル画像（図 2 のサムネイル画像 g ）を読み出して、サムネイル画像表示領域 3 1 3 に表示する。

【 0 0 4 3 】

タイムテーブル表示欄 3 1 4 は、素材コンテンツの時間軸上におけるセグメント位置を示すタイムテーブルを表示する欄である。認識結果表示制御手段 1 5 1 は、素材情報記憶手段 1 4 のセグメントの時間情報（図 2 の時間情報 t ）を参照して、タイムテーブルを生成し表示する。

スクロールバー表示欄 3 1 5 は、項目一覧が画面に収まらない場合に、どの部分のセグメントを表示しているのかを示すスクロールバーを表示する欄である。認識結果表示制御手段 1 5 1 は、スクロールバーの上下によって、画面上の項目一覧を更新する。

このように、項目一覧画面 3 1 を表示することで、操作者は、項目を確認することができ、一度に音声認識結果を表示する場合に比べて、音声認識結果を確認したいセグメントを容易に選択することができる。

【 0 0 4 4 】

この項目一覧画面 3 1 において、操作者が行う入力装置 2 のマウスのクリック、あるいは、表示装置 3 のタッチパネルへのタッチによる選択ボタン（図 4 中、「open」）3 1 1 の押下により、認識結果表示制御手段 1 5 1 は、項目一覧画面 3 1 において、セグメントの単語列の修正を行う編集領域 3 1 6（図 5 参照）を表示する。

【 0 0 4 5 】

図 5 は、編集領域 3 1 6 を表示した項目一覧画面 3 1 B を示す画面例である。

この項目一覧画面 3 1 B は、図 4 で説明した項目一覧画面 3 1 に対して、選択されたセグメントにおいて、動画表示領域 3 1 3 B と、編集領域 3 1 6 とが表示される。

【 0 0 4 6 】

動画表示領域 3 1 3 B は、セグメントに対応する素材コンテンツを再生する領域である。認識結果表示制御手段 1 5 1 は、当該セグメントが選択されたタイミングで、素材情報記憶手段 1 4 のセグメントの時間情報（図 2 の時間情報 t ）を参照して、対応する素材コンテンツの映像の先頭フレームを動画表示領域 3 1 3 B に表示する。この動画表示領域 3 1 3 B の画像領域をマウス等でクリック、あるいは再生開始ボタン s t を押下されることで、認識結果表示制御手段 1 5 1 は、映像 / 音声再生手段 1 5 3 に当該素材コンテンツの再生を指示する。

【 0 0 4 7 】

編集領域 3 1 6 は、セグメントに対応する単語列を表示し、編集対象となる領域である。認識結果表示制御手段 1 5 1 は、編集領域 3 1 6 に、素材情報記憶手段 1 4 に記憶されている当該セグメントに対応する単語列（図 2 の単語 w の列）を展開する。

なお、このとき、認識結果表示制御手段 1 5 1 は、選択ボタン 3 1 1 を、編集領域 3 1 6 を非表示とするボタン（図 4 中、「close」）とする。そして、選択ボタン（図 4 中、「close」）3 1 1 の押下により、認識結果表示制御手段 1 5 1 は、編集領域 3 1 6 を非表示とし、動画表示領域 3 1 3 B をサムネイル画像表示領域 3 1 3 として、図 4 の項目一覧画面 3 1 に表示を戻す。

図 1 に戻って、音声認識誤り修正支援装置 1 の構成について説明を続ける。

【 0 0 4 8 】

誤り修正手段 1 5 2 は、操作者の編集操作により、編集領域 3 1 6 (図 5) において、セグメントの単語列の誤りを修正するものである。この誤り修正手段 1 5 2 は、単語列を修正する編集動作においては、一般的なテキストエディタ (スクリーンエディタ) として機能する。ただし、誤り修正手段 1 5 2 は、単語列を修正する際に、音声を再生する機能を有する。

【 0 0 4 9 】

具体的には、誤り修正手段 1 5 2 は、編集領域 3 1 6 (図 5) において、マウスのクリック、あるいは、タッチパネルへのタッチにより、選択された単語から音声を再生する。また、音声再生中、再度、任意の位置を選択されることで、誤り修正手段 1 5 2 は、音声の再生を停止する。

10

【 0 0 5 0 】

図 6 は、編集領域における編集作業の一例を説明するための説明図である。

例えば、図 6 の編集領域 3 1 6 において、「 3 月 」が選択された場合、誤り修正手段 1 5 2 は、素材情報記憶手段 1 4 のセグメントの時間情報 (図 2 の時間情報 t) を参照して、対応する素材セグメントの位置から音声を再生するように、映像 / 音声再生手段 1 5 3 に指示する。なお、このとき、音声に連動して、動画表示領域 3 1 3 B (図 5) において、音声再生の時間に対応する映像を再生することとしてもよい。

ここで、操作者が、誤り (ここでは、「ハタ寒い」) を発見して修正箇所をマウスでクリック等することで、誤り修正手段 1 5 2 は、音声再生を停止してカーソル C を表示する。そして、誤り修正手段 1 5 2 は、操作者の編集操作により、誤りである「ハタ寒い」を「肌寒い」と修正する。そして、誤り修正手段 1 5 2 は、素材情報記憶手段 1 4 に記憶されている誤りのあった単語を、修正後の単語に置き換える。これによって、音声認識誤り修正支援装置 1 は、操作者による修正後の保存操作を省略することができる。

20

【 0 0 5 1 】

また、誤り修正手段 1 5 2 は、マウスクリック等で指定された単語位置から音声を再生する。

図 7 は、音声再生に連動して編集領域の単語の表示属性を変更する例を説明するための説明図である。例えば、図 7 に示すように、編集領域 3 1 6 において、音声の再生を開始したい箇所をマウス等で選択された場合、誤り修正手段 1 5 2 は、素材情報記憶手段 1 4 のセグメントの時間情報 (図 2 の時間情報 t) を参照し、選択した単語から再生停止の指示があるまで音声を再生するように、映像 / 音声再生手段 1 5 3 に指示する。

30

そして、誤り修正手段 1 5 2 は、図 7 に示すように、音声の再生位置とセグメント中の再生有無とを明示するように、音声の再生に連動して、再生される音声に対応する各単語の表示部分の表示属性を変更する。例えば、誤り修正手段 1 5 2 は、音声に対応する単語を、白黒反転または予め定めた色でカラー表示する。

【 0 0 5 2 】

このとき、誤り修正手段 1 5 2 は、操作者が行った操作のフィードバック情報を画面上に提示する。例えば、図 8 に示すように、誤り修正手段 1 5 2 は、選択された単語位置に音声の再生開始を示すポップアップメッセージ p o p 1 を表示し、音声が停止した単語位置に音声の再生終了を示すポップアップメッセージ p o p 2 を表示する。これによって、操作者が不慣れであっても、自身の操作内容を把握することができ、安心して操作を行うことができる。

40

【 0 0 5 3 】

また、誤り修正手段 1 5 2 は、指定された単語または単語列に対応する音声を繰り返し再生することもできる。

例えば、図 9 に示すように、編集領域 3 1 6 において、音声を再生したい単語または単語列をマウス等で選択 (図中、白黒反転領域) することで、誤り修正手段 1 5 2 は、ポップアップメニュー p m を表示し、「繰り返し再生」を選択されることで、対応する単語または単語列の音声を繰り返し再生する。

50

図 1 に戻って、音声認識誤り修正支援装置 1 の構成について説明を続ける。

【 0 0 5 4 】

映像 / 音声再生手段 1 5 3 は、素材コンテンツの映像および音声を再生するものである。この映像 / 音声再生手段 1 5 3 は、認識結果表示制御手段 1 5 1 または誤り修正手段 1 5 2 から指定された位置から、素材コンテンツ（映像・音声）を再生する。

【 0 0 5 5 】

書き起こし結果出力手段 1 6 は、編集手段 1 5 で修正された音声認識結果（書き起こし結果）を、外部に出力するものである。

この書き起こし結果出力手段 1 6 は、素材コンテンツのファイル名、または、素材コンテンツ内のセグメントの識別番号を指定されることで、素材情報記憶手段 1 4 に記憶されている該当する素材コンテンツまたはセグメントの単語列を読み出して出力する。

【 0 0 5 6 】

以上説明したように音声認識誤り修正支援装置 1 を構成することで、音声認識誤り修正支援装置 1 は、簡易なテキスト編集操作で、認識結果の単語とその元となった音声とを確認しながら、音声認識の誤りを修正することができる。また、音声認識誤り修正支援装置 1 は、素材コンテンツに対して、セグメント単位で部分的に誤り修正を行うことができる。

なお、音声認識誤り修正支援装置 1 は、コンピュータを、前記した各手段として機能させるための音声認識誤り修正支援プログラムで動作させることができる。

【 0 0 5 7 】

[音声認識誤り修正支援装置の動作]

次に、図 1 0 ~ 図 1 2 を参照して、本発明の実施形態に係る音声認識誤り修正支援装置 1 の動作について説明する。なお、ここでは、音声認識誤り修正支援装置 1 の動作として、素材コンテンツに対して音声認識による認識結果をセグメント単位で生成するセグメント情報生成動作と、認識結果をセグメント単位で表示装置 3 に提示するセグメント情報提示動作と、音声再生を行いながら認識結果を修正するセグメント修正動作と、について説明する。

【 0 0 5 8 】

（セグメント情報生成動作）

まず、図 1 0 を参照（適宜図 1 参照）して、音声認識誤り修正支援装置 1 のセグメント情報生成動作について説明する。

ステップ S 1 において、素材コンテンツ入力手段 1 0 は、音声認識を行う素材コンテンツを入力する。このとき、素材コンテンツ入力手段 1 0 は、入力した素材コンテンツを素材情報記憶手段 1 4 に書き込み記憶する。

【 0 0 5 9 】

ステップ S 2 において、音声認識手段 1 1 は、ステップ S 1 で入力した素材コンテンツの音声を認識し、テキストデータである認識結果と当該認識結果を構成する単語ごとの時間情報とを対応付けて生成する。

【 0 0 6 0 】

ステップ S 3 において、認識結果分割手段 1 2 は、ステップ S 2 で認識された認識結果を、予め定めた基準、例えば、映像のカット点、音声の無音区間等によりセグメントに分割する。このとき、認識結果分割手段 1 2 は、セグメント単位で、認識結果の単語と時間情報とを対応付けて、素材コンテンツを素材情報記憶手段 1 4 に書き込み記憶する。

【 0 0 6 1 】

ステップ S 4 において、項目情報抽出手段 1 3 は、ステップ S 3 で分割されたセグメントごとに、セグメントに含まれる特徴単語を項目として抽出するとともに、セグメントに対応する映像からサムネイル画像を抽出する。このとき、項目情報抽出手段 1 3 は、抽出した項目およびサムネイル画像を、セグメントに対応付けて素材情報記憶手段 1 4 に書き込み記憶する。

以上の動作によって、音声認識誤り修正支援装置 1 は、図 2 に示すように、素材情報記

10

20

30

40

50

憶手段 14 に、素材コンテンツと、素材コンテンツをセグメントに分割した各種情報とを記憶する。

【0062】

(セグメント情報提示動作)

次に、図 11 を参照 (適宜図 1 参照) して、音声認識誤り修正支援装置 1 のセグメント情報提示動作について説明する。

ステップ S10 において、素材コンテンツ選択手段 150 は、素材情報記憶手段 14 に記憶されている素材コンテンツのいずれかを選択するための選択ボタンを含んだ素材コンテンツ選択画面 30 (図 3 参照) を表示装置 3 に表示する。

【0063】

ステップ S11 において、素材コンテンツ選択手段 150 は、画面上で選択ボタンが押下されるまで待機し (ステップ S11 で No)、選択ボタンが押下された場合 (ステップ S11 で Yes)、ステップ S12 以降の制御を行う認識結果表示制御手段 151 に制御を移す。

【0064】

ステップ S12 において、認識結果表示制御手段 151 は、素材情報記憶手段 14 に記憶されている各種の情報に基づいて、セグメントごとに、項目と当該セグメントに含まれる単語列を表示するか否かを指定する選択ボタンとを含んだ項目一覧画面 31 (図 4 参照) を表示装置 3 に表示する。

【0065】

ステップ S13 において、認識結果表示制御手段 151 は、項目一覧画面で選択ボタン (open) が押下されるまで待機する (ステップ S13 で No)。

一方、選択ボタン (open) が押下された場合 (ステップ S13 で Yes)、ステップ S14 において、認識結果表示制御手段 151 は、図 5 に示すように、選択されたセグメントに対応して編集領域 316 を表示し、素材情報記憶手段 14 に記憶されている当該セグメントに対応する認識結果である単語列を編集領域 316 に展開する。

【0066】

この動作以降、音声認識誤り修正支援装置 1 は、操作者が画面上で編集結果を修正可能な状態に移行する。なお、選択ボタン (open) の押下により編集領域 316 を表示した場合、認識結果表示制御手段 151 は、任意のタイミングで、選択ボタン (close) の押下により編集領域 316 を非表示とすることができるが、この非表示の動作については図示を省略した。また、項目一覧画面 31B (図 5 参照) の動画表示領域 313B における素材コンテンツの再生動作についてもここでは説明を省略する。

以上の動作によって、音声認識誤り修正支援装置 1 は、素材コンテンツをセグメント単位で、音声認識の誤りを修正することが可能になる。

【0067】

(セグメント修正動作)

次に、図 12 を参照 (適宜図 1 参照) して、音声認識誤り修正支援装置 1 のセグメント修正動作について説明する。なお、セグメント修正動作は、操作者が行う任意の手順であるため、ここでは、音声再生と修正動作とを併せて行う動作の一例で説明する。

【0068】

ステップ S20 において、誤り修正手段 152 は、操作者のマウスのクリック、あるいは、タッチパネルへのタッチにより、編集領域 316 (図 5) 内の音声再生したい単語または単語列を選択する。このとき、誤り修正手段 152 は、映像 / 音声再生手段 153 を介して、素材情報記憶手段 14 のセグメントの時間情報を参照して、単語または単語列に対応する時間の音声再生する。これによって、操作者は、音声と音声認識された単語列とを対比して確認することができる。

【0069】

ステップ S21 において、誤り修正手段 152 は、操作者のマウスのクリック、あるいは、タッチパネルへのタッチにより、修正箇所の位置の指定を受け付ける。このとき、誤

10

20

30

40

50

り修正手段１５２は、音声が生語列の末尾まで再生されていない、あるいは、繰り返し再生中で、音声が生再生中であれば、音声の再生を停止する。

【００７０】

ステップＳ２２において、誤り修正手段１５２は、編集領域の指定された位置にカーソルを表示して、文字削除、文字挿入等の操作者の編集作業により、認識誤りを修正する。ここで、誤り修正手段１５２は、素材情報記憶手段１４の生語を修正結果で更新する。

【００７１】

ステップＳ２３において、誤り修正手段１５２は、操作者のマウスのクリック、あるいは、タッチパネルへのタッチにより、修正を行った箇所の位置の指定を受け付ける。このとき、誤り修正手段１５２は、映像／音声再生手段１５３を介して、素材情報記憶手段１４のセグメントの時間情報を参照して、生語または生語列に対応する時間の音声を生再生する。これによって、操作者は、修正結果が正しいか否かを確認することができる。

10

【００７２】

なお、図示を省略しているが、ステップＳ２３における操作者の確認で、修正箇所がまだ正しく修正されていない場合、ステップＳ２１に戻って、動作を繰り返す。

以上の動作によって、音声認識誤り修正支援装置１は、音声認識の誤りを修正する際に、修正対象箇所の音声を素早く再生し、簡易な操作で音声認識の誤り修正することができる。

【００７３】

以上、本発明の実施形態について説明したが、本発明は、この実施形態に限定されるものではない。

20

ここでは、素材コンテンツを、映像および音声を含んだものとして説明したが、音声のみの素材コンテンツであっても構わない。

その場合、項目情報抽出手段１３は、項目のみを抽出し、サムネイル画像を抽出しないこととすればよい。また、映像／音声再生手段１５３は、音声のみを再生する音声再生手段とすればよい。

【００７４】

また、ここでは、音声認識誤り修正支援装置１に、直接、修正端末（入力装置２、表示装置３、スピーカ４）を接続する構成としたが、これらは、ネットワークを介して接続する形態であっても構わない。

30

【００７５】

また、音声認識誤り修正支援装置１は、修正端末を複数備える構成であっても構わない。その場合、認識結果表示制御手段１５１は、ある修正端末が修正を行っているセグメントについて、他の修正端末が修正対象として選択しないように排他制御し、例えば、他の修正端末において、選択ボタンを表示しないようにする。

【００７６】

また、音声認識誤り修正支援装置１の編集手段１５は、認識結果を修正するサーバとして、画面制御を行うユーザインタフェースを提供し、ネットワークを介して接続された複数の修正端末が、当該ユーザインタフェースを介して動作するクライアントとして機能させることとしてもよい。これによって、ネットワークを介して、複数の地点で、音声認識の誤りを修正することができる。

40

【００７７】

また、音声認識誤り修正支援装置１は、音声認識手段１１を外部に備えてもよい。

例えば、図１３に示す音声認識誤り修正支援装置１Ｂの構成としてもよい。音声認識誤り修正支援装置１Ｂは、音声認識誤り修正支援装置１（図１）の音声認識手段１１を音声認識装置として外部に備える。この場合、認識結果分割手段１２は、音声認識手段１１から出力される音声の認識結果と当該認識結果を構成する生語ごとの時間情報とを、入力インタフェースである認識結果入力手段１７を介して入力すればよい。

なお、音声認識誤り修正支援装置１Ｂも、コンピュータを、前記した各手段として機能させるための音声認識誤り修正支援プログラムで動作させることができる。

50

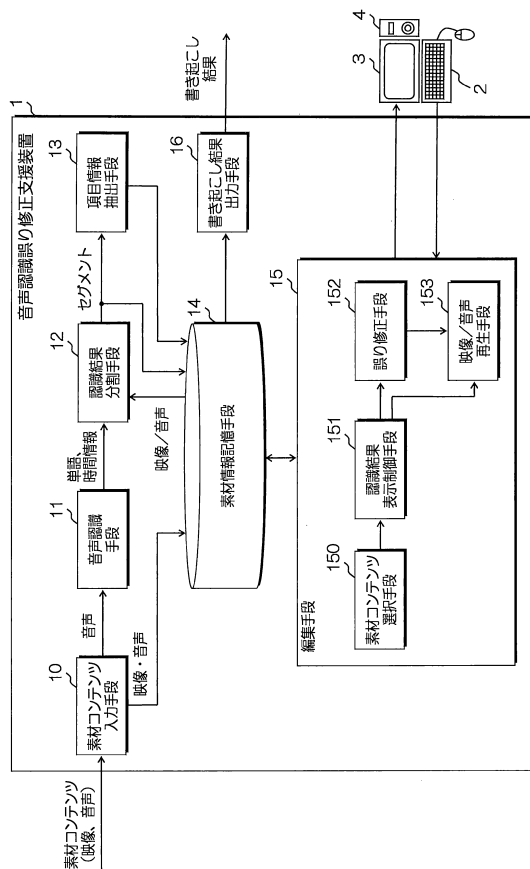
【符号の説明】

【0078】

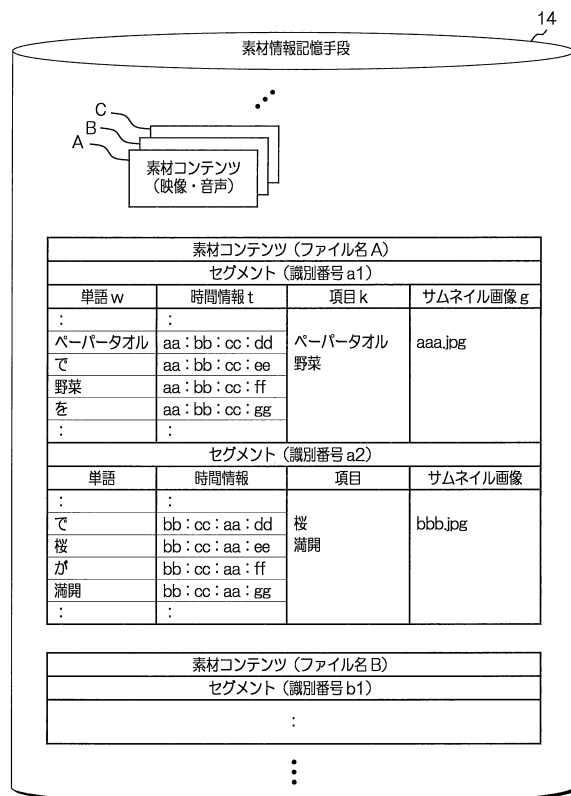
- 1, 1B 音声認識誤り修正支援装置
 10 素材コンテンツ入力手段
 11 音声認識手段
 12 認識結果分割手段
 13 項目情報抽出手段
 14 素材情報記憶手段（記憶手段）
 15 編集手段
 150 素材コンテンツ選択手段
 151 認識結果表示制御手段
 152 誤り修正手段
 153 映像／音声再生手段
 16 書き起こし結果出力手段
 17 認識結果入力手段

10

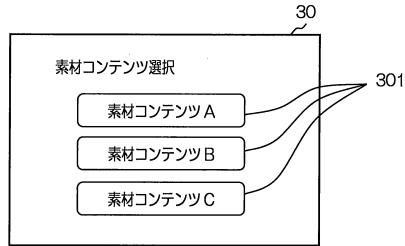
【図1】



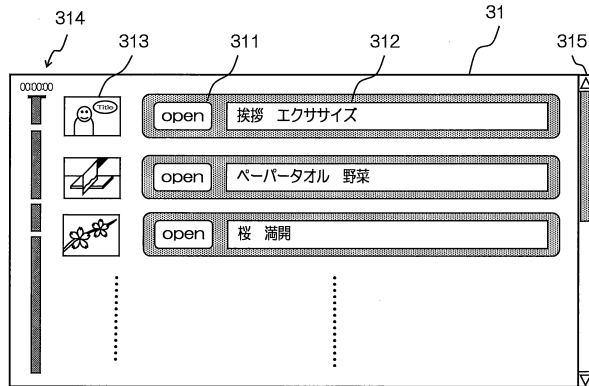
【図2】



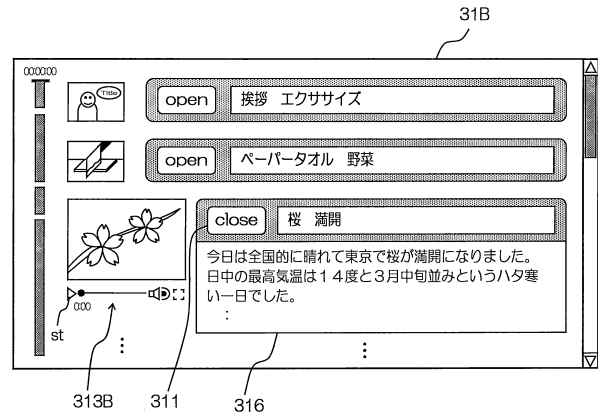
【図 3】



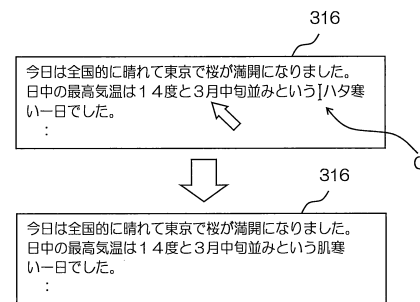
【図 4】



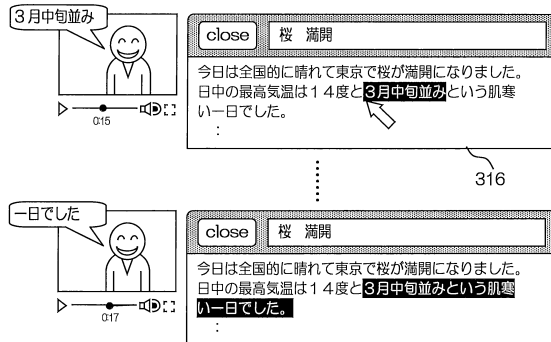
【図 5】



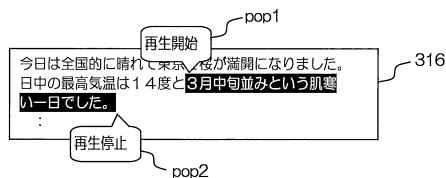
【図 6】



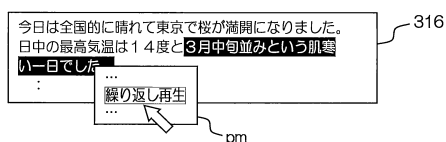
【図 7】



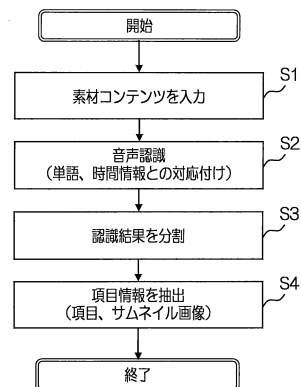
【図 8】



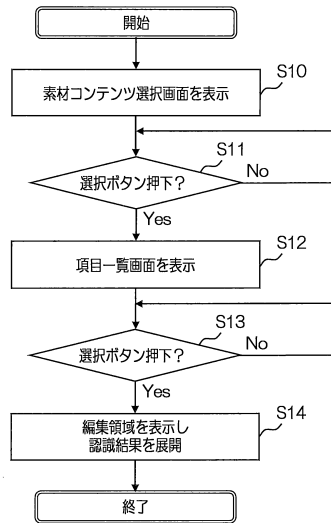
【図 9】



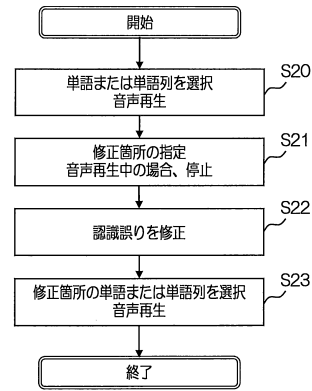
【図 10】



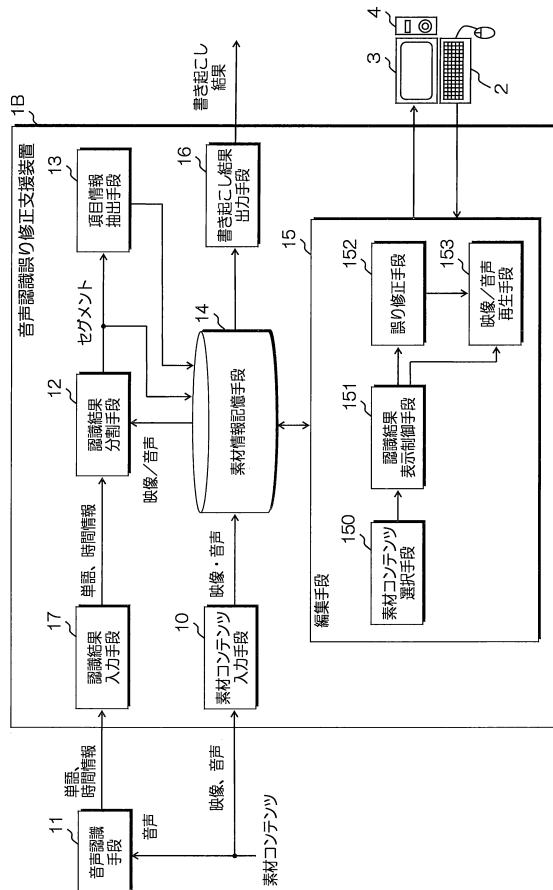
【図 1 1】



【図 1 2】



【図 1 3】



フロントページの続き

- (72)発明者 一木 麻乃
東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内
- (72)発明者 伊藤 均
東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内
- (72)発明者 所澤 愛子
東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内
- (72)発明者 小林 彰夫
東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内

審査官 渡部 幸和

- (56)参考文献 特開2006-330170(JP,A)
特開2013-020411(JP,A)
特開2014-044363(JP,A)
米国特許出願公開第2015/0088505(US,A1)
特開2006-202321(JP,A)
特開2001-282291(JP,A)
荒井 孝, "字幕放送 ニュース字幕用音声認識システムの整備", 放送技術 第65巻
第12号, 西村 瓊江 兼六館出版株式会社, 2012年11月29日, 第65巻, 139~142

- (58)調査した分野(Int.Cl., DB名)
G10L 15/00-19/26