

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3945778号

(P3945778)

(45) 発行日 平成19年7月18日(2007.7.18)

(24) 登録日 平成19年4月20日(2007.4.20)

(51) Int. Cl.			F I		
G 1 0 L	15/22	(2006.01)	G 1 0 L	15/22	4 6 0 Z
G 1 0 L	15/00	(2006.01)	G 1 0 L	15/00	2 0 0 G
G 1 0 L	15/10	(2006.01)	G 1 0 L	15/10	3 0 0 G
G 1 0 L	15/08	(2006.01)	G 1 0 L	15/08	3 0 0 Z

請求項の数 16 (全 27 頁)

(21) 出願番号	特願2004-71229 (P2004-71229)	(73) 特許権者	390009531
(22) 出願日	平成16年3月12日(2004.3.12)		インターナショナル・ビジネス・マシー ズ・コーポレーション
(65) 公開番号	特開2005-258198 (P2005-258198A)		INTERNATIONAL BUSIN ESS MACHINES CORPO RATION
(43) 公開日	平成17年9月22日(2005.9.22)		アメリカ合衆国10504 ニューヨーク 州 アーモンク ニュー オーチャード ロード
審査請求日	平成17年1月14日(2005.1.14)	(74) 代理人	100086243 弁理士 坂口 博
		(74) 代理人	100091568 弁理士 市位 嘉宏
		(74) 代理人	100108501 弁理士 上野 剛史

最終頁に続く

(54) 【発明の名称】 設定装置、プログラム、記録媒体、及び設定方法

(57) 【特許請求の範囲】

【請求項1】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置であって、

前記音声の内容を示す内容データを取得する内容データ取得部と、

再生される前記音声をもとに音声認識した文字データを分割して複数の認識データを生成する音声認識部と、

前記複数の認識データの各々に一致する文字列を前記内容データから検出する文字列検出部と、

前記文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字列検出部と、

前記認識データのうち前記文字列検出部により一致する文字が検出されなかった各文字について、当該文字の読みに含まれる音素に一致する音素を、前記内容データの読みの中から検出する音素検出部とを備え、

前記文字列検出部は、前記認識データのうち前記音素検出部により一致する音素が検出された文字に一致する文字として、前記内容データにおいて当該音素を含む文字を更に検出し、

前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部を更に備え

10

20

る設定装置。

【請求項 2】

前記内容データの読み方の候補である読み候補を複数生成する読付部を更に備え、
前記音素検出部は、前記認識データのうち前記文字検出部により前記内容データに一致する文字が検出されなかった各文字について、当該文字の読みに含まれる音素に一致する音素を、前記読付部により生成された複数の前記読み候補の何れかの中から検出する
請求項 1 記載の設定装置。

【請求項 3】

前記読付部は、前記内容データにおける複数の前記読み候補の各々を、当該読み候補により読まれる可能性を示す情報に対応付けて生成し、
前記音素検出部は、前記認識データに含まれる文字の読みに含まれる音素を、読まれる可能性が高い順に、前記複数の読み候補の各々と比較する
請求項 2 記載の設定装置。

10

【請求項 4】

前記複数の認識データの各々が文字列に一致する確度である信頼度を算出する信頼度算出部を更に備え、
前記表示設定部は、前記内容データにおける複数の前記文字列のうち連続した 2 つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場合に、先に表示すべき前記文字列の末尾に後に表示すべき前記文字列を連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う

20

請求項 1 記載の設定装置。

【請求項 5】

前記信頼度算出部は、前記文字列検出部により一致する文字列が検出された認識データに対応付けて、前記文字検出部により一致する文字列が検出された認識データと比較して、より高い信頼度を算出する
請求項 4 記載の設定装置。

【請求項 6】

前記認識データのうち前記文字検出部により一致する文字が検出されなかった各文字について、当該文字の読みに含まれる音素に一致する音素を、前記内容データの読みの中から検出する音素検出部を更に備え、
前記文字検出部は、前記認識データのうち前記音素検出部により一致する音素が検出された文字に一致する文字として、前記内容データにおいて当該音素を含む文字を検出し、
前記信頼度算出部は、前記音素検出部により一致する音素が検出された文字を含む認識データに対応付けて、前記音素検出部により一致する音素が検出されることなく前記文字検出部により一致する文字が検出された認識データと比較して、より低い信頼度を生成する

30

請求項 4 記載の設定装置。

【請求項 7】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置であって、
前記音声の内容を示す内容データを取得する内容データ取得部と、
再生される前記音声を音声認識した文字データを分割して複数の認識データを生成する音声認識部と、
前記複数の認識データの各々が文字列に一致する確度である信頼度を算出する信頼度算出部と、
前記複数の認識データの各々に一致する文字列を前記内容データから検出すると共に、予め定められた基準信頼度未満の信頼度の認識データである低信頼データについては、さらに、当該低信頼データに後続する認識データに一致する文字列を検出できなかった場合に、当該低信頼データに一致する文字列は検出できないと判断する文字列検出部と、

40

50

前記文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出部と、前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部とを備える設定装置。

【請求項 8】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置であって、

前記音声の内容を示す内容データを取得する内容データ取得部と、

再生される前記音声を音声認識した文字データを分割して複数の認識データを生成すると共に、さらに、音声認識した前記複数の認識データが、再生される音声の内容と一致する可能性を示す音声認識確信度を、認識データ毎に生成する音声認識部と、

前記複数の認識データの各々について、音声認識確信度がより高い認識データに一致する文字列を、当該認識データと比較して音声認識確信度が低い認識データに先立って前記内容データから検出すると共に、第 1 の前記認識データに一致する第 1 の文字列及び第 2 の前記認識データに一致する第 2 の文字列を検出した場合に、前記第 1 の認識データに後続しかつ前記第 2 の認識データに先行する認識データに一致する文字列として、前記第 1 の文字列に後続し前記第 2 の文字列に先行する文字列を検出する文字列検出部と、

前記文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出部と、

前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部とを備える設定装置。

【請求項 9】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置であって、

再生される前記音声の内容を示す内容データに含まれる複数の文字列の各々に対応付けて、当該文字列を表示すべき時刻、及び、当該時刻に当該文字列を内容とする音声の再生される確度である信頼度を取得する信頼度取得部と、

前記複数の文字列のうち連続した 2 つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場合に、先に表示すべき前記文字列の末尾に後に表示すべき前記文字列を連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う表示設定部とを備える設定装置。

【請求項 10】

前記表示設定部は、先に表示すべき前記文字列に対応する信頼度が、後に表示すべき前記文字列に後続する後続文字列に対応する信頼度より高い場合に、連結した前記文字列の末尾に前記後続文字列を更に連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う

請求項 9 記載の設定装置。

【請求項 11】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置として、コンピュータを機能させるプログラムであって、

前記コンピュータを、

前記音声の内容を示す内容データを取得する内容データ取得部と、

再生される前記音声を音声認識した文字データを分割して複数の認識データを生成する音声認識部と、

前記複数の認識データの各々に一致する文字列を前記内容データから検出する文字列検

10

20

30

40

50

出部と、

前記文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出部と、

前記認識データのうち前記文字検出部により一致する文字が検出されなかった各文字について、当該文字の読みに含まれる音素に一致する音素を、前記内容データの読みの中から検出する音素検出部として機能させ、

前記文字検出部は、前記認識データのうち前記音素検出部により一致する音素が検出された文字に一致する文字として、前記内容データにおいて当該音素を含む文字を更に検出し、

10

前記コンピュータを、更に、

前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部として機能させるプログラム。

【請求項 1 2】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置として、コンピュータを機能させるプログラムであって、

前記コンピュータを、

再生される前記音声の内容を示す内容データに含まれる複数の文字列の各々に対応付けて、当該文字列を表示すべき時刻、及び、当該時刻に当該文字列を内容とする音声再生される確度である信頼度を取得する信頼度取得部と、

20

前記複数の文字列のうち連続した2つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場合に、先に表示すべき前記文字列の末尾に後に表示すべき前記文字列を連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う表示設定部として機能させるプログラム。

【請求項 1 3】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定方法であって、

コンピュータにより、

前記音声の内容を示す内容データを取得する内容データ取得段階と、

再生される前記音声を音声認識した文字データを分割して複数の認識データを生成する音声認識段階と、

30

前記複数の認識データの各々に一致する文字列を前記内容データから検出する文字列検出段階と、

前記文字列検出段階において一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出段階と、

前記認識データのうち前記文字検出段階において一致する文字が検出されなかった各文字について、当該文字の読みに含まれる音素に一致する音素を、前記内容データの読みの中から検出する音素検出段階部と、

40

前記認識データのうち前記音素検出段階において一致する音素が検出された文字に一致する文字として、前記内容データにおいて当該音素を含む文字を更に検出する段階と、

前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定段階とを備える設定方法。

【請求項 1 4】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定方法であって、

コンピュータにより、

50

再生される前記音声の内容を示す内容データに含まれる複数の文字列の各々に対応付けて、当該文字列を表示すべき時刻、及び、当該時刻に当該文字列を内容とする音声再生される確度である信頼度を取得する信頼度取得段階と、

前記複数の文字列のうち連続した2つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場合に、先に表示すべき前記文字列の末尾に後に表示すべき前記文字列を連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う表示設定段階とを備える設定方法。

【請求項15】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置として、コンピュータを機能させるプログラムであって、

前記コンピュータを、

前記音声の内容を示す内容データを取得する内容データ取得部と、

再生される前記音声を音声認識した文字データを分割して複数の認識データを生成する音声認識部と、

前記複数の認識データの各々が文字列に一致する確度である信頼度を算出する信頼度算出部と、

前記複数の認識データの各々に一致する文字列を前記内容データから検出すると共に、予め定められた基準信頼度未満の信頼度の認識データである低信頼データについては、さらに、当該低信頼データに後続する認識データに一致する文字列を検出できなかった場合に、当該低信頼データに一致する文字列は検出できないと判断する文字列検出部と、

前記文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出部と、

前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部として機能させるプログラム。

【請求項16】

内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置として、コンピュータを機能させるプログラムであって、

前記コンピュータを、

前記音声の内容を示す内容データを取得する内容データ取得部と、

再生される前記音声を音声認識した文字データを分割して複数の認識データを生成すると共に、さらに、音声認識した前記複数の認識データが、再生される音声の内容と一致する可能性を示す音声認識確信度を、認識データ毎に生成する音声認識部と、

前記複数の認識データの各々について、音声認識確信度がより高い認識データに一致する文字列を、当該認識データと比較して音声認識確信度が低い認識データに先立って前記内容データから検出すると共に、第1の前記認識データに一致する第1の文字列及び第2の前記認識データに一致する第2の文字列を検出した場合に、前記第1の認識データに後続しかつ前記第2の認識データに先行する認識データに一致する文字列として、前記第1の文字列に後続し前記第2の文字列に先行する文字列を検出する文字列検出部と、

前記文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出部と、

前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部と

として機能させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

10

20

30

40

50

本発明は、設定装置、プログラム、記録媒体、及び設定方法に関する。特に本発明は、音声に同期して音声の内容を表示する処理を制御する設定装置、プログラム、記録媒体、及び設定方法に関する。

【背景技術】

【0002】

近年、IT機器の性能が飛躍的に向上し、ブロードバンド等のコンピュータネットワークが整備されるのに伴って、動画などのデジタルコンテンツが盛んに配信されるようになってきている。動画は、文字などの静的なコンテンツと比較して多くの利用者にとって分かり易く、かつ説得力が高い。更に、ケーブルテレビ及びCSテレビ等の普及により、テレビ番組のチャンネル数が増加しており、動画コンテンツは様々な分野において更に広く用いられていくことが予想される。

10

【0003】

動画により提供される情報をより多くの利用者に適切に提供するには、動画に対応付けて音声の内容を示す字幕を表示することが必要である。更に、2007年には、放送等される全ての動画に字幕を付与することが目標として掲げられている。このため、動画に対して適切な字幕を表示する技術の進歩が社会的に要請されている。

【0004】

従来、音声を認識して音声の内容を示す文字列を生成する音声認識技術により、字幕を生成する方法が提案されている。しかしながら、音声認識技術は、音声を誤認識して誤った文字列を生成する場合がある。また、句読点又は記号等は、音声として表現されないもので、音声認識技術によってこれらの記号を適切に表示させることはできない。このため、音声認識技術をそのまま字幕生成に適用することはできず、音声認識結果を修正して字幕を作成していた(特許文献2参照。)

20

【0005】

他の方法として、動画の音声の台本を、適切な長さの文字列毎に分割して、適切なタイミングで表示する方法が提案されている。しかしながら、高機能な動画編集ソフトウェアを用いた場合であっても、手作業で適切なタイミングを決定するのは困難であった。このため、従来、再生される音声と台本とを比較して、台本中の各文字列を表示すべきタイミングを決定する技術が提案されている(特許文献1及び3参照。)

【0006】

非特許文献1については後述する。

【特許文献1】特開平10-254478号公報

【特許文献2】特開2000-89786号公報

【特許文献3】特開平10-136260号公報

【非特許文献1】「テレビドラマのシナリオと音声トラックの自動対応付け」、谷村正剛ら著、自然言語情報処理26-4、1999年5月28日発行

30

【発明の開示】

【発明が解決しようとする課題】

【0007】

特許文献1及び3の技術は、まず、音声を分析することにより、音声所定期間発せられなかった部分を文の切れ目と判断する。そして、音声を分析した結果得られた文の文頭の音素と、台本における各文の文頭に含まれる音素とを比較することにより、音声と台本との対応付けを生成する。これにより、台本中の各文を、その文に一致した音声が発せられる時に表示すべきであることが分かる。

40

【0008】

しかしながら、音声が発せられない部分は、文の切れ目とは限らない。例えば、話者は、迷ったり困ったとき、息継ぎをするとき、一瞬考えたとき、又は強調したいとき等、様々な場合に間をおいて話す。従って、この技術により文の切れ目を適切に判断することは難しく、台本中の各文と一致する音声を検出するのは困難である。音声と字幕の表示タイミングが相違すると、例えば、話者が話し始めているにもかかわらず字幕が表示されない

50

という不都合、又は、話者が話し始めていないのにも関わらずクイズの答えなどが先に表示されてしまう不都合が発生してしまう。

【0009】

更に、この技術は、台本中の各文をそのまま字幕として表示するため、利用者の読みやすさ又は表示装置の画面のサイズ等を考慮して文を分割又は結合することはできない。また、この技術は、音声認識の精度に関わらず同様の字幕を生成するため、音声認識の認識率が将来向上したとしても、字幕の表示精度を向上させることはできない。

【0010】

そこで本発明は、上記の課題を解決することのできる設定装置、プログラム、記録媒体、及び設定方法を提供することを目的とする。この目的は特許請求の範囲における独立項に記載の特徴の組み合わせにより達成される。また従属項は本発明の更なる有利な具体例を規定する。

10

【課題を解決するための手段】

【0011】

上記課題を解決するために、本発明の第1の形態においては、内容が予め定められた音声の再生に同期してその内容を表示する表示タイミングを設定する設定装置であって、音声の内容を示す内容データを取得する内容データ取得部と、再生される音声を音声認識した文字データを分割して複数の認識データを生成する音声認識部と、複数の認識データの各々に一致する文字列を内容データから検出する文字列検出部と、文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を内容データから検出することにより、当該認識データに一致する文字列を内容データから検出する文字検出部と、前記認識データのうち前記文字検出部により一致する文字が検出されなかった各文字について、当該文字の読みに含まれる音素に一致する音素を、前記内容データの読みの中から検出する音素検出部とを備え、前記文字検出部は、前記認識データのうち前記音素検出部により一致する音素が検出された文字に一致する文字として、前記内容データにおいて当該音素を含む文字を更に検出し、内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部を更に備える設定装置、当該設定装置を用いた設定方法、コンピュータを当該設定装置として機能させるプログラム、当該プログラムを記録した記録媒体を提供する。

20

30

【0012】

なお、上記の発明の概要は、本発明の必要な特徴の全てを列挙したものではなく、これらの特徴群のサブコンビネーションもまた、発明となりうる。

【発明の効果】

【0013】

本発明によれば、音声の再生に同期して、その音声の内容を示す文字列を表示することができる。

【発明を実施するための最良の形態】

【0014】

以下、発明の実施の形態を通じて本発明を説明するが、以下の実施形態は特許請求の範囲にかかる発明を限定するものではなく、また実施形態の中で説明されている特徴の組み合わせの全てが発明の解決手段に必須であるとは限らない。例えば、明快な方法として全体的にまたは部分的に重み付けなどを利用して幾つかの処理を一括処理させるような手段も想定しうる事はいうまでもない。

40

【0015】

図1は、字幕設定装置10の機能ブロック図である。字幕設定装置10は、内容が予め定められた音声を含む動画の再生に同期して、その内容を字幕として表示するタイミングを字幕行表示装置15に設定することを目的とする。字幕設定装置10は、内容データ取得部100と、音声認識部110と、文字列検出部120と、文字検出部130と、読付部140と、音素検出部150と、信頼度算出部160と、信頼度取得部165と、表示

50

設定部 170 とを備える。

【0016】

内容データ取得部 100 は、音声の内容を示す内容データ 20 を取得する。音声認識部 110 は、話者により話された音声をマイクなどで取得することにより、再生される音声を音声認識する。音声認識には、既存の多様な技術を適用可能である。そして、音声認識部 110 は、音声認識した文字データを分割して複数の認識データを生成する。例えば、生成した複数の認識データを、認識データ 30 - 1 ~ N とする。ここで、認識データ 30 - 1 ~ N の各々は、音声認識された時刻に対応付けられていることが望ましい。音声認識部 110 は、生成した認識データ 30 - 1 ~ N を文字列検出部 120 に送る。

【0017】

文字列検出部 120 は、認識データ 30 - 1 ~ N の各々に一致する文字列を内容データ 20 から検出し、検出結果を信頼度算出部 160 に送る。文字列検出部 120 は、文字検出部 130 から受け取った検出結果に基づいて一致する文字列を検出してもよい。文字検出部 130 は、文字列検出部 120 により一致する文字列が検出されなかった各認識データについて、その認識データに含まれる各文字に一致する文字を内容データ 20 から検出することにより、その認識データに一致する文字列を内容データ 20 から検出する。また、文字検出部 130 は、音素検出部 150 により一致する音素が検出された文字に一致する文字として、内容データ 20 においてその音素を含む文字を検出してもよい。そして、文字検出部 130 は、検出結果を文字列検出部 120 に送る。

【0018】

読付部 140 は、内容データ 20 の読み方の候補である読み候補を複数生成する。更に、読付部 140 は、これらの読み候補の各々を、その読み候補により読まれる可能性を示す情報に対応付けて生成してもよい。音素検出部 150 は、認識データ 30 - 1 ~ N のうち文字検出部 130 により一致する文字が検出されなかった各文字を、その文字の読みを含む音素に展開する。例えば、音素検出部 150 は、漢字をその漢字の読みを示す平仮名に変換してもよい。

【0019】

そして、音素検出部 150 は、認識データ 30 - 1 ~ N のうち文字検出部 130 により一致する文字が検出されなかった各文字について、その文字の読みに含まれる音素に一致する音素を、読付部 140 により生成された複数の読み候補の何れかの中から検出する。更に、音素検出部 150 は、一致する音素を検出できなかった文字については、その文字に含まれる音素が発せられる時間の長さ等に基づいて、一致する文字を検出してもよい。そして、音素検出部 150 は、検出結果を文字検出部 130 に送る。

【0020】

信頼度算出部 160 は、認識データ 30 - 1 ~ N の各々が文字列に一致する確度である信頼度を、認識データ毎に算出する。ここで、信頼度とは、各認識データが音声認識された時刻に、その認識データに一致する文字列を内容とする音声再生される確度をいう。例えば、信頼度算出部 160 は、文字列検出部 120 のみにより一致する文字列が検出された認識データに対応付けて、文字列検出部 120 及び文字検出部 130 により一致する文字列が検出された認識データと比較して、より高い信頼度を算出してもよい。そして、信頼度算出部 160 は、内容データ 20 の各文字列に信頼度を対応付けて信頼度取得部 165 に送る。

【0021】

信頼度取得部 165 は、各文字列に対応付けて、その文字列を表示すべき時刻、即ちその文字列に一致する認識データとして音声認識された音声の再生時刻を、文字列検出部 120 を介して音声認識部 110 から取得する。更に、信頼度取得部 165 は、その再生時刻にその文字列を内容とする音声再生される確度である信頼度を、信頼度算出部 160 から取得する。

【0022】

表示設定部 170 は、内容データ 20 に含まれる文字列の各々を表示させる表示タイミ

10

20

30

40

50

ングを、その文字列に一致する認識データとして音声認識された音声の再生時に設定する。例えば、表示設定部 170 は、文字列を表示すべき時刻をその文字列に対応付けた表示タイミング情報 40 を、字幕行表示装置 15 に出力してもよいし、字幕行表示装置 15 から参照可能なデータベース等に格納してもよい。更に、表示設定部 170 は、表示すべき字幕の設定情報を字幕行表示装置 15 から取得し、取得したその設定情報に基づいて文字列を連結してもよい。設定情報とは、例えば、字幕行表示装置 15 の表示部において 1 行に表示可能な文字数、又は、句点若しくは読点において字幕を改行するか否かを示す情報をいう。

【0023】

図 2 は、内容データ 20 の一例を示す。内容データ 20 は、音声の内容として、漢字、仮名、アルファベット、及び句読点により構成される文字列を含む。内容データ 20 を用いることにより、音声認識の結果をそのまま表示する技術とは異なり、発音されない記号等も含めて適切な字幕を生成できる。以降、内容データ 20 が、「近年 IT を取り巻く環境は著しく変化いたしております。ここで流れをおさらいします。」という文章である場合について説明する。

10

【0024】

図 3 は、認識データ 30 - 1 ~ N の一例を示す。音声認識部 110 は、認識データ 30 - 1 ~ N の各々を、その認識データが認識された認識時刻に対応付けて生成する。ここで、認識時刻とは、例えば、音声の再生開始時から、その認識データとして認識される部分の再生時点までに経過した時間である。一例として、音声認識部 110 は、認識データ 30 - 1 として「金ねうん」の文字列を生成し、認識データ 30 - 1 に対応付けて「02.103」の認識時刻を生成する。

20

【0025】

図 3 を図 2 と比較することにより、音声認識部 110 が、音声認識の結果、内容データ 20 とは異なる誤った文字列を生成していることが分かる。本実施例における字幕設定装置 10 は、音声認識の結果がこのように誤っている場合であっても、適切な表示タイミングを設定することができる。

【0026】

図 4 は、表示タイミング情報 40 の一例を示す。表示設定部 170 は、表示タイミング情報 40 として、表示すべき文字列に、その文字列を表示すべき表示時刻に対応付けた情報を生成し、字幕行表示装置 15 に設定する。ここで、表示時刻とは、音声の再生開始時から文字列を表示すべき時刻までの時間である。一例として、表示設定部 170 は、「近年 IT を取り巻く環境は」の文字列を、表示時刻の一例であるタイムスタンプ「02.103」に対応付けた情報を生成する。

30

【0027】

ここで、「近年 IT を取り巻く環境は」の文字列は、認識データ 30 - 1 ~ 4 の組に対応する。このように、表示設定部 170 は、表示すべき時刻を設定するのみならず、必要に応じて文字列を連結する処理を行う。

【0028】

図 5 は、字幕設定装置 10 が表示タイミングを設定する処理の動作フローを示す。内容データ取得部 100 は、音声の内容を示す内容データ 20 を取得する (S500)。ここで、内容データ取得部 100 は、取得すべき内容データを利用者からの指示に基づいて選択してもよい。音声認識部 110 は、再生される音声を音声認識し (S510)、音声認識した文字データを分割して認識データ 30 - 1 ~ N を生成する (S515)。例えば、音声認識部 110 は、予め定められた単語数、例えば 1 から 3 単語毎に分割してもよいし、予め定められた再生時間、例えば 1.5 秒毎に分割してもよい。

40

【0029】

これに代えて、音声認識部 110 は、無音状態が所定の長さ、例えば 100 ミリ秒以上継続した部分を、認識データの境界と判断してもよい。好ましくは、音声認識部 110 は、文字データを、できるだけ短い意味のまとまり、例えば文より短い単語、句、又は節等

50

の文字列毎に分割する。これにより、タイムスタンプの生成精度を高めることができる。音声認識部 110 によるこの分割の処理を、細分音声認識行処理と呼ぶ。

【0030】

また、好ましくは、音声認識部 110 は、音声認識した認識データ 30 - 1 ~ N が、再生される音声の内容と一致する可能性を示す音声認識確信度を、音声認識データ毎に更に生成する。更に、音声認識部 110 は、音声認識した複数の文字データの中から、分割すべき文字データを、利用者からの指示に基づいて選択してもよい。

【0031】

続いて、文字列検出部 120 は、認識データ 30 - 1 ~ N の各々に一致する文字列を内容データ 20 から検出する (S520)。この処理の中で、必要に応じて、文字検出部 130 は、認識データに含まれる各文字に一致する文字を内容データ 20 から検出してもよい。また、認識データ 30 - 1 ~ N のうち文字検出部 130 により一致する文字が検出されなかった各文字について、その文字の読みに含まれる音素に一致する音素を内容データ 20 から検出してもよい。詳細は後述する。

10

【0032】

表示設定部 170 は、文字列検出部 120 により何れの文字列にも一致しない認識データがあるか否か判断する (S530)。そのような認識データがある場合に (S530: YES)、表示設定部 170 は、文字列検出部 120 により内容データ 20 において一致する文字列が検出されなかったその認識データを、その認識データとして音声認識された音声の再生時に表示させるべく、表示すべき文字列に追加する (S540)。

20

【0033】

更に、表示設定部 170 は、表示すべき字幕の設定情報を字幕行表示装置 15 から取得し、取得したその設定情報に基づいて文字列を連結する (S550)。表示設定部 170 は、この連結処理を、字幕表示に先立って予め行ってもよいし、字幕を表示すべき指示を字幕行表示装置 15 から受けた場合に動的に行ってもよい。これにより、表示設定部 170 は、字幕行表示装置 15 の表示部のサイズ等に応じて、適切な字幕を生成することができる。以降、表示設定部 170 によるこの連結処理を、最適字幕行表示処理と呼ぶ。そして、表示設定部 170 は、内容データ 20 に含まれる文字列の各々を表示させる表示タイミングを、その文字列に一致する認識データとして音声認識された音声の再生時に設定する (S560)。

30

【0034】

図 6 は、S520 における処理の詳細を示す。文字列検出部 120 は、まず、認識データ 30 - 1 ~ N の各々に含まれる各文字に一致する文字を内容データ 20 から検出する (S600)。従来、この検出処理として効率的な方法である DP マッチングが、提案されている (非特許文献 1 参照。)。文字列検出部 120 は、例えば、DP マッチングによりこの検出処理を行ってもよいし、他の手法によりこの検出処理を行ってもよい。文字列検出部 120 は、検出結果として、内容データ 20 中の文字を順次縦軸に配列し、認識データ 30 - 1 ~ N を順次横軸に配置した座標軸において、内容データ 20 中の各文字と認識データ 30 - 1 ~ N 中の各文字が一致する点を順次経路するグラフである最適マッチング路を生成してもよい。

40

【0035】

そして、文字列検出部 120 は、認識データ 30 - 1 ~ N の各々に一致する文字列を再度検出するべく、以下の処理を各認識データについて繰り返す (S610)。本実施例において、文字列検出部 120 は、複数の認識データに対して、認識された順に、即ち音声として再生された順に、以下の処理を行う。これに代えて、文字列検出部 120 は、音声認識確信度が高い順に、即ち、音声認識確信度がより高い認識データに対して、その認識データと比較して音声認識確信度が低い認識データに先立って、以下の処理を行ってもよい。

【0036】

文字列検出部 120 は、その認識データに一致する文字列を内容データ 20 から検出す

50

る (S 6 2 0)。続いて、文字列検出部 1 2 0 は、一致する文字列を検出できなかった認識データについて、その認識データに含まれる文字と一致する文字を検出する処理が必要か否か判断する (S 6 3 0)。処理が必要な場合に (S 6 3 0 : Y E S)、文字検出部 1 3 0 は、その認識データに含まれる文字に一致する文字を検出する (S 6 4 0)。なお、このように、一致する文字列が検出できなかったデータあるいは単純に一括した文字列に対して、そのデータに含まれる部分データについて一致する文字を順次再帰的にスケールレベルを遷移させながら検出する処理を、本実施例においては、アップスケーリング処理と呼ぶ。

【 0 0 3 7 】

続いて、信頼度算出部 1 6 0 は、認識データ 3 0 - 1 ~ N の各々が文字列に一致する確度である信頼度を、認識データ毎に算出する (S 6 5 0)。例えば、信頼度算出部 1 6 0 は、文字列検出部 1 2 0 のみにより一致する文字列が検出された認識データに対応付けて、文字列検出部 1 2 0 及び文字検出部 1 3 0 により一致する文字列が検出された認識データと比較して、より高い信頼度を算出してよい。

10

【 0 0 3 8 】

また、信頼度算出部 1 6 0 は、音素検出部 1 5 0 により一致する音素が検出された文字を含む認識データに対応付けて、音素検出部 1 5 0 により一致する音素が検出されることなく文字検出部 1 3 0 により一致する文字が検出された認識データと比較して、より低い信頼度を生成する。即ち、アップスケーリング処理の段階が増加するのに応じてより低い信頼度を生成する。これにより、異なる文字で音素が偶然一致したようなエラーを含み得る認識データに対しては、より低い信頼度を生成することができる。

20

【 0 0 3 9 】

字幕設定装置 1 0 は、以上の処理を各認識データについて繰り返す (S 6 6 0)。

本図に示すように、文字列検出部 1 2 0 は、各認識データに一致する文字列を内容データ 2 0 から検出する処理に先立って、検出精度の高い D P マッチング等により、文字単位の一一致を判断する。そして、文字列検出部 1 2 0 は、D P マッチングにより既に一致する文字が検出された認識データについて、その認識データと一致する文字列を内容データ 2 0 から再度検出する。これにより、一致する文字を検出する精度を高められると共に、その一致が認識データ単位で判断されるものであるか、又は文字単位で判断されるものであるか判断できる。これに代えて、文字列検出部 1 2 0 は、D P マッチング等の文字単位の一一致を判断しなくともよい。

30

【 0 0 4 0 】

図 7 は、S 6 2 0 における処理の詳細を示す。文字列検出部 1 2 0 は、まず、検出対象の認識データに一致する文字列を内容データ 2 0 から検出する。更に、文字列検出部 1 2 0 は、検出対象の認識データの検出結果のみならず、その前に検出する対象であった認識データ、及び更にその前に検出する対象であった認識データの検出結果に基づいて、更に以下の処理を行う。

【 0 0 4 1 】

本図における丸印は、認識データに一致する文字列が検出されたことを示す。一方、X 印は、認識データに一致する文字列が検出されていないことを示す。例えば文字列検出部 1 2 0 は、1 つ前の認識データ及び対象の認識データの各々に一致する文字列を検出した場合には、2 つ前の認識データの検出結果に関わらず、対象の認識データに一致する文字列を検出したと判断する。

40

【 0 0 4 2 】

文字列検出部 1 2 0 は、1 つ前の認識データに一致する文字列を検出し、かつ対象の認識データに一致する文字列を検出できなかった場合には、2 つ前の認識データの検出結果に関わらず、対象の認識データについての処理を保留し、S 6 2 0 の処理を終え、次の認識データについての処理に移る。

【 0 0 4 3 】

1 つ前の認識データに一致する文字列が検出できない場合において、文字列検出部 1 2

50

0 は、以下の処理を行う。

文字列検出部 120 は、2 つ前の認識データ及び対象の認識データの各々に一致する文字列を検出した場合には、当該 1 つ前の認識データが、2 つ前の認識データ及び対象の認識データの各々に一致する各文字列の間の文字列に一致すると判断する。

【0044】

文字列検出部 120 は、2 つ前の認識データに一致する文字列を検出し、対象の認識データに一致する文字列を検出できない場合に、対象の認識データについての処理を保留し、S620 の処理を終え、次の認識データについての処理に移る。但し、対象の認識データが、検出対象の最後の認識データである場合には、文字列検出部 120 は、1 つ前の認識データ及び対象の認識データを連結したデータを対象として、文字列検出部 130 により一致する文字を検出させる。即ちこの場合、S630 において、文字列検出部 120 は、一致する文字を検出する処理が必要と判断する。

10

【0045】

文字列検出部 120 は、2 つ前の認識データに一致する文字列を検出できず、対象の認識データに一致する文字列を検出できた場合に、2 つ前の認識データ及び 1 つ前の認識データを連結したデータを対象として、文字列検出部 130 により一致する文字を検出させる。文字列検出部 130 の検出結果に基づく認識データの信頼度が、予め定められた基準信頼度未満であれば、文字列検出部 120 は、その認識データに対応付けて、その認識データが信頼度の低い旨を示す低信頼データである旨の再評価フラグを付す。本図においてはこのフラグを三角印で表す。

20

【0046】

文字列検出部 120 は、2 つ前、1 つ前、及び対象の認識データの何れにも一致する文字列を検出できなかった場合に、1 つ前、2 つ前、及び更にその前の認識データを結合したデータを対象として、文字列検出部 130 により一致する文字を検出させる。

【0047】

一方、文字列検出部 120 は、2 つ前の認識データに一致する文字列を検出できず、1 つ前の認識データに再評価フラグが付されている場合においては、以下の処理を行う。

文字列検出部 120 は、対象の認識データに一致する文字列を検出した場合に、1 つ前の認識データから再評価フラグを取り除くことにより、当該 1 つ前の認識データに一致する文字列を検出したと判断する。即ち三角印を丸印に変更する。一方、対象の認識データに一致する文字列を検出できなかった場合に、文字列検出部 120 は、1 つ前の認識データから再評価フラグを取り除くことにより、当該 1 つ前の認識データに一致する文字列を検出できないと判断する。即ち、三角印をバツ印に変更する。

30

【0048】

このように、文字列検出部 120 は、複数の認識データの各々について、その認識データに完全に一致する文字列のみならず、その認識データの前後の認識データが一致した場合に、その認識データについても一致したと判断する。より正確には、文字列検出部 120 は、第 1 の認識データに一致する第 1 の文字列及び第 2 の認識データに一致する第 2 の文字列を検出した場合に、第 1 の認識データに後続しかつ第 2 の認識データに先行する認識データに一致する文字列として、第 1 の文字列に後続し第 2 の文字列に先行する文字列を検出する。即ち、一致とは、完全一致のみならず、前後の認識データが一致したことに基づくこの一致を含む。以降、この一致を、挟み打ち処理による一致と呼ぶ。

40

これにより、文字又は音素単位の一一致を検出する処理をできるだけ減少させて、処理の効率を高めることができる。更に、文字単位の一一致を検出する必要がある場合であっても、検出範囲を限定することができるので、効率がよい。

【0049】

図 8 は、S640 における処理の詳細を示す。文字列検出部 130 は、検出対象の認識データに含まれる各文字について、以下の処理を繰り返す (S800)。まず、文字列検出部 130 は、その文字に一致する文字を内容データ 20 から検出する (S810)。そして、文字列検出部 130 は、一致する文字を検出できなかった認識データについて、その認識

50

データに含まれる文字に含まれる音素と一致する音素を検出する処理が必要か否か判断する(S 8 2 0)。

【 0 0 5 0 】

処理が必要な場合に(S 8 2 0 : Y E S)、音素検出部 1 5 0 は、文字検出部 1 3 0 により一致する文字が検出されなかった各文字について、その文字の読みに含まれる音素に一致する音素を、内容データ 2 0 の読みの中から検出する(S 8 3 0)。文字検出部 1 3 0 は、以上の処理を各文字について繰り返す(S 8 4 0)。

【 0 0 5 1 】

図 9 は、S 8 1 0 における第 1 の処理の詳細を示す。文字検出部 1 3 0 は、検出対象の文字が認識データの末尾の文字でない場合に、本図の処理を行う。まず、文字検出部 1 3 0 は、検出対象の文字に一致する文字を内容データ 2 0 から検出する。更に、文字検出部 1 3 0 は、検出対象の文字の検出結果のみならず、その前に検出する対象であった文字、及び、検出対象の認識データ(認識データの組を含む)の先頭の文字の検出結果に基づいて、更に以下の処理を行う。

【 0 0 5 2 】

文字検出部 1 3 0 は、1 つ前の文字に一致する文字を検出し、かつ対象の文字に一致する文字を検出した場合には、S 8 1 0 における処理を終了し、次の文字に対する検出処理に移る。一方、文字検出部 1 3 0 は、1 つ前の文字に一致する文字を検出し、かつ対象の文字に一致する文字を検出できなかった場合には、対象の文字についての検出処理を保留して、次の文字についての処理に移る。

【 0 0 5 3 】

文字検出部 1 3 0 は、1 つ前の文字に一致する文字を検出できなかった場合においては、他の条件に応じて以下の処理を行う。

文字検出部 1 3 0 は、先頭の文字及び対象の文字の各々に一致する文字を検出した場合に、文字の一致を検出すべき対象の認識データ全体が、内容データ 2 0 の文字列に一致したと判断する。即ち、文字検出部 1 3 0 は、同一の認識データ内の文字については、複数の文字についても挟み打ち処理による一致の判断を行う。このように、文字の一致とは、文字の完全一致のみならず、前後の文字が一致したことに基づく一致を含む。

【 0 0 5 4 】

一方、文字検出部 1 3 0 は、対象の文字に一致する文字を検出できなかった場合には、対象の文字についての検出処理を保留して、次の文字についての処理に移る。文字検出部 1 3 0 は、先頭の文字に一致する文字を検出できず、対象の文字に一致する文字を検出できた場合には、S 8 1 0 における処理を終了し、次の文字に対する検出処理に移る。

【 0 0 5 5 】

図 1 0 は、S 8 1 0 における第 2 の処理の詳細を示す。文字検出部 1 3 0 は、検出対象の文字が認識データの末尾の文字である場合に、本図の処理を行う。具体的には、文字検出部 1 3 0 は、検出対象の認識データ(認識データの組を含む)の先頭の文字の検出結果、及び、検出対象の次の認識データの文字列検出部 1 2 0 による検出結果に基づいて、以下の処理を行う。

【 0 0 5 6 】

文字検出部 1 3 0 は、先頭の文字に一致する文字を検出し、かつ次の認識データに一致する文字列が検出されている場合に、対象の認識データ全体を一致と判断する。一方、文字検出部 1 3 0 は、先頭の文字に一致する文字を検出し、かつ次の認識データに一致する文字列が検出されていない場合に、先頭の文字に後続する文字列の音素を検出対象として、音素検出部 1 5 0 により一致する音素を検出させる。

【 0 0 5 7 】

文字検出部 1 3 0 は、先頭の文字に一致する文字を検出せず、かつ次の認識データに一致する文字列が検出されている場合に、次の認識データに先行する文字列の音素を検出対象として、音素検出部 1 5 0 により一致する音素を検出させる。一方、文字検出部 1 3 0 は、先頭の文字に一致する文字を検出せず、かつ次の認識データに一致する文字列が検出

10

20

30

40

50

されていない場合に、文字を検出する対象の認識データ全体を検出対象として、音素検出部150により一致する音素を検出させる。

【0058】

図11は、S830における第1の処理の詳細を示す。音素検出部150は、検出対象の音素が認識データの末尾の文字でない場合に、本図の処理を行う。まず、音素検出部150は、検出対象の音素に一致する音素を内容データ20の所定の読み候補から検出する。音素検出部150は、一致する音素が検出できなかった場合には、検出対象の音素を読まれる可能性が高い順に複数の読み候補の各々と比較するべく、次に可能性の高い読み候補と比較する。何れの読み候補にも一致しない場合には、音素検出部150は、その音素についての処理を保留して、次の音素の処理に移る。

10

【0059】

続いて、検出対象の音素に一致する音素を検出した場合には、音素検出部150は、検出対象の認識データ(認識データの組を含む)の先頭の文字の検出結果、及び、検出対象の1つ前の音素の検出結果に基づいて、以下の処理を行う。

先頭文字に一致する文字が検出されている場合に、文字検出部130は、検出対象の認識データ内において一致する文字が検出されていない各文字について、一致する文字が検出されたと判断する。このように、音素の一致に基づいて、文字についての挟み打ち処理による一致を判断してもよい。一方、音素検出部150は、対象の音素に一致する音素を検出したその他の場合には、対象の音素についての処理を終了して、次の音素についての処理に移る。

20

【0060】

図12は、S830における第2の処理の詳細を示す。音素検出部150は、検出対象の音素が認識データの末尾の文字である場合に、本図の処理を行う。具体的には、音素検出部150は、先頭文字に一致する文字の検出結果と、次の認識データに一致する文字列の検出結果、又は、対象の認識データが最後の認識データ(例えば認識データ30-N)であるか否かとに基づいて、以下の処理を行う。

【0061】

まず、次の認識データに一致する文字列が検出されていない場合と、検出対象の認識データが最後の認識データでない場合とにおいて、音素検出部150は、S830の処理、即ち一致する音素を検出する処理を終了する。この結果、文字列検出部120は、次の認識データについての処理に移る。一方、次の認識データに一致する文字列が検出されている場合、又は、検出対象の認識データが最後の認識データである場合においては、以下の処理を行う。

30

【0062】

音素検出部150は、先頭文字に一致する文字が検出されている場合には、検出対象の認識データ内の不一致文字を一致するものと判断する。一方、音素検出部150は、先頭文字に一致する文字が検出されていない場合には、音声の内容に関わらず音声の長さ又は文字の長さに基づいて一致する音素を検出する強制割り振り処理を行う。信頼度算出部160は、この強制割り振り処理により一致する音素が検出された認識データに対応付けて、この強制割り振り処理によらず一致する文字が検出された認識データと比較して、更に低い信頼度を算出する。

40

【0063】

図13は、S550における処理の詳細を示す。表示設定部170は、利用者から入力された、表示すべき字幕の設定情報を字幕行表示装置15等から取得する(S1310)。例えば、表示設定部170は、設定情報を示すコマンド等をパースすることにより、設定情報の内容を解析する(S1320)。設定情報とは、字幕行表示装置15の表示部の1行に表示させる文字数であってもよいし、字幕を句読点で改行するか否かの指示であってもよい。

【0064】

信頼度取得部165は、内容データ20に含まれる複数の文字列の各々に対応付けて、

50

その文字列を表示すべき時刻、及びその時刻にその文字列を内容とする音声再生される確度である信頼度を、信頼度算出部160から取得する(S1325)表示設定部170は、各認識データと一致する内容データ20内の各文字列について、その文字列が、設定情報の条件を満たすか否か判断する(S1330)。

【0065】

満たしていない場合に(S1330:NO)。表示設定部170は、複数の文字列を連結する(S1340)。具体的には、表示設定部170は、複数の文字列のうち連続した2つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場合に、先に表示すべきその文字列の末尾に後に表示すべきその文字列を連結した文字列を、先に表示すべき文字列を表示すべき時刻に表示させる設定を行う。

10

【0066】

そして、表示設定部170は、S1330に処理を戻して判断を繰り返す。この結果、表示設定部170は、設定情報の条件を満たすまで、文字列の連結を繰り返す。例えば、表示設定部170は、S1340の処理を行う直前において先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に後続する後続文字列に対応する信頼度より高い場合に、S1340において連結した文字列の末尾にその後続文字列を更に連結した文字列を、先に表示すべきその文字列を表示すべき時刻に表示させる設定を行ってもよい。

このように、表示設定部170は、表示デバイスの機能・特徴に応じて文字列を連結することにより、最適なユーザビリティ(可用性)を利用者に提供することができる。

20

【0067】

以上、図1から図13において説明したように、字幕設定装置10は、音声認識した文字データに対して表示タイミングを設定する場合には、細分認識行処理により文字データを1から3単語程度の長さの認識データ毎に分割して、内容データ20内の文字列との一致を判断する。これに対して、字幕行を表示する場合には、表示部等の特徴に基づいてこれらの認識データを適切に連結する。即ち、表示タイミングの設定と、字幕行生成とでは、異なるサイズのデータを処理対象とする。これにより、双方の処理を効率的に行うことができる。また、字幕設定装置10は、音声認識の結果を用いて表示タイミングの設定を行うので、音声認識技術の進歩に伴い、表示タイミングの設定精度を向上させることができる。

30

【0068】

本実施例における字幕設定装置10により行った実験結果を次に示す。本実験において、字幕設定装置10は、アドリブに基づく32行分の音声と、台本の定められた86行分の音声とを入力とする。DPマッチングによって、全ての行の一致が判断された結果、そのうち12%の文においてタイムスタンプに誤りが生じている。そして、文字列検出部120により66行分の文字列が検出され、文字検出部130により36行分の文字が検出された結果、タイムスタンプの誤りは一切生じていない。音素検出部150により6行分の文字の音素が検出された結果、2%の文字においてタイムスタンプに誤りが生じている。更に、音素検出部150により強制割り振りが処理された結果、4%の文字においてタイムスタンプに誤りが生じている。このように、従来効率が高いアルゴリズムとして知られているDPマッチングと比較して、より高い精度で表示タイミングを設定することができる。

40

【0069】

また、この実験において、表示設定部170は、文字検出部130により検出した文字を含む文字列のうち2つを、他の文字列に連結して表示させ、音素検出部150により音素が検出された文字を含む文字列のうち3つを、他の文字列に連結して表示させた。更に、強制割り振りにより検出された文字を含む文字列のうち10の文字列を、他の文字列に連結して表示させた。このように、アップスケーリング処理の段階が進むのに応じて低い信頼度を生成することにより、誤っている可能性の高いタイムスタンプを有する文字列を

50

、他の文字列に連結して表示させる。この結果、字幕行の表示タイミングの精度を高めることができる。

【0070】

図14は、字幕設定装置10として機能するコンピュータのハードウェア構成の一例を示す。字幕設定装置10は、ホストコントローラ1482により相互に接続されるCPU1400、RAM1420、グラフィックコントローラ1475、及び表示装置1480を有するCPU周辺部と、入出力コントローラ1484によりホストコントローラ1482に接続される通信インターフェイス1430、ハードディスクドライブ1440、及びCD-ROMドライブ1460を有する入出力部と、入出力コントローラ1484に接続されるROM1410、フレキシブルディスクドライブ1450、及び入出力チップ1470を有するレガシー入出力部とを備える。

10

【0071】

ホストコントローラ1482は、RAM1420と、高い転送レートでRAM1420をアクセスするCPU1400及びグラフィックコントローラ1475とを接続する。CPU1400は、ROM1410及びRAM1420に格納されたプログラムに基づいて動作し、各部の制御を行う。グラフィックコントローラ1475は、CPU1400等がRAM1420内に設けたフレームバッファ上に生成する画像データを取得し、表示装置1480上に表示させる。これに代えて、グラフィックコントローラ1475は、CPU1400等が生成する画像データを格納するフレームバッファを、内部に含んでもよい。

【0072】

20

入出力コントローラ1484は、ホストコントローラ1482と、比較的高速な入出力装置である通信インターフェイス1430、ハードディスクドライブ1440、及びCD-ROMドライブ1460を接続する。通信インターフェイス1430は、ネットワークを介して外部の装置と通信する。ハードディスクドライブ1440は、字幕設定装置10が使用するプログラム及びデータを格納する。CD-ROMドライブ1460は、CD-ROM1495からプログラム又はデータを読み取り、RAM1420を介して入出力チップ1470に提供する。

【0073】

また、入出力コントローラ1484には、ROM1410と、フレキシブルディスクドライブ1450や入出力チップ1470等の比較的低速な入出力装置とが接続される。ROM1410は、字幕設定装置10の起動時にCPU1400が実行するブートプログラムや、字幕設定装置10のハードウェアに依存するプログラム等を格納する。フレキシブルディスクドライブ1450は、フレキシブルディスク1490からプログラム又はデータを読み取り、RAM1420を介して入出力チップ1470に提供する。入出力チップ1470は、フレキシブルディスク1490や、例えばパラレルポート、シリアルポート、キーボードポート、マウスポート等を介して各種の入出力装置を接続する。

30

【0074】

字幕設定装置10に提供されるプログラムは、フレキシブルディスク1490、CD-ROM1495、又はICカード等の記録媒体に格納されて利用者によって提供される。プログラムは、入出力チップ1470及び/又は入出力コントローラ1484を介して、記録媒体から読み出され字幕設定装置10にインストールされて実行される。

40

【0075】

字幕設定装置10にインストールされて実行されるプログラムは、内容データ取得モジュールと、音声認識モジュールと、文字列検出モジュールと、文字検出モジュールと、読付モジュールと、音素検出モジュールと、信頼度算出モジュールと、信頼度取得モジュールと、表示設定モジュールとを含む。各モジュールが字幕設定装置10に働きかけて行わせる動作は、図1から図13において説明した字幕設定装置10における、対応する部材の動作と同一であるから、説明を省略する。

【0076】

以上に示したプログラム又はモジュールは、外部の記憶媒体に格納されてもよい。記憶

50

媒体としては、フレキシブルディスク1490、CD-ROM1495の他に、DVDやPD等の光学記録媒体、MD等の光磁気記録媒体、テープ媒体、ICカード等の半導体メモリ等を用いることができる。また、専用通信ネットワークやインターネットに接続されたサーバシステムに設けたハードディスク又はRAM等の記憶装置を記録媒体として使用し、ネットワークを介してプログラムを字幕設定装置10に提供してもよい。

【0077】

図15は、文字列検出部120による処理の一例を説明する図である。文字列検出部120は、第1の認識データである「取り巻く」及び第3の認識データである「変化いたして」の各々に一致する文字列を内容データ20から検出する。一方、文字列検出部120は、第1及び第3の認識データの間の認識データである「緩急は著しく」に完全に一致する文字列を検出できない。このような場合には、文字列検出部120は、挟み打ち処理により、認識データである「緩急は著しく」に一致する文字列として、内容データ20における「取り巻く」及び「変化いたして」の間の文字列である「環境は著しく」を検出することができる。同様に、文字列検出部120は、認識データである「流れをお洗い」に一致する文字列として、「流れをおさらい」を検出することができる。

10

【0078】

図16は、音素検出部150による第1の処理の一例を示す。音素検出部150は、文字検出部130により一致する文字が検出されなかった各文字を、その文字の読みに含まれる音素に展開する。本例において、検出対象の文字が日本語であるので、音素検出部150は、漢字及び仮名の混じった文字列を、その文字列の読みを示す平仮名に展開する。即ち、音素検出部150は、「金ねうん」及び「愛ティーを」を、「きんねうん」及び「あいていーを」に展開する。

20

【0079】

一方、読付部140は、内容データである「近年ITを」の読み方の候補の1つとして、「きんねんあいていーを」を生成する。この結果、文字列検出部120は、内容データにおいて、認識データである「愛ティーを」に一致する文字列として、「愛ティーを」の音素に一致する音素を含む文字列である「ITを」を検出することができる。更に、この検出結果に基づいて、文字列検出部120は、認識データである「金ねうん」に一致する文字列として、「近年」を検出してもよい。

【0080】

図17は、音素検出部150による第2の処理の一例を示す。音素検出部150は、文字検出部130により一致する文字が検出されなかった各文字を、その文字の読みに含まれる音素に展開する。本例において、検出対象の文字が日本語であるので、音素検出部150は、漢字及び仮名の混じった文字列を、その文字列の読みを示す平仮名に展開する。即ち、音素検出部150は、「逆しすせ」及び「五人を」を、「ぎゃくしすせ」及び「ごにんを」に展開する。なお、本例においては、図16と比較して音声認識処理による認識率が低い。

30

【0081】

図16と同様に、読付部140は、内容データである「近年ITを」の読み方の候補の1つとして、「きんねんあいていーを」を生成する。しかしながら、音素検出部150は、「ぎゃくしすせ」及び「ごにんを」に音素が一致する文字又は文字列を、「きんねんあいていーを」の中から検出することができない。

40

【0082】

この場合、音素検出部150は、音声の内容に関わらず音声の長さ又は文字の長さに基づいて一致する音素を検出する強制割り振り処理を行う。例えば、「ぎゃくしすせ」として認識された音声の再生時間と、「ごにんを」として認識された音声の再生時間との比率に基づいて、「ぎゃくしすせ」の音素が、「きんねん」の音素に一致すると判断してもよいし、「ごにんを」の音素が、「あいていーを」の音素に一意すると判断してもよい。

【0083】

図18は、本実施例による処理の概要をまとめた図である。内容データ取得部100は

50

、内容データ20、例えば、「アクセシビリティについて」という文字列を取得する。音声認識部110は、音声認識処理により認識データ30-1~3、例えば「汗しびれ」、「地位」、及び「について」を生成する。文字列検出部120は、まず、DPマッチングにより、「アクセシビリティについて」及び「汗しびれ地位について」を比較して、認識データ30-1~3の各々に含まれる各文字に一致する文字を内容データ20から検出する。

【0084】

この結果、文字列検出部120は、認識データ30-1に一致する文字列として「アクセシビ」を検出し、認識データ30-2に一致する文字列として「リティ」を検出する。この検出結果に基づいてそのまま字幕を作成した場合には、2つの問題がある。1つ目の問題は、「地位」として音声認識された時間に表示すべき文字列が「ティ」であるにも関わらず、「リティ」が表示されてしまうことである。即ち、文字列「リティ」を表示すべき時間を示すタイムスタンプが誤っている。

10

【0085】

2つ目の問題は、「アクセシビリティ」という1つの単語が、音声認識処理の辞書などに登録されていないので、2つの文字列「アクセシビ」及び「リティ」に分割されて検出されていることである。これにより、字幕行において「アクセシビ」及び「リティ」の間で改行される恐れがある。

【0086】

文字列検出部120は、認識データ30-1~3の各々に一致する文字列を内容データ20から再度検出する。そして、信頼度算出部160は、認識データ30-1~3の各々が内容データ20における文字列を検出する確度である信頼度を算出する。この結果、信頼度算出部160は、認識データ30-1である「汗しいびれ」に対応付けて、認識データ30-2である「地位」と比較して高い信頼度を算出する。

20

【0087】

表示設定部170は、字幕の1行に表示可能な文字数が10文字である旨の設定情報を取得する。この場合、表示設定部170は、「アクセシビ」及び「リティ」を連結して「アクセシビリティ」を生成するが、「アクセシビリティ」を「について」に連結しない。その結果、表示設定部170は、表示タイミング情報40として、「アクセシビリティ」及び「について」の各々を、「41.5」及び「50.5」等の所定の時刻に表示すべき旨の情報を生成することができる。

30

【0088】

以上、本発明を実施の形態を用いて説明したが、本発明の技術的範囲は上記実施の形態に記載の範囲には限定されない。上記実施の形態に、多様な変更または改良を加えることが可能であることが当業者に明らかである。その様な変更または改良を加えた形態も本発明の技術的範囲に含まれ得ることが、特許請求の範囲の記載から明らかである。

【0089】

以上に示す実施例によると、以下の各項目に示す設定装置、プログラム、記録媒体、及び設定方法が実現される。

(項目1) 内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置であって、前記音声の内容を示す内容データを取得する内容データ取得部と、再生される前記音声を音声認識した文字データを分割して複数の認識データを生成する音声認識部と、前記複数の認識データの各々に一致する文字列を前記内容データから検出する文字列検出部と、前記文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出部と、前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部とを備える設定装置。

40

(項目2) 前記認識データのうち前記文字検出部により一致する文字が検出されなかつ

50

た各文字について、当該文字の読みに含まれる音素に一致する音素を、前記内容データの読みの中から検出する音素検出部を更に備え、前記文字検出部は、前記認識データのうち前記音素検出部により一致する音素が検出された文字に一致する文字として、前記内容データにおいて当該音素を含む文字を検出する項目1記載の設定装置。

【0090】

(項目3) 前記内容データの読み方の候補である読み候補を複数生成する読付部を更に備え、前記音素検出部は、前記認識データのうち前記文字検出部により前記内容データに一致する文字が検出されなかった各文字について、当該文字の読みに含まれる音素に一致する音素を、前記読付部により生成された複数の前記読み候補の何れかの中から検出する項目2記載の設定装置。

10

(項目4) 前記読付部は、前記内容データにおける複数の前記読み候補の各々を、当該読み候補により読まれる可能性を示す情報に対応付けて生成し、前記音素検出部は、前記認識データに含まれる文字の読みに含まれる音素を、読まれる可能性が高い順に、前記複数の読み候補の各々と比較する項目3記載の設定装置。

(項目5) 前記複数の認識データの各々が文字列に一致する確度である信頼度を算出する信頼度算出部を更に備え、前記文字列検出部は、予め定められた基準信頼度未満の信頼度の認識データである低信頼データについて、当該低信頼データに後続する認識データに一致する文字列を検出できなかった場合に、当該低信頼データに一致する文字列は検出できないと判断する項目1記載の設定装置。

(項目6) 前記複数の認識データの各々が文字列に一致する確度である信頼度を算出する信頼度算出部を更に備え、前記表示設定部は、前記内容データにおける複数の前記文字列のうち連続した2つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場合に、先に表示すべき前記文字列の末尾に後に表示すべき前記文字列を連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う項目1記載の設定装置。

20

【0091】

(項目7) 前記信頼度算出部は、前記文字列検出部により一致する文字列が検出された認識データに対応付けて、前記文字列検出部により一致する文字列が検出された認識データと比較して、より高い信頼度を算出する項目6記載の設定装置。

(項目8) 前記認識データのうち前記文字列検出部により一致する文字が検出されなかった各文字について、当該文字の読みに含まれる音素に一致する音素を、前記内容データの読みの中から検出する音素検出部を更に備え、前記文字列検出部は、前記認識データのうち前記音素検出部により一致する音素が検出された文字に一致する文字として、前記内容データにおいて当該音素を含む文字を検出し、前記信頼度算出部は、前記音素検出部により一致する音素が検出された文字を含む認識データに対応付けて、前記音素検出部により一致する音素が検出されることなく前記文字列検出部により一致する文字が検出された認識データと比較して、より低い信頼度を生成する項目6記載の設定装置。

30

(項目9) 前記音声認識部は、音声認識した前記複数の認識データが、再生される音声の内容と一致する可能性を示す音声認識確信度を、認識データ毎に更に生成し、前記文字列検出部は、音声認識確信度がより高い認識データに一致する文字列を、当該認識データと比較して音声認識確信度が低い認識データに先立って検出し、第1の前記認識データに一致する第1の文字列及び第2の前記認識データに一致する第2の文字列を検出した場合に、前記第1の認識データに後続しかつ前記第2の認識データに先行する認識データに一致する文字列として、前記第1の文字列に後続し前記第2の文字列に先行する文字列を検出する項目1記載の設定装置。

40

【0092】

(項目10) 前記表示設定部は、前記文字列検出部により前記内容データにおいて一致する文字列が検出されなかった認識データを、当該認識データとして音声認識された音声の再生時に表示させる設定を行う項目1記載の設定装置。

(項目11) 内容が予め定められた音声の再生に同期して前記内容を表示する表示タイ

50

ミングを設定する設定装置であって、再生される前記音声の内容を示す内容データに含まれる複数の文字列の各々に対応付けて、当該文字列を表示すべき時刻、及び、当該時刻に当該文字列を内容とする音声再生される確度である信頼度を取得する信頼度取得部と、前記複数の文字列のうち連続した2つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場合に、先に表示すべき前記文字列の末尾に後に表示すべき前記文字列を連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う表示設定部とを備える設定装置。

(項目12) 前記表示設定部は、先に表示すべき前記文字列に対応する信頼度が、後に表示すべき前記文字列に後続する後続文字列に対応する信頼度より高い場合に、連結した前記文字列の末尾に前記後続文字列を更に連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う項目11記載の設定装置。

10

【0093】

(項目13) 内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置として、コンピュータを機能させるプログラムであって、前記コンピュータを、前記音声の内容を示す内容データを取得する内容データ取得部と、再生される前記音声を音声認識した文字データを分割して複数の認識データを生成する音声認識部と、前記複数の認識データの各々に一致する文字列を前記内容データから検出する文字列検出部と、前記文字列検出部により一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出部と、前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定部として機能させるプログラム。

20

(項目14) 内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定装置として、コンピュータを機能させるプログラムであって、前記コンピュータを、再生される前記音声の内容を示す内容データに含まれる複数の文字列の各々に対応付けて、当該文字列を表示すべき時刻、及び、当該時刻に当該文字列を内容とする音声再生される確度である信頼度を取得する信頼度取得部と、前記複数の文字列のうち連続した2つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場合に、先に表示すべき前記文字列の末尾に後に表示すべき前記文字列を連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う表示設定部として機能させるプログラム。

30

(項目15) 項目13又は項目14に記載のプログラムを記録した記録媒体。

【0094】

(項目16) 内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定方法であって、コンピュータにより、前記音声の内容を示す内容データを取得する内容データ取得段階と、再生される前記音声を音声認識した文字データを分割して複数の認識データを生成する音声認識段階と、前記複数の認識データの各々に一致する文字列を前記内容データから検出する文字列検出段階と、前記文字列検出段階において一致する文字列が検出されなかった各認識データについて、当該認識データに含まれる各文字に一致する文字を前記内容データから検出することにより、当該認識データに一致する文字列を前記内容データから検出する文字検出段階と、前記内容データに含まれる文字列の各々を表示させる表示タイミングを、当該文字列に一致する認識データとして音声認識された音声の再生時に設定する表示設定段階とを備える設定方法。

40

(項目17) 内容が予め定められた音声の再生に同期して前記内容を表示する表示タイミングを設定する設定方法であって、コンピュータにより、再生される前記音声の内容を示す内容データに含まれる複数の文字列の各々に対応付けて、当該文字列を表示すべき時刻、及び、当該時刻に当該文字列を内容とする音声再生される確度である信頼度を取得する信頼度取得段階と、前記複数の文字列のうち連続した2つの文字列について、先に表示すべき文字列に対応する信頼度が、後に表示すべき文字列に対応する信頼度より高い場

50

合に、先に表示すべき前記文字列の末尾に後に表示すべき前記文字列を連結した文字列を、先に表示すべき前記文字列を表示すべき時刻に表示させる設定を行う表示設定段階とを備える設定方法。

【図面の簡単な説明】

【0095】

【図1】図1は、字幕設定装置10の機能ブロック図である。

【図2】図2は、内容データ20の一例を示す。

【図3】図3は、認識データ30-1~Nの一例を示す。

【図4】図4は、表示タイミング情報40の一例を示す。

【図5】図5は、字幕設定装置10が表示タイミングを設定する処理の動作フローを示す 10

。

【図6】図6は、S520における処理の詳細を示す。

【図7】図7は、S620における処理の詳細を示す。

【図8】図8は、S640における処理の詳細を示す。

【図9】図9は、S810における第1の処理の詳細を示す。

【図10】図10は、S810における第2の処理の詳細を示す。

【図11】図11は、S830における第1の処理の詳細を示す。

【図12】図12は、S830における第2の処理の詳細を示す。

【図13】図13は、S550における処理の詳細を示す。

【図14】図14は、字幕設定装置10として機能するコンピュータのハードウェア構成 20の一例を示す。

【図15】図15は、文字列検出部120による処理の一例を説明する図である。

【図16】図16は、音素検出部150による第1の処理の一例を示す。

【図17】図17は、音素検出部150による第2の処理の一例を示す。

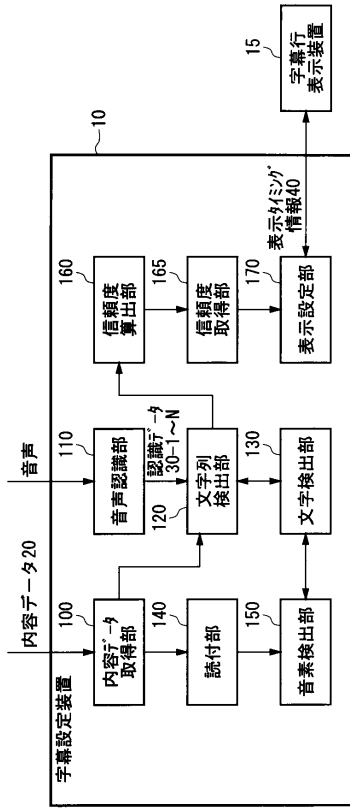
【図18】図18は、本実施例による処理の概要をまとめた図である。

【符号の説明】

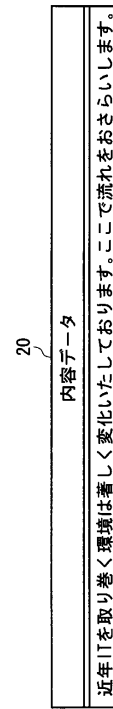
【0096】

10	字幕設定装置	
15	字幕行表示装置	
20	内容データ	30
30	認識データ	
40	表示タイミング情報	
100	内容データ取得部	
110	音声認識部	
120	文字列検出部	
130	文字検出部	
140	読付部	
150	音素検出部	
160	信頼度算出部	
165	信頼度取得部	40
170	表示設定部	

【 図 1 】



【 図 2 】



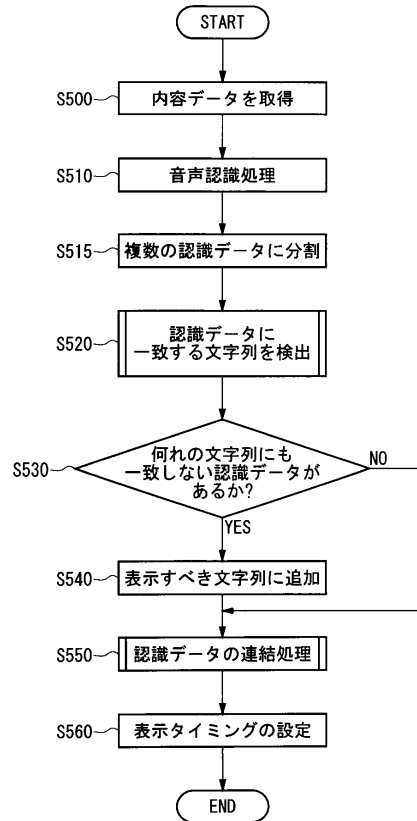
【 図 3 】

認識時刻	認識データ	
02.103	金ねうん	30-1
03.211	愛ティーを	30-2
05.235	取り巻く	30-3
06.456	繰急は	30-4
08.387	著しく変化いたして	30-5
11.125	おります。	30-6
16.887	ここで	30-7
17.001	流れをお洗い	30-8
18.051	致します。	30-9
⋮	⋮	
⋮	⋮	
XX	XX	30-N

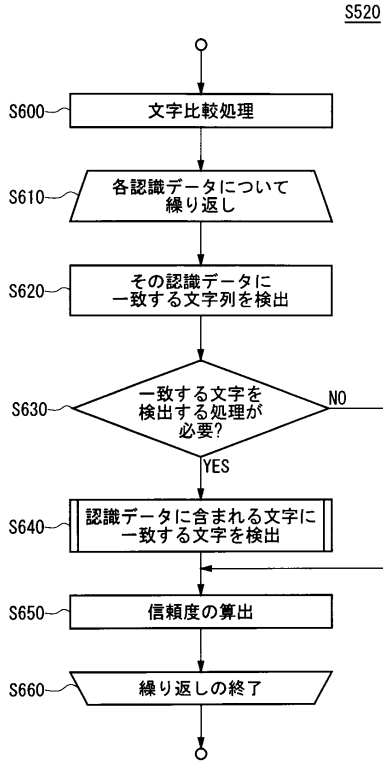
【 図 4 】

表示時刻	文字列
02.103	近年ITを取り巻く環境は
08.387	著しく変化いたしております。
16.887	ここで流れをおさらいします。

【 図 5 】



【 図 6 】

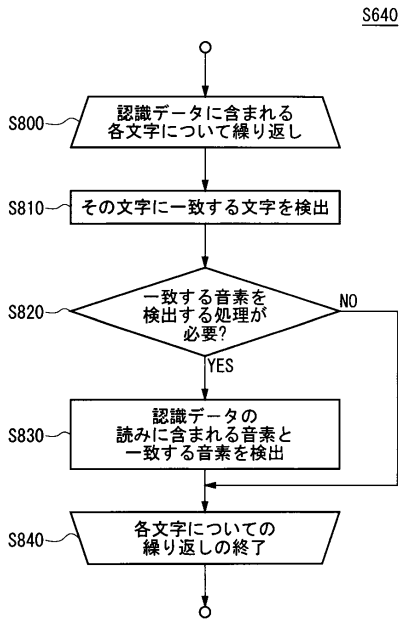


【 図 7 】

S620

2つ前の認識データ	1つ前の認識データ	対象の認識データ	処理内容
○/×	○	○	対象の認識データを一致と判断
		×	保留
		○	1つ前の認識データを一致と判断
		×	保留、但し最後の認識データならば、末尾確定文字比較
	×	○	2つ前、1つ前の認識データに含まれる文字と一致する文字を検出(範囲確定) 信頼度が基準信頼度未満ならば、再評価フラグ(△)
		×	2つ前、1つ前の認識データに含まれる文字と一致する文字を検出(範囲不確定)
		○	1つ前の認識データを一致と判断
	△	×	1つ前の認識データを一致しないと判断

【 図 8 】



【 図 9 】

S810

先頭文字	前の文字	対象文字	処理内容
○/×	○	○	対象文字を一致と判断
		×	保留
○		○	認識データ内の不一致文字を一致と判断
	×	×	保留
		○	対象文字を一致と判断
×		×	保留

【 図 1 0 】

S810

先頭文字	次の認識データ	処理内容
○	○	対象の認識データを一致と判断
	×	先頭文字確定、次の認識データ不確定音素検出処理
×	○	先頭文字不確定、次の認識データ確定音素検出処理
	×	先頭文字不確定、次の認識データ不確定音素検出処理

【 図 1 1 】

S830

先頭文字	前の音素	対象音素	処理内容
○/×	○	○	音素一致と判断
		×	保留、又は、次の読み候補
○		○	認識データ内の不一致文字を一致と判断
	×	×	保留、又は、次の読み候補
×		○	音素一致と判断
		×	保留、又は、次の読み候補

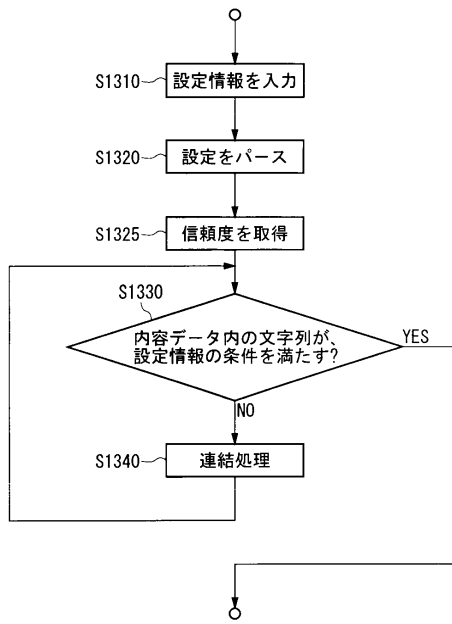
【 図 1 2 】

S830

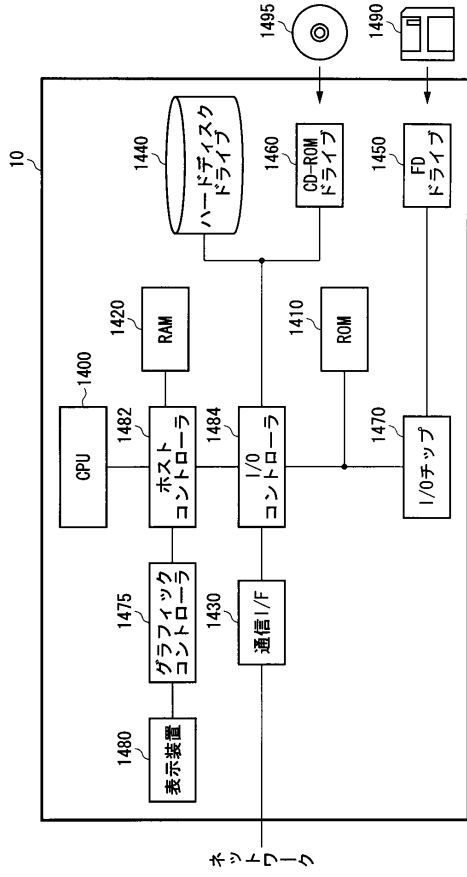
先頭文字	次の認識データ、又は最後の認識データ	処理内容
○	○	認識データ内の不一致文字を一致と判断
	×	保留
×	○	強制割り振り
	×	保留

【 図 1 3 】

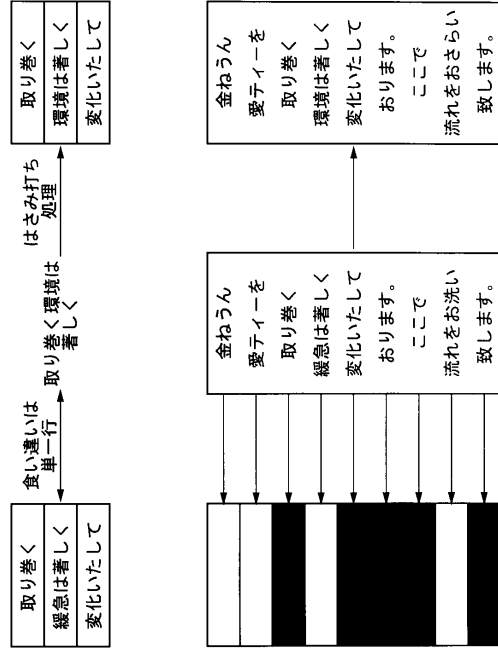
S550



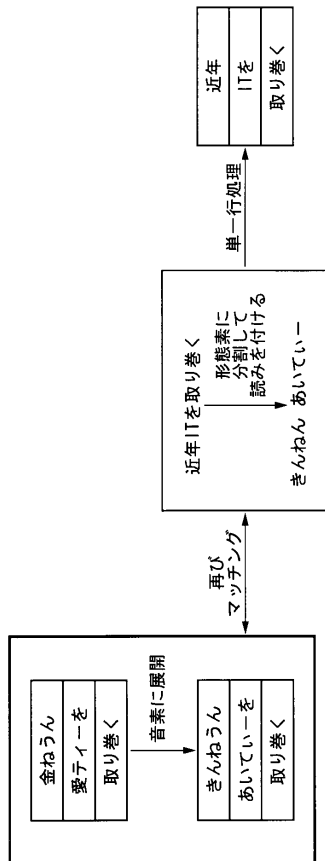
【図14】



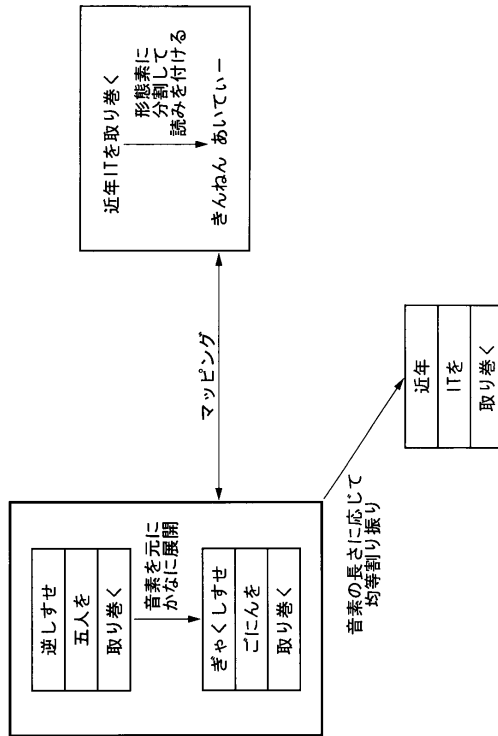
【図15】



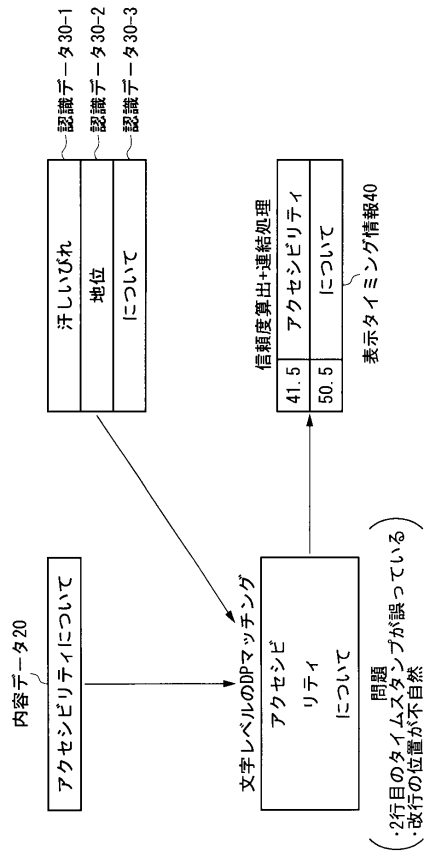
【図16】



【図17】



【 図 18 】



フロントページの続き

- (72)発明者 宮本 晃太郎
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内
- (72)発明者 東海林 みどり
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内

審査官 櫻本 剛

- (56)参考文献 特開平10-308887(JP,A)
特開2001-034151(JP,A)
特開2005-045503(JP,A)
丸山他, 字幕送出タイミング検出におけるワード列ペアモデルの構成検討, 日本音響学会平成10年度秋季研究発表会講演論文集, 日本, 1998年 6月24日, Vol.1, p.25-26

- (58)調査した分野(Int.Cl., DB名)
G10L 15/00-15/28
JSTPlus(JDream2)
IEEE