

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2010-55030  
(P2010-55030A)

(43) 公開日 平成22年3月11日(2010.3.11)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G10L 15/06 (2006.01)</b>	G10L 15/06 300Y	5D015
	G10L 15/06 300C	
	G10L 15/06 310Z	

審査請求 未請求 請求項の数 7 O L (全 33 頁)

(21) 出願番号	特願2008-222817 (P2008-222817)	(71) 出願人	000004352 日本放送協会 東京都渋谷区神南2丁目2番1号
(22) 出願日	平成20年8月29日 (2008. 8. 29)	(74) 代理人	100064908 弁理士 志賀 正武
		(74) 代理人	100108578 弁理士 高橋 詔男
		(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
		(72) 発明者	今井 亨 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
		Fターム(参考)	5D015 AA01 GG01

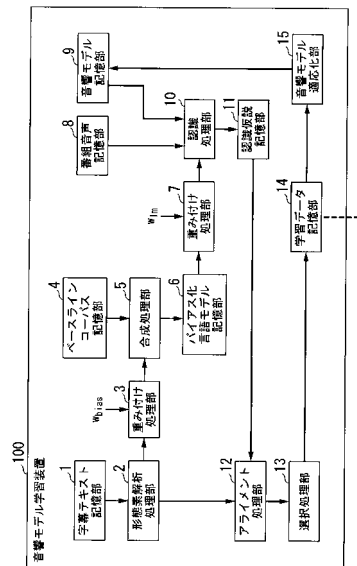
(54) 【発明の名称】 音響処理装置およびプログラム

(57) 【要約】

【課題】字幕と発話内容の一致率が低いオフライン字幕を利用した場合にも、高い認識率によって音響モデルの学習データを自動生成し、多様な発話スタイルに対応できる音響モデルを作成できる音響モデル学習装置を提供する。

【解決手段】一致区間のみを切り出して利用するのではなく、その他の区間から得られる情報も、音響モデル学習のために用いる。そのため、音声認識結果と書き起こし字幕テキストの一致区間を利用して学習データを得る際に、各形態素の信頼度を導入することにより、一致区間以外の音声も利用して学習データを自動生成させる。

【選択図】 図 1



## 【特許請求の範囲】

## 【請求項 1】

音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部と、

前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力する認識処理部と、

前記認識処理部が出力した認識結果データに基づき、教師系列データを取得し、前記教師系列データに含まれる各々の音素に信頼度データを付加し、前記信頼度データが付加された前記教師系列データと該教師系列データに対応する音響特徴量データとを、学習データとして出力する選択処理部と、

前記選択処理部が出力した学習データを読み込み、前記学習データに含まれる各々の前記音素の観測確率データを算出し、前記音素ごとに、前記学習データに含まれる前記音響特徴量データと算出された前記観測確率データとに基づき、且つ、前記信頼度データを重みとして乗じて得られる音響モデルを用いて、前記音響モデル記憶部を更新する音響モデル適応化部と、

を具備することを特徴とする音響処理装置。

## 【請求項 2】

請求項 1 に記載の音響処理装置であって、

前記選択処理部は、前記音声に対応するテキストデータを読み込み、前記認識結果データに含まれる最尤パスを前記教師系列データとして取得し、当該教師系列データと前記テキストデータが一致する一致区間に含まれる前記音素に対しては最高信頼度を表わす信頼度データを付加し、その他の区間に含まれる前記音素に対しては最低信頼度を表わす信頼度データを付加する、

ことを特徴とする音響処理装置。

## 【請求項 3】

請求項 1 に記載の音響処理装置であって、

前記選択処理部は、前記音声に対応するテキストデータを読み込み、前記認識結果データに含まれる最尤パスと前記テキストデータが一致する一致区間と時間軸上で対立する区間を枝刈りする処理を行ない、この枝刈り処理の結果得られる系列を前記教師系列データとして取得し、当該教師系列データに含まれる言語的単位ごとの事後確率を前記言語的単位に含まれる前記音素に対する信頼度データとして付加する、

ことを特徴とする音響処理装置。

## 【請求項 4】

音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部と、

前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力する認識処理部と、

前記認識処理部が出力した認識結果データと前記音声に対応するテキストデータとを読み込み、前記認識結果データに含まれる最尤パスと前記テキストデータが一致する一致区間と時間軸上で対立する区間を枝刈りする処理を行ない、この枝刈り後の認識結果データに含まれる言語的単位ごとの事後確率を前記言語的単位に含まれる各々の音素に対する信頼度データとして付加した教師系列データを求め、前記教師系列データとこれに対応する音響特徴量データとを、学習データとして出力する選択処理部と、

前記選択処理部が出力した学習データを読み込み、フォワードバックワードアルゴリズムを用いて前記学習データに含まれる各々の音素の観測確率を求め、前記音素ごとに、前記観測確率で重み付けられた前記音響特徴量データとして得られる音響モデルを用いて、前記音響モデル記憶部を更新する音響モデル適応化部と、

を具備することを特徴とする音響処理装置。

## 【請求項 5】

請求項 1 に記載の音響処理装置であって、

10

20

30

40

50

前記選択処理部は、前記音声に対応するテキストデータを読み込み、前記認識結果データから、前記テキストデータとの一致が最大となるようなパスを選択して前記教師系列データとして取得するとともに、当該教師系列データと前記テキストデータが一致する一致区間に含まれる前記音素に対しては最高信頼度を表わす信頼度データを付加する、ことを特徴とする音響処理装置。

【請求項 6】

音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部と、

前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力する認識処理部と、

前記認識処理部が出力した認識結果データと前記音声に対応するテキストデータとを読み込み、前記認識結果データから、前記テキストデータとの一致が最大となるようなパスを選択して前記教師系列データとして取得するとともに、当該教師系列データ内で前記認識結果データと前記テキストデータとが一致する一致区間を選択し、前記一致区間と時間軸上で対立する区間を枝刈りする処理を行ない、この枝刈り後の認識結果データに含まれる言語的単位ごとの事後確率を前記言語的単位に含まれる各々の音素に対する信頼度データとして付加した教師系列データを求め、前記教師系列データとこれに対応する音響特徴量データとを、学習データとして出力する選択処理部と、

前記選択処理部が出力した学習データを読み込み、フォワードバックワードアルゴリズムを用いて前記学習データに含まれる各々の音素の観測確率を求め、前記音素ごとに、前記観測確率で重み付けられた前記音響特徴量データとして得られる音響モデルを用いて、前記音響モデル記憶部を更新する音響モデル適応化部と、

を具備することを特徴とする音響処理装置。

【請求項 7】

音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部を備えるコンピュータに、

前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力する認識処理過程と、

前記認識処理過程で出力した認識結果データに基づき、教師系列データを取得し、前記教師系列データに含まれる各々の音素に信頼度データを付加し、前記信頼度データが付加された前記教師系列データと該教師系列データに対応する音響特徴量データとを、学習データとして出力する選択処理過程と、

前記選択処理過程で出力した学習データを読み込み、前記学習データに含まれる各々の前記音素の観測確率データを算出し、前記音素ごとに、前記学習データに含まれる前記音響特徴量データと算出された前記観測確率データとに基づき、且つ、前記信頼度データを重みとして乗じて得られる音響モデルを用いて、前記音響モデル記憶部を更新する音響モデル適応化過程と、

の処理を実行させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、統計的処理に基づく音声処理に関する。特に、各音素の音響特徴量の統計量を音響モデルとして処理を行なう音響処理装置およびそのプログラムに関する。また特に、音響モデルの学習に関する。

【背景技術】

【0002】

統計モデルを用いた音声認識では、各音素の統計量を得るために大量の学習データが必要になる。この学習データは、大量の音声データとその音声に対する書き起こしの対であり、従来人手により時間とコストをかけて作成されてきた。

【0003】

10

20

30

40

50

現在、放送番組の一部には、高齢者や聴覚障害者の番組視聴など音声からの情報が十分に得られない環境でも番組内容が理解できるように、番組音声に対応する字幕テキストが付与されているものがある。この字幕テキストにはオンラインで付与された字幕テキストとオフラインで付与された字幕テキストの2種類があり、前者は番組中の発話内容とほぼ一致している。但し、後者は話し言葉の冗長性を除去し、簡潔で読みやすい字幕を付与できるため、発話内容と字幕との一致率が低い。

【0004】

上記の一致率が高いオンライン字幕を利用できる場合には、音声区間の切り出しと、切り出された音声に対応する字幕区間の切りだしを行なうことで、音響モデルの学習データを容易に自動作成することができ、読み上げ音声を中心に比較的高い音声認識精度が得られるニュース番組を対象とした場合の認識精度の改善が報告されている（非特許文献1）。この方法では、字幕区間の切り出しに、字幕テキストで適応化した言語モデルを用いて得られた番組音声の認識結果を利用し、字幕テキストと認識結果の一致区間を音響モデルの学習データとしている。この番組音声の認識は、主に音声区間の切り出しと字幕と発話内容の不一致部分の検出を目的として行われ、言語モデルの適応化に用いる字幕テキストの重みは比較的小さいことが特徴である。

10

【0005】

図14は、従来技術による音響モデル学習装置の機能構成を示すブロック図である。この図に示す音響モデル学習装置200は、字幕テキストを形態素単位に分割した後、字幕テキストに $W_{bias}$ の重みを付けて、バイアスのある言語モデル（バイアス化言語モデル86）を得る。この言語モデルを用いて番組音声を認識する（認識処理部90）。本装置の主目的は、字幕と発話内容の不一致区間の検出と音声と字幕の切り出しであるため、バイアス重み $W_{bias}$ （4程度）と言語重み $W_{lm}$ （10程度）は比較的小さな値を用いる。次に、番組全体で認識仮説と字幕テキストのDPマッチングを行い（アライメント処理部92）、字幕テキストと認識結果が3単語以上連続して一致している区間を選択し、番組音声の中の音声区間と対応する単語列である学習データ（学習データ記憶部94）が得られる。

20

【非特許文献1】Long Nguyen, 外1名, "Light Supervision in Acoustic Model Training", IEEE International Conference on Acoustics, Speech, and Signal Processing, 2004. Proceedings. (ICASSP apos;04), 2004年, Volume 1, p. 17 - 21

30

【発明の開示】

【発明が解決しようとする課題】

【0006】

現在、放送番組の対談部分などの比較的自由なスタイルで発話された部分の認識精度は、読み上げ音声の認識精度に比べて低く、このような部分の認識精度を向上させるため、様々なスタイルで発話された音響モデルの学習データが必要とされている。

【0007】

放送番組の対談部分などの比較的自由なスタイルでは、オフラインで作成される字幕テキストが用いられるが、前述の通りこの字幕テキストの作成の際は簡潔さが優先されるため、字幕テキストと実際の発話内容との一致率は低い。従来技術による音響モデル学習装置は字幕テキストと実際の発話内容との一致率が高いことを前提としており、従来技術による方法ではこの一致率の低いデータからは高い認識精度が得られず、その結果として音響モデルの学習データを効率よく多量に生成することができないという問題があった。

40

【0008】

本発明のより具体的な課題は、次の通りである。対談などの自由発話の音声を学習する場合、議事録や放送番組の字幕テキストなど、一般的に入手容易な書き起こしテキストは不完全であることが多く、十分な学習効果が得られないことがある。例えば、対談の自由発話には、「あの」、「で」、「まあ」などの多くの不要語が含まれており、これらの不要語は書き起こしのテキストには盛り込まれないことが多い。よって、そのような不完全な書き起こしのテキストからは、これらの不要語が得られない。また、字幕テキストを作

50

成する目的に応じて、話し言葉に含まれる冗長な部分を読みやすいフレーズに置き換えて書き起こしのテキストが作成されることも多い。従って、その冗長な部分に対応するテキストが得られない。

【0009】

また、従来のように一致部分の音声区間を切り出して学習する技術では、不一致部分の音声は学習には用いられない。そのため、音声と書き起こしテキストの一致率が低い場合には、少量の学習データしか得られない。また、音声認識の入力となる発話切り出し単位と学習音声の切り出し単位とが異なり、細かく切り出された音声を学習データとしなければならない、学習時と認識時の音声切り出し単位の不接合を避けられないという問題もあった。また、従来技術を用いた場合には、認識結果（最尤仮説）と書き起こしテキストとの間の不一致には、例えば、「こと」と「事」のような表記の揺らぎによる不一致や、「ですから」と「です/から」といったように形態素分割が異なるために起こる不一致が見られ、これらの区間の音声学習データとしては用いられなかったという問題もあった。

10

【0010】

本発明は、上記のような事情を考慮して為されたものであり、字幕と発話内容の一致率が低いオフライン字幕を利用した場合にも、高い認識率によって音響モデルの学習データを自動生成し、多様な発話スタイルに対応できる音響モデルを作成することを目的とするものである。また、本発明は、字幕と発話内容の一致率が低いオフライン字幕を利用した場合にも、より多くの量の音響モデルの学習データを生成することのできる音響処理装置およびプログラムを提供することを目的とする。

20

【課題を解決するための手段】

【0011】

上記の課題を解決するために、本発明では、次のような手段を用いる。即ち、所定数の形態素が連続して一致している区間が一致区間であるが、そのような一致区間を切り出して利用するのではなく、その他の区間から得られる情報も、音響モデル学習のために用いる。より具体的には、音声認識結果と書き起こし字幕テキストの一致区間を利用して学習データを取得の際に、各形態素（単語）の信頼度を導入することにより、一致区間以外の音声も利用して学習データを自動生成させる。このように一致区間に対応する音声を切り出さずに学習できるようになるため、認識時と学習時の発話単位の不整合が解消されることにより認識精度の向上が期待される。

30

【0012】

より具体的には、本発明の特別な技術的特長（Special Technical Feature）は、次に述べる構成あるいはその部分集合による構成によるものである。その構成とは即ち、

- 1) 音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部を備え、
- 2) 認識処理部が、前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力し、
- 3) 選択処理部は、前記認識処理部が出力した認識結果データに基づき、各々の音素に対する信頼度データを付加した教師系列データを求め、前記教師系列データとこれに対応する音響特徴量データとを、学習データとして出力し、そして、
- 4) 音響モデル適応化部は、前記選択処理部が出力した学習データを読み込み、前記学習データに含まれる各々の音素の観測確率を求め、前記音素ごとに、音響モデルを用いて、前記音響モデル記憶部を更新する。

40

【0013】

そして、本発明の第1の態様による音響処理装置は、音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部と、前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力する認識処理部と、前記認識処理部が出力した認識結果データに基づき、教師系列データを取得し、前記教師系列データに含まれる各々の音素に信頼度データを付加し、前記信頼度データが付加された前記教師系列データと該教師系列データに対応する音響特徴

50

量データとを、学習データとして出力する選択処理部と、前記選択処理部が出力した学習データを読み込み、前記学習データに含まれる各々の前記音素の観測確率データを算出し、前記音素ごとに、前記学習データに含まれる前記音響特徴量データと算出された前記観測確率データとに基づき、且つ、前記信頼度データを重みとして乗じて得られる音響モデルを用いて、前記音響モデル記憶部を更新する音響モデル適応化部とを具備することを特徴とするものである。

ここで、言語的単位とは、言語的にまとまりのある単位であり、例えば、形態素、単語などである。また、認識仮説とは、認識処理の結果得られる認識の仮説であり、ある時刻でみたときには、複数の仮説が並立する場合もある。例えば形態素を単位としてみたとき、認識仮説は、形態素がアークに対応し、形態素間の接続点がノードに対応する形のラティス構造（時間方向の半順序構造）を有する。ある時刻において複数のアーク（相互に対立するアーク）が存在するとき、それらに対する確率を与えることができるが、本態様ではその確率を信頼度データとしている。また、教師系列データとは、信頼度データが付与されたラティスである。

本態様では、音響処理装置は、認識仮説の最尤パスと（書き起こしの字幕などの）テキストデータとの一致に基づき音声の区間を切り出す代わりに、教師音素列の信頼度を導入し、信頼度に基づいて学習する部分と学習しない部分を特定する。選択処理部が出力する学習データには、音声と教師音素列に加えて、それらに対応する教師音素列の信頼度を含んでいる。そして音響モデル適応化部は、音響モデルの統計量を推定する際に、例えばフォワードバックワードアルゴリズム（Forward-Backward Algorithm）やビタビアルゴリズム（Viterbi Algorithm）により得られる教師音素列の各音素の観測確率に、上記の信頼度による重み付けを行なって推定する。

このような構成により、認識結果データの中の一一致区間だけではなく、その他の部分の音声の音響特徴量データを用いて、音響モデルの学習（音響モデル記憶部の更新（適応化））を行なえる。また、信頼度データを用いて、それによって重み付けをしていることにより、信頼度に応じた適切な学習を行なうことが可能となる。

#### 【0014】

本発明の第2の態様による音響処理装置は、前記の音響処理装置において、前記選択処理部は、前記音声に対応するテキストデータを読み込み、前記認識結果データに含まれる最尤パスを前記教師系列データとして取得し、当該教師系列データと前記テキストデータが一致する一致区間に含まれる前記音素に対しては最高信頼度を表わす信頼度データを付加し、その他の区間に含まれる前記音素に対しては最低信頼度を表わす信頼度データを付加する、ことを特徴とする。

ここで、テキストデータとは、音声に対応するものであり、例えば書き起こしの字幕テキストのデータである。但し、音声（実際の発話内容）と字幕テキストの内容とは、必ずしも完全に一致しているとは限らない。また、信頼度データは、認識結果データに含まれるあるパス（パスは、言語的単位（形態素等）の系列に対応）が信頼できる度合いを表わす数値である。最高信頼度とは、その仮説の正しさを表わす確率が1である状態に対応する。最低信頼度とは、その仮説が正解に対して無情報である状態に対応する。信頼度データを0以上1以下の実数で表わしたとき、信頼度データが大きいほど信頼できる度合いが高く、最高信頼度は信頼度データ「1」に相当し、最低信頼度は信頼度データ「0」に相当する。

本態様では、音響処理装置は、最尤仮説とテキストデータの一致区間のみを学習するために、一致区間に含まれる形態素（単語）およびその形態素内の音素にのみ信頼度1を与え、それ以外の区間に含まれる形態素およびその形態素内の音素の信頼度を0とする。

これにより、一致区間に高い信頼度を与えて音響モデルの学習をすることができる。

#### 【0015】

本発明の第3の態様による音響処理装置は、前記の音響処理装置において、前記選択処理部は、前記音声に対応するテキストデータを読み込み、前記認識結果データに含まれる最尤パスと前記テキストデータが一致する一致区間と時間軸上で対立する区間を枝刈りす

る処理を行ない、この枝刈り処理の結果得られる系列を前記教師系列データとして取得し、当該教師系列データに含まれる言語的単位ごとの事後確率を前記言語的単位に含まれる前記音素に対する信頼度データとして付加する、ことを特徴とする。

本態様においては、一致区間に対立する区間を枝刈りしているため、一致区間の事後確率は1となる。つまり、この一致区間に含まれる音素には最高信頼度を表わす信頼度データ(=1)が付加される。そして、本態様では、音響処理装置は、認識結果データの形態素ラティスを利用し、最尤仮説とテキストデータとの一致区間以外の形態素にも非零の信頼度を付与する。これは、時間軸上で最尤仮説と対立する(重複する)仮説の枝刈り処理を行なった後の認識仮説のラティスを用いて、最尤仮説の各形態素の事後確率を算出し、この事後確率を信頼度として付与する方法である。

10

これにより、一致区間に最高信頼度を表わす信頼度データを付加し、その他の区間にも事後確率に応じた非零の信頼度データを付加し、その重みを用いて音響モデルの学習を行なえる。

#### 【0016】

上記の各態様の音響処理装置は、最尤仮説に含まれる各形態素を教師とし、これに基づいて信頼度を付与し、音響モデルの学習を行なっている。

これに対して、本発明の第4の態様による音響処理装置は、音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部と、前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力する認識処理部と、前記認識処理部が出力した認識結果データと前記音声に対応するテキストデータとを読み込み、前記認識結果データに含まれる最尤パスと前記テキストデータが一致する一致区間と時間軸上で対立する区間を枝刈りする処理を行ない、この枝刈り後の認識結果データに含まれる言語的単位ごとの事後確率を前記言語的単位に含まれる各々の音素に対する信頼度データとして付加した教師系列データを求め、前記教師系列データとこれに対応する音響特徴量データとを、学習データとして出力する選択処理部と、前記選択処理部が出力した学習データを読み込み、フォワードバックワードアルゴリズムを用いて前記学習データに含まれる各々の音素の観測確率を求め、前記音素ごとに、前記観測確率で重み付けられた前記音響特徴量データとして得られる音響モデルを用いて、前記音響モデル記憶部を更新する音響モデル適応化部とを具備することを特徴とするものである。

20

30

前記の第1の態様による音響処理装置と比較して、本態様の特徴は、音響処理装置が、観測確率で重み付けられた音響特徴量データとして得られる音響モデルを用いて音響モデル記憶部を更新する点である。また、選択処理部は、認識結果データに含まれる最尤パスとテキストデータが一致する一致区間と時間軸上で対立する区間を枝刈りする処理を行ない、この枝刈り後の認識結果データに含まれる言語的単位ごとの事後確率を前記言語的単位に含まれる各々の音素に対する信頼度データとして付加する。

本態様では、音響処理装置は、枝刈り処理後の認識仮説のラティス(このラティスには、下で述べる観測確率に信頼度が含まれている)を教師系列として、このラティスに直接フォワードバックワードアルゴリズムを適用して、音響モデルを学習する。この方法では、単語の信頼度(事後確率)は直接用いず、フォワードバックワードアルゴリズムにより得られる観測確率に信頼度が含まれている。

40

#### 【0017】

本発明の第5の態様による音響処理装置は、前記の音響処理装置において、前記選択処理部は、前記音声に対応するテキストデータを読み込み、前記認識結果データから、前記テキストデータとの一致が最大となるようなパスを選択して前記教師系列データとして取得するとともに、当該教師系列データと前記テキストデータが一致する一致区間に含まれる前記音素に対しては最高信頼度を表わす信頼度データを付加することを特徴とするものである。

本態様では、音響処理装置は、認識仮説のラティスとテキストデータの一致が最大になる形態素系列を教師系列とする。そして、一致区間には最高信頼度を表わす信頼度データ

50

を付加し、その他の区間にも適宜信頼度データを付加する。つまり、最尤仮説ではないパス内に一致区間が存在する場合にもそのような一致区間に最高信頼度が与えられる。これにより、最尤仮説とテキストデータとの間の表記の揺らぎや形態素分割の異なりに起因して不一致区間と判定されてしまう区間を一致区間として利用し、音響モデルを学習することができる。

【0018】

本発明の第6の態様による音響処理装置は、音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部と、前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力する認識処理部と、前記認識処理部が出力した認識結果データと前記音声に対応するテキストデータとを読み込み、前記認識結果データから、前記テキストデータとの一致が最大となるようなパスを選択して前記教師系列データとして取得するとともに、当該教師系列データ内で前記認識結果データと前記テキストデータとが一致する一致区間を選択し、前記一致区間と時間軸上で対立する区間を枝刈りする処理を行ない、この枝刈り後の認識結果データに含まれる言語的単位ごとの事後確率を前記言語的単位に含まれる各々の音素に対する信頼度データとして付加した教師系列データを求め、前記教師系列データとこれに対応する音響特徴量データとを、学習データとして出力する選択処理部と、

10

前記選択処理部が出力した学習データを読み込み、フォワードバックワードアルゴリズムを用いて前記学習データに含まれる各々の音素の観測確率を求め、前記音素ごとに、前記観測確率で重み付けられた前記音響特徴量データとして得られる音響モデルを用いて、前記音響モデル記憶部を更新する音響モデル適応化部とを具備することを特徴とする。

20

前記の第4の態様による音響処理装置が最尤パスとテキストデータとが一致する一致区間を基準として対立する区間を枝刈りするのに対して、この第6の態様による音響処理装置は、認識結果データの全体（最尤パス以外の仮説も含む）とテキストデータとが一致する一致区間を基準として対立する区間を枝刈りする点が特徴である。

【0019】

また、本発明の一態様は、音素と該音素に対応する音響特徴量とを関連付けた音響モデルを記憶する音響モデル記憶部を備えるコンピュータに、前記音響モデル記憶部から読み出した音響モデルを用いて音声の認識処理を行い、認識仮説を含んだ認識結果データを出力する認識処理過程と、前記認識処理過程で出力した認識結果データに基づき、教師系列データを取得し、前記教師系列データに含まれる各々の音素に信頼度データを付加し、前記信頼度データが付加された前記教師系列データと該教師系列データに対応する音響特徴量データとを、学習データとして出力する選択処理過程と、前記選択処理過程で出力した学習データを読み込み、前記学習データに含まれる各々の前記音素の観測確率データを算出し、前記音素ごとに、前記学習データに含まれる前記音響特徴量データと算出された前記観測確率データとに基づき、且つ、前記信頼度データを重みとして乗じて得られる音響モデルを用いて、前記音響モデル記憶部を更新する音響モデル適応化過程との処理を実行させるプログラム。

30

【0020】

また、さらに、次のA1からA5までの態様として上記課題を解決するようにしても良い。

40

【0021】

[A1] 本発明の一態様は、上記の音響処理装置において、前記音響モデル記憶部は、前記音響モデルを番組ごとに記憶するものであり、前記認識処理部は、前記番組の音声と、当該番組に対応した前記音響モデルを用いて音声の認識処理を行ない、前記選択処理部は、前記番組ごとに前記学習データを出力し、前記音響モデル適応化部は、当該番組用の前記音響モデルを更新することを特徴とする。

この構成によれば、認識対象の音声は放送等の番組の音声であり、番組ごとに音響モデルを持ち、番組ごとの認識結果を用いて音響モデルの適応化を行なえる。よって、番組ごとに特有の音響的特徴がある場合に、その特有の特徴に応じた適応化を行なうことができ

50



、認識精度が上がり、より多量の学習モデルを獲得できる。

【0022】

[A2] 本発明の一態様は、上記の音響処理装置において、言語モデルを記憶する言語モデル記憶部をさらに具備し、前記認識処理部は、前記言語モデル記憶部から読み出した前記言語モデルと前記音響モデル記憶部から読み出した前記音響モデルとを所定の重み比率値で重み付けして用いて、認識処理を行なうものであり、且つ、前記認識処理部は、音響モデル適応化部によって更新された前記音響モデルを用いて再度認識処理を行なうものであり、このとき使用する前記重み比率値は、前回の認識処理のときよりも前記言語モデルがより小さい重みで用いられる値とすることを特徴とする。

この構成によれば、認識処理と、認識処理結果を用いた音響モデルの適応化とを、繰り返し行なう場合に、徐々に、言語モデルを用いる重みが相対的に低下する。よって、当初の認識率を高めるために言語モデルの重みを大きめにとった場合にも、音響モデルの適応化の進展に合わせて、言語モデルの重みを減ずることができ、認識精度が上がり、より多量の学習モデルを獲得できる。

【0023】

[A3] 本発明の一態様は、上記の音響処理装置において、前記音声の中の発話部分の特徴量を表わす発話モデルを記憶する発話モデル記憶部と、前記発話モデル記憶部から読み出した前記発話モデルを用いて前記音声の中の発話部分の音声区間を切り出す音声区間切り出し部と、前記選択処理部が出力した前記学習データを用いて、前記発話モデル記憶部に記憶される前記発話モデルを更新する発話モデル適応化部とをさらに具備し、前記認識処理部は、前記音声区間切り出し部によって切り出された前記音声区間の音声を認識処理することを特徴とする。

この構成によれば、発話モデルを用いて音声区間の中の発話部分の音声区間を切り出す。そして、認識処理の結果得られた学習データを用いて、発話モデルを適応化する。よって、発話部分の音声区間の切り出しと、切り出された音声区間の認識処理と、認識結果に基づく学習データを用いた発話モデルの適応化の処理を繰り返すことができる。よって、音声区間の切り出しの精度が向上し、認識精度が上がり、より多量の学習モデルを獲得できる。

【0024】

[A4] 本発明の一態様は、上記の音響処理装置において、話者ごと又は話者属性ごとの統計量を表わす話者属性別発話モデルを記憶する話者属性別発話モデル記憶手段と、前記話者属性別発話モデル記憶手段から読み出した前記話者属性別発話モデルを用いて前記音声の中の所定音声区間における話者又は話者属性を推定する話者属性推定手段とをさらに具備し、前記音響モデル記憶部は、前記話者ごとまたは前記話者属性ごとに区別して前記音響モデルを記憶するものであり、前記認識処理部は、前記話者属性推定手段による推定結果に応じて、該当する前記話者用または前記話者属性用の前記音響モデルを前記音響モデル記憶部から読み出し、この読み出された前記音響モデルを用いて前記音声の認識処理を行なうものであり、さらに、前記話者属性推定手段による推定結果に対応する前記学習データを用いて、前記話者属性別発話モデル記憶手段に記憶された、当該推定結果に対応する前記話者ごとまたは前記話者属性ごとの話者属性別発話モデルを更新する話者属性別発話モデル適応化手段を備えることを特徴とする。

この構成によれば、話者属性別発話モデルを用いて話者又は話者属性を推定する。また、話者又は話者属性ごとの音響モデルを備えており、前記の推定結果に応じた音響モデルを用いた認識処理を行なう。この認識処理の結果得られる学習データを用いて話者属性別発話モデルの適応化を行なう。この処理を繰り返すことにより、話者属性別発話モデルの質が向上し、話者又は話者属性を推定する精度が向上し、即ち、話者又は話者属性ごとの音響モデルに適した音声区間を切り出す精度が向上する。よって、認識精度が上がり、より多量の学習モデルを獲得できる。

【0025】

[A5] 本発明の一態様は、上記の音響処理装置において、前記音響モデル適応化部は

、前記話者属性推定手段による推定結果に対応する前記学習データを用いて、当該推定結果に対応する前記話者ごとまたは前記話者属性ごとの前記音響モデルを更新するものであることを特徴とする。

この構成によれば、前記の推定結果に対応した学習データを用いて、当該推定結果に対応した音響モデルを適応化する。よって、音響モデルの質がより向上し、認識精度が上がり、より多量の学習モデルを獲得できる。

【発明の効果】

【0026】

本発明によると、字幕が付与されている様々な番組音声から、大量の音響モデルの学習データを効率的に得ることができる。また、字幕と実際の音声との一致率が低い番組音声からも、より効率的に、音響モデルの学習データを得ることが可能となる。特に、字幕と音声の一致しない区間の一部からも学習データを得ることができる。また、これまで認識できなかった番組や話者の認識が可能になったり、音声認識の認識精度の向上が可能になったりする。

10

【発明を実施するための最良の形態】

【0027】

以下、図面を参照しながら本発明の実施形態について説明する。

[第1の実施の形態]

図1は、第1の実施形態による音響モデル学習装置の機能構成を示すブロック図である。図示するように、音響モデル学習装置100は、字幕テキスト記憶部1と、形態素解析処理部2と、重み付け処理部3と、ベースラインコーパス記憶部4と、合成処理部5と、バイアス化言語モデル記憶部6（言語モデル記憶部）と、重み付け処理部7と、番組音声記憶部8と、音響モデル記憶部9と、認識処理部10と、認識仮説記憶部11と、アライメント処理部12と、選択処理部13と、学習データ記憶部14と、音響モデル適応化部15とを含んで構成される。

20

【0028】

字幕テキスト記憶部1は、放送番組の字幕テキストのデータを記憶する。形態素解析処理部2は、字幕テキスト記憶部1から読み出した字幕テキストについて、既存技術を用いて形態素解析処理を行い、その結果を、形態素に区切られた列として出力する。

【0029】

重み付け処理部3は、形態素解析された字幕テキストについて、ベースラインコーパスと合成するための重み付けを行なう。つまり、重み付け処理部3は、字幕テキストの出現頻度に重み値 $W_{bias}$ を乗ずるようなバイアスをかける。ベースラインコーパス記憶部4は、ベースラインコーパスを記憶する。ベースラインコーパスは、電子化された自然言語の大量の文章からなるデータベースである。例えば、蓄積された過去のニュース原稿のテキストをベースラインコーパスとして利用する。なお、ベースラインコーパスに対して予め統計処理を行い、後述する言語モデルに合う形式で記憶しておくようにしても良い。

30

【0030】

合成処理部5は、重み付け処理部3によって重み付けされた字幕テキストのデータと、ベースラインコーパス記憶部4から読み出したベースラインコーパスとを用いて、これらを合成し、出現する語に関する統計的な数値とともに表わしたバイアス化言語モデルを作成する。バイアス化言語モデル記憶部6は、合成処理部5によって作成されたバイアス付きの言語モデルを記憶する。言語モデルの詳細については、後で図面を参照しながら説明する。

40

【0031】

重み付け処理部7は、認識処理部10による認識処理のために、バイアス化言語モデル記憶部6に記憶されている言語モデルに対して重み値 $W_{lm}$ による重み付けを行なう。これにより、認識処理部10は、言語モデルと音響モデルとを $W_{lm} : 1$ の重みの比で用いることとなる。

【0032】

50

番組音声記憶部 8 は、番組音声を記憶している。この番組音声は、字幕テキスト記憶部 1 に記憶される字幕テキストのデータと対応するものである。

音響モデル記憶部 9 は、音素とその音素に対応する音響特徴量とを関連付けた音響モデルのデータを記憶する。音響モデルのデータについては後でも説明する。

【 0 0 3 3 】

認識処理部 10 は、音響モデル記憶部 9 から読み出した音響モデルと重み値  $W_{1m}$  による重み付けをつけた言語モデルとを用いて、番組音声記憶部 8 から読み出す音声の認識処理を行い、認識仮説（認識結果データ）を出力する。認識仮説記憶部 11 は、この認識仮説を記憶する。

【 0 0 3 4 】

アライメント処理部 12 は、形態素解析処理部 2 によって形態素解析処理済みの字幕テキストデータと、認識仮説記憶部 11 から読み出す認識仮説との一致部分を見つける処理を行なう。具体的には、アライメント処理部 12 は、これら両データをサーチし、所定数以上連続して語が一致しているか所定数以上連続して文字が一致している箇所を見つける処理を行なう。

【 0 0 3 5 】

選択処理部 13 は、アライメント処理部 12 によってアライメントされたデータに基づいて、認識仮説の各パスに信頼度データを付加する処理を行なう。

一例として、選択処理部 13 は、最尤パスと字幕テキストデータとが一致する区間（およびその区間に含まれる音素）には信頼度 1（最高信頼度を表わす）を与え、その他の区間（およびその区間に含まれる音素）には信頼度 0（最低信頼度を表わす）を与える。そして、選択処理部 13 は、信頼度データが付加された認識結果データ（これが教師系列データであり、この系列中に音素が含まれる）と、これに対応する音響特徴量データとを、対応付けて出力する。具体的には、選択処理部 13 は、このデータを学習データ記憶部 14 に書き込む。

【 0 0 3 6 】

学習データ記憶部 14 は、選択処理部 13 によって出力される学習データを記憶する。音響モデル適応化部 15 は、学習データ記憶部 14 から読み出した学習データを用いて、音響モデル記憶部 9 に記憶される音響モデルを適応化する（更新する）処理を行なう。

【 0 0 3 7 】

この音響モデル学習装置 100 の特徴は、音響モデル適応化部 15 が存在し、選択処理部 13 によって出力される学習データを用いて、最大事後推定法（MAP）や最尤線形回帰分析法（MLLR）等の適応化技術により、音響モデル記憶部 9 に記憶されている音響モデルを適応化（更新）し、さらに適応化された音響モデルを用いて認識処理を繰り返し行なうようにした点である。

【 0 0 3 8 】

また、従来技術と異なり、比較的大きな重み値  $W_{bias}$ （1000 程度）と重み値  $W_{1m}$ （16 程度）を用いて認識処理を行なう点も、特徴的である。字幕と発話の一致率が低く、自由発話の比率が高い番組では、上記のような重み値を用いることが、認識精度の向上を図りながらより多くの学習データを得られることにつながる。

これは、実験結果からも得られた適切な重み値である。具体的には、本願発明者らが、NHK（日本放送協会）の実際に番組の音声および字幕テキストを用いて行なった実験では、 $W_{bias} < 5000$  の領域において、 $W_{1m}$  の変化による単語誤認識率（WER）の差が小さい。また、 $W_{1m} < 18$  の領域において、 $W_{bias}$  の変化による単語誤認識率の差が小さい。そして、 $W_{bias} < 5000$  且つ  $W_{1m} < 18$  の場合に、単語誤認識率がそれほど上がらずに、且つ得られる学習データの量が多くなる。そして、 $W_{bias}$  が 1000 程度で  $W_{1m}$  が 16 程度のときに、特に、単語誤認識率がそれほど上がらずに、且つ得られる学習データの量が多くなる。

【 0 0 3 9 】

次に、言語モデルについて説明する。

10

20

30

40

50

図 2 は、バイアス化言語モデル記憶部 6 によって記憶され認識処理部 10 によって認識処理の際に使用される言語モデルのデータの構成を示す概略図である。

図 2 ( a ) は、形態素単体での出現確率を表わすテーブルを示す。このテーブルは、表形式のデータであり、形態素と出現確率の項目 ( 列 ) を有する。テーブルの行は、形態素の種類ごとに存在する。各形態素と、その形態素の出現確率とが関連付けられてこのテーブルに保持される。

図 2 ( b ) は、形態素と形態素の接続確率を表わすテーブルである。このテーブルは、表形式のデータであり、形態素 1 と形態素 2 と接続確率の項目 ( 列 ) を有する。テーブルの行は、形態素 1 と形態素 2 の組み合わせごとに存在する。各行は、形態素 1 が出現したときその形態素 1 に直ちに続いて形態素 2 が出現する ( 接続する ) 確率である接続確率の情報を表わしている。

10

#### 【 0 0 4 0 】

次に、音響モデルについて説明する。

図 3 は、音響モデル記憶部 9 によって記憶され認識処理部 10 によって認識処理の際に使用される音響モデルのデータの構成を示す概略図である。図示するように、音響モデルは、音素と、その音素が発話されるときの音響的特徴量とを関連付けて保持するデータである。音響的特徴量としては、例えば、10 ミリ秒ごとの間隔でそれぞれ開始する 25 ミリ秒のフレーム内の周波数パワー分布を基にした M F C C (メル周波数ケプストラム係数, Mel-Frequency Cepstrum Coefficient) や、P L P (Perceptual Linear Predictive) などを用いる。

20

#### 【 0 0 4 1 】

次に、認識処理部 10 による認識結果である認識仮説について説明する。

図 4 は、認識処理部 10 によって出力され認識仮説記憶部 11 によって記憶される認識仮説のデータ構成およびデータ例を示す概略図である。認識仮説は、論理的には、言語的単位 ( 形態素 ) をアークとするラティス構造 ( 半順序 ( partial order ) 構造 ) のグラフで表現される。

#### 【 0 0 4 2 】

図 4 ( a ) は、表形式で表現される認識仮説のデータ構成を示す。図示するように、表形式で表わした認識仮説は、各行がラティス上の各アークに対応し、始点ノード、終点ノード、言語的単位、信頼度の各項目 ( 列 ) を有する。始点ノードおよび終点ノードは、ラティス上のノードに便宜的に与えられたノード識別情報を値として持つ。言語的単位は、そのアークに対応する形態素である。信頼度は、認識結果におけるそのアークの信頼度 ( 言語モデルおよび音響モデルを基にした統計的処理で得られる認識結果の信頼度 ) を表わす数値で表わされるものであり、その数値の範囲は 0 以上 1 以下である。図示しているデータ例は、例えば 1 行目のデータに相当するアークの始点ノードは「 1 」であり、終点ノードは「 2 」であり、言語的単位は「 F 」である。また例えば 2 行目のデータに相当するアークの始点ノードは「 1 」であり、終点ノードは「 2 」であり、言語的単位は「 H 」である。3 行目以下のデータも同様である。同図では信頼度数値の記載を省略している。

30

#### 【 0 0 4 3 】

図 4 ( b ) は、図 4 ( a ) の表形式のデータが表現するラティスを絵的に表現した概略図である。同図における丸印がそれぞれノードに相当し、アークはノード間を結ぶ矢印付きの線で表わされている。また、各アークに対応する言語的要素が付記されている。例えば一番左のノード ( ノード「 1 」 ) を始点とするアークは 3 本あり、それらはそれぞれ、言語的要素が「 F 」で終点がノード「 2 」 ( 左から二番目のノード ) 、言語的要素が「 H 」で終点がノード「 2 」、言語的要素が「 M 」で終点がノード「 3 」 ( 左から三番目のノード ) である。

40

#### 【 0 0 4 4 】

次に、音響モデル学習装置 100 の動作および処理手順について説明する。

図 5 は、音響モデル学習装置 100 全体の処理手順を示すフローチャートである。

図示するように、ステップ S 0 1 において、まず形態素解析処理部 2 が、字幕テキスト

50

記憶部 1 から読み出した字幕テキストの形態素解析処理を行なう。形態素解析処理は、既存の技術により可能である。

【 0 0 4 5 】

次に、ステップ S 0 2 において重み値  $W_{bias}$  を用いて形態素解析処理結果に重み付けをするとともに、ステップ S 0 3 において重み付けされた形態素解析処理結果とベースラインコーパスとを合成する処理を行なう。具体的には、合成処理部 5 は、形態素解析処理部 2 による処理結果を統計処理し、各形態素の出現確率を算出するとともに、2 つの連続する形態素間の接続確率を算出し、図 2 ( a ) に示した形式の出現確率データおよび図 2 ( b ) に示した形式の接続確率データを得る。また、合成処理部 5 は、ベースラインコーパス記憶部 4 に記憶されているベースラインコーパスを基に、必要に応じて同様の統計

10

【 0 0 4 6 】

次に、ステップ S 0 4 において重み値  $W_{lm}$  を用いてバイアス化言語モデル記憶部 6 から読み出される言語モデルに重み付けするとともに、ステップ 0 5 において認識処理部 1 0 は、重み付けされたバイアス化言語モデルと音響モデル記憶部 9 から読み出した音響モデルを用いて、番組音声記憶部 8 に記憶されている番組音声の認識処理を行い、その結果

20

【 0 0 4 7 】

次に、ステップ S 0 6 において、アライメント処理部 1 2 は、形態素解析処理された字幕テキストと認識処理の結果得られた認識仮説とのアライメント処理を行なう。このアライメント処理は、両データが単語単位で一致する部分を探索することによって行ない、或いは両データ中の複数の単語が連続して一致する区間を探索することによって行なう。アライメント処理の結果、字幕テキストの中の区間と番組音声の中の区間がデータ的に対応付けられる。

30

【 0 0 4 8 】

次に、ステップ S 0 7 において、選択処理部 1 3 は、認識仮説に信頼度を付加する処理を行なう。具体的な処理の一例は次の通りである。選択処理部 1 3 は、まずアライメント処理部 1 2 によるアライメント結果を用いて、元の字幕テキストと認識仮説に含まれる最尤パスとが 3 単語以上連続して一致している一致区間を選択し、その一致区間（およびその区間に含まれる各音素）に対しては、信頼度 1 を与える。そして、選択処理部 1 3 は、その他の区間に対しては信頼度 0 を与える。そして、選択処理部 1 3 は、そのようにして得られた学習データを、学習データ記憶部 1 4 に書き込む。

【 0 0 4 9 】

次に、ステップ S 0 8 において、繰り返し処理を終了するか否かを判定する。この判定は、例えば、認識結果の精度が十分かどうかにより行なう。より具体的には、この判定は ( a ) ステップ S 0 5 ~ S 0 9 のループを繰り返した回数（例えば、この回数を 3 ~ 4 回として良い場合もある）、( b ) 番組音声全体のうちの選択処理部 1 3 によって選択された一致区間の長さの比率、( c ) 番組音声全体のうちの選択処理部 1 3 によって選択された一致区間の前回からの増分比率、などのいずれかによって行なう。

40

判定の結果、繰り返し処理を終了しない場合（ステップ S 0 8 : N O ）には、次のステップ S 0 9 に進む。

判定の結果、繰り返しを終了する場合（ステップ S 0 8 : Y E S ）には、このフローチャート全体の処理を終了する。このとき、学習データ記憶部 1 4 に累積的に書き込まれたデータが、本装置による学習処理の結果として得られた音響モデルである。

50

## 【 0 0 5 0 】

ステップ S 0 9 においては、音響モデル適応化部 1 5 は、得られた学習データを用いて、音響モデル記憶部 9 に記憶されている音響モデルを適応化する（更新する）処理を行なう。なお、音響モデル適応化部 1 5 が如何に学習データを用いて具体的に音響モデルの更新を行うかについては、後で詳述する。

## 【 0 0 5 1 】

上記のように、音響モデル学習装置 1 0 0 は、選択処理部 1 3 によって出力される学習データを用いて、音響モデル適応化部 1 5 が、音響モデル記憶部 9 に記憶されている音響モデルを適応化し、さらに適応化された音響モデルを用いて認識処理を繰り返し行なうようにしている。

10

## 【 0 0 5 2 】

次に、上述した音響モデル適応化部 1 5 による音響モデルの適応化処理について詳しく説明する。本実施形態では、音声認識による音声認識区間  $r$  をそのまま用いて音響モデルを学習するため、認識結果の各形態素や音素  $j$  の信頼度を導入する。音響モデル適応化部 1 5 は、フォワードバックワードアルゴリズムを用いて、各時刻における音素  $j$  の観測確率に、信頼度による重みを与えて最尤推定を行ない、これにより得られる平均ベクトルを用いて音響モデル記憶部 9 に記憶されている音響モデルの更新を行なう。最尤推定により得られる音素（状態） $j$  の平均ベクトルは、下の式（1）により計算される。

## 【 0 0 5 3 】

【数 1】

20

$$\hat{\mu}_j = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} C_j^r L_j^r(t) o_t^r}{\sum_{r=1}^R \sum_{t=1}^{T_r} C_j^r L_j^r(t)} \quad \dots \quad (1)$$

## 【 0 0 5 4 】

【数 2】

30

$$o_t^r$$

## 【 0 0 5 5 】

は、音声認識による音声切り出し区間  $r$  の時刻  $t$  における特徴量の観測ベクトルである。また、

## 【 0 0 5 6 】

【数 3】

40

$$L_j^r(t)$$

## 【 0 0 5 7 】

は、音声切り出し区間  $r$  の時刻  $t$  における音素  $j$  の観測確率である。なお、音響モデル適応化部 1 5 が、動的計画法の一種であるフォワードバックワードアルゴリズムによる手順を実行することにより、この観測確率を得ることができる。

また、

50

【 0 0 5 8 】

【 数 4 】

$$C_j^r$$

【 0 0 5 9 】

は、音声切り出し区間  $r$  における音素  $j$  の信頼度である。

本実施形態では、選択処理部 13 は、次の式 (2) により信頼度を与えている。

10

【 0 0 6 0 】

【 数 5 】

$$C_j^r = \begin{cases} 1 & \text{(最尤仮説と字幕の一致区間)} \\ 0 & \text{(その他の区間)} \end{cases} \quad \dots (2)$$

【 0 0 6 1 】

つまり、認識結果の最尤仮説のパスと元の字幕テキストとが一致する区間においては信頼度は 1 (信頼度最大を表わす) であり、その他の区間においては信頼度 0 (信頼度最低を表わす) とする。

20

【 0 0 6 2 】

図 6 は、本実施形態において選択処理部 13 が付加する信頼度の例を表わす概略図であり、最尤仮説のパスと書き起こしによる字幕テキストの一致を利用する方法による信頼度を表わしている。図 6 (a) は、書き起こしによる字幕テキストに対応するグラフを示す。図示するように、字幕テキストに対応するグラフは、形態素 a、b、c、d、e にそれぞれ対応するアークが直列につながって構成されている。図 6 (b) は、最尤仮説の教師単語列に対応するグラフを示す。図示するように、最尤仮説に対応するグラフは、形態素 a、f、c0、c1、q、e' にそれぞれ対応するアークが直列につながって構成されている。また、図 6 (b) に付記されている括弧内の数値は、信頼度である。本例において両者を比較すると、形態素 a および d が一致しており、その他の形態素は不一致である。したがって、選択処理部 13 は、形態素 a および d の区間にはそれぞれ信頼度 1.0 を与え、その他の区間には信頼度 0.0 を与えている。

30

そして、選択処理部 13 は、各々の区間の信頼度を、その区間に含まれる音素の信頼度として付加する。

【 0 0 6 3 】

本実施形態による音響処理装置の特徴は、音声認識による音声切り出し区間  $r$  をそのまま用いて音響モデルの学習を行なうため、認識結果に含まれる状態 (単語や音素)  $j$  の信頼度を導入した点にある。そして、音響モデル適応化部 15 が音響モデルの統計量の推定を行なう際に、フォワードバックワードアルゴリズムにより得られる音素  $j$  の観測確率に、信頼度による重みを与えている。そして、音響モデル適応化部 15 は、上記の式 (1) による計算で得られる音素  $j$  の平均ベクトルを用いて、音響モデル記憶部 9 に記憶されている音響モデルを更新する。

40

【 0 0 6 4 】

また、式 (2) で表わされる信頼度を付与することにより、認識結果の最尤パスと字幕テキストの一致区間に含まれる音素に対応付けられた特徴量のみを、音響モデルの学習に利用することができる。

【 0 0 6 5 】

なお、式 (1) による計算では最尤推定に基づく特徴量の平均ベクトルを得てこれを音

50

響モデルの適応化に用いたが、その代わりに、最尤推定に基づく分散や混合重みなどを用いて音響モデルを更新するようにしても良い。さらに、最尤推定に限らず、MAP推定 (maximum a posteriori estimation、最大事後確率推定) やMLLR (Maximum Likelihood Linear Regression、最尤線形回帰) 推定などを用いるようにしても良い。

【0066】

ここで、仮に、

【0067】

【数6】

$$C_j^r$$

10

【0068】

をすべての区間において等しく1にすると、これは、認識結果と字幕テキストの一致に関する情報を全く利用しない教師なし学習となる。

【0069】

また、上記の方法では、フォワードバックワードアルゴリズムにより得られる観測確率を用いたが、その代わりに、例えばビタビアルゴリズムを用いて得られる最尤パス (ビタビパス) を利用するようにしても良い。この場合は、最尤パス上の観測確率には1を与え、その他のパス上の観測確率には0を与えるようにする。

20

【0070】

図7は、音響モデル学習装置100による音響モデルの学習の概略を補足的に説明するための図である。

図7において、形態素解析結果は字幕テキストを形態素解析して得られたデータであり、このデータは形態素 (単語) の一次元的な列である (図中のそれぞれの四角が形態素に相当)。また、認識結果は、認識処理部10による音声認識処理の結果得られる認識仮説のラティスに含まれるパスの一部に相当するデータである (図中のそれぞれの四角が形態素に相当)。そして、この図においては形態素解析結果と最尤パスとが一致する一致区間をハッチングで表わしている。また、これらの各々の区間には、音響特徴量および信頼度が対応している。本実施形態では、上記の一致区間の信頼度は1.0であり、その他の区間の信頼度は0.0である。そして、各区間に含まれる音素と、対応する音響特徴量と、対応する信頼度とを関連付けたものが、選択処理部13によって得られる学習データである。

30

【0071】

本実施形態では、認識処理の結果である学習データを用いて、音響モデルの適応化を行っている。そして、適応化された音響モデルを用いて再度認識を行い、学習データを生成する。この音響モデルの適応化と、適応化された音響モデルを用いた認識処理を繰り返すことにより、学習データの正確性が増し、より多くの量の学習データを効率的に得ることができる。

40

【0072】

[第2の実施の形態]

次に第2の実施形態について説明する。本実施形態では、認識結果の単語ラティスを用いて、最尤パスと字幕テキストとの一致区間以外の形態素をも学習データとして利用する。

なお、以下では、本実施形態特有の部分のみを説明し、その他の部分については前記の実施形態1と同様であるため説明を省略する。

【0073】

本実施形態においては、選択処理部13は、認識仮説のラティスの枝刈り処理を行い、

50



枝刈りの後のラティス内の最尤パスの各形態素の事後確率を求め、この事後確率を各形態素および形態素内の各音素の信頼度として付与する。ここで、枝刈りの対称となるのは、ある時間において最尤パスと対立するすべてのパスである。このような枝刈り処理の結果として、最尤パスと字幕テキストが一致する区間の各形態素の信頼度は1（信頼度最大を表わす）となる。つまり、最尤パスと字幕テキストが一致する区間の形態素については、第1の実施形態と同様の学習を行なうこととなる。

【0074】

次に、選択処理部13による認識仮説の枝刈りの処理についてより詳細に説明する。

図8は、選択処理部13によって枝刈りされた後の認識仮説のデータ例を示す概略図である。そして、図4に示したデータが、選択処理部13によって枝刈りされる前の認識仮説のデータである。

まず、選択処理部13は、図4に示した認識仮説のデータと字幕テキストのデータとの一致区間を選択する。ここでは選択処理部13は1形態素以上連続して両者が一致する区間を一致区間として選択するものとする。その結果、形態素B（始点ノードが「2」で終点ノードが「3」）と形態素C（始点ノードが「3」で終点ノードが「4」）と形態素E（始点ノードが「5」で終点ノードが「6」）が一致区間として選択される。

【0075】

その後、選択処理部13は、上で選択された一致区間の時間を含み且つ字幕テキストのデータとは一致しない認識仮説を枝刈りする。具体的には、図4に示した形態素Bの区間の時間を含み且つ形態素Bとは一致しない認識仮説（言い換えれば、形態素Bと対立する認識仮説）である形態素IとMとNは、選択処理部13による枝刈りの対象となる。同様に、一致区間として選択された形態素Cと対立する認識仮説である形態素JとNとOは、選択処理部13による枝刈りの対象となる。同様に、一致区間として選択された形態素Eと対立する認識仮説である形態素LとPは、選択処理部13による枝刈りの対象となる。

つまり、図4に示すラティスのうち、形態素IとJとLとMとNとOとPのそれぞれのアークが選択処理部13によって枝刈りされる。その枝刈りの結果として残るのは、形態素BとCとEとFとGとHとKのそれぞれのアークであり、これらのアークに対応するデータが、図8に示すデータである。

【0076】

図8(a)は、そのようなラティス構造のグラフを表わす表形式のデータであり、この表は、前記の形態素BとCとEとFとGとHとKに対応する行のデータを有している。またこの表において、形態素BとCとEは、それぞれ一致区間に含まれていたものであるため、その結果として、選択処理部13は、これらの形態素に対応する信頼度がいずれも最高の1となるように認識仮説のデータを更新している。また、その他の、形態素FとHとKとGのアークについては、選択処理部13は、枝刈り後の事後確率を算出し、その算出結果をそれぞれの信頼度データとして認識仮説のデータを更新する。図示する例では、形態素Fの信頼度は0.1、形態素Hの信頼度は0.9、形態素Kの信頼度は0.2、形態素Gの信頼度は0.8である。

図8(b)は、図8(a)の表形式のデータが表現するラティスを絵的に示した概略図である。

【0077】

なお、選択処理部13が事後確率を算出する方法は次の通りである。選択処理部13は、動的計画法の一種であるフォワードバックワードアルゴリズムを用いることで、デコーダー（認識処理部10）の出力から得られる形態素mのスコア $L(m)$ を基に、事後確率を算出する。ここで、 $L(m)$ としては、音響モデルの尤度、言語モデルの尤度、あるいはそれらの和を用いる。まず、形態素mの前向き確率 $P(m)$ と後ろ向き確率 $P(m)$ とを以下の式(3)~(7)の通り定義する。

【0078】

10

20

30

40

【数 7】

$$\alpha(b) = 1.0 \quad \dots (3)$$

【0079】

【数 8】

$$\beta(e) = 1.0 \quad \dots (4)$$

10

【0080】

【数 9】

$$\alpha(m) = \sum_{m' \in LEFT(m)} \alpha(m') L(m) \quad \dots (5)$$

20

【0081】

【数 10】

$$\beta(m) = \sum_{m' \in RIGHT(m)} \beta(m') L(m) \quad \dots (6)$$

【0082】

【数 11】

30

$$Z = \alpha(e) = \beta(b) \quad \dots (7)$$

【0083】

ここでは、事後確率を計算するため、図4(b)などに示したラティスの両端に、それぞれ、始端単語と終端単語に相当するアークを追加する。始端単語はその図のラティスの左端に追加され、「b」と表わされる。終端単語はその図のラティスの右端に追加され、「e」と表わされる。ここで前向きとは、始端から終端の方向であり、図のラティスでは左から右の方向である。また後ろ向きとは、終端から始端の方向であり、図のラティスでは右から左の方向である。また、LEFT(m)は、形態素mに左から(つまり始端側から)接続する形態素の集合である。また、RIGHT(m)は、形態素mに右から(つまり終端側から)接続する形態素の集合である。

40

【0084】

式(3)に表わすように、始端単語(b)の前向き確率は1.0である。また、式(4)に表わすように、終端単語(e)の後ろ向き確率は1.0である。

式(5)に表わすように、形態素mの前向き確率は、形態素mに左から接続する各形態素m'の前向き確率(m')に当該形態素mのスコアL(m)を乗じた値の、左から接続する全ての形態素についての総和である。このように前向き確率は再帰的な定義となっているが、始端単語から前向きに順次計算していくことにより、ラティス中の全ての形態

50

素の前向き確率を算出できる。

式(6)に表わすように、形態素 $m$ の後ろ向き確率は、形態素 $m$ に右から接続する各形態素 $m'$ の後ろ向き確率 $\beta(m')$ に当該形態素 $m$ のスコア $L(m)$ を乗じた値の、右から接続する全ての形態素についての総和である。このように後ろ向き確率は再帰的な定義となっているが、終端単語から後ろ向きに順次計算していくことにより、ラティス中の全ての形態素の後ろ向き確率を算出できる。

式(7)に表わすように、終端単語の前向き確率および始端単語の後ろ向き確率を $Z$ とする。

【0085】

そして、形態素 $m$ の事後確率 $\gamma(m)$ は、 $\alpha(m)$ 、 $\beta(m)$ 、 $Z$ を用いて、式(8)により算出することができる。

【0086】

【数12】

$$\gamma(m) = \frac{\alpha(m)\beta(m)L(m)}{Z} \quad \dots \quad (8)$$

【0087】

以上のように、選択処理部13は、認識結果に含まれる最尤パスと字幕テキストとの一致区間と時間軸上で対立する区間を枝刈りする処理を行ない、この枝刈り処理の結果を教師系列データとして取得し、この教師系列データに含まれる形態素ごとの事後確率をその形態素に含まれる音素の信頼度データとして付加する。

そして、音響モデル適応化部15は、選択処理部13が得た信頼度データを用いて、第1の実施形態と同様に、例えば前記の式(1)を用いて音素に対応するベクトルの更新式を得て、これにより音響モデル記憶部9に記憶された音響モデルを更新する。

【0088】

[第3の実施の形態]

次に第3の実施形態について説明する。本実施形態では、選択処理部13が、第2の実施形態における枝刈り処理と同様の処理を行なう。そして、枝刈り後の認識結果のラティスに対して、直接フォワードバックワードアルゴリズムによる処理を行い、各音素の統計量を学習する。この場合、フォワードバックワードアルゴリズムの処理より得られる観測確率に、既に信頼度が含まれる。

【0089】

第2の実施形態の場合と同様に本実施形態においても、選択処理部13は、枝刈り処理を行なう結果、最尤パスと字幕テキストとの一致区間では信頼度1(信頼度最大を表わす)を付与する。一方で、最尤パスと字幕テキストが一致しない区間においては、選択処理部13は、各々のパスの事後確率に応じた信頼度を付与する。

【0090】

そして、音響モデル適応化部15は、得られている観測ベクトルと、選択処理部13により求められた観測確率とに基づき、下の式(9)を用いて音素 $j$ の平均ベクトルを計算する。

【0091】

10

20

30

40

【数 1 3】

$$\hat{\mu}_j = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} L_j^r(t) o_t^r}{\sum_{r=1}^R \sum_{t=1}^{T_r} L_j^r(t)} \quad \dots \quad (9)$$

【0092】

この式(9)においては、

【0093】

【数 1 4】

$$L_j^r(t)$$

【0094】

は、区間  $r$  の時刻  $t$  における音素  $j$  の観測確率であり、この観測確率は上述したフォワードバックワードアルゴリズムによる処理を用いた方法で選択処理部 13 によって算出されたものであり、信頼度が含まれている。

【0095】

本実施形態で上記の平均ベクトルを計算する場合、式(9)からも明らかなように、ある時刻において対立する複数のパス(候補)間で、共通して出現する度合いの高い音素については、対立候補それぞれにおける当該音素の観測確率がすべて加算されるため、当該音素に対応する観測ベクトルの影響が強く効いた平均ベクトルが算出される。一方で、同じく式(9)からわかるように、対立する複数のパスのいずれかにしか現れない音素については、そのパス中の音素に対応する観測ベクトルが平均ベクトルには強い影響を与えない。

【0096】

図9は、対立する複数のパス間で共通して出現する音素の例を示す概略図である。この図に沿って本実施形態の特徴を補足説明する。この図における時間軸は左から右への方向である。また、この図に示す認識結果のラティスは対立する2系列のパスを含んでいる。図における上側のパスは、音素  $j_1$ 、 $j_2$ 、 $j_3$  からなる形態素と、音素  $j_4$ 、 $j_5$ 、 $j_6$  からなる形態素を含んでいる。また、下側のパスは、音素  $j_7$ 、 $j_2$ 、 $j_3$  からなる形態素と、音素  $j_8$ 、 $j_9$ 、 $j_{10}$  からなる形態素を含んでいる。ここで、各々の音素の観測確率が時刻  $t$  をパラメタとする関数になることは式(9)にも表わされているとおりであるが、上記の音素  $j_1 \sim j_{10}$  のうち  $j_2$  と  $j_3$  だけは上側のパスと下側のパスの両方に共通して出現している。従って、式(9)によれば、このような音素  $j_2$  および  $j_3$  は、認識結果のラティス全体の中の図示している部分から学習データを得るのに、他の音素( $j_1$  と、 $j_4 \sim j_{10}$ ) に比べてより大きな重みで作用する。

【0097】

以上述べたように、本実施形態では、選択処理部 13 は、認識結果データに含まれる最尤パスと字幕テキストのデータが一致する一致区間と時間軸上で対立する区間を枝刈りする処理を行ない、この枝刈り後の認識結果データに含まれる形態素ごとの事後確率を当該形態素に含まれる各々の音素に対する信頼度データとして実質的に付加し、この信頼度データが付加された教師系列データ(ラティス)とそれに対応する音響特徴量データとを、学習データとして出力する処理を行なう。

そして、音響モデル適用化部 15 は、学習データ記憶部 14 からこの学習データを読み

10

20

30

40

50

込み、フォワードバックワードアルゴリズムを用いて、教師系列データ（ラティス）に含まれる各音素の観測確率を求め、この観測確率で重み付けられた観測ベクトルの平均ベクトルを音素ごとに算出し、得られた平均ベクトルを音響モデルとして用いて音響モデル記憶部に記憶されている音響モデルを適応化する。

【 0 0 9 8 】

このように、本実施形態による方法では、対立する複数の候補間で、含まれる音素のバリエーションに応じた結果が得られることとなる。

【 0 0 9 9 】

[ 第 4 の実施形態 ]

次に、第 4 の実施形態を説明する。第 1 ~ 第 3 の実施形態が最尤パスと字幕テキストが一致する区間を基準として教師系列を求めていたのに対し、本実施形態の特徴は、認識結果のラティスと字幕テキストが一致する区間を基準として教師系列を作成する点である。

本実施形態では、選択処理部 13 は、認識結果のラティスと字幕テキストとのビタピアライメントを取ることにより、認識結果のラティスと字幕テキストが一致する区間を基に教師系列を得る。

【 0 1 0 0 】

図 6 が最尤パスと字幕テキストとが一致する区間を用いた方法を説明しているのに対して、図 10 は、本実施形態による選択処理部 13 が、認識結果のラティスと字幕テキストの一致区間を抽出する処理を概念的に説明する概略図である。図示する例では、( a ) の書き起こしの字幕テキストデータは、a - b - c - d - e という 5 つの形態素に対応するアークが直列する形のグラフによって表わされる。一方で、( b ) の認識結果のラティスは、複数の対立候補のアーク（各々のアークは形態素に対応する）を含んだ構造を有しており、字幕テキストデータに対応する形態素 a、b、c、e の他にも、形態素 c<sub>0</sub> や c<sub>1</sub> や e' や、その他（f、g、h、i、j、k、l、m、n、・・・、以下記載省略）の形態素に対応するアークを含んでいる。そして、( c ) は、( a ) の字幕テキストと ( b ) の認識結果のラティスを基に、これらの一致区間を抽出して得られた教師形態素列を示している。( c ) の例は、a - b - c - g - e という 5 つの形態素に対応するアークが直列する形のグラフを表わしている。ここで、形態素 a、b、c、e は字幕テキストと認識結果のラティスが一致する区間であるため、選択処理部 13 は、それらに信頼度 1.0 を付与している。また、字幕テキストにおける形態素 d に一致するパスが認識結果のラティス内に存在しないため、選択処理部 13 は、それに対応する時間の区間からは形態素 g を選択するとともに、それに信頼度 0.1 を付与している。ここで、( b ) に示した認識結果のラティスにおける最尤パスは、形態素 c<sub>0</sub> - c<sub>1</sub> の系列を含み、また形態素 e' の系列を含んでいるが、これら c<sub>0</sub>、c<sub>1</sub>、e' はいずれも字幕テキストと一致するものではない。そこで、選択処理部 13 は、c<sub>0</sub> - c<sub>1</sub> の系列に対立する区間から c を選び、また、e' に対立する区間である e を選んでいる。これは、前述の通り、c および e がそれぞれ字幕テキストと一致するためである。

【 0 1 0 1 】

上記の処理に具体例を当てはめると次の通りである。

例えば、形態素 c は「ですから」であり、形態素 c<sub>0</sub> - c<sub>1</sub> の系列は「です / から」に対応する。本実施形態の方法では、認識結果における最尤パスが「です / から」を含むものであっても、字幕テキストと一致する「ですから」を選択して教師形態素列を作る。つまり、最尤パスにおける形態素分割の結果が字幕テキストと異なることに起因して不一致区間と判定されてしまう区間を、一致区間と同等に扱うことができるようになる。

また例えば、形態素 e は「こと」に対応し、形態素 e' は「事」に対応する。本実施形態の方法では、認識結果における最尤パスが「事」を含むものであっても、字幕テキストと一致する「こと」を選択して教師形態素列を作る。つまり、表記の揺らぎに起因して不一致区間と判定されてしまう区間を、一致区間と同等に扱うことができるようになる。

【 0 1 0 2 】

このように、本実施形態の方法では、選択処理部 13 は、最尤仮説と字幕テキストとの

間では不一致となるものの、ラティス内の、その最尤仮説に対立する仮説の中に字幕テキストと一致するテキストがあると期待される場合に、認識結果のラティスが含む情報を最大限に活かすことが可能となる。

【 0 1 0 3 】

[ 第 5 の実施形態 ]

次に、第 5 の実施形態を説明する。第 5 の実施形態は、第 4 の実施形態と同様に、認識結果のラティスと字幕テキストが一致する区間を基準として教師系列を作成する点である。

本実施形態では、選択処理部 1 3 は、認識結果のラティスと字幕テキストとのビタピアライメントを取ることにより、認識結果のラティスと字幕テキストが一致する区間を得る。

そして、選択処理部 1 3 は、得られた一致区間と対立する区間の枝刈り処理を行なう。ここで、枝刈り処理の手法そのものは第 3 の実施形態におけるそれと同様であるが、第 3 の実施形態が最尤パスと字幕テキストとの一致区間を基準としていたのに対し、本実施形態は、認識結果のラティスと字幕テキストが最大に一致する区間を基準とする。

そのように枝刈り後のラティスが得られた後の処理は、第 3 の実施形態における処理と同様である。つまり、音響モデル適用化部 1 5 は、選択処理部 1 3 が出力した学習データを読み込み、フォワードバックワードアルゴリズムを用いて学習データに含まれる各々の音素の観測確率を求め、音素ごとに、観測確率で重み付けられた音響特徴量データとして得られる音響モデルを用いて、前記音響モデル記憶部を更新する。

【 0 1 0 4 】

次に、第 1 ~ 第 5 のいずれかの実施形態を用いたさらなるバリエーションとして、複数の追加実施形態を説明する。要約すると、追加実施形態の 1 は、各々の番組用に、学習データを抽出する。また、追加実施形態の 2 は、言語モデルの重みを適宜変える。また、追加実施形態の 3 は、話者の交代を検出し、話者ごと、又は話者属性ごとの音響モデルを抽出する。

【 0 1 0 5 】

[ 追加実施形態の 1 ]

本発明の追加実施形態の 1 について説明する。

図 1 1 は、同実施形態による音響モデル学習装置 1 0 1 の機能構成を示すブロック図である。図示するように、音響モデル学習装置 1 0 1 は、内部に、音響モデル学習装置（番組 A 用）1 0 1 A と音響モデル学習装置（番組 B 用）1 0 1 B と音響モデル学習装置（番組 C 用）1 0 1 C と音響モデル学習装置（番組 D 用）1 0 1 D とを含んで構成されている。これら音響モデル学習装置 1 0 1 A ~ 1 0 1 D の各々は、それぞれ単一の番組用のものである。

【 0 1 0 6 】

ここで、番組とは、典型的にはテレビやラジオの放送番組であるが、公衆によって直接受信されることを目的とする無線通信の送信であるところの「放送」の番組に限らず、有線通信によるテレビ放送やラジオ放送の番組、インターネットを介してパケットデータとして配信される動画（音声を含む）による放送や音声放送の番組、通信衛星から多数の受信者に向けて送信される番組、視聴者或いは聴取者のリクエストに応じてオン・デマンド的に配信される番組など、音声を含むコンテンツを含んでいる。

なお、音響モデル学習装置 1 0 1 が内部に備える番組個別用の音響モデル学習装置の数は、4 に限定されず、任意の自然数として構成しても良い。

【 0 1 0 7 】

音響モデル学習装置 1 0 1 A ~ 1 0 1 D の各々は、図示するように、第 1 の実施形態で説明した音響モデル学習装置 1 0 0 と同様の構成・機能を有し、同様の動作をする。

本実施形態においても、選択処理部 1 3 は、アライメント処理部 1 2 によってアライメントされたデータを用いて、認識結果データに信頼度を付加し、学習データとして出力する。

10

20

30

40

50

## 【0108】

本実施形態のポイントは、音響モデル学習装置101A～101Dの各々が専用の番組音声記憶部8と音響モデル記憶部9を備え、個々の音響モデル記憶部9に各々専用の音響モデルを記憶し、この各々専用の音響モデルを用いて認識処理部10が認識処理を行い、その結果得られる学習データが各々専用の学習データ記憶部14に書き込まれ、この学習データが蓄積されて出力されるとともに、この得られた学習データを用いて音響モデル適応化部15が当該番組用の音響モデル記憶部9を適応化する点であり、そのように番組ごとの音響モデルを用いて番組ごとの学習データを得る処理を繰り返す点である。

言い換えれば、音響モデル記憶部9は音響モデルを番組ごとに記憶するものであり、認識処理部10は番組の音声当該番組に対応した音響モデルを用いて音声の認識処理を行なうものであり、選択処理部13は番組ごとに学習データを出力し、音響モデル適応化部15は、当該番組用の音響モデルを更新する。

10

## 【0109】

なお、音響モデル学習装置101A～101Dは、そのすべての構成要素をそれぞれが専用に備えるようにしても良いし、一部の構成要素を共用にしても良い。例えば、認識処理部10とアライメント処理部12と選択処理部13と音響モデル適応化部15の処理機能自体は共通のハードウェアを用いて、音響モデル記憶部9に記憶される音響データと認識仮説記憶部11に記憶される認識仮説と学習データ記憶部14に記憶される学習データとが、それぞれの番組専用の領域に分けて管理されるように装置を構成しても良い。

## 【0110】

本実施形態の利点は次の通りである。例えば、対談形式のテレビ放送番組の音声を認識させて学習データを得ることを考えた場合、そしてその番組の形式がレギュラー話者（毎回出演する話者）とゲスト話者（特定回のみ出演する話者）の2人による複数回にわたるシリーズ番組を対象とした場合（ここでは放送の各回が番組A～Dに相当）、異なるゲスト話者の学習データが重要であるならば、番組回ごとの音響モデルの適応化処理を繰り返して行なったほうが、複数回に共通の音響モデルを適応化処理するよりも、話者適応の効果が得られると言える。この場合には、各回から得られた学習データを集めて最終的な音響モデルの学習データを得る。

20

## 【0111】

さらに、上記のような対談形式の番組に限らず、放送番組の中には、限られた数人の話者しかいない番組や、特定の話者が継続してレギュラー出演している番組などがある。本実施形態では、音響モデルの適応化を番組ごとに行なうため、他の番組の学習データが混ざらず、その結果として高い認識精度が得られる。

30

なお、この実施形態は、前述した第1から第5のいずれかの実施形態と組み合わせて実現してよい。言い換えれば、第1から第5のいずれかの実施形態で述べた、選択処理部13および音響モデル適応化部15の処理を用いて、番組ごとに音響モデルの学習を行なうようにする。

## 【0112】

[追加実施形態の2]

次に、本発明の追加実施形態の2について説明する。上記の実施形態では、認識処理部10による認識処理の際の言語モデルと音響モデルそれぞれの重みに影響する重み値 $W_{1m}$ として固定的な値を用いていた。本実施形態では、音響モデルの適応度合いに応じて、この重み値 $W_{1m}$ を変化させる。

40

## 【0113】

なお、本実施形態においても、選択処理部13は、アライメント処理部12によってアライメントされたデータを用いて、認識結果データに信頼度を付加し、学習データとして出力する。

## 【0114】

図12は、本実施形態による音響モデル学習装置の処理手順を示すフローチャートである。本フローチャートのステップS31からS39までは、図5に示したS01からS0

50

9までに対応し、同様の処理手順であるのでここでは説明を省略する。そして、本フローチャートのステップS39の処理に続いて、ステップS40では、バイアス化言語モデルの重み付け係数を更新する。一例としては、ステップS34からS40までの繰り返し処理の第n回目( $n = 1, 2, 3, \dots$ )における重み値(この重み値はnの関数であるため $W_{1m}(n)$ と表わす)を、

$$W_{1m}(n) = 13 - 0.5n \quad \dots \quad (10)$$

とする。つまり、上の式(10)に従えば、1回目の認識処理に用いる重み値 $W_{1m}(1)$ は12.5であり、2回目の認識処理に用いる重み値 $W_{1m}(2)$ は12.0であり、3回目の認識処理に用いる重み値 $W_{1m}(3)$ は11.5であり、以下同様に、前回の認識処理のときよりも小さい重み値 $W_{1m}(n)$ を用いる。これは、認識処理において、前回の認識処理のときよりも言語モデルによる制約の影響をより小さくすることを意味する。

ステップS40の処理が終わると、ステップS34の処理に戻って新たな重み付け係数 $W_{1m}$ による言語モデルの重み付けを行い、以下、ステップS35以降の処理に続く。

#### 【0115】

ここでは言語モデルに対する重み値 $W_{1m}$ を用いて認識処理を行なっているが、これは、言語モデルと音響モデルの重みの比率を $W_{1m} : 1$ の比としていることになる。そしてこれは、言語モデルの重み値を1に固定して音響モデルの重み値 $W_{am}$ を用いて(つまり言語モデルと音響モデルの重みの比率を $1 : W_{am}$ の比として)も相対的に同じことであり、この場合、本実施形態では認識処理の繰り返しごとに音響モデルの重み値 $W_{am}$ が徐々に大きくなるように変化させても、同様に、前回の認識処理のときよりも言語モデルによる制約の影響をより小さくすることを意味する。また、重み値 $W_{1m}$ と $W_{am}$ とを用いて言語モデルと音響モデルの重みの比率を $W_{1m} : W_{am}$ の比としても同様であり、本実施形態では、認識処理において前回の認識処理のときよりも言語モデルによる制約の影響がより小さくなるようにする。

#### 【0116】

また、本実施例の変形例として、重み値 $W_{1m}$ を固定したままで、当該番組音声に対応する字幕テキストのための重み値 $W_{bias}$ を認識処理の繰り返しに伴って徐々に小さくしていても、今回認識処理においては前回の認識処理のときよりも言語モデルが相対的に音響モデルよりもより小さい重みで用いられるという同様の作用が得られる。

さらにまた、本実施例の変形例として、重み値 $W_{1m}$ と重み値 $W_{bias}$ の両方を認識処理の繰り返しに伴って徐々に小さくしていても、今回認識処理においては前回の認識処理のときよりも言語モデルが相対的に音響モデルよりもより小さい重みで用いられるという同様の作用が得られる。

#### 【0117】

言い換えれば、本実施形態あるいはその変形例においては、認識処理部10は、バイアス化言語モデル記憶部6から読み出した言語モデルと音響モデル記憶部9から読み出した前記音響モデルとを所定の重み比率値で重み付けして用いて、認識処理を行なうものである。また、認識処理部10は、音響モデル適応化部15によって更新された音響モデルを用いて再度認識処理を行なうものであり、このとき使用する前記重み比率値は、前回の認識処理のときよりも言語モデルがより小さい重みで用いられる値としている。

#### 【0118】

これにより、音響モデルの適応度合いに応じて言語的な制約を減じることとなり、これによってより一層、音響モデルの学習の精度が向上する。

さらに詳細に述べると、前記の各実施形態では、認識精度を高めるため、従来技術による方法よりも強い言語的な制約を用いた認識処理を行なうようにしている。つまり、言語モデルの重み付けの度合いを比較的大きくしている。しかしながら、音響モデルを適応化する処理を繰り返すことにより、比較的小さな言語的な制約で認識精度を確保できるようになる。従って、本実施形態では、繰り返し処理による音響モデルの精度向上に応じて言語的な制約を減ずるようにしている。この言語的な制約の調整により、字幕と発話の不一致区

10

20

30

40

50



間の検出制度の向上が可能となる。

なお、この実施形態は、前述した第1から第5のいずれかの実施形態と組み合わせて実現してよい。言い換えれば、第1から第5のいずれかの実施形態で述べた、選択処理部13および音響モデル適応化部15の処理を用いながら、重み値 $W_{1m}$ を徐々に変化させる。

#### 【0119】

[追加実施形態の3]

次に、本発明の追加実施形態の3について説明する。本実施形態では、話者の交代あるいは話者の属性の交代を検出し、話者に依存した音響モデル或いは話者の属性に依存した音響モデルを用いて認識処理を行なう点が特徴的である。

10

#### 【0120】

図13は、本実施形態による音響モデル学習装置の機能構成を示すブロック図である。図示する構成において、音響モデル学習装置102が有する字幕テキスト記憶部1と形態素解析処理部2と重み付け処理部3とベースラインコーパス処理部4と合成処理部5とバイアス化言語モデル記憶部6と重み付け処理部7と番組音声記憶部8は、前述した実施形態におけるそれらと同様のものである。さらに、本実施形態の音響モデル学習装置102は、その特有の構成として、音響モデル記憶部9Fおよび9Mと、認識処理部10MFと、認識仮説記憶部11Fおよび11Mと、アライメント処理部12Fおよび12Mと、選択処理部13Fおよび13Mと、学習データ記憶部14Fおよび14Mと、音響モデル適応化部15Fおよび15Mと、女声発話モデル記憶部16F（発話モデル記憶部、話者属性性別発話モデル記憶手段）と、男声発話モデル記憶部16M（発話モデル記憶部、話者属性性別発話モデル記憶手段）と、発話モデル適応化部17Fおよび17Mと、音声区間切り出し部18（話者属性推定手段）とを含んで構成される。

20

#### 【0121】

女声発話モデル記憶部16Fは、女性の声の統計的な音響特徴量を含む女声発話モデルを記憶している。また、男声発話モデル記憶部16Mは、男性の声の統計的な音響特徴量を含む男声発話モデルを記憶している。つまり、女声発話モデル記憶部16Fと男声発話モデル記憶部16Mは、話者ごと又は話者属性ごとの統計量を表わす話者属性性別発話モデルを記憶するための話者属性性別発話モデル記憶手段としても機能する。なお、女声および男声の統計的音響特徴量としてはMFCFを用いている。この統計的音響特徴量としては、他にも、PLPやHMM（隠れマルコフモデル）やホルマント周波数の分布などを用いることができる。

30

音声区間切り出し部18は、女声発話モデル記憶部16Fから読み出した女声発話モデルと男声発話モデル記憶部16Mから読み出した男声発話モデルとを用いて、番組音声の中の、女声で発話されている部分と、男声で発話されている部分と、発話以外の部分（番組中の、例えば無音部分や、背景雑音のみの部分や、音楽の部分や、効果音の部分など）とを区別する。つまり音声区間切り出し部18は、入力される音声が入力される発話であるか否かを検出するとともに、発話である場合にはその話者属性（ここでは、話者の性別）を検知する。つまり、音声区間切り出し部18は、話者属性性別発話モデルを用いて音声の中の所定音声区間における話者又は話者属性を推定する話者属性推定手段としても機能する。そして、音声区間切り出し部18は、女声による音声区間と男声による音声区間とをそれぞれ切り出した形で認識処理部10MFに渡す。

40

#### 【0122】

また、音響モデル学習装置102は、女声用と男声用の音響モデルを区別して管理し、記憶している。具体的には、音響モデル記憶部9Fには女声用の音響モデルが記憶され、音響モデル記憶部9Mには男声用の音響モデルが記憶されている。つまり、音響モデル記憶部9Fと9Mは、話者ごとまたは話者属性ごとに区別して音響モデルを記憶している。

#### 【0123】

そして、認識処理部10MFは、音声区間切り出し部18から渡された女声音声区間については音響モデル記憶部9Fから読み出した女声用音響モデルを用いて、また音声区間

50

切り出し部 1 8 から渡された男声音声区間については音響モデル記憶部 9 M から読み出した男声用音響モデルを用いて、それぞれ認識処理を行なう。つまり、認識処理部 1 0 M F は、音声区間切り出し部 1 8 によって切り出された音声区間の音声を認識処理する。またつまり、認識処理部 1 0 M F は、話者属性推定手段による推定結果に応じて、該当する話者用または話者属性用の音響モデルを用いて前記音声の認識処理を行なう。そして、認識処理部 1 0 M F は、女声音声区間から得られた認識結果仮説を認識仮説記憶部 1 1 F に書き込み、男声音声区間から得られた認識結果仮説を認識仮説記憶部 1 1 M に書き込む。

【 0 1 2 4 】

なお、男女別の音声の統計量 (Male speech model, Female speech model) を用いて音声区間の検出と男女の話者交代を検出し、性別に依存する音響モデルを用いて認識を行なうには、Toru IMAI, Shoei SATO, Shinichi HOMMA, Kazuo ONOE, Akio KOBAYASHI 「Online Speech Detection and Dual-Gender Speech Recognition for Captioning Broadcast News」 (IEICE Transactions on Information and Systems 2007 E90-D(8):1286-1291) に記載された方法を利用可能である。

10

【 0 1 2 5 】

アライメント処理部 1 2 F は、認識仮説記憶部 1 1 F から読み出した女声音声区間の認識仮説のみを用いて、前述した実施形態と同様のアライメント処理を行なう。選択処理部 1 3 F は、アライメント処理部 1 2 F による処理結果に基づき、女声音声区間の認識仮説に、前述の手法を用いて信頼度データを付加し、得られた学習データを学習データ記憶部 1 4 F に書き込む。

20

これと同様に、アライメント処理部 1 2 M は、認識仮説記憶部 1 1 M から読み出した男声音声区間の認識仮説のみを用いて、前述した実施形態と同様のアライメント処理を行なう。選択処理部 1 3 M は、アライメント処理部 1 2 M による処理結果に基づき、男声音声区間の認識仮説に、前述の手法を用いて信頼度データを付加し、得られた学習データを学習データ記憶部 1 4 M に書き込む。

このように、音声区間切り出し部 1 8 によって切り出された女声音声区間および男声音声区間それぞれに基づいて、女性用および男性用のそれぞれ専用の学習データが得られる。

【 0 1 2 6 】

音響モデル適応化部 1 5 F は、学習データ記憶部 1 4 F から読み出した女声用学習データのみを用いて、音響モデル記憶部 9 F に記憶されている女声用音響モデルを適応化 (更新) する。また、音響モデル適応化部 1 5 M は、学習データ記憶部 1 4 M から読み出した男声用学習データのみを用いて、音響モデル記憶部 9 M に記憶されている男声用音響モデルを適応化 (更新) する。つまり、音響モデル適応化部 1 5 F と 1 5 M は、話者属性推定手段による推定結果に対応する学習データを用いて、当該推定結果に対応する話者ごとまたは話者属性ごとの音響モデルを更新するものである。

30

【 0 1 2 7 】

発話モデル適応化部 1 7 F は、学習データ記憶部 1 4 F から読み出した女声用学習データのみを用いて、女声発話モデル記憶部 1 6 F に記憶されている女声発話モデルを適応化 (更新) する。また、発話モデル適応化部 1 7 M は、学習データ記憶部 1 4 M から読み出した男声用学習データのみを用いて、男声発話モデル記憶部 1 6 F に記憶されている男声発話モデルを適応化 (更新) する。つまり、発話モデル適応化部 1 7 F と 1 7 M は、話者属性推定手段による推定結果に対応する学習データを用いて、話者属性別発話モデル記憶手段に記憶された、当該推定結果に対応する話者ごとまたは話者属性ごとの話者属性別発話モデルを更新するものである。

40

【 0 1 2 8 】

また、音響モデル学習装置 1 0 2 は、上記の一連の認識処理、アライメント処理、選択処理、そして音響モデル適応化処理と発話モデル適応化処理を、繰り返し行なう。

【 0 1 2 9 】

上記のような構成および作用により、入力音声の統計的音響特徴量 (女声発話モデルお

50

よび男声発話モデル)を利用して高精度な音声区間の切り出しを行なう方法を併用する場合において、得られた学習データを用いて切り出しに使用する上記の統計量も適応化することを繰り返すことができる。これにより、例えば雑音区間を発話区間の音声として認識してしまうような音声区間の切り出し誤りが減少し、学習データの質のさらなる向上が可能となる。

また、話者属性(性別)に依存した音響モデルを使用して認識処理を行なうことにより、話者属性非依存の音響モデルを用いた場合に比べ、高い認識精度を得ることができ、より多くの学習データを作成することができる。

また、入力音声の統計量を利用して話者或いは話者属性の推定手段(音声区間切り出し部18)を備えた音声アルゴリズムを併用する場合に、得られた学習データを用いて話者推定用の統計量(女声発話モデルおよび男声発話モデル)を適応化するとともに、話者(或いは話者属性)ごとに音響モデルを設けて(音響モデル記憶部9Fおよび9M)、これらをそれぞれ適応化しているため、さらに、認識精度の向上が図れる。

#### 【0130】

なお、この実施形態は、前述した第1から第5のいずれかの実施形態と組み合わせて実現してよい。言い換えれば、選択処理部13Fおよび13Mはそれぞれ、第1から第5のいずれかの実施形態で述べた選択処理部13の処理を用いる。また、音響モデル適応化部15Fおよび15Mはそれぞれ、第1から第5のいずれかの実施形態で述べた音響モデル適応化部15の処理を用いる。このようにして、話者ごと、あるいは、話者属性ごとの音響モデルの学習を行なう。

#### 【0131】

なお、本実施形態の更なる変形例として、追加実施形態の1で説明した番組ごとの音響モデルを管理する構成(このとき、適宜、発話モデルも番組ごとに管理するようにしても良い)や、追加実施形態の2で説明した処理の繰り返しに伴って言語モデルの重み付けを徐々に変化させる構成を併用しても良い。

#### 【0132】

また、性別ごとに音声区間切り出しのための発話モデルおよび認識処理のための音響モデルを設ける代わりに、或いは性別に加えて、他の話者属性ごと或いは話者個人ごとに、発話モデルや音響モデルを設けて、その話者属性ごと或いは話者個人ごとの音声区間切り出し処理や認識処理を行なうようにしても良い。「他の話者属性」とは、例えば、年齢層などである。このとき、話者の年齢層を例えば、少年期(5歳から14歳)、青年期(15歳から24歳)、壮年期(25歳から44歳)、中年期(45歳から64歳)、高年期(65歳以上)などに分類する。

#### 【0133】

また、発話環境ごとに、発話モデルや音響モデルを設けるようにしても良い。ここで「発話環境ごと」とは、例えば、話者が原稿を読み上げている形式の場合、対談あるいは座談形式の場合、雑談の場合などである。

#### 【0134】

また、本実施形態では、発話モデルを利用して音声区間を切り出す構成と、話者属性ごとに音響モデルを設けるとともに話者属性ごとに音声区間を切り出してそれぞれ専用の音響モデルを用いて認識を行い、音響モデルを適応化する構成との両方を用いているが、これらのいずれか一方のみの構成を含むようにしても良い。

#### 【0135】

また、本実施形態では、認識仮説を、論理的には言語的単位(形態素)をアークとするラティス構造のグラフで表現して処理を行なったが、その代わりに、言語的単位として音素をアークとするラティス構造のグラフで表現して、同様の処理を行なうようにしても良い。

#### 【0136】

<コンピュータシステムを用いた実施形態>

なお、上述した各実施形態における音響モデル学習装置の機能をコンピュータで実現す

10

20

30

40

50

るようにしても良い。その場合、この音響モデル学習の機能を実現するためのプログラムをコンピュータ読み取り可能な記録媒体に記録して、この記録媒体に記録されたプログラムをコンピュータシステムに読み込ませ、実行することによって実現しても良い。なお、ここでいう「コンピュータシステム」とは、OSや周辺機器等のハードウェアを含むものとする。また、「コンピュータ読み取り可能な記録媒体」とは、フレキシブルディスク、光磁気ディスク、ROM、CD-ROM等の可搬媒体、コンピュータシステムに内蔵されるハードディスク等の記憶装置のことをいう。さらに「コンピュータ読み取り可能な記録媒体」とは、インターネット等のネットワークや電話回線等の通信回線を介してプログラムを送信する場合の通信線のように、短時刻の間、動的にプログラムを保持するもの、その場合のサーバやクライアントとなるコンピュータシステム内部の揮発性メモリのように、一定時刻プログラムを保持しているものも含んでも良い。また上記プログラムは、前述した機能の一部を実現するためのものであっても良く、さらに前述した機能をコンピュータシステムにすでに記録されているプログラムとの組み合わせで実現できるものであっても良い。

10

## 【0137】

以上、この発明の実施形態について図面を参照して詳述してきたが、具体的な構成はこの実施形態に限られるものではなく、この発明の要旨を逸脱しない範囲の設計等も含まれる。

例えば、一致区間を選択する際に、選択処理部13は字幕テキストと認識仮説が3単語以上連続して一致している区間を選択することを説明したが、「3」に限らず、適宜所定の単語数以上連続して位置している区間を選択するようにしても良い。また、語数を基準とする代わりに、所定の文字数（例えば、5文字）以上連続して一致している区間を選択するようにしても良い。

20

## 【0138】

## &lt; 認識実験の結果 &gt;

ここで、本願発明者らが実際に行なった認識実験について説明する。

2004年10月14日、および2007年6月25日、27日、28日（本放送と再放送を含む）にNHKで放送された5回分の番組「きょうの健康」を評価音声として認識実験を行なった。各放送回には約2600単語（形態素）の発話があり総計12807単語を評価した。

30

## 【0139】

音響モデルの学習用には、2004年から2007年に放送されたNHK「きょうの健康」から、評価対象と同一のコンテンツを除く93回分の音声と字幕を用いて、放送回ごとに下記の4手法（LS0～LS3）を適用して約20時間分の学習データを生成し、性別依存不特定話者音響モデルをMLLRとMAPで適応化して認識実験を行なった。

LS0：これは従来技術による手法であり、認識結果（最尤仮説）と書き起こしのテキストデータとの一致部分の音声を切り出して、学習音声とした。

LS1：信頼度を全て1として、一致区間情報を利用しない教師なし学習とした。

LS2：最尤仮説と書き起こしのテキストデータとの一致部分の信頼度を1とし、それ以外の部分の信頼度を0として、音声を切り出さずに一区間のみを用いた学習とした。

40

LS3：最尤仮説と書き起こしのテキストデータとの一致部分の信頼度を1とし、最尤仮説の単語の事後確率から不一致部分の信頼度を得て、学習した。

## 【0140】

このときの言語的制約の強さは、 $W_{bias} = 1000$ 、 $W_{lm} = 16$ を用いた。こうして作成されたこの適応化言語モデルでのテストセットパープレキシティーは22程度であった。また、上記の20時間分の学習データのうち、認識結果と字幕の一致部分、つまり従来のLS法で得られる学習音声の量は5～8時間程度であった。

言語モデルは、情報番組の書き起こしの字幕テキストデータを基に、番組ホームページから得られる各放送回の番組内容紹介テキストに20倍の重みを付けて学習した番組依存言語モデルである。この言語モデルのテストセットパープレキシティーは59であり、未

50

知語率は 0.8%であった。

【0141】

認識実験および学習データの生成には、男女別のモノフォンHMMを用いて音声区間を検出し、男女の自動判別を行ないながら性別依存トライフォンHMMの両方を適応化する。したがって、適応化モデルにより発話の切り出し精度と男女の判別精度の向上も期待される。ベースラインの音響モデルは、340時間(男性)と250時間(女性)のニュース音声から学習され、モノフォンHMM(126状態32混合分布)、状態共有トライフォンHMM(4000状態、16混合分布)が男女別に作成された。

【0142】

実験結果は次の通りである。最尤仮説を教師単語列とした場合のLS0~LS3でのWERを考察する。まず、最尤仮説と書き起こしの字幕テキストの一致部分のみを学習するLS0とLS2を比較すると、LS2では、信頼度を用いることにより認識と同一の発話区間検出を利用できたため、より大きな改善が得られた。次に、LS0以外の3手法を比較すると、収束の速さに若干の差があるものの、最終的にはほぼ同等の改善が得られた。本実験条件では、教師単語列作成時の言語的制約( $W_{bias}$ ,  $W_{lm}$ )が事前に最適化されていたため、書きお越しの一致の有無にかかわらず、全ての音声を学習(LS1)しても改善が得られた。一方、ラティスと字幕テキストの一致が最大になる単語列を教師とした場合、提案法LS3以外のLS1とLS2のWERの増加が見られた。本実験条件では、仮説の信頼度に関わらず教師単語列を選択したため、学習データ生成時の認識誤りを教師とすることが多くなり、評価音声の認識精度が低下したと思われる。これに対して提案法LS3では、教師単語列の制度が低い場合でも頑健な学習結果が得られ、WERの増加は小さかった。

10

20

【0143】

本発明は、不完全な書き起こしの字幕データから効率的に音響モデルを学習するため、教師単語系列に信頼度を与え、さらに認識ラティスの事後確率を信頼度とする方法を用いた。字幕テキストを用いた音響モデルの学習に本発明の方法を適用した結果、信頼度の導入により、従来法よりも低いWERが得られることが認識実験を通して実証された。さらに事後確率を用いることで、頑健な学習が可能であることを示すことができた。

【図面の簡単な説明】

【0144】

30

【図1】本発明の第1の実施形態による音響モデル学習装置の機能構成を示したブロック図である。

【図2】同実施形態で用いる言語モデルのデータ構成を示した概略図である。

【図3】同実施形態で用いる音響モデルのデータ構成を示した概略図である。

【図4】同実施形態で用いる認識仮説記憶部のデータ構成およびデータ例、並びにそのデータが表わすラティス構造を示した概略図である。

【図5】同実施形態による音響モデル学習の処理手順を示したフローチャートである。

【図6】同実施形態により最尤仮説のパスと書き起こしによる字幕テキストの一致を利用して信頼度を付加する例を表わす概略図である。

【図7】同実施形態の音響モデル学習装置による音響モデルの学習の概略を説明するための概略図である。

40

【図8】第2の実施形態において枝刈りされた後の認識仮説のデータ例、並びにそのデータが表わすラティス構造を示した概略図である。

【図9】第3の実施形態において対立する複数のパス間で共通して出現する音素の例を説明するための概略図である。

【図10】第4の実施形態において、選択処理部13が認識結果のラティスと字幕テキストの一致区間を抽出する処理を概念的に説明する概略図である。

【図11】本発明の追加実施形態の1による音響モデル学習装置の機能構成を示したブロック図である。

【図12】本発明の追加実施形態の2による音響モデル学習装置の処理手順を示したフロ

50

ーチャートである。

【図 1 3】本発明の追加実施形態の 3 による音響モデル学習装置の機能構成を示したブロック図である。

【図 1 4】従来技術による音響モデル学習装置の機能構成を示したブロック図である。

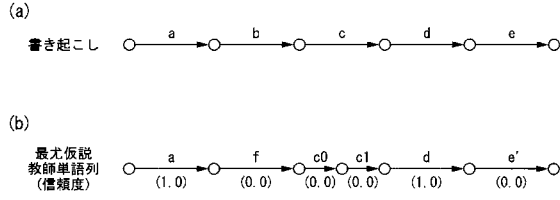
【符号の説明】

【 0 1 4 5 】

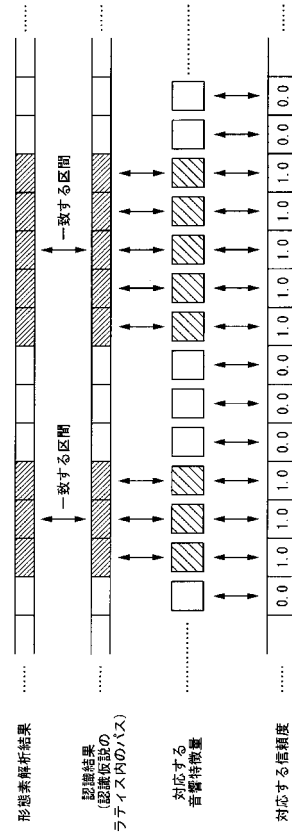
- |                              |  |    |
|------------------------------|--|----|
| 1                            | 字幕テキスト (Caption text) 記憶部                          |    |
| 2                            | 形態素解析 (Morphological analysis) 処理部                 |    |
| 3                            | 重み付け処理部  |    |
| 4                            | ベースラインコーパス (Baseline corpus) 記憶部                   | 10 |
| 5                            | 合成処理部  |    |
| 6                            | バイアス化言語モデル (Biased LM) 記憶部 (言語モデル記憶部)              |    |
| 7                            | 重み付け処理部  |    |
| 8                            | 番組音声 (Program audio) 記憶部                           |    |
| 9, 9 F, 9 M                  | 音響モデル (AM) 記憶部                                     |    |
| 10, 10 M F                   | 認識 (Recognition) 処理部                               |    |
| 11, 11 F, 11 M               | 認識仮説記憶部  |    |
| 12, 12 F, 12 M               | アライメント (Alignment) 処理部                             |    |
| 13, 13 F, 13 M               | 選択 (Selection) 処理部                                 |    |
| 14, 14 F, 14 M               | 学習データ (Transcripts) 記憶部                            | 20 |
| 15, 15 F, 15 M               | 音響モデル適応化部 (Adaptation)                             |    |
| 16 F                         | 女声発話モデル (Female speech model) 記憶部 (話者属性別発話モデル記憶手段) |    |
| 16 M                         | 男声発話モデル (Male speech model) 記憶部 (話者属性別発話モデル記憶手段)   |    |
| 17 F, 17 M                   | 発話モデル適応化部 (Adaptation)                             |    |
| 18                           | 音声区間切り出し部 (話者属性推定手段)                               |    |
| 100, 101, 101 A ~ 101 D, 102 | 音響モデル学習装置  |    |



【 図 6 】



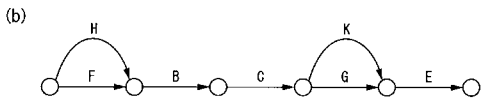
【 図 7 】



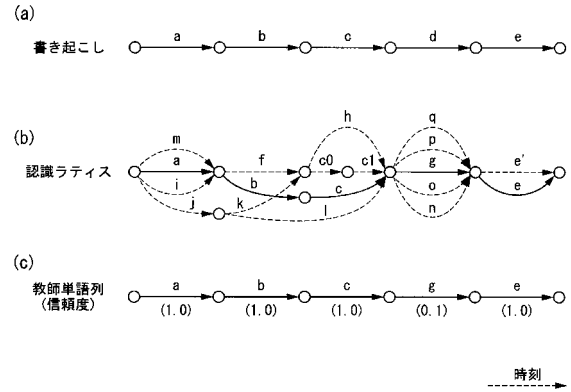
【 図 8 】

(a)

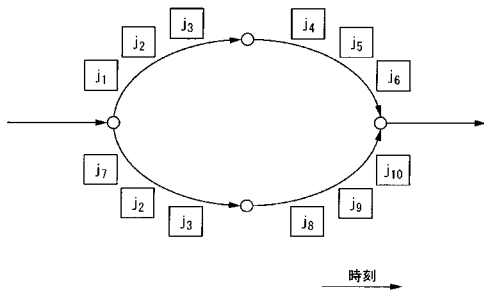
始点ノード	終点ノード	言語的単位 (形態素)	信頼度
1	2	F	0.1
1	2	H	0.9
2	3	B	1
3	4	C	1
4	5	K	0.2
4	5	G	0.8
5	6	E	1
⋮	⋮	⋮	⋮



【 図 10 】

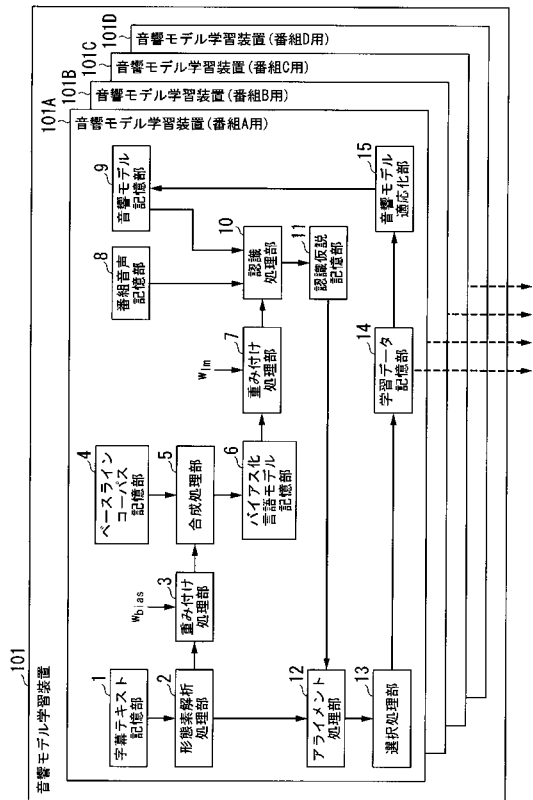


【 図 9 】

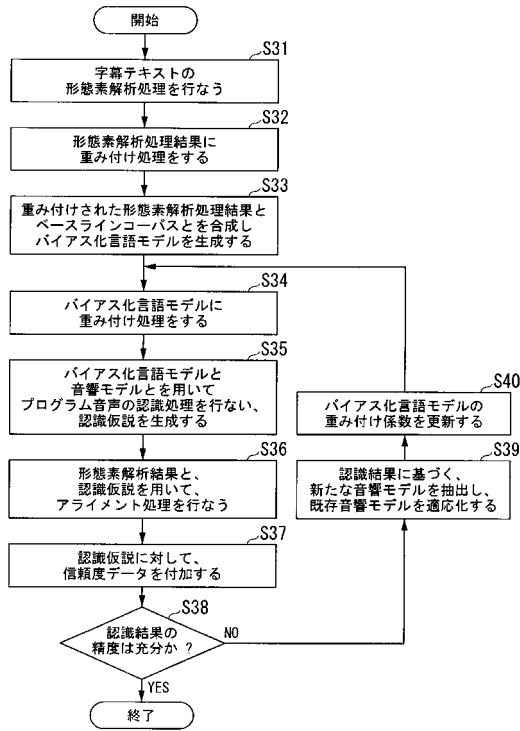




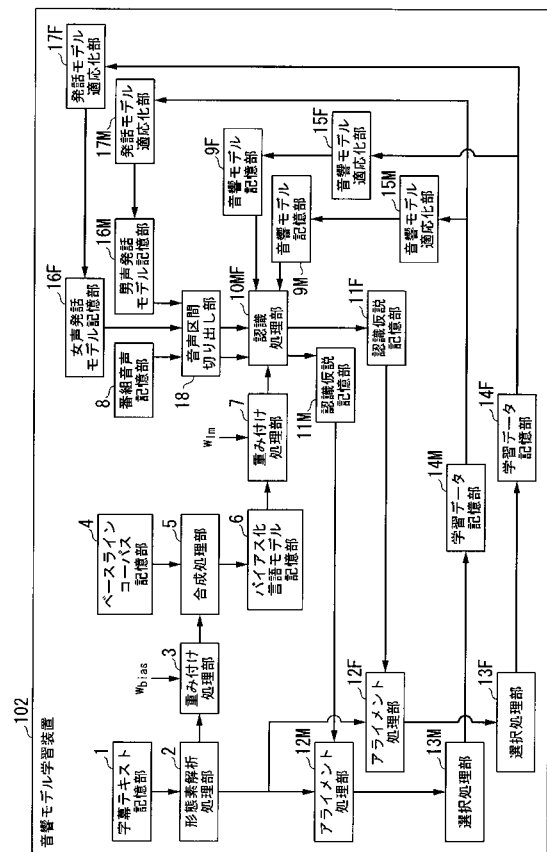
【図 1 1】



【図 1 2】



【図 1 3】



【図 1 4】

