

(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2019-12095
(P2019-12095A)

(43) 公開日 平成31年1月24日(2019.1.24)

(51) Int.Cl.		F I		テーマコード (参考)
G 1 0 L 15/06 (2013.01)		G 1 0 L 15/06	3 0 0 C	
G 1 0 L 15/18 (2013.01)		G 1 0 L 15/18	3 0 0 H	
G 1 0 L 15/187 (2013.01)		G 1 0 L 15/187		

審査請求 未請求 請求項の数 7 O L (全 13 頁)

(21) 出願番号	特願2017-126929 (P2017-126929)	(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29年6月29日 (2017.6.29)	(74) 代理人	110001807 特許業務法人磯野国際特許商標事務所
		(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	尾上 和穂 北海道室蘭市山手町1-3-50 日本放送協会 室蘭放送局内

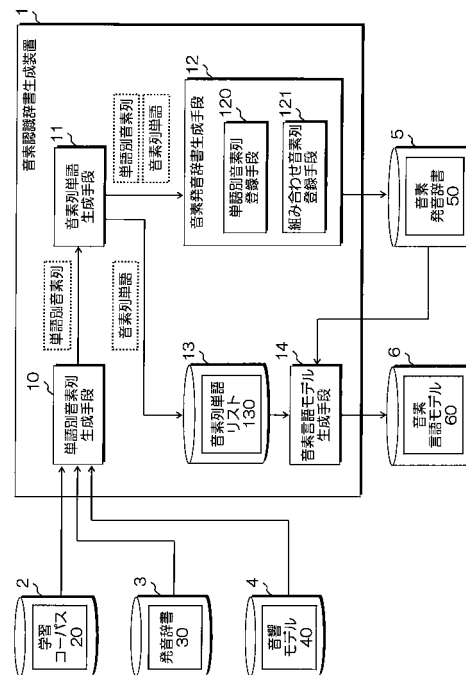
(54) 【発明の名称】 音素認識辞書生成装置および音素認識装置ならびにそれらのプログラム

(57) 【要約】

【課題】音素認識に用いる音素発音辞書および音素言語モデルを生成する音素認識辞書生成装置を提供する。

【解決手段】音素認識辞書生成装置1は、学習コーパスから、強制アライメントにより単語別音素列を生成する単語別音素列生成手段10と、単語別音素列を1単語のテキストデータ形式に変換して音素列単語を生成する音素列単語生成手段11と、音素列単語を見出し語とし、当該音素列単語に対応する単語別音素列を発音表記とすることで、音素発音辞書50を生成する音素発音辞書生成手段12と、音素列単語のリストから、N-gram言語モデルを学習し、音素言語モデル60を生成する音素言語モデル生成手段14と、を備える。

【選択図】図1



【特許請求の範囲】**【請求項 1】**

音響モデルと発音辞書と学習コーパスとを用いて、音素認識に用いる音素発音辞書および音素言語モデルを生成する音素認識辞書生成装置であって、

前記学習コーパスの音声を、前記音響モデルと前記発音辞書とに基づいて音声認識し、前記発音辞書に登録されている見出し語に対応する単語ごとの音素列である単語別音素列を生成する単語別音素列生成手段と、

前記単語別音素列を 1 単語のテキストデータ形式に変換して音素列単語を生成する音素列単語生成手段と、

前記音素列単語を見出し語とし、当該音素列単語に対応する前記単語別音素列を発音表記とすることで、前記音素発音辞書を生成する音素発音辞書生成手段と、

前記音素列単語生成手段で生成される前記音素列単語のリストから前記音素列単語の連鎖として N - g r a m 言語モデルを学習することにより、前記音素言語モデルを生成する音素言語モデル生成手段と、

を備えることを特徴とする音素認識辞書生成装置。

10

【請求項 2】

前記単語別音素列生成手段は、前記単語別音素列の音素間に音素以外の予め定めた文字を挿入することで、前記音素列単語を生成することを特徴とする請求項 1 に記載の音素認識辞書生成装置。

【請求項 3】

前記音素発音辞書生成手段は、予め定めた数の音素を組み合わせた音素列を前記テキストデータ形式に変換した見出し語とし、当該見出し語に対応する音素列を発音表記として前記音素発音辞書に登録することを特徴とする請求項 1 または請求項 2 に記載の音素認識辞書生成装置。

20

【請求項 4】

前記音素言語モデル生成手段は、前記音素列単語生成手段で生成される前記音素列単語のリストに存在しない音素列単語の連鎖に対して、スムージングにより N グラム確率を与え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素認識辞書生成装置。

【請求項 5】

コンピュータを、請求項 1 から請求項 4 のいずれか一項に記載の音素認識辞書生成装置として機能させるための音素認識辞書生成プログラム。

30

【請求項 6】

音響モデルと、請求項 1 から請求項 4 のいずれか一項に記載の音素認識辞書生成装置により生成された音素発音辞書および音素言語モデルとを用いて、音声の音素を認識する音素認識装置であって、

前記音響モデルと前記音素発音辞書と前記音素言語モデルとにより、前記音声を音素列単語単位で認識する認識手段と、

この認識手段で認識された 1 単語のテキストデータ形式である音素列単語を、個々の音素に分離して音素列を生成する音素列生成手段と、

を備えることを特徴とする音素認識装置。

40

【請求項 7】

コンピュータを、請求項 6 に記載の音素認識装置として機能させるための音素認識プログラム。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、発話音声の音素認識に用いる音素発音辞書および音素言語モデルを生成する音素認識辞書生成装置およびそのプログラム、ならびに、音素発音辞書および音素言語モデルを用いた音素認識装置およびそのプログラムに関する。

50

【背景技術】

【0002】

通常、音声認識では、単語と当該単語の発音系列（音素列）とを対応付けた発音辞書を用いている。この発音辞書には、一般的な辞書に記載されているような読みが発音として登録されている。

しかし、表記上の読みと実際に発話された発音とでは異なることが多い。例えば、放送番組では、ニュース番組のアナウンサの正確な（発音辞書の発音と近い）発音に比べ、情報番組の出演者の発話は曖昧な発音が多い。

そこで、統計的機械翻訳モデルを利用して、アナウンサ等の正確な発音を前提とした音素列から、発音が不明瞭な発話の音素列の単語を推定して、発音辞書を拡張する技術が開示されている（特許文献1参照）。

10

【0003】

特許文献1の技術（以下、従来技術という）では、認識対象音素の前後の音素に対する依存性（環境依存）を考慮して音素認識を行う。

この従来技術は、学習コーパスから、トライフォンを1つの単語として発音辞書を学習するとともに、トライフォンの接続確率を与える言語モデルを学習する。ここで、トライフォンは、例えば、「警察」の発音では、「（けー）k - e : + s」, 「（さ）e : - s + a」, 「（つ）s - a + t s」のように、中心音素を含めた前後の発音を含めて表現したものである。

【0004】

20

そして、従来技術は、音声と書き起こしテキストとを対応付けた学習コーパスから強制音素アライメントを行った音素列（標準音素列）と、音素のトライフォンの言語モデルおよび発音辞書を用いて学習コーパスの音声を音素認識した音素列（実発話音素列）とを用いて、統計的機械翻訳モデルを学習する。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2016-161765号公報

【発明の概要】

【発明が解決しようとする課題】

30

【0006】

前記した従来技術は、統計的機械翻訳モデルを学習するために、強制音素アライメントを行った音素列（標準音素列）と、音素認識した音素列（実発話音素列）とを用いる。この統計的機械翻訳モデルの精度を高めるには、標準音素列と実発話音素列の質が重要になる。

従来技術で、アナウンサ等の正確な発音の音声とその書き起こしテキストとを学習コーパスとして用いて標準音素列と実発話音素列とを生成した場合、理想的には、それぞれの音素列がほぼ同じであることが望ましい。

しかし、従来技術では、標準音素列と実発話音素列とをDP（Dynamic Programming）マッチングした結果、音素が異なる割合（音素異なり率）が、22.8%あり、さらなる音素認識の精度改善が望まれている。

40

【0007】

本発明は、このような問題に鑑みてなされたものであり、音素認識の精度を高める音素認識辞書（音素発音辞書および音素言語モデル）を生成する音素認識辞書生成装置およびそのプログラム、ならびに、音素発音辞書および音素言語モデルを用いた音素認識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】

【0008】

前記課題を解決するため、本発明に係る音素認識辞書生成装置は、音響モデルと発音辞書と学習コーパスとを用いて、音素認識に用いる音素発音辞書および音素言語モデルを生

50

成する音素認識辞書生成装置であって、単語別音素列生成手段と、音素列単語生成手段と、音素発音辞書生成手段と、音素言語モデル生成手段と、を備える。

【0009】

かかる構成において、音素認識辞書生成装置は、単語別音素列生成手段によって、学習コーパスの音声を音響モデルと発音辞書とに基づいて音声認識し、発音辞書に登録されている見出し語に対応する単語ごとの音素列である単語別音素列を生成する。

そして、音素認識辞書生成装置は、音素列単語生成手段によって、単語別音素列を1単語のテキストデータ形式に変換して音素列単語を生成する。例えば、音素列単語生成手段は、単語別音素列の音素間のスペースに音素以外の予め定めた文字（例えば、“+”）を挿入することで、音素列単語を生成する。これによって、音素認識辞書生成装置は、音素列単語を1単語として扱うことが可能になる。

10

【0010】

そして、音素認識辞書生成装置は、音素発音辞書生成手段によって、音素列単語を見出し語とし、当該音素列単語に対応する単語別音素列を発音表記とすることで、音素発音辞書を生成する。これによって、音素発音辞書生成手段は、単語単位で音素列の発音を音素発音辞書に登録する。

さらに、音素認識辞書生成装置は、音素言語モデル生成手段によって、音素列単語生成手段で生成される音素列単語のリストから音素列単語の連鎖としてN-gram言語モデルを学習することにより、音素言語モデルを生成する。これによって、音素言語モデル生成手段は、音素認識を行う際の音素列単語の接続確率を計算するため音素列単語の出現確率をモデル化する。

20

【0011】

なお、音素認識辞書生成装置は、コンピュータを、単語別音素列生成手段、音素列単語生成手段、音素発音辞書生成手段、音素言語モデル生成手段として機能させるための音素認識辞書生成プログラムで動作させることができる。

【0012】

また、前記課題を解決するため、本発明に係る音素認識装置は、音響モデルと、音素認識辞書生成装置により生成された音素発音辞書および音素言語モデルとを用いて、音声の音素を認識する音素認識装置であって、認識手段と、音素列生成手段と、を備える。

【0013】

かかる構成において、音素認識装置は、認識手段によって、音響モデルと音素発音辞書と音素言語モデルとにより、音声を音素列単語単位で認識する。これによって、認識手段は、単語の繋がりに依存した音素列を認識することが可能になる。

30

そして、音素認識装置は、音素列生成手段によって、認識手段で認識された1単語のテキストデータ形式である音素列単語を、個々の音素に分離して音素列を生成する。例えば、音素列生成手段は、単語別音素列の音素間に挿入されている予め定めた文字（例えば、“+”）をスペースに置き換えることで、個々の音素に分離する。

【0014】

なお、音素認識装置は、コンピュータを、認識手段、音素列生成手段として機能させるための音素認識プログラムで動作させることができる。

40

【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。

本発明によれば、音素列を単語単位とした音素発音辞書および音素言語モデルを生成することができる。

この音素発音辞書および音素言語モデルを用いることで、音素認識する際の音素の連結確率を、単に音素の前後の依存性だけではなく、音素の単語内および単語間における依存性も加味して算出することが可能になり、音声から音素を認識する際の認識精度を高めることができる。

【図面の簡単な説明】

50

【0016】

【図1】本発明の第1実施形態に係る音素認識辞書生成装置の構成を示すブロック構成図である。

【図2】図1の単語別音素列生成手段における単語別音素列の生成例を説明するための説明図であって、(a)は学習コーパスの音声の書き起こし例、(b)は発音辞書の一部、(c)は生成した単語別音素列の例を示す。

【図3】音素の表記例を示す図である。

【図4】図1の音素発音辞書生成手段が生成する音素発音辞書の例を示す図である。

【図5】図1の音素列単語生成手段が生成する音素列単語リストの例を示す図である。

【図6】図1の音素言語モデル生成手段が生成する音素言語モデルの例を示す図である。

【図7】本発明の第1実施形態に係る音素認識辞書生成装置の動作を示すフローチャートである。

【図8】本発明の第2実施形態に係る音素認識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。

<第1実施形態>

〔音素認識辞書生成装置の構成〕

まず、図1を参照して、本発明の第1実施形態に係る音素認識辞書生成装置1の構成について説明する。

【0018】

音素認識辞書生成装置1は、音声データから音素を認識するための辞書として、音素発音辞書および音素言語モデルを生成するものである。この音素認識辞書生成装置1は、学習コーパス記憶装置2、発音辞書記憶装置3および音響モデル記憶装置4にそれぞれ記憶されている学習コーパス20、発音辞書30および音響モデル40から、音素発音辞書50および音素言語モデル60を生成する。

【0019】

具体的には、音素認識辞書生成装置1は、学習コーパスから強制アライメントにより単語別音素列を生成し、生成した単語別音素列を1単語のテキストデータ形式に変換して音素列単語を生成する。そして、音素認識辞書生成装置1は、生成した音素列単語を見出し語とし、当該音素列単語に対応する単語別音素列を発音表記とすることで、音素発音辞書50を生成する。さらに、音素認識辞書生成装置1は、生成した音素列単語のリストから、N-gram言語モデルを学習し、音素言語モデル60を生成する。

【0020】

学習コーパス20は、予め大量の音声データ(音声コーパス)と、音声データの書き起こしテキスト(テキストコーパス)とを対応付けたデータである。この学習コーパス20は、例えば、ニュース番組、情報番組等におけるアナウンサ、リポータ等の約1000時間程度の音声(音声コーパス)と、その音声を書き起こしたテキスト(テキストコーパス)である。

【0021】

発音辞書30は、所定の文字列である見出し語(ここでは、単語とする)ごとに、その発音表記(音素列)を示した辞書である。

この発音辞書30は、一般的な発音辞書であって、例えば、人手を介して見出し語(単語)とその発音表記(音素列)とを対応付けた辞書である。

【0022】

音響モデル40は、大量の音声データから予め学習したディープニューラルネットワーク(DNN:Deep Neural Network)音響モデルである。例えば、DNNの入力には、メルフィルタバンク対数パワーの40次元に時間変化()を加えて11フレーム分の特徴量を連結(スプライス)した特徴量を用い、DNNの隠れ層を8層とする。

なお、音響モデル40における音響特徴量の尤度計算は、隠れマルコフモデル(HMM

10

20

30

40

50

: Hidden Markov Model) や、ガウス混合モデル (GMM: Gaussian mixture model) 音響モデルであっても構わない。

以下、音素認識辞書生成装置 1 の構成について詳細に説明する。

【0023】

音素認識辞書生成装置 1 は、図 1 に示すように、単語別音素列生成手段 10 と、音素列単語生成手段 11 と、音素発音辞書生成手段 12 と、音素列単語リスト記憶手段 13 と、音素言語モデル生成手段 14 と、を備える。

また、音素認識辞書生成装置 1 は、生成した音素発音辞書 50 を記憶する音素発音辞書記憶装置 5 と、生成した音素言語モデル 60 を記憶する音素言語モデル記憶装置 6 と、を外部に接続している。もちろん、音素発音辞書記憶装置 5 および音素言語モデル記憶装置 6 は、音素認識辞書生成装置 1 の内部に備える構成としてもよい。また、音素発音辞書記憶装置 5 および音素言語モデル記憶装置 6 は、1 つの記憶装置で構成してもよい。

【0024】

単語別音素列生成手段 10 は、発音辞書 30 と音響モデル 40 とに基づいて、学習コーパス 20 の音声 (音声コーパス) を強制アライメントすることで、発音辞書 30 に登録されている見出し語に対応する単語ごとに、音声の音素列を切り分けて単語別音素列を生成するものである。

【0025】

この単語別音素列生成手段 10 は、学習コーパス 20 の音声から、音響モデル 40 に対応する音響特徴量 (メル周波数ケプストラム係数等) を抽出する。そして、単語別音素列生成手段 10 は、発音辞書 30 と音響モデル 40 とを用いて、音声の書き起こしテキスト (テキストコーパス) を事前知識とする音声認識を行い、発音辞書 30 に登録されている文字列 (見出し語) に対応して強制アライメントする。これにより、単語別音素列生成手段 10 は、図 2 (b) に示されているように、発音辞書 30 に登録されている単語に複数存在する発音の音素列に対し、尤も音声に近い発音の音素列を選択し、単語別音素列を生成する。

【0026】

図 2 は、単語別音素列生成手段 10 における単語別音素列の生成例を示す。例えば、単語別音素列生成手段 10 は、学習コーパス 20 として、「世界一短い東京の橋でイベントが開かれました」の音声データを入力した場合、音響モデル 40 に対応する音響特徴量を抽出する。

そして、単語別音素列生成手段 10 は、音声データに対応する図 2 (a) に示す学習コーパス 20 の書き起こしテキスト「世界一短い東京...」を事前知識として、図 2 (b) に示す発音辞書 30 と、音響モデル 40 と、を用いて音声認識を行う。

【0027】

これによって、単語別音素列生成手段 10 は、図 2 (c) に示すように、単語ごとの音素列 (単語別音素列) 「s e k a i i c h i / m i j i k a i / t o : k y o : / ...」 (ここで、“ ” はスペースを示す) を生成する。

単語別音素列生成手段 10 は、生成した単語別音素列を音素列単語生成手段 11 に出力する。

【0028】

音素列単語生成手段 11 は、単語別音素列生成手段 10 で生成された単語別音素列を、単語ごとに 1 単語のテキストデータ形式に変換した音素列単語を生成するものである。

この音素列単語生成手段 11 は、単語別音素列の音素間に音素以外の予め定めた文字を挿入することで、個々に分離した音素列を、1 単語のテキストデータ形式に変換する。

【0029】

具体的には、音素列単語生成手段 11 は、音素ごとにスペースを含んだ単語別音素列のスペースを、音素以外の予め定めた文字に置き換えて 1 つの単語テキストとする。例えば、音素列単語生成手段 11 は、単語別音素列のスペースを“+”に置き換え、“s e k a i i c h i”を“s + e + k + a + i + i + c h + i”等に変換する。

10

20

30

40

50

【 0 0 3 0 】

音素列単語生成手段 1 1 は、スペースを含んだ単語別音素列と、テキスト置換した音素列単語とを対にして、順次、音素発音辞書生成手段 1 2 に出力する。また、音素列単語生成手段 1 1 は、テキスト置換した音素列単語のみを、順次、音素列単語リスト記憶手段 1 3 に書き込む。

【 0 0 3 1 】

音素発音辞書生成手段 1 2 は、音素列を単語とみなした音素列単語の発音辞書である音素発音辞書を生成するものである。音素発音辞書生成手段 1 2 は、図 1 に示すように、単語別音素列登録手段 1 2 0 と、組み合わせ音素列登録手段 1 2 1 と、を備える。

【 0 0 3 2 】

単語別音素列登録手段 1 2 0 は、単語別音素列と音素列単語とを対として登録した音素発音辞書を生成するものである。単語別音素列登録手段 1 2 0 は、音素列単語生成手段 1 1 で生成された音素列単語を見出し語とし、音素列単語と対となる単語別音素列をその見出し語の発音として、音素発音辞書記憶装置 5 の音素発音辞書 5 0 に登録する。

【 0 0 3 3 】

なお、単語別音素列登録手段 1 2 0 は、同じ見出し語となる音素列単語に対して、異なる発音の単語別音素列が入力された場合、見出し語に複数の発音を登録する。また、単語別音素列登録手段 1 2 0 は、同じ見出し語となる音素列単語に対して、同じ発音の単語別音素列が入力された場合、登録を行わないこととする。

【 0 0 3 4 】

組み合わせ音素列登録手段 1 2 1 は、任意の音素の組み合わせで構成される音素列を単語とみなした見出し語と、その音素列とを対として、音素発音辞書に登録するものである。

具体的には、組み合わせ音素列登録手段 1 2 1 は、図 3 に示す音素の例において、すべての音素（図 3 の例では、4 0 音素）に対して、予め定めた最大音素数（ここでは、“4”とする）の音素の組み合わせ（ $4 0^1 + 4 0^2 + 4 0^3 + 4 0^4$ 通り）の音素列を、音素発音辞書記憶装置 5 の音素発音辞書 5 0 に登録する

【 0 0 3 5 】

この組み合わせ音素列登録手段 1 2 1 は、音素列単語生成手段 1 1 と同様に、音素を組み合わせさせた音素列を、1 つのテキストデータ形式に変換する。具体的には、組み合わせ音素列登録手段 1 2 1 は、音素を組み合わせさせた音素列のスペースを音素以外の予め定めた 1 つのテキスト（ここでは、“+”）に置き換えた単語に変換し、見出し語とする。

【 0 0 3 6 】

ここで、図 4 を参照して、音素発音辞書生成手段 1 2 が音素発音辞書記憶装置 5 に登録する音素発音辞書 5 0 の例について説明する。

図 4 に示すように、音素発音辞書 5 0 は、単語別音素列登録手段 1 2 0 で登録される辞書 A と、組み合わせ音素列登録手段 1 2 1 で登録される辞書 B とで構成される。

辞書 A は、学習コーパス 2 0 の書き起こしに含まれる単語の発音を示す単語音素列のスペース部分を“+”に置き換えた単語別音素列を見出し語とし、スペースを含んだ音素列（単語別音素列）を見出し語に対応する発音表記とする。

【 0 0 3 7 】

辞書 B は、すべての音素の予め定めた最大音素数の組み合わせにおいて、音素列のスペース部分を“+”に置き換えた組み合わせ音素列を見出し語とし、スペースを含んだ音素列を見出し語に対応する発音表記とする。これによって、学習コーパス 2 0 に含まれていない音素の組み合わせであっても、音素発音辞書 5 0 内に見出し語と発音表記とが登録される。

図 1 に戻って、音素認識辞書生成装置 1 の構成について説明を続ける。

【 0 0 3 8 】

音素列単語リスト記憶手段 1 3 は、音素列単語生成手段 1 1 で生成される音素列単語を、音素列単語リストとして記憶するものである。音素列単語リスト記憶手段 1 3 は、半導

10

20

30

40

50

体メモリ、ハードディスク等の一般的な記憶装置で構成することができる。

【0039】

図5に、音素列単語リスト記憶手段13に記憶される音素列単語リスト130の例を示す。図5に示すように、音素列単語リスト130は、音素列単語生成手段11で生成した単語別音素列のスペースを“+”に置き換えた音素列単語を逐次記憶したものである。

この音素列単語リスト130には、学習コーパス20の書き起こしに含まれる単語の音素列を1つの単語として順次書き込まれる。

【0040】

音素言語モデル生成手段14は、音素列単語リスト記憶手段13に記憶されている音素列単語リスト130から、音素言語モデルを学習により生成するものである。

音素言語モデルは、任意の音素列単語の単語列において、それが文である確率（尤度）を付与する確率モデル（統計的言語モデル）である。この音素言語モデルは、例えば、N-gram言語モデルであって、以下の式（1）に示すように、音素列単語の列 $w_1 w_2 \dots w_{i-1}$ の後にi番目の音素列単語 w_i が出現する条件付き確率（Nグラム確率）を与えるモデルである。なお、桁あふれを防止するため、式（1）の尤度を対数とし、対数尤度とすることが好ましい。

【0041】

【数1】

$$P(w_i | w_{i-N+1} \dots w_{i-1}) \quad \dots \text{式 (1)}$$

【0042】

例えば、学習コーパスの書き起こしで「東京の橋で」という単語列が存在する場合、音素言語モデル生成手段14は、音素列単語リスト130として生成される「t + o : + k y + o :」、「n + o」、「h + a + s h + i」、「d + e」の音素列単語からなる「t + o : + k y + o : n + o h + a + s h + i d + e」という学習テキストでN-gram言語モデルを学習する。

【0043】

なお、音素言語モデル生成手段14は、学習テキストとして音素列単語リスト130に現れない音素列単語の連鎖には、一般的なスムージング手法によってNグラム確率を与える。音素言語モデル生成手段14は、スムージング手法として、例えば、バックオフスムージング（back-off smoothing）を用いることができる。バックオフスムージングは、学習テキストに出現しない音素列単語の連鎖のNグラム確率を、連鎖数の少ない音素列単語の連鎖に与えられているNグラム確率から推定するものである。

【0044】

これによって、音素言語モデル生成手段14は、すべての音素の組み合わせを含んだ音素発音辞書50に登録されている見出し語の音素列単語の連鎖に、Nグラム確率を付与することができる。

音素言語モデル生成手段14は、生成した音素言語モデルを音素言語モデル記憶装置6に書き込み記憶する。

【0045】

図6に、音素言語モデル記憶装置6に記憶される音素言語モデル60の例を示す。ここでは、N-gram言語モデルとして、2-gram言語モデルの例を示す。

図6に示すように、音素言語モデル60は、2つの音素列単語 w_1, w_2 に対して、Nグラム確率（ $\log P(w_2 | w_1)$ ）を対応付けたものである。

【0046】

以上説明したように音素認識辞書生成装置1を構成することで、音素認識辞書生成装置1は、発話音声から音素を認識するための辞書として、音素発音辞書および音素言語モデルを生成することができる。このように生成された音素発音辞書および音素言語モデルは、音素認識を行う際に、単に音素の前後の依存性だけではなく、音素の単語内および単語

10

20

30

40

50

間における音素列の依存性を加味して、音素認識の精度を高めることができる。

なお、音素認識辞書生成装置 1 は、図示を省略したコンピュータを、前記した各手段として機能させるプログラム（音素認識辞書生成プログラム）で動作させることができる。

【0047】

〔音素認識辞書生成装置の動作〕

次に、図 7 を参照（構成については適宜図 1 参照）して、本発明の第 1 実施形態に係る音素認識辞書生成装置 1 の動作について説明する。

【0048】

ステップ S 1 において、単語別音素列生成手段 10 は、学習コーパス 20 の音声から音響特徴量を抽出し、発音辞書 30 と音響モデル 40 を用いて、学習コーパス 20 の音書の書き起こしテキストを事前知識とする音声認識を行い、発音辞書 30 に登録されている見出し語に対応して強制アライメントした単語別音素列を生成する。

10

【0049】

ステップ S 2 において、音素列単語生成手段 11 は、ステップ S 1 で生成した単語別音素列の音素間のスペースを音素以外の予め定めた 1 つのテキスト（例えば、“+”）に置き換えて、音素列単語を生成する。これによって、以降の動作において、単語別音素列を、スペースのない、1 つの単語テキストとして扱うことが可能になる。

【0050】

ステップ S 3 において、音素列単語生成手段 11 は、ステップ S 2 で生成した音素列単語を、順次、音素列単語リスト記憶手段 13 に書き込み記憶する。これによって、音素列単語リスト記憶手段 13 には、学習コーパス 20 の音声に対応する音素列を単語ごとにテキスト化した音素列単語リスト 130 が記録される。

20

【0051】

ステップ S 4 において、音素発音辞書生成手段 12 は、単語別音素列登録手段 120 によって、ステップ S 2 で生成した音素列単語を見出し語とし、ステップ S 1 で生成した単語別音素列をその見出し語に対応する発音表記として、音素発音辞書記憶装置 5 の音素発音辞書 50 に登録する（図 4 の辞書 A 参照）。

【0052】

ステップ S 5 において、単語別音素列生成手段 10 は、学習コーパス 20 の音声についてすべて入力終了したか否かを判定する。ここで、学習コーパス 20 の入力がない場合（ステップ S 5 で No）、音素認識辞書生成装置 1 は、ステップ S 1 に動作を戻す。

30

一方、学習コーパス 20 の入力終了した場合（ステップ S 5 で Yes）、音素認識辞書生成装置 1 は、ステップ S 6 に動作を進める。

【0053】

ステップ S 6 において、音素発音辞書生成手段 12 は、組み合わせ音素列登録手段 121 によって、任意の音素の組み合わせで構成される音素列を単語とみなした見出し語と、その音素列とを対として、音素発音辞書記憶装置 5 の音素発音辞書 50 に登録する（図 4 の辞書 B 参照）。これによって、学習コーパス 20 からは抽出することができない音素の並びに対して、見出し語と発音表記とを割り当てることができる。

40

【0054】

ステップ S 7 において、音素言語モデル生成手段 14 は、ステップ S 3 で順次、音素列単語リスト記憶手段 13 に記憶された音素列単語リスト 130 から、N-gram 言語モデルの音素言語モデル 60 を生成し、音素言語モデル記憶装置 6 に記憶する。

【0055】

さらに、ステップ S 8 において、音素言語モデル生成手段 14 は、音素発音辞書 50 に登録されている音素の組み合わせから生成された見出し語を含めて、学習コーパスとして音素列単語リスト 130 に現れない音素列単語の連鎖に対して、スムージング手法によって N グラム確率を与える。これによって、音素言語モデル 60 を用いて音素認識を行う際に、音素列単語の連結確率が“0”になることを防止することができる。

50

以上の動作によって、音素認識辞書生成装置 1 は、音声から音素を認識するための辞書として、音素発音辞書および音素言語モデルを生成する。

【0056】

<第2実施形態>

〔音素認識装置〕

次に、図8を参照して、本発明の第2実施形態に係る音素認識装置200について説明する。

【0057】

音素認識装置200は、音響モデルと、音素認識辞書生成装置1で生成した音素発音辞書および音素言語モデルとを用いて、音声データから音素を認識するものである。この音素認識装置200は、音響モデル記憶装置4、音素発音辞書記憶装置5および音素言語モデル記憶装置6にそれぞれ記憶されている音響モデル40、音素発音辞書50および音素言語モデル60を用いて、音声データから音素を認識する。

10

【0058】

音響モデル40は、図1で説明した音響モデルと同じであって、大量の音声データから予め学習した音素ごとの音響特徴量をディープニューラルネットワーク(DNN)によってモデル化したものである。

【0059】

音素発音辞書50は、図1で説明した音素認識辞書生成装置1で生成されたものである(図4参照)。

20

音素言語モデル60は、図1で説明した音素認識辞書生成装置1で生成されたものである(図6参照)。

【0060】

音素認識装置200は、図8に示すように、認識手段201と、音素列生成手段202と、を備える。

【0061】

認識手段201は、音響モデル40と、音素発音辞書50と、音素言語モデル60とを用いて、音声データから音素列を認識するものである。

この認識手段201は、外部から入力される音声データから音響特徴量を抽出し、音響モデル40と音素発音辞書50とから音素列単語の候補をリストアップする。そして、認識手段201は、その候補の中で、音素言語モデル60に基づく接続確率が最大となる音素列単語を認識結果とする。

30

【0062】

具体的には、認識手段201は、音素列単語列 w_1, w_2, \dots, w_n で、以下の式(2)に示す、 w_{n-1} の次に w_n が出現する確率(事後確率) $P(w_n | w_{n-1})$ の接続確率が最大となる音素列単語列を認識する。

【0063】

【数2】

$$P(w_1, w_2, \dots, w_n) = \prod_{k=1}^n P(w_k | w_{k-1}) \quad \dots \text{式(2)}$$

40

【0064】

このように、認識手段201は、一般的な音声認識が発音辞書に登録されている単語単位で音声を認識するのに対し、音素発音辞書50に登録されている単語とみなした音素列単語単位で音声を認識する。

認識手段201は、認識した音素列単語を、順次、音素列生成手段202に出力する。

【0065】

音素列生成手段202は、認識手段201で認識された1単語のテキストデータ形式である音素列単語から音素列を生成するものである。

50

具体的には、音素列生成手段 202 は、音素列単語から、音素以外の予め定めた文字（ここでは、“+”）をスペースに置き換えて、音素列を生成する。例えば、音素列生成手段 202 は、音素列単語 “s + e + k + a + i + i + c h + i” を音素列 “s e k a i i c h i” に変換して出力する。

この音素列生成手段 202 が行う変換処理は、図 1 で説明した音素列単語生成手段 11 の変換処理の逆変換に相当する。

【0066】

以上説明したように音素認識装置 200 を構成することで、従来、音響モデルにおけるトライフォン HMM により文脈として前後の音素の依存性で認識をしていた音素認識に対し、音素認識装置 200 は、単語の繋がりをを用いた、より長い文脈の依存性を考慮して音素認識を行う。

10

【0067】

これによって、音素認識装置 200 は、従来よりも精度よく音素認識を行うことができる。具体的には、従来技術の課題で説明したように、従来の音素認識の音素異なり率が 22.8% であったのに対し、音素認識装置 200 は、音素異なり率を 1.2% に改善することができた。

なお、音素認識装置 200 は、図示を省略したコンピュータを、前記した各手段として機能させるプログラム（音素認識プログラム）で動作させることができる。

【0068】

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

20

ここでは、音素発音辞書 50 の見出し語と音素言語モデル 60 の接続対象とを、音素列単語生成手段 11（図 1 参照）が生成した単語別音素列のスペースを “+” とした音素列単語とすることで、1 単語分の音素列を 1 つの単語として扱うこととした。

【0069】

しかし、音素列を 1 単語とみなす手法は、これに限定されるものではない。例えば、1 単語分の音素列の末尾に音素以外の予め定めた文字（例えば、“¥”）を付加することとしてもよいし、単語分の音素列の前後に音素以外の予め定めた文字（例えば、“<”, “>”）を付加することとしてもよい。

この場合も、音素列生成手段 202（図 8 参照）は、音素列単語生成手段 11（図 1 参照）が行った処理の逆変換を行えばよい。

30

【0070】

また、ここでは、音素言語モデル生成手段 14 が生成する音素言語モデル 60 として、2 - g r a m 言語モデルを例示した。

しかし、音素言語モデル生成手段 14 は、N - g r a m 言語モデルであれば、1 - g r a m 言語モデル、3 - g r a m 言語モデル等であっても構わない。

【符号の説明】

【0071】

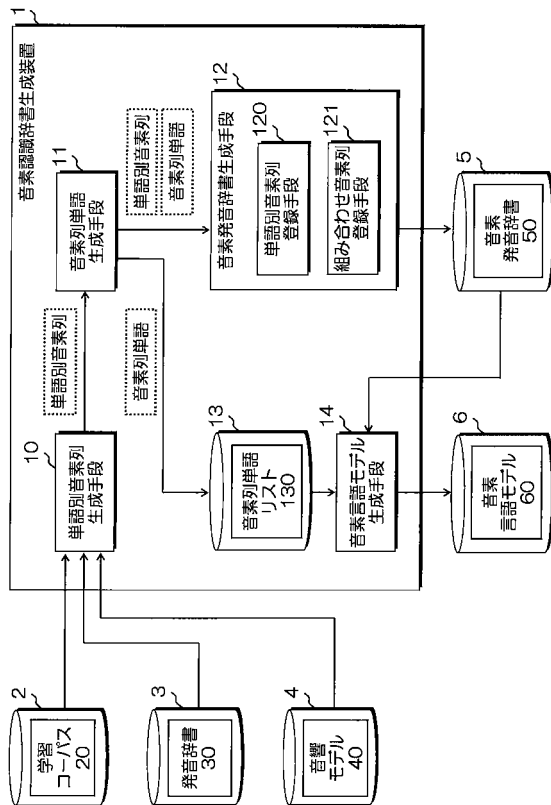
- 1 音素認識辞書生成装置
- 10 単語別音素列生成手段
- 11 音素列単語生成手段
- 12 音素発音辞書生成手段
- 120 単語別音素列登録手段
- 121 組み合わせ音素列登録手段
- 13 音素列単語リスト記憶手段
- 130 音素列単語リスト
- 14 音素言語モデル生成手段
- 2 学習コーパス記憶装置
- 20 学習コーパス
- 3 発音辞書記憶装置

40

50

- 3 0 発音辞書
- 4 音響モデル記憶装置
- 4 0 音響モデル
- 5 音素発音辞書記憶装置
- 5 0 音素発音辞書
- 6 音素言語モデル記憶装置
- 6 0 音素言語モデル

【 図 1 】



【 図 2 】

(a) 20

世界ー 短い 東京 の 橋 で イベント が 開か れ ました

(b) 30

見出し語 (単語)	発音表記 (音素列)
世界ー	s _Δ e _Δ k _Δ a _Δ i _Δ i _Δ ch _Δ i / s _Δ e _Δ k _Δ a _Δ i _Δ ch
短い	m _Δ i _Δ j _Δ i _Δ k _Δ a _Δ i
東京	t _Δ o _Δ ky _Δ o _Δ / t _Δ o _Δ u _Δ ky _Δ o _Δ u / t _Δ o _Δ o _Δ ky _Δ o _Δ u
⋮	⋮

(c)

s_Δe_Δk_Δa_Δi_Δi_Δch_Δi / m_Δi_Δj_Δi_Δk_Δa_Δi / t_Δo_Δky_Δo_Δ / . . .

【 図 3 】

a, a:, b, by, ch, d, dy, e, e:, f, g, gy, h, hy, i, i:, j, k, ky, m, my, n, ny, o, o:, p, py, r, ry, s, sh, t, ts, u, u:, w, y, z, N, Q

【 図 4 】

見出し語 (単語)	発音表記 (音素列)
s+e+k+a+i+i+ch+i	s _Δ e _Δ k _Δ a _Δ i _Δ i _Δ ch _Δ i
m+i+j+i+k+a+i	m _Δ i _Δ j _Δ i _Δ k _Δ a _Δ i
t+o:+ky+o:	t _Δ o _Δ : _Δ ky _Δ o _Δ :
n+o	n _Δ o
h+a+sh+i	h _Δ a _Δ sh _Δ i
:	:
a	a
a:	a:
b	b
:	:
a+a	a _Δ a
a+a:	a _Δ a:
a+b	a _Δ b
:	:
a+a+a+a	a _Δ a _Δ a _Δ a
a+a+a+a:	a _Δ a _Δ a _Δ a:
a+a+a+b	a _Δ a _Δ a _Δ b
:	:

【 図 5 】

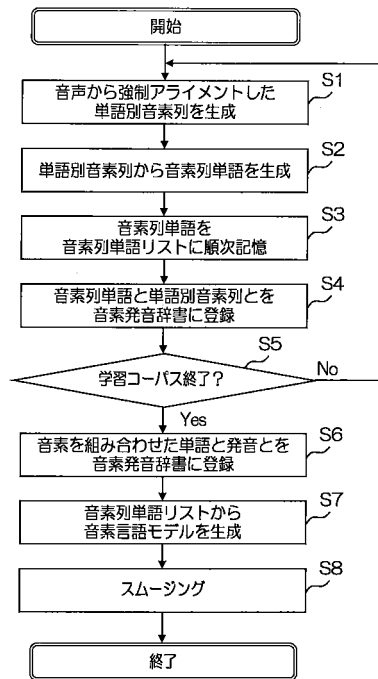
```

s+e+k+a+i+i+ch+i
m+i+j+i+k+a+i
t+o:+ky+o:
n+o
h+a+sh+i
d+e
i+b+e+N+t+o
g+a
h+i+r+a+k+a
r+e
m+a+sh+i
t+a
:
    
```

【 図 6 】

$\log P(w_2 w_1)$	w ₁	w ₂
-2.198956	s+o+r+r+e+d+e+w+a	h+i+r+o+sh+m+a
-1.394523	s+o+r+r+e+d+e+w+a	h+i+t+o+m+a+z+u
-8.217644	s+o+r+r+e+d+e+w+a	h+u+sh+i+N
:	:	:

【 図 7 】



【 図 8 】

