

(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2005-92682

(P2005-92682A)

(43) 公開日 平成17年4月7日(2005.4.7)

(51) Int. Cl.⁷

G06F 17/28

F I

G06F 17/28

P

テーマコード(参考)

5B091

審査請求 未請求 請求項の数 3 O L (全 15 頁)

(21) 出願番号 特願2003-327491 (P2003-327491)
 (22) 出願日 平成15年9月19日(2003.9.19)

(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
 (74) 代理人 100070150
 弁理士 伊東 忠彦
 (72) 発明者 後藤 功雄
 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 加藤 直人
 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) 5B091 AB04 CA21 CC03 CC16 EA01

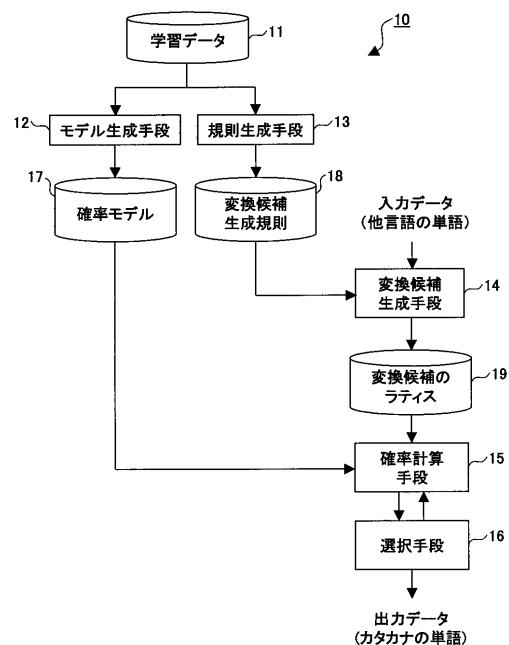
(54) 【発明の名称】 翻字装置、及び翻字プログラム

(57) 【要約】

【課題】 他言語の単語からカタカナの単語への高精度な翻字を実現する。

【解決手段】 他言語の単語とカタカナの単語とにおける部分文字列が対応付けられたデータに基づいて、変換候補の規則を生成する変換候補規則生成手段と、前記他言語の単語を文脈情報に基づいて変換単位に分割するための分割確率を取得するモデルと、前記他言語とカタカナとの部分文字列の対応確率を計算するモデルとを生成するモデル生成手段と、前記他言語で入力される単語を前記変換候補規則生成手段により得られる変換規則に基づいて、前記カタカナの変換候補と前記他言語での変換単位とを生成する変換候補生成手段と、前記モデル生成手段により得られるモデルと、前記他言語と前記カタカナとの文脈情報とに基づいて、変換候補の生起確率を計算する確率計算手段と、前記確率計算手段により得られる生起確率が最大となる変換候補を選択する変換候補選択手段とを有する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

他言語の単語からカタカナの単語を生成するための翻字装置において、

前記他言語の単語と前記カタカナの単語とにおける部分文字列が対応付けられたデータに基づいて、変換候補の規則を生成する変換候補規則生成手段と、

前記他言語の単語を単語内の文脈情報に基づいて変換単位に分割するための分割確率を取得するモデルと、前記他言語とカタカナとの部分文字列の対応確率を単語内の文脈情報に基づいて計算するモデルとを生成するモデル生成手段と、

前記他言語で入力される単語を前記変換候補規則生成手段により得られる変換規則に基づいて、前記カタカナの変換候補と前記他言語での変換単位とを生成する変換候補生成手段と、

前記モデル生成手段により得られるモデルと、前記他言語と前記カタカナとの文脈情報とに基づいて、変換候補の生起確率を計算する確率計算手段と、

前記確率計算手段により得られる生起確率が最大となる変換候補を選択する変換候補選択手段とを有することを特徴とする翻字装置。

【請求項 2】

前記変換候補選択手段は、

予め設定される評価式に基づいて、前記変換候補生成手段にて得られる変換候補からカタカナの変換候補を選択することを特徴とする請求項 1 に記載の翻字装置。

【請求項 3】

他言語の単語からカタカナの単語を生成するための処理をコンピュータに実行させるための翻字プログラムにおいて、

前記他言語の単語と前記カタカナの単語とにおける部分文字列が対応付けられたデータに基づいて、変換候補の規則を生成する変換候補規則生成処理と、

前記他言語の単語を単語内の文脈情報に基づいて変換単位に分割するための分割確率を取得するモデルと、前記他言語とカタカナとの部分文字列の対応確率を単語内の文脈情報に基づいて計算するモデルとを生成するモデル生成処理と、

前記他言語で入力される単語を前記変換候補規則生成処理により得られる変換規則に基づいて、前記カタカナの変換候補と前記他言語での変換単位とを生成する変換候補生成処理と、

前記モデル生成処理により得られるモデルと、前記他言語と前記カタカナとの文脈情報とに基づいて、変換候補の生起確率を計算する確率計算処理と、

前記確率計算処理により得られる生起確率が最大となる変換候補を選択する変換候補選択処理とをコンピュータに実行させるための翻字プログラム。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、翻字装置、及び翻字プログラムに係り、特に、他言語の単語からカタカナの単語へ変換するための翻字装置、及び翻字プログラムに関する。

【背景技術】**【0002】**

従来、英語や日本語等の文字が異なる言語間において、固有名詞は多くの場合に元の単語の発音を表す外来語に翻訳される。特に、日本語では、カタカナを用いた単語に翻訳される場合が多い。

【0003】

ここで、他言語からカタカナへの翻字処理に関する技術は、すでに開示されており、例えば、他言語を発音記号（音韻体系）へ変換してから、発音記号をカタカナへ変換する方式（例えば、非特許文献 1 参照。）や他言語からカタカナへ直接変換する小規模な変換テーブルを用いて変換する方式（例えば、非特許文献 2 参照。）がある。

10

20

30

40

50

【0004】

また、カタカナ以外の文字への翻字も提案されている。例えば、「発音を考慮した変換単位」を用いる英語から韓国語への翻字処理の手法がある（例えば、非特許文献3参照）。また、英語から韓国語への翻字において、決定木を用いて変換単位の曖昧性を解消する手法がある（例えば、非特許文献4参照。）。更に、英語から韓国語への翻字において、決定木を用いて英語の文脈を考慮して変換候補の部分文字列を決定する手法がある（例えば、非特許文献5参照。）。

【非特許文献1】堀内 雄一，山崎 一生．1990．英単語のアルファベット表記から仮名表記への変換．情報処理学会自然言語処理研究会報告，No．79-1，pp．1-8．

【非特許文献2】住吉 英樹，相沢 輝昭．英語固有名詞の片カナ変換．1994．情報処理学会論文誌，Vol．35，No．1，pp．35-45．

【非特許文献3】Byung-Ju Kang and Key-Sun Choi．2000．Automatic Transliteration and Back-Transliteration by Decision Tree Learning．International Conference on Language Resources and Evaluation，pp．1135-1411．

【非特許文献4】In-Ho Kang and GilChang Kim．2000．English-to-Korean Transliteration using Multiple Unbounded Overlapping Phoneme Chunks．The 18th International Conference on Computational Linguistics，Vol．1，pp．418-424．

【非特許文献5】Jong-Hoon Oh and Key-Sun Choi．2002．An English-Korean Transliteration Model using Pronunciation and Contextual rules．The 19th International Conference on Computational Linguistics．

【発明の開示】

【発明が解決しようとする課題】

【0005】

しかしながら、非特許文献1に記載された技術は、他言語を発音記号へ変換することが困難であり、また非特許文献2に記載された技術は、小規模な変換テーブルを用いた方式の場合に詳細な文脈利用ができないため精度に問題があり、高精度な翻字処理を行うことはできない。

【0006】

また、非特許文献3に記載された技術は、変換単位に複数の長さの部分文字列を用いる場合における変換元の部分文字列の選択の曖昧性を考慮しているが、更に高精度な変換を行うためには、翻字元と翻字先との文脈情報を考慮する必要がある。

【0007】

更に、非特許文献4に記載された技術は、変換先の候補との対応関係を考慮せずに元の英語の単語の情報のみで一意に決定しているため、高精度な変換を行っているとはいえない。また、非特許文献5に記載された技術は、英語の1文字を変換の単位をした「発音を考慮しない変換単位」を用いているため、「発音を考慮した変換単位」と比べると精度が低下する。

【0008】

本発明は、上述した問題点に鑑みなされたものであり、他言語からカタカナへの高精度な翻字を行うための翻字装置、及び翻字プログラムを提供することを目的とする。

【課題を解決するための手段】

【0009】

10

20

30

40

50

上記課題を解決するために、本件発明は、以下の特徴を有する課題を解決するための手段を採用している。

【0010】

請求項1に記載された発明は、他言語の単語からカタカナの単語を生成するための翻字装置において、前記他言語の単語と前記カタカナの単語とにおける部分文字列が対応付けられたデータに基づいて、変換候補の規則を生成する変換候補規則生成手段と、前記他言語の単語を単語内の文脈情報に基づいて変換単位に分割するための分割確率を取得するモデルと、前記他言語とカタカナとの部分文字列の対応確率を単語内の文脈情報に基づいて計算するモデルとを生成するモデル生成手段と、前記他言語で入力される単語を前記変換候補規則生成手段により得られる変換規則に基づいて、前記カタカナの変換候補と前記他言語での変換単位とを生成する変換候補生成手段と、前記モデル生成手段により得られるモデルと、前記他言語と前記カタカナとの文脈情報とに基づいて、変換候補の生起確率を計算する確率計算手段と、前記確率計算手段により得られる生起確率が最大となる変換候補を選択する変換候補選択手段とを有することを特徴とする。

10

【0011】

請求項1記載の発明によれば、部分文字列が対応付けられたデータを利用して、変換規則を適用する単位となる変換単位への分割確率に基づいてカタカナの変換候補を選択することにより、他言語からカタカナへの翻字を高精度に行うことができる。また、モデル生成手段により得られるモデルと前記他言語と前記カタカナの文字との文脈情報とにより生起確率を算出することで、より高精度にカタカナの変換候補の選択を高精度に行うことができ、他言語からカタカナへの翻字を高精度に行うことができる。

20

【0012】

請求項2に記載された発明は、前記変換候補選択手段は、予め設定される評価式に基づいて、前記変換候補生成手段にて得られる変換候補からカタカナの変換候補を選択することを特徴とする。

【0013】

請求項2記載の発明によれば、予め設定される評価式を用いることで、一定の評価基準により容易に変換候補を選択することができる。これにより、他言語からカタカナへの翻字を高精度に行うことができる。

【0014】

請求項3に記載された発明は、他言語の単語からカタカナの単語を生成するための処理をコンピュータに実行させるための翻字プログラムにおいて、前記他言語の単語と前記カタカナの単語とにおける部分文字列が対応付けられたデータに基づいて、変換候補の規則を生成する変換候補規則生成処理と、前記他言語の単語を単語内の文脈情報に基づいて変換単位に分割するための分割確率を取得するモデルと、前記他言語とカタカナとの部分文字列の対応確率を単語内の文脈情報に基づいて計算するモデルとを生成するモデル生成処理と、前記他言語で入力される単語を前記変換候補規則生成処理により得られる変換規則に基づいて、前記カタカナの変換候補と前記他言語での変換単位とを生成する変換候補生成処理と、前記モデル生成処理により得られるモデルと、前記他言語と前記カタカナとの文脈情報とに基づいて、変換候補の生起確率を計算する確率計算処理と、前記確率計算処理により得られる生起確率が最大となる変換候補を選択する変換候補選択処理とをコンピュータに実行させる。

30

40

【0015】

請求項3記載の発明によれば、部分文字列が対応付けられたデータを利用して、変換規則を適用する単位となる変換単位への分割確率に基づいてカタカナの変換候補を選択することにより、他言語からカタカナへの翻字を高精度に行うことができる。また、モデル生成処理により得られるモデルと前記他言語と前記カタカナの文字との文脈情報とにより生起確率を算出することで、より高精度にカタカナの変換候補の選択を高精度に行うことができ、他言語からカタカナへの翻字を高精度に行うことができる。また、実行プログラムをコンピュータにインストールすることにより、容易に他言語からカタカナへの翻字を実

50

現することができる。

【発明の効果】

【0016】

本発明によれば、他言語の単語からカタカナの単語への高精度な翻字を実現する。

【発明を実施するための最良の形態】

【0017】

<本発明の概要>

本発明は、外来語のカタカナで表現される単語を、元の外国語（以後、他言語という）から生成するものである。そのために、カタカナの単語とその対訳の他言語の単語との両方の単語内で発音的に類似している部分に対応付けたデータベースに基づいて変換候補を生成し、その変換候補の適用スコアを統計的に学習し、学習結果を利用して翻字（音訳）を行う。

10

【0018】

以下に、上記のような特徴を有する本発明における翻字装置、及び翻字プログラムを好適に実施した形態について、図面を用いて詳細に説明する。なお、本実施例では、他言語の例として、英語の場合について説明するが、韓国語等の言語においても本発明を適用することができる。

【0019】

<機能構成図>

図1は、本発明における翻字装置の機能構成の一例を示す図である。図1の翻字装置10は、学習データ11と、モデル生成手段12と、規則生成手段13と、変換候補生成手段14と、変換候補の確率計算手段15と、最適な変換候補を選択する選択手段16とを有するよう構成されている。

20

【0020】

まず、モデル生成手段12は、学習データ（コーパス）11を入力して確率モデル17を出力する。ここでは、最大エントロピー法に基づく学習を例として扱う。つまり、最大エントロピー法で利用する素性関数を定義して確率モデル17を作成して出力する。なお、素性関数の定義内容については後述する。

【0021】

規則生成手段13は、単語内の部分文字列の対応がつけられた学習データ11を用いて、英語の各部分文字列から変換されているカタカナの部分の文字列の変換候補の規則を生成する。

30

【0022】

例えば、「シ/ソー/ラ/ス：the/sau/ru/s」のデータからは、“the” “シ”，“sau” “ソー”，“ru” “ラ”，“s” “ス”という変換候補生成規則を得る。このような変換候補を学習データ11中の全ての英語とカタカナの部分文字列の対応付けされたデータから生成し、変換候補生成規則18を作成する。

【0023】

なお、上述したモデル生成手段12及び規則生成手段13により確率モデル17及び変換候補生成規則18を作成するまでが学習フェーズとなる。つまり、学習データ11に基づいて、入力される他言語の単語を翻字する処理の前に実行される。また、以下に説明する変換候補生成手段14、確率計算手段15、及び選択手段16が、翻字を実行する実行フェーズとなる。

40

【0024】

次に、実際の翻字を行う際には、例えば、他言語として英語の単語からなる入力データが変換候補生成手段14に入力される。変換候補生成手段14は、英語の単語からカタカナの単語を直接推定する。ここで、変換候補生成手段14における候補生成の内容について具体的に説明する。

【0025】

翻字を行う英語の単語からカタカナの単語の先頭に“^”、単語の末尾に“\$”等の識

50

別子を追加して、英語の単語 E を以下に示す (1) 式のように表現する。

【 0 0 2 6 】

【 数 1 】

$$E = e_0^{m+1} = e_0 e_1 \dots e_{m+1} \quad \dots(1)$$

$$e_0 = \wedge \quad e_{m+1} = \$ \quad \dots(2)$$

10

ここで、 e_j は、英語の単語の j 番目の文字であり、 m は、英語の単語の “ \wedge ” と “ $\$$ ” 以外の文字数である。また、 e_0^{m+1} は、 e_0 から e_{m+1} までの文字列であることを示している。

【 0 0 2 7 】

この英語の単語の各部分に対する対応付けされた英語の部分文字列 eu (English Unit) と、カタカナの部分文字列 ku (Katakana unit) とからなる変換候補生成規則の適用方法は、 E の文字列中に一致する変換候補生成規則の eu を全て適用し、その eu に対応する全ての ku により、ラティス $L\{K\}$ を作成する。

【 0 0 2 8 】

20

ここで、一例として図 2 に英語の単語「actinium」の変換候補のラティス $L\{K\}$ の例を示す。 $L\{K\}$ 中の “ \wedge ” から “ $\$$ ” までの各経路 P_d (P_1, P_2, \dots, P_q) 中の部分文字列を繋いだ文字列が変換先の単語の候補となる。例えば、図 2 において、「c」には、「キ (ki)」、「ク (ku)」、及び「ック (kku)」の 3 つの候補があることを示している。なお、 q は、 $L\{K\}$ 中の “ \wedge ” から “ $\$$ ” までの経路数を示している。

【 0 0 2 9 】

ここで、 $L\{K\}$ 中のある経路 P_d を選択した場合について説明する。この場合の P_d 中の “ \wedge ” 及び “ $\$$ ” 以外の部分文字列の数を $n(P_d)$ とする。また、 P_d 中の部分文字列に、先頭から順番に番号を付与する。上述の条件により、 P_d に対する英語の単語 E とその変換結果のカタカナの単語 K は、次のようになる。

30

【 0 0 3 0 】

【 数 2 】

$$E = e_0^{m+1} = e_0 e_1 \dots e_{m+1} = eu_0^{n(P_d)+1} = eu_0 eu_1 \dots eu_{n(P_d)+1} \quad \dots(3)$$

$$K = k_0^{l(P_d)+1} = k_0 k_1 \dots k_{l(P_d)+1} = ku_0^{n(P_d)+1} = ku_0 ku_1 \dots ku_{n(P_d)+1} \quad \dots(4)$$

$$e_0 = k_0 = eu_0 = ku_0 = \wedge, \quad e_{m+1} = k_{l(P_d)+1} = eu_{m+1} = ku_{n(P_d)} = \$ \quad \dots(5)$$

40

ここで、 k_j はカタカナの単語の j 番目の文字であり、 $m(P_d)$ はカタカナの単語の “ \wedge ” 及び “ $\$$ ” 以外の文字数である。なお、(3) 式における $eu_0^{n(P_d)+1}$ は、 eu_0 から $eu_{n(P_d)+1}$ までの文字列を示し、(4) 式における $ku_0^{n(P_d)+1}$ は、 ku_0 から $ku_{n(P_d)+1}$ までの文字列を示している。

【 0 0 3 1 】

$L\{K\}$ 中の各 P_d における (4) 式の $ku_0^{n(P_d)+1}$ が変換候補のカタカナ単語となる。また、(3) 式の $eu_0^{n(P_d)+1}$ が (4) 式の変換候補を出力する際の

50

英語の単語中の変換単位を示している。

【0032】

変換候補生成手段14は、他言語の単語等が格納されている変換候補生成規則18を入力し、変換候補のラティス19を出力する。出力された変換候補のラティス19は、確率計算手段15に入力される。

【0033】

次に、確率計算手段15は、入力された変換候補のラティス19と、確率モデル17とに基づいて、変換候補の生起確率を計算して選択手段16に出力する。選択手段16は、生起確率に基づいて変換候補を選択して出力する。また、確率計算手段15の処理と選択手段16の処理は交互に繰り返しながら少しずつ処理を行い、最適な変換候補として生起確率が最大となる変換候補を選択する。なお、選択手段16は、確率計算手段15による文脈情報を用いた変換候補の評価に基づいて変換候補の選択を行う。

10

【0034】

ここで、変換候補の評価手法について説明する。まず、英語の単語を入力して対応するカタカナの単語E^を推定するためには、以下に示す(6)式を満たすKを求めればよい。

【0035】

【数3】

$$E^{\wedge} = \arg \max_K P(K | E) \quad \dots(6)$$

20

ここで、P(K | E)は、Eが与えられた場合の、Kの条件付き確率分布を表す。しかしながら、(6)式を直接求めることは未知の単語に対して難しい。そこで、(3)式、(4)式により(6)式中の単語を部分文字列に分解する。分解した式を(7)式に示す。

【0036】

【数4】

$$E^{\wedge} = \arg \max_K \sum_{eu_0^{n(P_d)+1}} \sum_{ku_0^{n(P_d)+1}} P(K | ku_0^{n(P_d)+1}) P(ku_0^{n(P_d)+1} | eu_0^{n(P_d)+1}) P(eu_0^{n(P_d)+1} | E) \quad \dots(7)$$

30

(7)式では、ラティス上の同じKを示す全ての変換候補の確率を合計することで、結果が得られることを示している。

【0037】

また、(7)式のP(eu_0^{n(P_d)+1} | E)は、英語の単語から生成される部分文字列の確率分布であり、変換単位推定モデルと呼ぶ。また、P(ku_0^{n(P_d)+1} | eu_0^{n(P_d)+1})は、英語の部分文字列から生成されたカタカナの部分文字列の確率分布であり翻訳モデルと呼ぶ。更に、P(K | ku_0^{n(P_d)+1})は、カタカナの部分文字列からカタカナの単語が生成される確率分布である。

40

【0038】

ここで、上述の変換単位推定モデル、翻訳モデル、及び確率分布である式、P(K | ku_0^{n(P_d)+1}) P(ku_0^{n(P_d)+1} | eu_0^{n(P_d)+1}) P(eu_0^{n(P_d)+1} | E)に実際の値を入力した例を図に示す。図3は、変換単位推定モデル、翻訳モデル、及び確率分布に実際の値を適用した一例の図である。なお、図3では、“ア

50

クチニウム (actinium) ” を変換単位推定モデル、翻訳モデル、及び確率分布を示す式に適用し、「変換単位推定モデル×翻訳モデル×確率分布」を示している。

【0039】

ここで、(7)式の $P(ku_0^{n(p_d)+1} | eu_0^{n(p_d)+1})$ を、単語単位の処理から部分文字列単位の処理に分解する。これにより、下記に示す(8)式のようになる。

【0040】

【数5】

$$P(ku_0^{n(p_d)+1} | eu_0^{n(p_d)+1}) = \prod_{i=0}^{n(p_d)+1} P(ku_i | ku_0^{i-1}, eu_0^{n(p_d)+1}) \quad \dots(8)$$

10

更に、 $P(ku_i | ku_0^{i-1}, eu_0^{n(p_d)+1})$ の条件の英語の文字列を eu_i と、 eu_i の前 a 文字、 eu_i の後 b 文字だけに近似し、カタカナの文字列を ku_i の前 c 文字だけに近似する。

【0041】

【数6】

20

$$\prod_{i=0}^{n(p_d)+1} P(ku_i | ku_0^{i-1}, eu_0^{n(p_d)+1}) \approx \prod_{i=0}^{n(p_d)+1} P(ku_i | k_{start_ku(i)-c}^{start_ku(i)-1}, e_{start_eu(i)-a}^{start_eu(i)-1}, eu_i, e_{start_eu(i)+1}^{start_eu(i)+b}) \quad \dots(9)$$

ここで、 $start_eu(i)$ は、i 番目の部分文字列 eu_i の初めの文字の位置を示し、 $start_ku(i)$ は i 番目の部分文字列 ku_i の初めの文字の位置を示している。また、上述の a, b, c は定数を示している。

30

【0042】

(7)式の変換単位推定モデル $P(eu_0^{n(p_d)+1} | E)$ の確率は、 $E = eu_0^{m+1}$ の単語を部分文字列に分割する確率(分割確率)であるので、各文字の間が分割点にあるかどうかで全ての分割パターンを表現することができる。分割可能な部分は $m+1$ 個あり、それらが分割点かそうでないかの2値を取ることで、全ての部分文字列への分割を表現することができる。ここで、 e_j と e_{j+1} との間が部分文字列の分割になるかどうかを Z_j で表現する。

【0043】

【数7】

40

$$Z_j = \begin{cases} \text{分割点の場合} \\ \text{分割点ではない場合} \end{cases} \quad \dots(10)$$

Z_j を用いて、 $P(eu_0^{n(p_d)+1} | E)$ を単語単位の処理から文字単位の処理に分解する。

50

【 0 0 4 4 】

【 数 8 】

$$P(eu_0^{n(p_d)+1} | E) = \prod_{j=0}^m P(z_j | z_0^{j-1}, e_0^{m+1}) \quad \dots(11)$$

更に、 Z_j の前 a ' 文字と後 b ' 文字と、 Z_j の前の c ' の分割情報（分割点か、又は 10
分割点ではないかの情報）とを考慮するように近似する。

【 0 0 4 5 】

【 数 9 】

$$\prod_{j=0}^m P(z_j | z_0^{j-1}, e_0^{m+1}) \approx \prod_{j=0}^m P(z_j | z_{j-c}^{j-1}, e_{j-a'+1}^{j+b'}) \quad \dots(12)$$

20

(9) 式、(1 2) 式を用いると、(7) 式は次のようになる。

【 0 0 4 6 】

【 数 1 0 】

$$k^{\wedge} \approx \arg \max_K \sum_{ku_0^{n(p_d)+1}} \sum_{eu_0^{n(p_d)+1}} P(K | ku_0^{n(p_d)+1}) \prod_{i=1}^{n(p_d)+1} P(ku_i | k_{start_ku(i)-c}^{start_ku(i)-1}, e_{start_eu(i)-a}^{start_eu(i)-1}, eu_i, e_{start_eu(i)+1}^{start_eu(i)+b}) \prod_{j=0}^m P(z_j | z_{j-c}^{j-1}, e_{j-a'+1}^{j+b'}) \quad \dots(13)$$

30

この (1 3) 式が、本発明における第 1 の評価式である。

【 0 0 4 7 】

また、(1 3) 式とは別の方法について説明する。(1 3) 式に示すように同じ K を出
力する $ku_0^{n(p_d)+1}$ と $eu_0^{n(p_d)+1}$ についての合計を取らずに、部分
文字列の組み合わせを 1 つだけ選択するように近似する。ここでは、K は、 $ku_0^{n(p_d)+1}$ としている。

【 0 0 4 8 】

【 数 1 1 】

40

$$k^{\wedge} \approx \arg \max_{K=ku_i^{(N)+1}} \prod_{i=1}^{n(p_d)+1} P(ku_i | k_{start_ku(i)-c}^{start_ku(i)-1}, e_{start_eu(i)-a}^{start_eu(i)-1}, eu_i, e_{start_eu(i)+1}^{start_eu(i)+b}) \prod_{j=0}^m P(z_j | z_{j-c}^{j-1}, e_{j-a'+1}^{j+b'}) \quad \dots(14)$$

この (1 4) 式が本発明における第 2 の評価式である。

【 0 0 4 9 】

なお、(1 3) 式、(1 4) 式に示す評価式では、 eu_i に対応する ku_i の確率を求
める際に、 eu_i の前 a 文字及び eu_i の後 b 文字の英語の文脈と、 ku_i の前 c 文字の
日本語の文脈情報を考慮している。これによって、 eu_i の発音を示す ku_i の推定精度 50

を向上させることができる。また、英単語を部分文字列に分割する際に、分割候補の部分の前 a ' 文字と b ' 文字、前 c ' の分割情報という文脈情報を考慮している。

【 0 0 5 0 】

このように、文脈情報を用いて、元の単語を部分文字列へ分割する確率と、元の単語の部分文字列をカタカナの部分文字列へ変換する確率とからカタカナの単語の生起確率を計算し、ビタビアルゴリズム (V i t e r b i a l g o r i t h m) を利用して効率的に確率が最大となるカタカナを選択することで、他言語の単語から高精度にカタカナの単語へ変換することができる。つまり、ラティス L { K } 中の経路から、本発明における第 1 の評価式又は第 2 の評価式を満たすカタカナの文字列を選択し、カタカナの単語として出力する。

10

【 0 0 5 1 】

< 確率モデルの生成における素性関数の定義内容 >

次に、モデル生成手段 1 2 にて行う確率モデルの生成における素性関数の定義内容について説明する。

【 0 0 5 2 】

本発明における評価式 ((1 3) 式 , (1 4) 式) で文脈を考慮する場合には、最大エントロピー法に基づいて構築した確率モデルを利用する。この確率モデルを利用するとモデルが対応できるデータが過疎になることを避けながら文脈情報を全て考慮して確率を求めることができる。なお、本発明におけるモデルの生成においてはこの限りではなく、他の統計的手法を用いてもよい。また、以下の説明では、最大エントロピー法による学習を例として扱う。

20

【 0 0 5 3 】

まず、(1 3) 式と (1 4) 式中の翻訳モデルである (1 5) 式で利用する素性関数について説明する。

【 0 0 5 4 】

【 数 1 2 】

$$P(ku_i | k_{start_ku(i)-1}^{start_ku(i)-1}, e_{start_eu(i)-a}^{start_eu(i)-1}, eu_i, e_{start_eu(i)+1}^{start_eu(i)+b}) \quad \dots(15)$$

30

ここでは、データが過疎になりにくいように文字情報だけでなく、子音、母音、半母音の区別の情報も利用する。そこで、 e_j の子音、母音、半母音の区別の情報を $G(e_j)$ と表す。

【 0 0 5 5 】

【 数 1 3 】

$$G(e_j) = \begin{cases} \text{vowel} & (e_j = \{a, i, u, e, o\}) \\ \text{semi-vowel} & (e_j = \{h, y\}) \\ \text{consonant} & (\text{else}) \end{cases} \quad \dots(16)$$

40

この eu_i , e_j , $G(e_j)$, ku_i を夫々 1 つの属性として、それらの属性の組み合わせにより、素性関数を定義する条件を作成する。

【 0 0 5 6 】

最大エントロピー法に基づいてモデルを構築する際に最も重要なことは、素性関数をど

50

のように定義するかという点にある。そこで、本発明では変換対象の部分文字列に距離が近いことと、連続していることが重要であると考え、以下の属性の組み合わせにより素性関数を定義する。

【0057】

「 $k u_i$ と $e u_i$ 」、「 $k u_i$ と、 $e u_i$ と、 $e u_i$ の前あるいは後、又は前後のいくつかの e 」、「 $k u_i$ と、 $e u_i$ と、 $e u_i$ の前あるいは後、又は前後のいくつかの $G(e)$ 」、「 $k u_i$ と $k_{start_ku}(i) - 1$ 」

また、分割モデルである(17)式では、次の組み合わせにより、素性関数を定義する。

【0058】

【数14】

$$P(z_j | z_{j-c}^{j-1}, e_{j-a+1}^{j+b}) \quad \dots(17)$$

「 Z_j と e_j と e_{j+1} 」、「 Z_j と e_j と e_{j+1} と e_j の前のいくつかの Z と e 」、「 Z_j と e_j と e_{j+1} と e_j の後のいくつかの Z 」、「 Z_j と e_j と e_{j+1} と e_j の前のいくつかの Z と e と後のいくつかの e 」

これにより、素性関数を定義することができるため、この素性関数を用いて最大エントロピー法による確率モデルを生成することができる。

【0059】

なお、最大エントロピー法によるモデルの学習では、例えば、Berger(1996)の確率モデルの構築手法等を用いて(13)式と(14)式で用いる(15)式、又は(17)式の確率モデルを、学習データ11を用いて構築することができる(Adam L. Berger, Stephen A. Della Pietra, and Vincent J. Della Pietra. 1996. A Maximum Entropy Approach to Natural Language Processing. Association for Computational Linguistics, Vol. 22, No. 1, pp. 39-71.)。

【0060】

これにより、本発明における素性関数を用いて最大エントロピー法による確率モデルを生成することができる。

【0061】

ここで、上述したように翻字装置における専用の装置構成により本発明における翻字を行うこともできるが、上述した内容をコンピュータに実行させることができる実行プログラム(翻字プログラム)を生成し、例えば、汎用のパーソナルコンピュータ、ワークステーション等に翻字プログラムをインストールすることにより、本発明における翻字が実現可能となる。

【0062】

ここで、本発明における実行可能なコンピュータのハードウェア構成例について図を用いて説明する。図4は、本発明における翻字処理が実現可能なハードウェア構成の一例を示す図である。

【0063】

図4におけるコンピュータ本体には、入力装置31と、出力装置32と、ドライブ装置33と、補助記憶装置34と、メモリ装置35と、各種制御を行うCPU(Central Processing Unit)36と、ネットワーク接続装置37とを有するよう構成されており、これらはシステムバスBで相互に接続されている。

【0064】

10

20

30

40

50

入力装置 31 は、使用者が操作するキーボード及びマウス等のポインティングデバイスを有しており、使用者からのプログラムの実行等、各種操作信号を入力する。出力装置 32 は、本発明における翻字処理を行うためのコンピュータ本体を操作するのに必要な各種ウィンドウやデータ等を表示するモニタを有し、CPU 36 が有する制御プログラムに基づいて実行結果等を表示することができる。

【0065】

ここで、本発明において、コンピュータ本体にインストールされる実行プログラムは、例えば、CD-ROM等の記録媒体 38 等により提供される。プログラムを記録した記録媒体 38 は、ドライブ装置 33 にセット可能であり、記録媒体 38 に含まれる実行プログラムが、記録媒体 38 からドライブ装置 33 を介して補助記憶装置 34 にインストールされる。

10

【0066】

補助記憶装置 34 は、ハードディスク等のストレージ手段であり、本発明における実行プログラムや、コンピュータに設けられた制御プログラムの他に、ドライブ装置 33 から読み取ることができる学習データや、学習フェーズにおいて作成された確率モデル 17 や変換候補生成規則 18 を蓄積し必要に応じて入出力を行うことができる。

【0067】

CPU 36 は、OS (Operating System) 等の制御プログラム、メモリ装置 35 により読み出され格納されている実行プログラムに基づいて、各種演算や各ハードウェア構成部とのデータの入出力等、コンピュータ全体の処理を制御して、上述した

20

【0068】

ネットワーク接続装置 37 は、通信ネットワーク等と接続することにより、実行プログラムを通信ネットワークに接続されている他の端末等から取得したり、翻字手順を規定したプログラムを実行することで得られた実行結果又は本発明における実行プログラム自体を他の端末等に提供することができる。

【0069】

上述したようなハードウェア構成により、特別な装置構成を必要とせず、低コストで高精度な翻字処理を実現できる。

30

【0070】

次に、上述したようなハードウェア構成により実行される翻字プログラムにおける処理手順について、フローチャートを用いて説明する。なお、フローチャートは学習フェーズと実行フェーズとに分けて示しており、翻字プログラムは、CPU 36 により図 4 に示す各構成部を用いた後述の各処理手順を実行する。つまり、ユーザから入力装置 31 を用いて翻字プログラムの実行指示が入力されると、補助記憶装置 34 に格納されている翻字プログラムをメモリ装置 35 に格納する。CPU 36 は、メモリ装置 35 に格納された翻字プログラムにしたがって本発明における翻字処理に係る機能を実行する。

【0071】

図 5 は、本発明の学習フェーズにおけるモデル生成手順を示す一例のフローチャートである。まず、学習データを入力する (S01)。この学習データは、対訳の単語内において対応付けがされているデータとなる。次に、確率モデルを生成する (S02)。ここでは、上述した素性関数を用い、(13) 式、(14) 式中の確率モデルを統計的に求める。その後、S02 により生成した確率モデルを出力する (S03)。

40

【0072】

次に、図 6 に、本発明の学習フェーズにおける変換候補生成規則作成手順の一例のフローチャートを示す。図 6 に示す変換候補生成規則作成処理では、まず、学習データを入力する (S11)。この学習データは、上述した確率モデル生成手順にて使用される学習データと同様であり、対訳の単語内において対応付けがされているデータである。次に、入力した学習データに基づいて変換候補生成規則を作成する (S12)。ここでは、部分対

50

応付けされたカタカナと英語の単語対を用いてカタカナへの翻字処理のための変換候補生成規則を作成する。その後、S 1 2にて生成された変換候補生成規則を出力する(S 1 3)。

【0073】

次に、実行フェーズについて図を用いて説明する。図7は、本発明の実行フェーズにおける翻字手順を示す一例のフローチャートである。

【0074】

図7において、まず、翻字を行うために他言語データが入力されると(S 2 1)、変換候補生成規則手順にて生成された変換候補生成規則を入力する(S 2 2)。次に、他言語データと変換候補生成規則とから変換候補を生成する(S 2 3)。具体的には、英語からカタカナへ変換する場合は、変換元となる英単語から変換候補生成規則を用いて変換先のカタカナの部分文字列からなるカタカナの変換候補のラティス $L\{K\}$ を生成する。

10

【0075】

次に、上述したモデル生成手順にて生成した確率モデルを入力する(S 2 4)。確率モデルを入力後、S 2 3にて生成した変換候補のラティス $L\{K\}$ を対象に、文脈情報を用いて元の単語を部分文字列へ分割する確率、及び元の単語の部分文字列をカタカナの部分文字列へ変換する確率から変換候補となるカタカナの単語の生起確率を計算する(S 2 5)。

【0076】

次に、最適な変換候補として、S 2 5にて計算された生起確率が最大となる変換候補を選択して出力する(S 2 6)。具体的には、(13)式、(14)式に示した評価式を満たす最適なカタカナの文字列を選択し、その文字列をカタカナの単語として出力する。

20

【0077】

ここで、本発明における第2の評価式である(14)式の上位解は、ダイナミックプログラミング(動的計画法)に基づく、最適な状態遷移が生じた場合の出力確率を求めるアルゴリズムであるビタビアルゴリズムによって、効率的に求めることができる。

【0078】

また、第1の評価式の(13)式を満たす解は、(14)式の上位解となる経路のみを取り扱うことにより、高精度な近似解を効率よく求めることができる。

【0079】

これにより、他言語の単語から高精度にカタカナの単語へ翻字することができる。また、実行プログラムを用いることで、特別な装置構成を必要とせず、汎用のコンピュータで本発明における翻字処理を実行できるため、低コストで高精度なカタカナへの翻字を実現することができる。

30

【0080】

上述したように本発明によれば、他言語の単語からカタカナの単語への高精度な翻字を実現することができる。これにより、例えば、辞書に登録がない場合でも翻字処理によって外国語の単語からカタカナの単語を生成することができるため、外国から日本語へ機械翻訳する際の翻訳率を向上することができる。

【0081】

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

40

【図面の簡単な説明】

【0082】

【図1】本発明における翻字装置の機能構成の一例を示す図である。

【図2】「actinium」の変換候補のラティス $L\{K\}$ の例を示す。

【図3】変換単位推定モデル、翻訳モデル、及び確率分布に実際の値を適用した一例の図である。

【図4】本発明における翻字処理が実現可能なハードウェア構成の一例を示す図である。

50

【図5】モデル生成手順を示す一例のフローチャートである。

【図6】変換候補生成規則作成手順を示す一例のフローチャートである。

【図7】翻字手順を示す一例のフローチャートである。

【符号の説明】

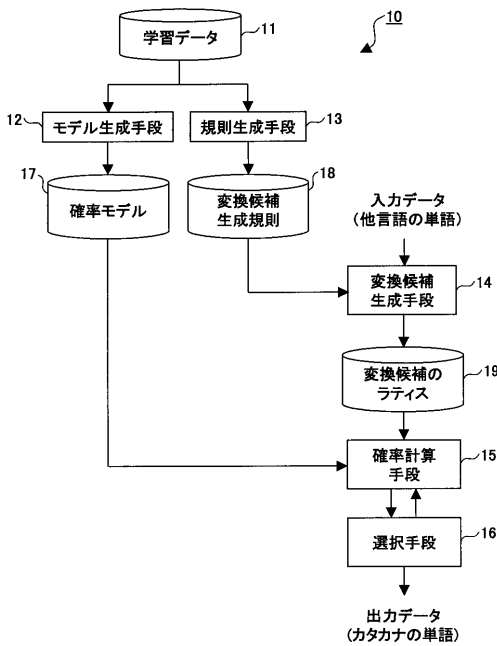
【0083】

- 10 翻字装置
- 11 学習データ
- 12 モデル生成手段
- 13 規則生成手段
- 14 変換候補生成手段
- 15 確率計算手段
- 16 選択手段
- 17 確率モデル
- 18 変換候補生成規則
- 19 変換候補のラティス
- 31 入力装置
- 32 出力装置
- 33 ドライブ装置
- 34 補助記憶装置
- 35 メモリ装置
- 36 CPU
- 37 ネットワーク接続装置
- 38 記録媒体

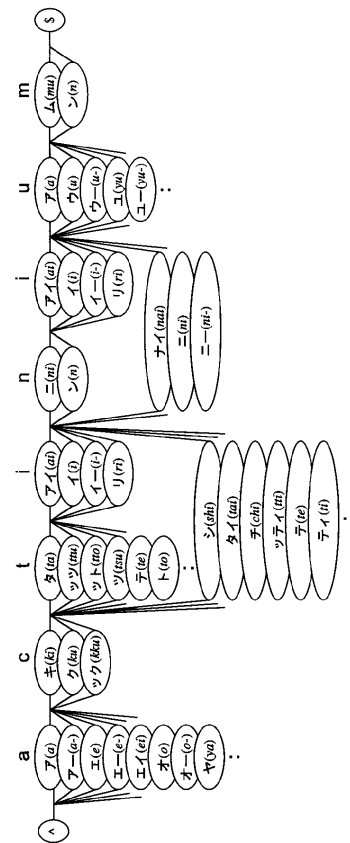
10

20

【図1】



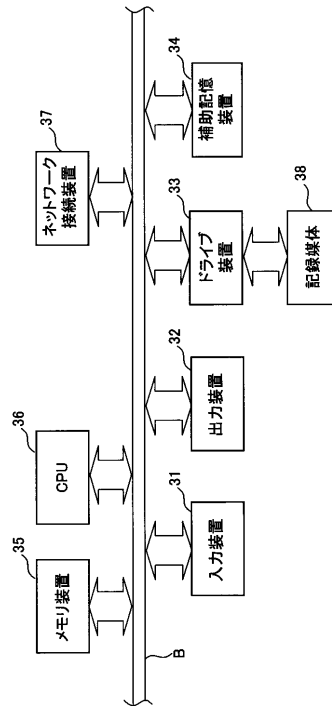
【図2】



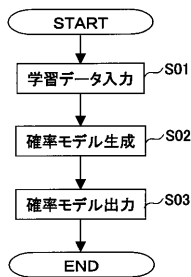
【 図 3 】

$P(\sim\text{アクチニウム}\$|\sim\text{/ア/ク/チ/ニ/ウ/ム/\$}) \times$
 (a ku chi ni u mu) (a/ ku/ chi/ ni/ u/ mu)
 $P(\sim\text{/ア/ク/チ/ニ/ウ/ム/\$}|\sim\text{/a/ku/chi/ni/u/mu/\$}) \times$
 (a/ ku/ chi/ ni/ u/ mu)
 $P(\sim\text{/a/c/ti/ni/u/m/\$}|\sim\text{actinium}\$)$

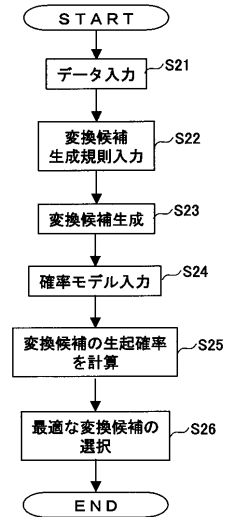
【 図 4 】



【 図 5 】



【 図 7 】



【 図 6 】

