

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3921543号

(P3921543)

(45) 発行日 平成19年5月30日(2007.5.30)

(24) 登録日 平成19年3月2日(2007.3.2)

(51) Int. Cl.

G06F 17/28 (2006.01)

F I

G06F 17/28

Q

請求項の数 3 (全 20 頁)

(21) 出願番号	特願2004-23914 (P2004-23914)	(73) 特許権者	301022471
(22) 出願日	平成16年1月30日(2004.1.30)		独立行政法人情報通信研究機構
(65) 公開番号	特開2005-216127 (P2005-216127A)		東京都小金井市貫井北町4-2-1
(43) 公開日	平成17年8月11日(2005.8.11)	(74) 代理人	100130111
審査請求日	平成16年1月30日(2004.1.30)		弁理士 新保 斉
		(74) 代理人	100090893
			弁理士 渡邊 敏
		(72) 発明者	内元 清貴
			東京都小金井市貫井北町4-2-1 独立 行政法人通信総合研究所内
		(72) 発明者	井佐原 均
			東京都小金井市貫井北町4-2-1 独立 行政法人通信総合研究所内
		審査官	和田 財太
			最終頁に続く

(54) 【発明の名称】 機械翻訳装置

(57) 【特許請求の範囲】

【請求項1】

所定の翻訳元言語で記述される翻訳元テキストを翻訳先言語で記述される翻訳先テキストに機械翻訳する機械翻訳装置であって、

翻訳元テキストを受理する受理手段、

該翻訳元テキストを形態素解析してその結果得られた各形態素を順に着目タームとし、翻訳元言語で記述された複数の文書を含む翻訳元言語文書データベースを用いて、該文書データベースの中から該着目タームが含まれる文書と、該文書データベースに含まれる全文書とのそれぞれにおける着目タームの分布間の距離を算出し、該距離が所定の閾値以上の着目タームを特徴的な意味を有する特徴語として抽出する特徴語抽出手段、

10

該特徴語を、翻訳先言語で表現される訳語に翻訳する訳語選択手段、

該訳語を含む文又は語句の少なくともいずれかを、翻訳先言語で記述された複数の文章を含む翻訳先言語データベースから抽出し、当該訳語と、文又は語句の少なくともいずれかとの関係を自動獲得した生成規則に基づいて文字単位候補を生成する文字単位候補生成手段、

生成された全ての文字単位候補間で依存関係が成立しうる文字単位候補対を全ての文字単位候補について抽出することを繰り返し、異なる依存関係で構成された翻訳先テキスト候補を生成する翻訳先テキスト候補生成手段、

各翻訳先テキスト候補の評価値を算出する評価手段、

該評価値に関連して、少なくとも翻訳先テキスト候補のうち1つを出力する出力手段

20

を備えたことを特徴とする機械翻訳装置。

【請求項 2】

所定の翻訳元言語で記述される翻訳元テキストを翻訳先言語で記述される翻訳先テキストに機械翻訳する機械翻訳装置であって、

翻訳元テキストを受理する受理手段、

該翻訳元テキストを形態素解析してその結果得られた各形態素から、単語列の主辞となる形態素のうち、品詞が動詞、形容詞、名詞、指示詞、副詞、接続詞、連体詞、感動詞のいずれかである語を特徴語として抽出する特徴語抽出手段、

該特徴語を、翻訳先言語で表現される訳語に翻訳する訳語選択手段、

該訳語を含む文又は語句の少なくともいずれかを、翻訳先言語で記述された複数の文章を含む翻訳先言語データベースから抽出し、当該訳語と、該文又は語句の少なくともいずれかとの関係を自動獲得した生成規則に基づいて文字単位候補を生成する文字単位候補生成手段、

生成された全ての文字単位候補間で依存関係が成立しうる文字単位候補対を全ての文字単位候補について抽出することを繰り返し、異なる依存関係で構成された翻訳先テキスト候補を生成する翻訳先テキスト候補生成手段、

各翻訳先テキスト候補の評価値を算出する評価手段、

該評価値に関連して、少なくとも翻訳先テキスト候補のうち 1 つを出力する出力手段を備えたことを特徴とする機械翻訳装置。

【請求項 3】

前記機械翻訳装置において、

前記特徴語抽出手段で抽出された特徴語のうち、該特徴語間の依存関係情報を抽出する依存関係抽出手段を備え、

該依存関係情報を有する特徴語から生成された文字単位候補については、前記翻訳先テキスト候補生成手段において、該依存関係情報を用いて翻訳先テキスト候補を生成することを特徴とする請求項 1 又は 2 に記載の機械翻訳装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は機械翻訳装置に関する。特に、翻訳元言語のテキストを入力し、翻訳先言語のテキストを出力するまでの処理手法に係る技術である。

【背景技術】

【0002】

従来の機械翻訳方法においては、例えば特許文献 1 に開示される機械翻訳の手法が知られている。該特許文献でも開示されているように従来の機械翻訳方法は、例えば日本語から英語への翻訳を行う場合に、日本語の解析を行って、文の意味構造を導出する。

この意味構造は木構造や意味ネットワークで表現できる。そして、生成された意味ネットワークを解析し、適切な訳語を選択して目的とする訳文を生成する処理を行っているものがある。

【0003】

ここで、意味ネットワークは意味記号を付加したノードと、ノード間の関係情報を付加したアークとによって記述されており、各ノードには訳文を生成する際に必要な文生成規則が付加される。生成時には文生成規則に基づき、英語の単語とその文法属性とを検索し、次にアークの情報と整合する単語、品詞、文型をそれぞれ決定する処理を行う。

【0004】

【特許文献 1】特公平 8 - 3 3 8 9 5 号公報

【0005】

また、近年、大量のコーパスが利用可能となり、自然な表層文を生成する目的にもコーパスが利用されるようになってきている。その典型例の 1 つが機械翻訳に用いられる言語モデルであり、例えば非特許文献 1 で開示されている。

10

20

30

40

50

ここで示される統計的機械翻訳では、原言語で書かれたテキストをS、目的言語で書かれたテキストをTとして、Sが与えられたときにTを生成する確率Pが最大になるようなテキストを最適な翻訳として出力する。

この時用いられるモデルとしては、単語や句を原言語から目的言語に置き換えるための翻訳モデルと、置き換えた単語や句を目的言語側で尤もらしい順序に並べ替えるための言語モデルがある。

【0006】

【非特許文献1】Brown,P.F., Cocke,J., Piera,S.A.D.,Pietra,V.J.D., Jelinek,F., La Ifferty,J.D., Mercer,R.L., and Roossin,P.S.(1990). "A Statistical Approach to Machine Translation." Computational Linguistics,16(2),79-85

10

【0007】

このような従来 of 統計的機械翻訳の手法では、与えられた語の集合を並べ換えると自然な文を生成できるという仮定がある。つまり、自然な文を生成するための語の集合は翻訳モデルにより過不足なく生成されることが前提となっている。

しかし、この前提のためには大規模な対訳コーパスが必要であり、日本語などの比較的コーパスが整備された言語が原言語であっても、対象言語との対訳コーパスの状況、対象言語におけるコーパスの状況によっては、上記従来 of 手法では十分な翻訳精度の向上が望めない問題があった。

【0008】

また、特に話し言葉や、文法的な誤り、記述上の誤りなど原言語の表現が不適切な場合に、これら従来 of 方法では正しい意味構造の解析や翻訳モデルの作用が不可能であり、全く誤った翻訳結果を出力してしまうことがある。

20

【発明の開示】

【発明が解決しようとする課題】

【0009】

本発明は、上記従来技術の有する問題点に鑑みて創出されたものであり、その目的は、翻訳元言語又は翻訳先言語に十分なコーパス等の環境が整っていない場合や、入力するテキストが不正確な場合にも、好適な翻訳先言語のテキストを出力可能な機械翻訳装置を提供することである。

【課題を解決するための手段】

30

【0010】

本発明は、上記の課題を解決するために、次のような機械翻訳装置を提供する。

すなわち、所定の翻訳元言語で記述される翻訳元テキストを翻訳先言語で記述される翻訳先テキストに機械翻訳する機械翻訳装置である。

そして、請求項1に記載の機械翻訳装置は、(1)翻訳元テキストを受理する受理手段、(2)該翻訳元テキストを形態素解析してその結果得られた各形態素を順に着目タームとし、翻訳元言語で記述された複数の文書を含む翻訳元言語文書データベースを用いて、該文書データベースの中から該着目タームが含まれる文書と、該文書データベースに含まれる全文書とのそれぞれにおける着目タームの分布間の距離を算出し、該距離が所定の閾値以上の着目タームを特徴的な意味を有する特徴語として抽出する特徴語抽出手段、(3)該特徴語を、翻訳先言語で表現される訳語に翻訳する訳語選択手段、(4)該訳語を含む文又は語句の少なくともいずれかを、翻訳先言語で記述された複数の文章を含む翻訳先言語データベースから抽出し、当該訳語と、該文又は語句の少なくともいずれかとの関係を自動獲得した生成規則に基づいて文字単位候補を生成する文字単位候補生成手段、(5)生成された全ての文字単位候補間で依存関係が成立しうる文字単位候補対を全ての文字単位候補について抽出することを繰り返し、異なる依存関係で構成された翻訳先テキスト候補を生成する翻訳先テキスト候補生成手段、(6)各翻訳先テキスト候補の評価値を算出する評価手段、(7)該評価値に関連して、少なくとも翻訳先テキスト候補のうち1つを出力する出力手段を備えたことを特徴とする。

40

【0011】

50

請求項 2 に記載の機械翻訳装置は、上記請求項 1 の機械翻訳装置において、特徴語抽出手段が翻訳元テキストを形態素解析してその結果得られた各形態素から、単語列の主辞となる形態素のうち、品詞が動詞、形容詞、名詞、指示詞、副詞、接続詞、連体詞、感動詞のいずれかである語を特徴語として抽出する構成である。

【 0 0 1 2 】

請求項 3 に係る機械翻訳装置は、前記特徴語抽出手段で抽出された特徴語のうち、該特徴語間の依存関係情報を抽出する依存関係抽出手段を備え、該依存関係情報を有する特徴語から生成された文字単位候補については、前記翻訳先テキスト候補生成手段において、該依存関係情報を用いて翻訳先テキスト候補を生成することができる。

【 発明の効果 】

【 0 0 1 3 】

以上の発明により次の効果を奏する。

すなわち請求項 1 ないし 3 に記載の機械翻訳装置によれば、新しい機械翻訳手法を導入することで、翻訳元言語と翻訳先言語に十分なコーパス（特に対訳コーパス）が整っていない場合でも、高精度な機械翻訳装置を提供することができる。

また、入力するテキストから特徴語を抽出するため、文法的に不正確なテキストを入力しても、自然で理解のしやすい翻訳を実現することができる。

【 0 0 1 4 】

特に、請求項 3 に記載の機械翻訳装置では、翻訳元言語テキストにおける依存関係を、翻訳先言語の生成ステップで用いることができるので、高精度な翻訳に寄与することができる。

【 0 0 1 5 】

このような機械翻訳装置は、特に口語を入力したものや、記事のタイトルなど文法的に正確でない翻訳元言語の文章を翻訳する際に適している。得られる結果は大意が分かりやすい翻訳先テキストである。

また、本実施例で示した日本語・英語間のようなコーパス環境の整った言語間はもとより、地域性の強い言語や、発展途上国の固有言語など、コーパスが未整備な言語に翻訳する際にも、実際の機械翻訳結果を生成することができる。

【 発明を実施するための最良の形態 】

【 0 0 1 6 】

以下、本発明の最良と考えられる実施形態を、図面に示す実施例を基に説明する。なお、実施形態は下記に限定されるものではない。

図 1 には本発明による第 1 実施例の機械翻訳方法のフローチャートを示す。図のように、翻訳元言語である日本語のテキスト（1）を入力受理し、該テキスト（1）中の特徴的単語を抽出（2）する。抽出された特徴的単語は特徴的単語データ（6）として記録する。なお、本発明はいかなる言語をも対象とするので、本発明は単語として分類できるものに限らず、単語、語句など広く含む特徴語を抽出し、後段の処理を行ってもよい。

【 0 0 1 7 】

次に、翻訳元言語と翻訳先言語の辞書データ（7）を用いて特徴的単語の最適な訳語を選択（3）する。これにより、翻訳先言語の特徴的単語が生成され特徴的単語訳語データ（8）として記録する

さらに、特徴的単語訳語データ（8）を用い、翻訳先言語のテキストを生成（4）する。このように特徴的単語訳語から、自然な英語テキスト（5）を生成することで機械翻訳を行うプロセスは、従来創出されていない、本発明の最大の特徴である。

【 0 0 1 8 】

テキスト生成（4）は、特徴的単語訳語から所定の生成規則に基づいて文字単位候補を生成し、その各文字単位間の依存関係の組み合わせをかえて翻訳先テキスト候補を生成する。さらに、翻訳先テキスト候補を評価して最も適切と判断される英文テキスト（5）を出力する。

【 0 0 1 9 】

10

20

30

40

50

また、本発明では図2に示すような第2実施例の機械翻訳方法を提供することができる。翻訳元言語である日本語のテキスト(10)を入力し、該テキスト(10)中の特徴的な単語を抽出(11)する。抽出された特徴的単語は特徴的単語データ(16)として記録する。

【0020】

特徴的単語抽出時に、各単語間の依存関係を抽出(12)する。該単語間依存関係データ(17)を記録する。

そして、翻訳元言語と翻訳先言語の辞書データ(18)を用いて特徴的単語の最適な訳語を選択(13)する。これにより、翻訳先言語の特徴的単語が生成され特徴的単語訳語及び前述の依存関係のデータ(19)を記録する

10

さらに、該データ(19)を用い、翻訳先言語のテキストを生成(14)する。このように特徴的単語訳語から、自然な英語テキスト(15)を生成する。

【0021】

テキスト生成(14)は、特徴的単語訳語から所定の生成規則に基づいて文字単位候補を生成し、その各文字単位間の依存関係を仮定して単数又は複数の翻訳先テキスト候補を生成する。上記で特徴的単語の依存関係が記録されている場合には、その依存関係を元に翻訳先テキスト候補を生成することで、依存関係に誤りのない高精度な翻訳先テキストが生成できる。

さらに、翻訳先テキスト候補を評価して最も適当と判断される英文テキスト(5)を出力する。

20

【0022】

本発明の第3実施例としては、上記第2実施例のテキスト生成(14)における処理を一部変更し、記録した単語間依存関係(17)に関わらず、特徴的単語訳語から所定の生成規則に基づいて文字単位候補を生成し、その各文字単位間の依存関係を仮定して単数又は複数の翻訳先テキスト候補を生成する。

その上で、翻訳先テキスト候補を評価する際に、上記単語間依存関係との同一性を調べ、その値が高い英語テキストを出力する。

【0023】

以上に述べたような機械翻訳方法を実現する機械翻訳装置の構成を図3に示す。本装置(30)は、例えばある新聞記事「官邸前などでドゥダエフ政権部隊と激しい市街戦を展開している。」という日本語テキスト(31)を入力すると、受理手段である入力受理部(40)で装置(30)内への取り込み処理を行い、特徴的単語抽出部(50)において「官邸」「ドゥダエフ」「部隊」「激しい」「戦」「展開」を日本語テキスト(31)から抽出する。

30

【0024】

さらに訳語選択部(60)で各特徴的単語の最適な訳語、ここでは「palace」「Dudayev」「troops」「fierce」「battle」「engage」を選択する。

選択された特徴的単語訳語を、テキスト生成部(70)において適切な単語間の補完を行いながらテキスト生成(35)し、出力部(80)から英語テキスト(32)を出力する。

40

次に各部(40)ないし(80)の詳細を説述する。

【0025】

入力受理部(40)は図4に示すようにCPU(41)とそれに接続されたスキャナ(42)や、CDドライブ、ハードディスクドライブ、MOドライブ、フロッピー(登録商標)ディスクドライブなどの外部記憶装置(43)等から構成される。また、CPU(41)の動作に伴い、必要に応じて公知のメモリを用いることもできる。

スキャナ(42)を備える場合にはCPU(41)において文字認識処理を行い、テキストデータに変換して外部記憶装置(43)に記録する。外部記憶装置(43)から直接日本語テキスト(31)のデータを読み出す場合にも、CPU(41)において本装置(30)で処理可能な形式にデータ変換を行うこともできる。

50

【 0 0 2 6 】

また、本発明はインターネットやイントラネットのネットワーク(44)を介して他のコンピュータサーバー等からテキストデータを受信することも可能である。

入力受理部(40)により日本語テキスト(31)は図5に示される特徴的単語抽出部(50)に送られる。

【 0 0 2 7 】

特徴的単語抽出部(50)の構成を図5に示す。ここでもCPU及びメモリが協働して各処理を行う。本特徴的単語抽出部(50)では、入力された日本語テキスト(31)からそのテキストの内容を特徴的に表す特徴語を抽出する。

このような技術は、言語処理において文書を要約する技術や、文書検索などの要素技術として公知の多数の手法が知られており、それらを適宜用いることができるが、ここでは一例として非特許文献2に記載の方法を用いる。

【 0 0 2 8 】

【非特許文献2】情報処理学会自然言語処理研究会 1999-NL-33, 1999「タームのrepresentativeness」を測る」久光徹、丹羽芳樹、辻井潤一

【 0 0 2 9 】

本方法によると、特徴語を選ぶために文書中の単語の話題性もしくは分野代表性(representativeness、本明細書ではこれを特徴性と呼ぶ。)を測ることが可能であり、かつ数値的な評価によるため、本発明の実施に好適である。以下に、簡単に説述する。

まず、本特徴的単語抽出部(50)では、公知の形態素解析技術を用いて、日本語テキスト(31)を形態素解析部(51)において形態素解析する。解析された形態素はメモリ又は図示しない外部記憶装置などに形態素テーブルとして記録する。

【 0 0 3 0 】

そして、形態素テーブルから形態素を順次読み出し、その形態素(以下、これを着目タームと呼ぶ)毎に特徴性を測る。

まず文書抽出部(52)において、着目タームWについて、Wを含む文書すべてを任意の文書データベース(56)から抽出する。文書データベース(56)は複数の日本語(翻訳元言語)の文書が含まれたものであり、外部記憶装置などに記憶されている。日本語単言語のコーパスや日英の対訳コーパスの日本語部分を用いてもよい。

【 0 0 3 1 】

次に、着目タームWが抽出された文書すべての集合における単語分布と、文書データベース(56)に含まれる全文書の単語分布とを、単語分布算出部(53)において算出し、各単語分布間の異なり度の度合いを測る。

具体的には異なり度合算出部(54)において次のような計算処理を行う。

【 0 0 3 2 】

すなわち、着目タームW、Wを含む文書すべての集合D(W)、全文書の集合D₀、D(W)における単語分布P_{D(W)}、D₀における単語分布P₀として、Wの特徴性Rep(W)を、2つの分布{P_{D(W)}, P₀}の距離Dist{P_{D(W)}, P₀}に基づいて定義する。

単語分布間の距離計測方法として、本実施例では対数尤度比を用いている。すなわち、全単語を{W₁, ..., W_n}、単語w_iがD(W)、D₀に出現する頻度をそれぞれk_i、K_iとするとき、P_{D(W)}、P₀の距離Dist{P_{D(W)}, P₀}を、次のように定義する。

【 0 0 3 3 】

【数1】

$$Dist(P_{D(W)}, P_0) = \sum_{i=1}^n k_i \log \frac{k_i}{\#D(W)} - \sum_{i=1}^n k_i \log \frac{K_i}{\#D_0}$$

ここで、#D(W)は着目タームWについてD(W)の含む単語数、#D₀は同様に全文

10

20

30

40

50

書の含む単語数である。

【0034】

数1の定義によると、 $\#D(W)$ が離れた着目ターム同士の特徴性を有効に比較することが難しいため、数2のように正規化を行った特徴性 $Rep(W)$ を定義する。なお $B(\cdot)$ は $\#D(W)$ が適切な数となる範囲内(例えば1000 $\#D(W)$ 20000)で特徴性が精度よく求められるような指数関数を用いた近似関数である。

【0035】

(数2)

$$Rep(W) = Dist\{P_{D(W)}, P_0\} / B(\#D(W))$$

【0036】

ここで、「する」などのように著しく $\#D(W)$ が大きい場合には、 $D(W)$ の抽出数を限定し、 $\#D(W)$ 20000を満たすようにすることで、上記近似関数を有効に用いることができると共に計算量を削減できる。

【0037】

特徴的単語抽出部(50)では以上の方法により特徴性を算出すると共に、所定の閾値に従って、特徴的単語決定部(55)により入力した日本語テキスト(31)の特徴的単語を抽出する。この結果は外部記憶装置に特徴的単語テーブル(33)として保存される。

【0038】

次に、図6に示す訳語選択部(60)では、特徴的単語(33)毎に最適な訳語、特徴的単語訳語(34)を選択する処理を行う。単語の翻訳は、通常辞書データベースを参照することにより可能であるが、異なる言語間ではしばしば単語の多義性の問題が生じる。すなわち、ある単語を入力しただけでは、複数の訳語のいずれを選択すべきかが不明であり、最適な訳語の選択は本機械翻訳装置(30)の翻訳精度に影響する極めて重要な技術である。

【0039】

このような多義性解消については、公知の優れた技術が多数提案されており、本発明の実施には任意の技術を用いることができるが、ここでは本件出願人らによって非特許文献3及び非特許文献4に開示される次のような訳語選択モデルを用いる例を示す。

【0040】

【非特許文献3】電子情報通信学会 言語理解とコミュニケーション研究会 NLC2001-41「翻訳メモリとコーパスを用いた学習に基づく訳語選択」内元清貴、関根聡、村田真樹、井佐原均 2001年

【非特許文献4】「Word Translation based on Machine Learning Models Using Translation Memory and Corpora」Kiyotaka Uchimoto, Satoshi Sekine, Masaki Murata, and Hitoshi Isahara. SENSEVAL-2, pp.155-158, 2001

【0041】

訳語選択部(60)は特徴的単語テーブル(33)を入力として、まず文字列の類似性判定部(61)において、文字列の類似性に基づく方法で訳語の決定を試みる。

すなわち、該テーブル(33)に含まれる特徴的単語が、対訳辞書データベース(63)において一義的に訳語が選択可能な場合には、該訳語を特徴的単語訳語テーブル(34)に出力する。

【0042】

また、日本語と英語の対訳コーパスデータベース(本実施例では対訳辞書との明確な区別を設けずに1個のデータベースとする)(63)を参照し、特徴的単語テーブル(33)との類似度を求める。

例えば、特徴的単語として「母」「遠慮」が含まれるとき、対訳コーパスデータベース(63)に含まれる「母に遠慮」「母への遠慮」「献金を遠慮」などとの類似度を求める。

【0043】

類似度を求める方法は、非特許文献5で開示されるUNIX(登録商標)のdiffコマンドを用いた方法によってもよく、以下の数3で求められる。

10

20

30

40

50

【 0 0 4 4 】

【非特許文献5】情報処理学会自然言語処理研究会 2001-NL-44, 2001 「diffと言語処理」 村田真樹、井佐原均

【 0 0 4 5 】

(数3)

類似度 = (特徴的単語テーブルと対訳コーパスとのdiffをとったときに一致した文字数) / (対訳コーパスの文字数)

【 0 0 4 6 】

ここで対訳コーパスについても、比較する前に機能語・動詞・形容詞の活用部分、サ変動詞をすべて削除する、あるいは該対訳コーパスによっても前述の特徴的単語の抽出を行い、それらの結果と類似度を比較するのが望ましい。

10

これにより、比較の対象として不適切な文字を多く含む対訳コーパスとの比較を避け、効果的な類似度の算出が可能となる。

【 0 0 4 7 】

類似度が所定の閾値(対訳コーパス及び特徴的単語テーブルの内容により適宜設定することができる)を超える対訳コーパスを特徴的単語の訳語を含む対訳コーパスとして選択する。

そして、該対訳コーパス中の訳語を特徴的単語訳語テーブル(34)に記録する。対訳コーパス(63)を用いるので、各単語の対訳関係はあらかじめ分かっており、特徴的単語と訳語の関係は機械的に決定できる。

20

【 0 0 4 8 】

ところで、対訳辞書・対訳コーパスデータベース(63)に入力される全ての特徴的単語テーブルに類似する文例を備えることは難しく、類似度が閾値以上とならないものが多数残ってしまう場合が考えられる。

そこでさらに、機械学習モデルを用いた類似性判定部(62)を備え、学習したデータ(64)を用いながら最適な訳語の選択を行う。

【 0 0 4 9 】

機械学習のモデルとしては、SVM(Support Vector Machine)を用いるが、ME(MaximumEntropy)、DL(Decision List)、SB(Simple Bayes)を用いてもよい。各分類クラスの確信度は基本的に文脈の集合をB、クラスの集合をAとするとき、文脈b(B)でクラスa(A)となる事象(a,b)の確率分布p(a,b)として求められる。SVMではこのような確率分布は得られないが、便宜的に最適のクラスに対して確率値を1、その他のクラスに対して0とする。

30

【 0 0 5 0 】

文脈bの素性としては、例えば次のものを用いることができる。

- (1) 形態素情報
- (2) 文字n-gram
- (3) 最大一致となる用例に関する情報
- (4) 内容語とその訳語候補の出現頻度

SVMによって、対訳コーパスを用いて機械学習を行い、学習データ(64)として記録する。

40

【 0 0 5 1 】

特徴的単語について対訳辞書データベース(63)から訳語候補を抽出する。そして、訳語候補の素性は英主辞単語、文集合、単語、出現頻度の組み合わせによって表される。機械学習による学習データ(64)を用い、訳語候補から特徴的単語訳語を選択し、特徴的単語訳語テーブル(34)に記録する。

このような訳語選択モデルには公知の手法を任意に用いることが可能であり、上記のような機械学習モデルによるものでなくともよい。

【 0 0 5 2 】

ここで、特徴的単語の訳語を選択する際に、固有表現などが抽出されて対訳辞書データベ

50

ースに訳語がない場合が想定される。このとき、訳語選択モデルに問い合わせ処理を行う処理部を設けてユーザーに問い合わせを行ってユーザが与えるようにしてもよい。

これを自動化する方法としては、非特許文献6に開示されるように、単に翻訳元言語の発音に従って、一定のルールにより翻訳先言語の文字に置き換える(音訳する)こともできるし、提案されているモデルを用いることもできる。ここで提案されているのは、人名や組織名などで、まず音訳を行ったり、適当な訳語で翻訳して訳語候補を作成し、その候補の中から所定のテキストデータベースに出現する頻度の高いものを訳語として選択するものである。テキストデータベースとしては、例えばインターネットのWeb情報などを利用することができる。

【0053】

次に、以上により形成された特徴的単語訳語テーブル(34)を、テキスト生成部(70)に入力し、英語テキストを生成する。

いくつかの単語を入力し、その単語を含むテキストを生成する方法としては次のような手法がある。すなわち、本件出願人が特許文献2で開示するテキスト生成方法を、翻訳先言語である英語に適用して用いる。

【0054】

【特許文献2】特開2003-271592号公報

【0055】

本テキスト生成部(70)の具体的な構成例として図7に示す各部を備える。テキスト生成部(70)は、例えばCPUとメモリ、ハードディスクなどの外部記憶媒体を備えるパーソナルコンピュータなどにより構成することができ、主な処理をCPUにおいて行い、処理の結果を随時メモリ、外部記憶媒体に記録する。

【0056】

本実施例で、入力された特徴的単語訳語(34)は2つの処理に用いられる。その1つは単語列生成規則獲得部(71)であり、もう1つは単語列候補生成部(72)である。

このとき、特徴的単語訳語(34)は単語列の主辞となる内容語であると定義する。また、内容語は、その語の品詞が、動詞、形容詞、名詞、指示詞、副詞、接続詞、連体詞、感動詞、未定義語である形態素の見出し語であるとし、それ以外の形態素の見出し語を機能語とする。

【0057】

単語列生成規則獲得部(71)では、特徴的単語訳語(34)が与えられたとき、それぞれを含む文を英語コーパス(75)から検索し、形態素解析、構文解析(係り受け解析)をする。そして、そこから特徴的単語訳語(34)を含む単語列を抽出して、特徴的単語訳語から単語列を生成する単語列生成規則(76)を獲得し、記録する。例えば、「palace」「at the palace」、「palace」「in the palace」、「battle」「battles」、「engage」「They are engaged in」、「engage」「I engaged」などの単語列生成規則(76)を獲得し、記録する。

【0058】

ここで、生成規則の自動獲得には次の手法を用いる。特徴的単語訳語の集合をVとし、特徴的単語訳語k(V)から単語列を生成する規則の集合をRkとするとき、規則rk(Rk)は次の形式で表現されるものと定義する。

$k \quad hk \quad m^*$

hkは特徴的単語訳語を含む主辞形態素、 m^* は同じ単語列内でhkに連続する任意個の形態素とする。特徴的単語訳語が与えられると、この形式を満たす規則を翻訳先言語のコーパス(75)から自動獲得する。

【0059】

一方、単語列候補生成部(72)では、単語列生成規則(76)を参照しながら、入力された特徴的単語訳語(34)から出力する英語テキスト(32)を構成する単語列の候補を生成する。

例えば、「palace」では自然なテキストを構成する単語列とはなりにくい、「at the p

10

20

30

40

50

「alace」あるいは「inthe palace」のように「palace」という単語と極めて密接な関連性を有する語句を付加し、後段の処理によるテキスト生成に備える。

【0060】

本実施例のように、単語列生成規則獲得部(71)によりコーパス(75)から入力する特徴的単語訳語(34)の単語列規則を生成することで、最小限の計算量で効果的に単語列生成規則を得ることができ、処理速度の向上に寄与する。

【0061】

もっとも、必ずしも特徴的単語訳語(34)に関連する語句をコーパスから抽出する構成を取る必然性はなく、計算能力に応じて任意の語句を入力された特徴的単語訳語(34)の前後に付加してもよい。あるいは、本装置(30)内に備える対訳辞書データベース(63)に含まれる慣用表現の情報から単語列を生成することもできる。上記「palace」「at the palace」などは対訳辞書データベースに記載される表現であり、単語列の候補として生成することができる。

【0062】

また、日本語など主格を多く省略する言語を入力した場合には、「engage」「They are engaged in」などのように主語を補って単語列候補を生成することができる。このとき、日本語などの多くの言語では主格が明らかな時や、形式主語であるときに省略されることに着目し、入力に主格が何であるかの情報だけでなく、主格がないという情報を用いることで、「engage」「He is engaged in」を生成せず、「engage」「They are engaged in」を生成するようにすることもできる。

【0063】

次に、テキスト候補生成部(73)においてテキスト候補を生成する。テキスト候補はグラフあるいは木の形で表現する。ここでは特徴的単語訳語(34)のうち、「palace」「troops」「engage」の3語の関係を例として説述する。

すなわち、図12のように、各単語列候補(34 aないし34 f)の間に係り受けの関係を仮定して、テキスト候補1(35')、テキスト候補2(35'')のような単語列を単位とした依存構造木の形でテキスト候補を生成する。このとき、3語の場合に全ての係り受け関係は $R! \times 2 = 12$ 通りであるが、翻訳先言語の文法・特性に合わせて語順の固定などにより候補の数を削減することができる。

【0064】

生成されたテキスト候補(35'など)は、評価部(74)でコーパスから学習した特徴的単語訳語生成モデル(77)や言語モデル(78)を用いて順序付けされる。

以下、特徴的単語訳語生成モデル(77)と、言語モデル(78)として形態素モデル及び係り受けモデルについて説述する。

【0065】

特徴的単語訳語生成モデルでは、次の5種類の情報を素性として用いたモデル(KM1ないし5)を考える。以下で、特徴的単語訳語の集合Vは、ある回数以上コーパスに出現した主辞単語の集合とし、単語列は前記で表現されるものと仮定する。また、各特徴的単語訳語は独立であり、与えられたテキストが単語列 $w_1 \cdots w_m$ からなるとき、特徴的単語訳語 k_i は単語 w_j ($1 \leq j \leq m$)に対応していると仮定する。図13にモデルの説明図を示す。

【0066】

[KM1]

前方の二単語を考慮(trigram)

k_i は前方の二単語 w_{j-1} と w_{j-2} のみに依存すると仮定する。

【数4】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_{j-1}, w_{j-2})$$

10

20

30

40

50

【 0 0 6 7 】

[K M 2]

後方の二単語を考慮(後方 trigram)

ki は後方の二単語wj+1 とwj+2 のみに依存すると仮定する。

【 数 5 】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_{j-1}, w_{j-2})$$

10

【 0 0 6 8 】

[K M 3]

係り単語列を考慮(係り単語列)

ki を含む単語列に係る単語列がある場合、ki はそのうち最も文末側の単語列の末尾から二単語wi とwi-1 のみに依存すると仮定する(図 1 3 参照)。

【 数 6 】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_i, w_{i-1})$$

20

【 0 0 6 9 】

[K M 4]

受け単語列を考慮(受け単語列)

ki を含む単語列を受ける単語列がある場合、ki はその単語列内の主辞単語から二単語ws とws+1 のみに依存すると仮定する(図 1 3 参照)。

【 数 7 】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_s, w_{s+1})$$

30

【 0 0 7 0 】

[K M 5]

係り単語列を最大二単語列考慮(係り二単語列)

ki を含む単語列に係る単語列がある場合、ki は、そのうち最も文末側の単語列の末尾から二単語wi 、wi-1 と、最も文頭側の単語列の末尾から二単語wh 、wh-1のみに依存すると仮定する(図 1 3 参照)。

40

【 数 8 】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_s, w_{s+1})$$

【 0 0 7 1 】

50

次に、形態素モデル (MM) について示す。形態素に付与すべき文法的属性が l 個あると仮定する。テキストつまり文字列が与えられたとき、その文字列が形態素であり、かつ $j(1 \leq j \leq l)$ 番目の文法的属性を持つとしたときの尤もらしさを確率値として求めるモデルを用いる。

テキスト T が与えられたとき、順序付き形態素集合 M が得られる確率は、各形態素 $m_i (1 \leq i \leq n)$ が独立であると仮定し、

【数 9】

$$P(M|T) = \prod_{i=1}^n P(m_i | m_{i-1}^{-1}, T)$$

10

と表す。ここで、 m_i は 1 から l までのいずれかの文法的属性を表わす。

【0072】

一方、係り受けモデル (DM) は、テキスト T と順序付き形態素集合 M が与えられたとき、各単語列に対する係り受けの順序付き集合 D が得られる確率は、各々の係り受け $d_1 \cdots d_n$ が独立であると仮定し、

【数 10】

$$P(M|T) = \prod_{i=1}^n P(m_i | m_{i-1}^{-1}, T)$$

20

と表わす。

【0073】

例えば、「palace troops fierce battle engage」の 5 つの特徴的単語訳語から「They are engaged in fierce battles with troops at the palace」と「I engaged the palace's troops in fierce battles.」の 2 つの候補が生成されたとする。係り受けモデルにより、このうち尤もらしい係り受け構造を持つ候補が優先される。

30

【0074】

以上に示すような各モデルを用い、本発明では評価部 (74) においてテキスト候補 (35' など) に評価付けを行う。

そして、評価値が最大あるいは閾値を超えるテキスト候補、あるいは評価値の上位 N 個を表層文に変換して出力する。

【0075】

出力部 (80) における出力方法としては、モニタによる表示の他、音声合成を用いた発声、翻訳システムなど他の言語処理システムへのデータ出力などが可能である。また、ネットワーク接続された他のコンピュータなどにテキストデータを送出してもよい。

40

【0076】

以上に示した実施例では、特徴的単語訳語の前後に語句を付加する構成を主としているが、本発明の実施においては特徴的単語訳語 (主辞単語に相当するもの) そのものを補完する構成をとることもできる。前述の通り、「engage」を補完するとき、語句を付加して「They are engaged」と補完する構成を示したが、新たに「They」を補完することができる。

【0077】

具体的には、図 7 の構成に図 14 の要素を追加する。すなわち、特徴的単語訳語 (34) を係り受け関係語抽出部 (79) にも入力し、該部 (79) ではコーパス (75) から該特徴的単語訳語 (34) と係り受け関係にある単語を抽出する。

50

そして、単語を新たな特徴的単語訳語として加え、もともと入力された特徴的単語訳語(34)と合わせて単語列候補生成部(72)における処理を行う。

【0078】

例えば、「They are engaged at the palace」そのものがコーパス(75)に無くとも、「they are engaged」と「they are at the palace」という係り受け関係がそれぞれコーパス(75)にあれば、それらに共通する単語「they」を新たに特徴的単語訳語として追加することによって、単語列候補生成部(72)によって「They are engaged at the palace」が生成できるようになる。

【0079】

本構成は、計量が少なく高速な特徴的単語訳語の追加が可能であるが、本発明では必ずしもコーパスから係り受け関係にある単語を抽出することに限らず、任意の特徴的単語訳語の候補を追加し、その中から評価部(74)における評価が結果的に最も高くなるものを出力してもよい。

これによって、特徴的単語訳語にテキストの意味を決定する重要な単語が、翻訳元言語の表現特性などにより欠落していたとしても、有意な翻訳先言語のテキストが出力できるようになる。

【0080】

本発明の機械翻訳方法において、さらに高精度な翻訳を実現する方法として、翻訳元言語テキスト(日本語テキスト(31))の単語の依存関係を、生成にも利用することが考えられる。以下、別実施例として説述する。

【0081】

具体的には図3の機械翻訳装置における特徴的単語抽出部(50)を、図8に示す依存関係を解析し、それをテーブルとして記録可能な特徴的単語抽出部(50')に置き換える。同様に、訳語選択部(60)は特徴的単語訳語の依存関係に変換可能な図9に示す訳語選択部(60')に、テキスト生成部(70)は単語訳語の依存関係を用いてテキスト生成が可能な図11又は図12のテキスト生成部(70')(70'')にそれぞれ置き換える。ここに述べた新たな機能以外の構成は上記実施例と同様であって、ここでは省略する。

【0082】

図8において、「官邸前などでドゥダエフ政権部隊と激しい市街戦を展開している。」という日本語テキスト(31)を入力する。依存関係解析部(57)では、公知の依存関係(係り受け関係)の解析方法を用いて、図15のように「官邸前などで」(80)「ドゥダエフ政権部隊と」(81)「激しい」(82)「市街戦を」(83)「展開している」(84)が依存関係(85)を有していることがわかる。ここで一般に依存関係解析部(57)における構文解析処理にはあらかじめ形態素解析を行うことが必要であり、形態素解析部(51)における解析結果を用いる。

【0083】

前述の処理により特徴語決定部(55)で決定された特徴的単語間の依存関係のみをとると、単語間依存関係抽出部(58)では、図16のように「官邸」(86)「部隊」(87)「激しい」(89)「戦」(90)の依存関係(91)が抽出される。

抽出された特徴的単語間の依存関係は単語間依存関係テーブル(36)として外部記憶装置やメモリなどに記録する。

【0084】

次に、図9の訳語選択部(60')では特徴的単語(33)と特徴的単語訳語(34)の関係から、日本語の単語間依存関係テーブル(36)を、英語の単語訳語間依存関係テーブル(37)に変換する。

具体的には、依存関係変換部(65)を設け、対訳辞書データベース(63)で一義的に決定できる単語、例えば部隊とtroopsのように対訳関係が明確な単語はそのまま依存関係テーブルの語を置き換える。一方、機械学習モデルを用いて訳語を選択したもの、例えば展開とengageなどのように多義性が生じるものは機械学習モデルによる類似性判定部(6

10

20

30

40

50

2) の判定結果を用いて単語訳語間依存関係テーブル(37)に記録する。

【0085】

これにより生成された単語訳語間依存関係テーブル(37)は図17のように、「palace」(92)「troops」(93)「fierce」(94)「battle」(95)「engage」(96)に置き換えられた上で、依存関係(97)が記録される。

ここで、固有名詞も特徴的単語とすると、図15でドゥダエフが抽出され、図18のように「Dudayev」(98)も追加することができる。

【0086】

図10のテキスト生成部(70')ではテキスト候補生成部(73)に抽出した単語訳語間依存関係テーブル(37)を入力して、上述の係り受け関係の仮定を行わず、係り受け

10

関係を決定する。
もちろん、単語訳語間依存関係テーブル(37)に依存関係の情報がない訳語間についてはテキスト候補生成部(73)において係り受け関係の仮定を行うこともできる。

【0087】

本構成によると、テキスト候補が高精度に生成できるため、評価部(74)では単語列候補生成部(72)・テキスト候補生成部(73)の生成結果をより厳密に評価することができる。

【0088】

一方、図11のテキスト生成部(70'')では評価部(74)に抽出した単語訳語間依存関係テーブル(37)を入力して、テキスト候補の評価に用いることができる。

20

ここで、評価部(74)は上述の言語モデル(78)により、各単語列に対する係り受けの順序に係る確率を用いるが、単語訳語間で単語訳語間依存関係テーブル(37)に記録された係り受け関係については、確率値を最大に設定し、当該テキスト候補の評価に用いる。

【0089】

本構成によると、仮に単語訳語間の依存関係の一部に解析誤差が生じた場合にも、テキスト候補生成部(73)ではそれにとらわれない候補の生成が可能のため、著しく不自然な係り受け関係は評価部(74)において相対的に低い評価とすることができ、本発明の特徴である自然な英語テキスト(32)の生成に寄与する。

【0090】

本発明は以上の構成により実現されるものであるが、特徴的単語の抽出処理、訳語の選択処理、いくつかのキーワードからテキストを生成する処理はいずれも公知のあるいは今後提供される言語処理技術を用いることができる。そして、各処理の高精度化に伴ってさらに翻訳精度の向上が期待されるものである。

30

【0091】

上記では説明の便宜のために、各部(40)(50)(60)(70)(80)を別個に説述したが、これらは一体的に例えば1台のパーソナルコンピュータによって提供することができる。特に、CPU、メモリ、入出力装置、ネットワークに接続するためのネットワークアダプタ(図示していない)、外部記憶装置などは共用することが望ましく、装置の簡略化に寄与することができる。

40

【0092】

外部記憶装置に記録される文書データベース(56)、対訳辞書・対訳コーパスデータベース(63)、コーパス(75)はいずれも同一のデータベースの一部又は全部を用いることが可能である。

また、これらは外部記憶装置上に記録される場合にとどまらず、ネットワーク上の複数のサーバーに記録されたものを収集するように構成してもよい。

【図面の簡単な説明】

【0093】

【図1】本発明による機械翻訳方法のフローチャートである。

【図2】本発明による機械翻訳方法(別実施例)のフローチャートである。

50

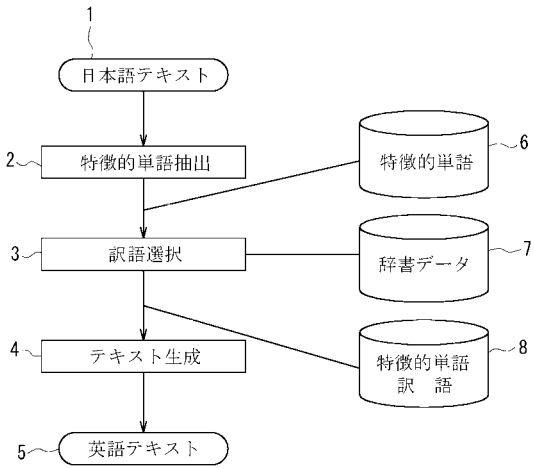
- 【図3】本発明による機械翻訳装置の全体構成図である。
- 【図4】本発明における入力部の構成図である。
- 【図5】本発明における特徴的単語抽出部の構成図である。
- 【図6】本発明における訳語選択部の構成図である。
- 【図7】本発明におけるテキスト生成部の構成図である。
- 【図8】本発明における特徴的単語抽出部（別実施例）の構成図である。
- 【図9】本発明における訳語選択部（別実施例）の構成図である。
- 【図10】本発明におけるテキスト生成部（別実施例）の構成図である。
- 【図11】本発明におけるテキスト生成部（別実施例）の構成図である。
- 【図12】特徴的単語訳語からのテキスト生成の例を示す説明図である。 10
- 【図13】特徴的単語訳語と単語列との関係を示す説明図である。
- 【図14】本発明に係る係り受け関係語抽出部の構成図である。
- 【図15】依存関係解析部における依存関係の解析結果を示す構造木である。
- 【図16】単語間依存関係抽出部における単語間依存関係テーブルの内容である。
- 【図17】依存関係変換部において変換された単語間依存関係テーブルの内容である。
- 【図18】同、固有名詞を特徴的単語とした時の単語間依存関係テーブルの内容である。

【符号の説明】

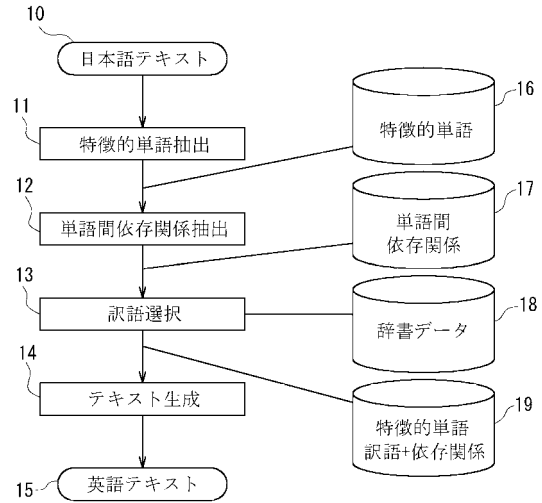
【0094】

- 30 機械翻訳装置
- 31 日本語テキスト 20
- 32 英語テキスト
- 33 特徴的単語
- 34 特徴的単語訳語
- 35 生成テキスト
- 40 入力受理部
- 50 特徴的単語抽出部
- 60 訳語選択部
- 70 テキスト生成部
- 80 出力部

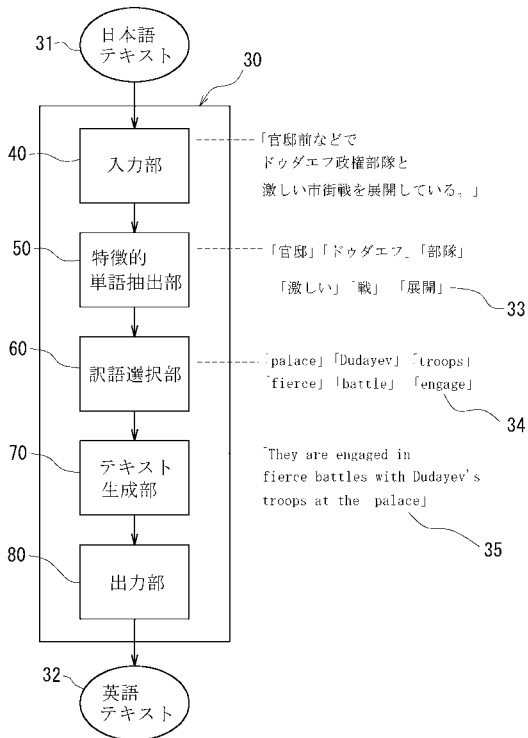
【図1】



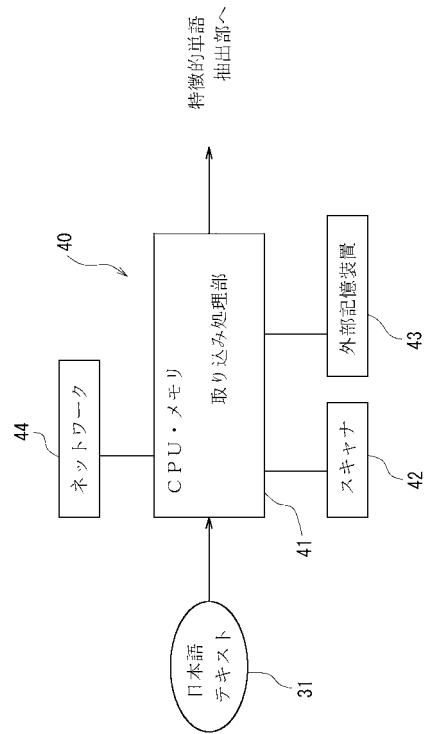
【図2】



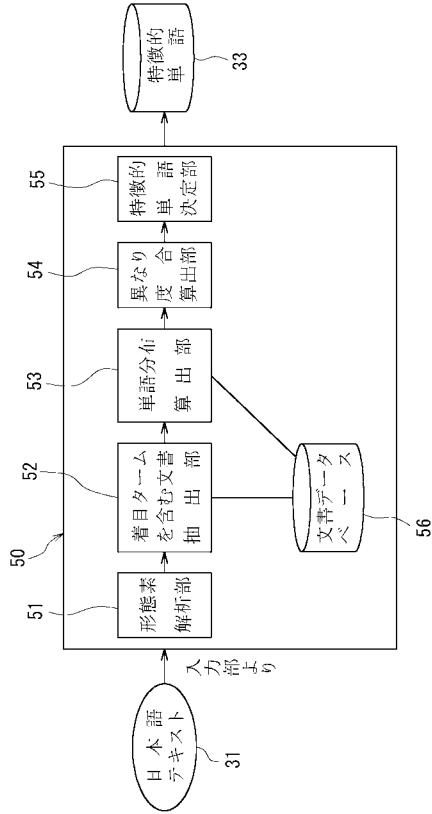
【図3】



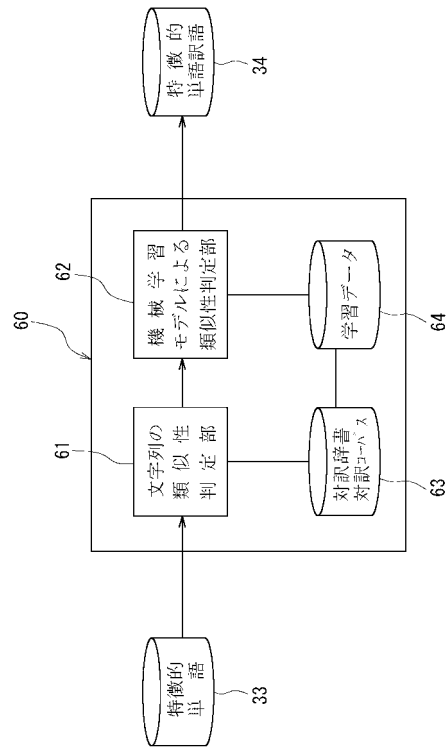
【図4】



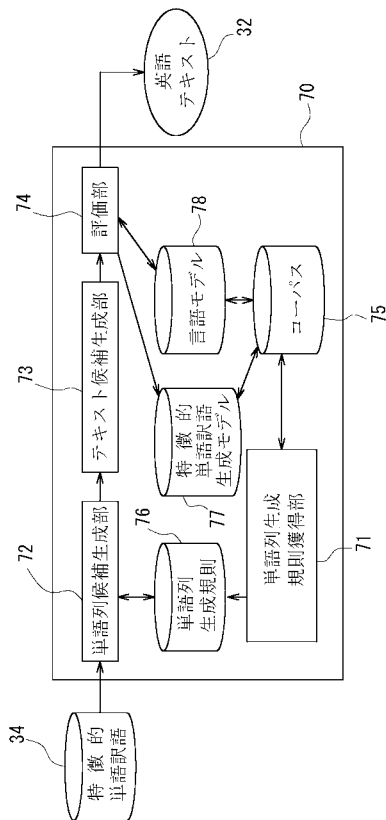
【 図 5 】



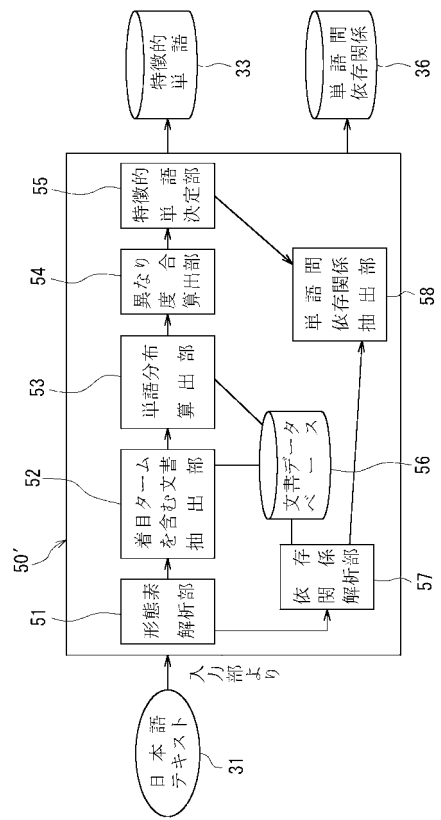
【 図 6 】



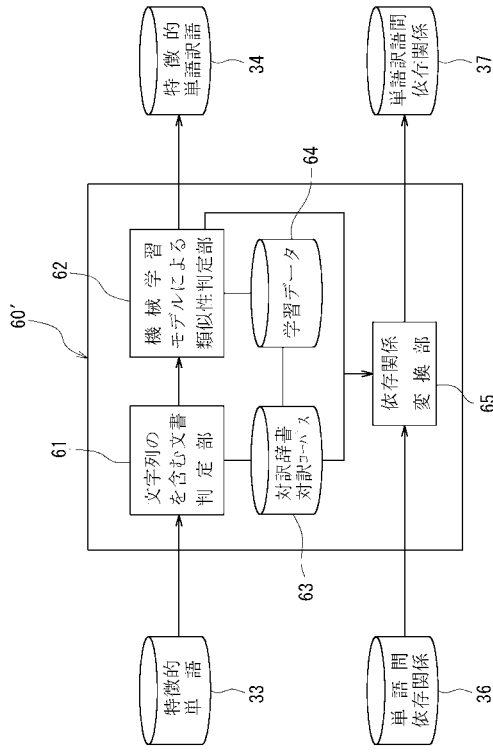
【 図 7 】



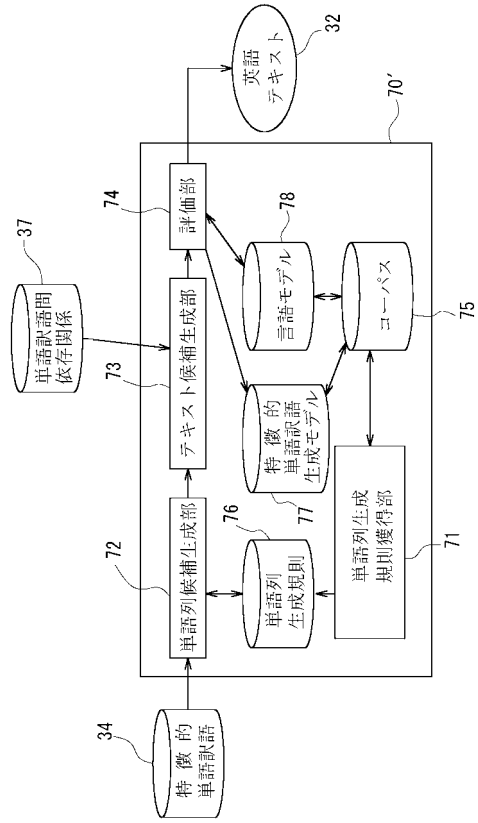
【 図 8 】



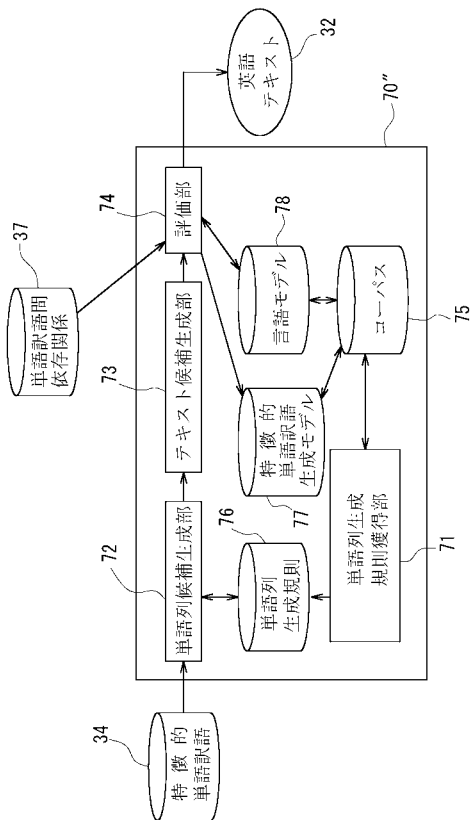
【 図 9 】



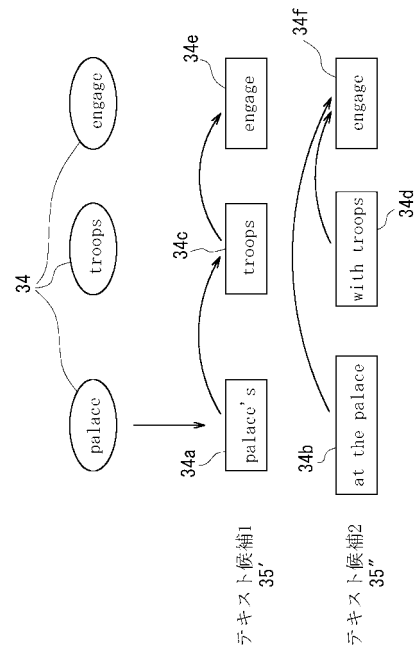
【 図 10 】



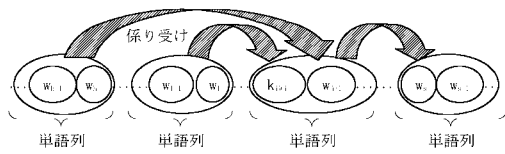
【 図 11 】



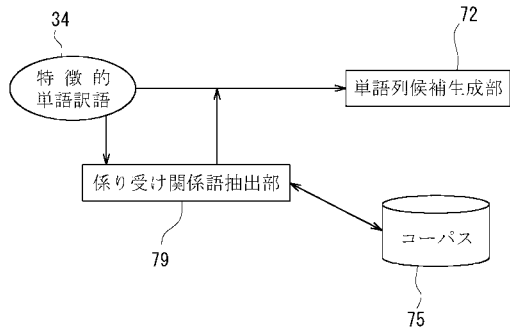
【 図 12 】



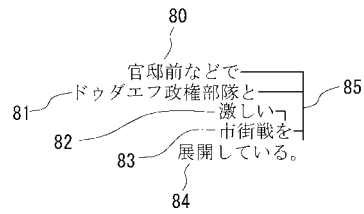
【図13】



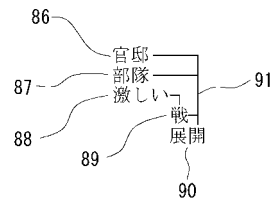
【図14】



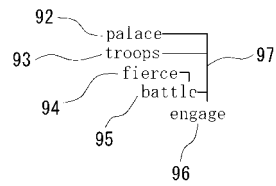
【図15】



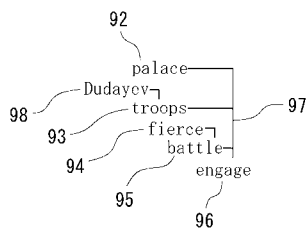
【図16】



【図17】



【図18】



フロントページの続き

(56)参考文献 特公平08-033895(JP, B2)

久光徹他, タームのrepresentativenessを測る, 情報処理学会研究報告99
-NL-133, 日本, 社団法人情報処理学会, 1999年 9月10日, Vol.99, No
.73, p.115 - p.122

(58)調査した分野(Int.Cl., DB名)

G06F 17/27 - 17/28