

(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2005-234800

(P2005-234800A)

(43) 公開日 平成17年9月2日(2005.9.2)

(51) Int.Cl.⁷

G06F 17/28

F I

G06F 17/28

Z

テーマコード (参考)

5B091

審査請求 未請求 請求項の数 16 O L (全 20 頁)

(21) 出願番号 特願2004-41751 (P2004-41751)

(22) 出願日 平成16年2月18日 (2004.2.18)

(出願人による申告) 平成15年度通信・放送機構、研究テーマ「大規模コーパス音声対話翻訳技術の研究開発」に関する委託研究、産業活力再生特別措置法第30条の適用を受ける特許出願

(71) 出願人 393031586
株式会社国際電気通信基礎技術研究所
京都府相楽郡精華町光台二丁目2番地2

(74) 代理人 100099933

弁理士 清水 敏

(72) 発明者 土居 蒼生
京都府相楽郡精華町光台二丁目2番地2
株式会社国際電気通信基礎技術研究所内

(72) 発明者 山本 博史
京都府相楽郡精華町光台二丁目2番地2
株式会社国際電気通信基礎技術研究所内

(72) 発明者 隅田 英一郎
京都府相楽郡精華町光台二丁目2番地2
株式会社国際電気通信基礎技術研究所内

Fターム(参考) 5B091 AA03 AB17 BA04 CA21 CC01
CC05 CC15 EA24

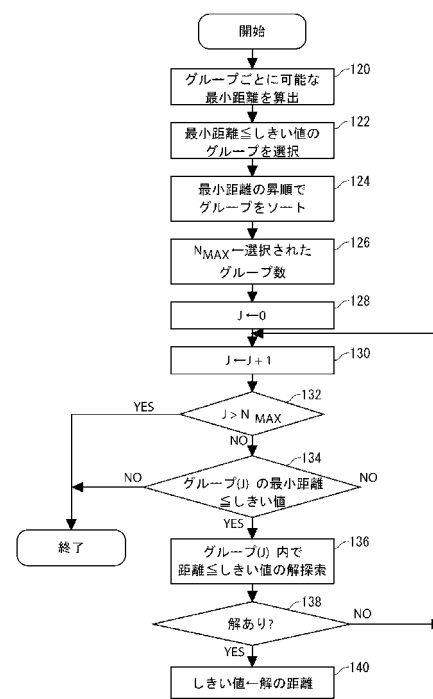
(54) 【発明の名称】 用例機械翻訳装置及び用例翻訳コンピュータプログラム、並びに用例検索装置及び用例検索コンピュータプログラム

(57) 【要約】

【課題】、高速で単語列編集距離を基準にして入力文と所定の関係を持つ用例を高速に検索可能にする。

【解決手段】 用例機械翻訳装置で用いられる用例検索装置は、各用例に含まれる第1の言語の単語列に含まれる内容語数及び機能語数に基づいて複数のグループに分割されたコーパスから、第1の言語の入力文に最も近い用例を検索するために、各グループに含まれる第1の言語の単語列と入力文との間で可能な最小の単語列編集距離を算出する最小距離算出部120と、距離の下限值がしきい値以下のものを選択するグループ選択部122と、選択されたグループに含まれる用例の中で、単語列編集距離により定義される入力文との距離が最小の解を探索する探索部124-140とを含む。探索にはA*アルゴリズムを使用しても良い。

【選択図】 図4



【特許請求の範囲】

【請求項 1】

コンピュータ読取可能な、第 1 及び第 2 の言語の対訳関係にある用例からなる用例コーパスを記憶するための第 1 の記憶手段と、

それぞれコンピュータ読取可能な、前記第 1 及び第 2 の言語の対訳辞書並びに前記第 1 及び第 2 の言語のシソーラスを記憶するための第 2 の記憶手段と、

前記第 1 の言語の入力文を受け、前記シソーラスを用いて前記用例コーパスから前記入力文と所定の関係を有する用例を検索するための用例検索手段と、

前記用例検索手段により検索された用例と、前記対訳辞書及び前記シソーラスとを用いて前記第 1 の言語の入力文を前記第 2 の言語の文に翻訳するための翻訳手段とを含み、

前記用例コーパスは、各用例に含まれる前記第 1 の言語の単語列に含まれる内容語数及び機能語数に基づいて複数のグループに分割され、

前記用例検索手段は、

前記複数のグループの各々について、予め定義された単語列編集距離を用い、当該グループに含まれる前記第 1 の言語の単語列と前記入力文との間で距離の下限値を算出するための最小距離算出手段と、

前記複数のグループのうち、前記最小距離算出手段により算出された距離の下限値が所定のしきい値以下のもののみを選択するためのグループ選択手段と、

前記グループ選択手段により選択されたグループに含まれる用例の中で、前記単語列編集距離により定義される入力文との距離が最小の解を探索するための探索手段とを含む、
用例機械翻訳装置。

【請求項 2】

前記用例検索手段はさらに、前記探索手段により解が見出されたことに応答して、前記探索手段による探索と並行して、前記所定のしきい値を前記見出された解の距離で置換するための手段を含む、請求項 1 に記載の用例機械翻訳装置。

【請求項 3】

前記複数のグループの各々に含まれる用例の、前記第 1 の言語の単語列は、一つの単語グラフ形式で表現され、

前記探索手段は、前記単語グラフの先頭ノードから最終ノードまでの可能な全経路について、当該経路に現れる単語列と入力単語列との照合を行なうことによって、前記入力単語列との間の単語列編集距離が最小となる経路を選択するための経路選択手段を含む、請求項 1 又は請求項 2 に記載の用例機械翻訳装置。

【請求項 4】

前記経路選択手段は、各グループに含まれる単語グラフの各経路に現れる単語列と、入力文との間の照合の途中経過を状態とする問題状態集合の中から、A*アルゴリズムを用いて単語列編集距離を最小にする目標状態を探索するための手段を含む、請求項 3 に記載の用例機械翻訳装置。

【請求項 5】

コンピュータ読取可能な、第 1 及び第 2 の言語の対訳関係にある用例からなる用例コーパスを記憶するための第 1 の記憶装置と、それぞれコンピュータ読取可能な、前記第 1 及び第 2 の言語の対訳辞書並びに前記第 1 及び第 2 の言語のシソーラスを記憶するための第 2 の記憶装置とを備えたコンピュータ上で実行されると、当該コンピュータを、前記用例コーパスを用いた用例機械翻訳装置として動作させる、用例翻訳コンピュータプログラムであって、

前記用例機械翻訳装置は、

前記第 1 の言語の入力文を受け、前記シソーラスを用いて前記用例コーパスから前記入力文と所定の関係を有する用例を検索するための用例検索手段と、

前記用例検索手段により検索された用例と、前記対訳辞書及び前記シソーラスとを用いて前記第 1 の言語の入力文を前記第 2 の言語の文に翻訳するための翻訳手段とを含み、

前記用例コーパスは、各用例に含まれる前記第 1 の言語の単語列に含まれる内容語数及

10

20

30

40

50

び機能語数に基づいて複数のグループに分割され、

前記用例検索手段は、

前記複数のグループの各々について、予め定義された単語列編集距離を用い、当該グループに含まれる前記第 1 の言語の単語列と前記入力文との間で距離の下限値を算出するための最小距離算出手段と、

前記複数のグループのうち、前記最小距離算出手段により算出された距離の下限値が所定のしきい値以下のもののみを選択するためのグループ選択手段と、

前記グループ選択手段により選択されたグループに含まれる用例の中で、前記単語列編集距離により定義される入力文との距離が最小の解を探索するための探索手段とを含む、用例翻訳コンピュータプログラム。

10

【請求項 6】

前記用例検索手段はさらに、前記探索手段により解が見出されたことに応答して、前記探索手段による探索と並行して、前記所定のしきい値を前記見出された解の距離で置換するための手段を含む、請求項 5 に記載の用例翻訳コンピュータプログラム。

【請求項 7】

前記複数のグループの各々に含まれる用例の、前記第 1 の言語の単語列は、一つの単語グラフ形式で表現され、

前記探索手段は、前記単語グラフの先頭ノードから最終ノードまでの可能な全経路について、当該経路に現れる単語列と入力単語列との照合を行なうことによって、前記入力単語列との間の単語列編集距離が最小となる経路を選択するための経路選択手段を含む、請求項 5 又は請求項 6 に記載の用例翻訳コンピュータプログラム。

20

【請求項 8】

前記経路選択手段は、各グループに含まれる単語グラフの各経路に現れる単語列と、入力文との間の照合の途中経過を状態とする問題状態集合の中から、A*アルゴリズムを用いて単語列編集距離を最小にする目標状態を探索するための手段を含む、請求項 7 に記載の用例翻訳コンピュータプログラム。

【請求項 9】

コンピュータ読取可能な、第 1 及び第 2 の言語の対訳関係にある用例からなる用例コーパスを記憶するための第 1 の記憶手段と、それぞれコンピュータ読取可能な、前記第 1 及び第 2 の言語の対訳辞書並びに前記第 1 及び第 2 の言語のシソーラスを記憶するための第 2 の記憶手段と、前記用例コーパスから検索された用例と、前記対訳辞書及び前記シソーラスとを用いて前記第 1 の言語の入力文を前記第 2 の言語の文に翻訳するための翻訳手段とを含む用例機械翻訳装置で使用され、前記第 1 の言語の入力文を受け、前記シソーラスを用いて前記用例コーパスから前記入力文と所定の関係を有する用例を検索するための用例検索装置であって、

30

前記用例コーパスは、各用例に含まれる前記第 1 の言語の単語列に含まれる内容語数及び機能語数に基づいて複数のグループに分割され、

前記用例検索装置は、

前記複数のグループの各々について、予め定義された単語列編集距離を用い、当該グループに含まれる前記第 1 の言語の単語列と前記入力文との間で距離の下限値を算出するための最小距離算出手段と、

40

前記複数のグループのうち、前記最小距離算出手段により算出された距離の下限値が所定のしきい値以下のもののみを選択するためのグループ選択手段と、

前記グループ選択手段により選択されたグループに含まれる用例の中で、前記単語列編集距離により定義される入力文との距離が最小の解を探索するための探索手段とを含む、用例検索装置。

【請求項 10】

さらに、前記探索手段により解が見出されたことに応答して、前記探索手段による探索と並行して、前記所定のしきい値を前記見出された解の距離で置換するための手段を含む、請求項 9 に記載の用例検索装置。

50

【請求項 1 1】

前記複数のグループの各々に含まれる用例の、前記第 1 の言語の単語列は、一つの単語グラフ形式で表現され、

前記探索手段は、前記単語グラフの先頭ノードから最終ノードまでの可能な全経路について、当該経路に現れる単語列と入力単語列との照合を行なうことによって、前記入力単語列との間の単語列編集距離が最小となる経路を選択するための経路選択手段を含む、請求項 9 又は請求項 1 0 に記載の用例検索装置。

【請求項 1 2】

前記経路選択手段は、各グループに含まれる単語グラフの各経路に現れる単語列と、入力文との間の照合の途中経過を状態とする問題状態集合の中から、A * アルゴリズムを用いて単語列編集距離を最小にする目標状態を探索するための手段を含む、請求項 1 1 に記載の用例検索装置。

【請求項 1 3】

コンピュータ読取可能な、第 1 及び第 2 の言語の対訳関係にある用例からなる用例コーパスを記憶するための第 1 の記憶装置と、それぞれコンピュータ読取可能な、前記第 1 及び第 2 の言語の対訳辞書並びに前記第 1 及び第 2 の言語のシソーラスを記憶するための第 2 の記憶装置と、前記対訳コーパスから検索された用例を用いて、前記第 1 の言語の入力文を前記第 2 の言語の文に翻訳する用例機械翻訳手段とを備えたコンピュータ上で実行されると、当該コンピュータを、前記第 1 の言語の入力文を受け、前記シソーラスを用いて前記用例コーパスから前記入力文と所定の関係を有する用例を検索する用例検索装置として動作させる、用例検索コンピュータプログラムであって、

前記用例コーパスは、各用例に含まれる前記第 1 の言語の単語列に含まれる内容語数及び機能語数に基づいて複数のグループに分割され、

前記用例検索装置は、

前記複数のグループの各々について、予め定義された単語列編集距離を用い、当該グループに含まれる前記第 1 の言語の単語列と前記入力文との間で距離の下限値を算出するための最小距離算出手段と、

前記複数のグループのうち、前記最小距離算出手段により算出された距離の下限値が所定のしきい値以下のもののみを選択するためのグループ選択手段と、

前記グループ選択手段により選択されたグループに含まれる用例の中で、前記単語列編集距離により定義される入力文との距離が最小の解を探索するための探索手段とを含む、用例検索コンピュータプログラム。

【請求項 1 4】

前記用例検索装置はさらに、前記探索手段により解が見出されたことに応答して、前記探索手段による探索と並行して、前記所定のしきい値を前記見出された解の距離で置換するための手段を含む、請求項 1 3 に記載の用例検索コンピュータプログラム。

【請求項 1 5】

前記複数のグループの各々に含まれる用例の、前記第 1 の言語の単語列は、一つの単語グラフ形式で表現され、

前記探索手段は、前記単語グラフの先頭ノードから最終ノードまでの可能な全経路について、当該経路に現れる単語列と入力単語列との照合を行なうことによって、前記入力単語列との間の単語列編集距離が最小となる経路を選択するための経路選択手段を含む、請求項 1 3 又は請求項 1 4 に記載の用例検索コンピュータプログラム。

【請求項 1 6】

前記経路選択手段は、各グループに含まれる単語グラフの各経路に現れる単語列と、入力文との間の照合の途中経過を状態とする問題状態集合の中から、A * アルゴリズムを用いて単語列編集距離を最小にする目標状態を探索するための手段を含む、請求項 1 5 に記載の用例検索コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は機械翻訳装置に関し、特に、用例を用いて機械翻訳を行なう用例機械翻訳装置並びに当該装置において用例を高速に検索するための用例検索装置に関する。

【背景技術】

【0002】

アナロジーに基づく機械翻訳の概念が提唱されて以来、このアイデアを具体化した用例翻訳が数多く提案されてきた。用例翻訳では、あらかじめ二つの言語で同じ意味を表す対訳からなる対訳コーパスを準備する。そしてこの対訳コーパスから対訳表現を抽出し、翻訳すべき文を構成する各部分に適合する対訳表現を見つけ、それらを組み合わせて翻訳文を生成する。

10

【0003】

一方、用例翻訳以前の代表的な翻訳方式としてルールベース翻訳がある。ルールベース翻訳では、言語現象について人間が内省することにより作成した翻訳のためのルールに基づいて翻訳を実行する。

【0004】

ルールベース翻訳では、翻訳ルールの精密な構築が翻訳システムの能力を決める。これに対し、用例翻訳では対訳コーパスからの学習によって能力が決定する。用例翻訳の大きな利点は、人手による翻訳ルールの記述が不要であり、翻訳システムの保守を含めた開発の効率が高いことである。この特長により、異なるドメインへの移植や新たな言語対への適用が容易になる。

20

【0005】

特許文献1において、用例翻訳の一つとしてDP(Dynamic Programming)マッチに基づいた翻訳方式が提案されている。

【0006】

この方式では、用例は原言語文と目的言語文との対である。各文は単語列として表現される。翻訳実行時には入力文との類似度の最も高いと判定される原言語文を持つ用例を検索する。検索された用例の原言語文と入力文との差異、及び差異部分と対応する目的言語文の部分を求め、翻訳パターンを生成する。

【0007】

30

他の用例翻訳方式と比べたときのこの方式の特徴は、以下の通りである。

【0008】

(1)多くの方式ではあらかじめ翻訳パターンを作成しているが、この方式では翻訳実行時に翻訳パターンを作成する。

【0009】

(2)多くの方式では構文解析又はツリーバンクの利用を仮定しているが、この方式ではそれらを利用しない。

【0010】

このように特許文献1において提案された翻訳方式では、用例を抽象化せず単語列の形のまま保持し、検索する。さらに訳文を生成する際にも用例の目的言語文の変更を最低限にとどめる。そのため、入力文に近い用例が存在すれば自然な表現の翻訳結果が得られるという特徴がある。また、多くの言語に高精度パーザを期待することができない現状を考えると、この方式は解析知識を使わないため用例翻訳方式の中でも特に多言語に適用することが容易であり、その点で他方式と比較して優れるものと考えられる。

40

【0011】

【特許文献1】特開2003-6193号公報

【特許文献2】特開平8-185482号公報

【非特許文献1】ラップ R.、「翻訳メモリのための品詞に基づく検索アルゴリズム」、LREC2002予稿集、pp.466-472、2002年(Rapp, R.: A Part-of-Speech-Based Search Algorithm for

50

or Translation Memories, Proc. of LREC 2002, pp. 466 - 472, 2002)

【発明の開示】

【発明が解決しようとする課題】

【0012】

特許文献1に記載の用例翻訳方式では、記憶した文の集合の中から最も入力文に類似したものを選び出す仕組みがその中核にある。特に大量の用例を利用する場合、その仕組みの効率的な実装が必須となる。

【0013】

用例検索処理の効率的な実装は、翻訳メモリの課題と共通である。翻訳メモリに関する非特許文献1に記載の研究では、品詞レベルでの文の完全一致に基づいた検索が提案されている。しかし、この用例翻訳方式では単語列の間の編集距離を基準にして用例を検索している。そのために、用例の検索では単語の挿入及び削除をも考慮しなければならない。そのため非特許文献1において提案された方式は用例翻訳方式には採用できない。

【0014】

特許文献2では、編集距離最小文を検索する方法を提示しているが、個々の候補に対する計算を繰返すので大量の用例を使った場合の高速化には限界があると考えられる。

【0015】

それゆえにこの発明の目的は、高速で単語列編集距離を基準にして用例を検索することができる用例検索装置及びそれを用いた用例機械翻訳装置、ならびにそれらのコンピュータプログラムを提供することである。

【課題を解決するための手段】

【0016】

本発明の第1の局面に係る用例機械翻訳装置は、コンピュータ読取可能な、第1及び第2の言語の対訳関係にある用例からなる用例コーパスを記憶するための第1の記憶手段と、それぞれコンピュータ読取可能な、第1及び第2の言語の対訳辞書並びに第1及び第2の言語のシソーラスを記憶するための第2の記憶手段と、第1の言語の入力文を受け、シソーラスを用いて用例コーパスから入力文と所定の関係を有する用例を検索するための用例検索手段と、用例検索手段により検索された用例と、対訳辞書及びシソーラスとを用いて第1の言語の入力文を第2の言語の文に翻訳するための翻訳手段とを含み、用例コーパスは、各用例に含まれる第1の言語の単語列に含まれる内容語数及び機能語数に基づいて複数のグループに分割され、用例検索手段は、複数のグループの各々について、予め定義された単語列編集距離を用い、当該グループに含まれる第1の言語の単語列と入力文との間で距離の下限値を算出するための最小距離算出手段と、複数のグループのうち、最小距離算出手段により算出された距離の下限値が所定のしきい値以下のもののみを選択するためのグループ選択手段と、グループ選択手段により選択されたグループに含まれる用例の中で、単語列編集距離により定義される入力文との距離が最小の解を探索するための探索手段とを含む。

【0017】

好ましくは、用例検索手段はさらに、探索手段により解が見出されたことに応答して、探索手段による探索と並行して、所定のしきい値を見出された解の距離で置換するための手段を含む。

【0018】

さらに好ましくは、複数のグループの各々に含まれる用例の、第1の言語の単語列は、一つの単語グラフ形式で表現されており、探索手段は、単語グラフの先頭ノードから最終ノードまでの可能な全経路について、当該経路に現れる単語列と入力単語列との照合を行なうことによって、入力単語列との間の単語列編集距離が最小となる経路を選択するための経路選択手段を含む。

【0019】

経路選択手段は、各グループに含まれる単語グラフの各経路に現れる単語列と、入力文

10

20

30

40

50

との間の照合の途中経過を状態とする問題状態集合の中から、A * アルゴリズムを用いて単語列編集距離を最小にする目標状態を探索するための手段を含んでもよい。

【0020】

本発明の第2の局面に係る用例翻訳コンピュータプログラムは、コンピュータ読取可能な、第1及び第2の言語の対訳関係にある用例からなる用例コーパスを記憶するための第1の記憶装置と、それぞれコンピュータ読取可能な、第1及び第2の言語の対訳辞書並びに第1及び第2の言語のシソーラスを記憶するための第2の記憶装置とを備えたコンピュータ上で実行されると、当該コンピュータを、用例コーパスを用いた用例機械翻訳装置として動作させる、用例翻訳コンピュータプログラムである。このプログラムにより実現される用例機械翻訳装置は、第1の言語の入力文を受け、シソーラスを用いて用例コーパスから入力文と所定の関係を有する用例を検索するための用例検索手段と、用例検索手段により検索された用例と、対訳辞書及びシソーラスとを用いて第1の言語の入力文を第2の言語の文に翻訳するための翻訳手段とを含み、用例コーパスは、各用例に含まれる第1の言語の単語列に含まれる内容語数及び機能語数に基づいて複数のグループに分割され、用例検索手段は、複数のグループの各々について、予め定義された単語列編集距離を用い、当該グループに含まれる第1の言語の単語列と入力文との間で距離の下限値を算出するための最小距離算出手段と、複数のグループのうち、最小距離算出手段により算出された距離の下限値が所定のしきい値以下のもののみを選択するためのグループ選択手段と、グループ選択手段により選択されたグループに含まれる用例の中で、単語列編集距離により定義される入力文との距離が最小の解を探索するための探索手段とを含む。

10

20

【0021】

好ましくは、用例検索手段はさらに、探索手段により解が見出されたことに応答して、探索手段による探索と並行して、所定のしきい値を見出された解の距離で置換するための手段を含む。

【0022】

さらに好ましくは、複数のグループの各々に含まれる用例の、第1の言語の単語列は、一つの単語グラフ形式で表現され、探索手段は、単語グラフの先頭ノードから最終ノードまでの可能な全経路について、当該経路に現れる単語列と入力単語列との照合を行なうことによって、入力単語列との間の単語列編集距離が最小となる経路を選択するための経路選択手段を含む。

30

【0023】

経路選択手段は、各グループに含まれる単語グラフの各経路に現れる単語列と、入力文との間の照合の途中経過を状態とする問題状態集合の中から、A * アルゴリズムを用いて単語列編集距離を最小にする目標状態を探索するための手段を含んでもよい。

【0024】

本発明の第3の局面に係る用例検索装置は、コンピュータ読取可能な、第1及び第2の言語の対訳関係にある用例からなる用例コーパスを記憶するための第1の記憶手段と、それぞれコンピュータ読取可能な、第1及び第2の言語の対訳辞書並びに第1及び第2の言語のシソーラスを記憶するための第2の記憶手段と、用例コーパスから検索された用例と、対訳辞書及びシソーラスとを用いて第1の言語の入力文を第2の言語の文に翻訳するための翻訳手段とを含む用例機械翻訳装置で使用され、第1の言語の入力文を受け、シソーラスを用いて用例コーパスから入力文と所定の関係を有する用例を検索するための用例検索装置である。用例コーパスは、各用例に含まれる第1の言語の単語列に含まれる内容語数及び機能語数に基づいて複数のグループに分割されている。用例検索装置は、複数のグループの各々について、予め定義された単語列編集距離を用い、当該グループに含まれる第1の言語の単語列と入力文との間で距離の下限値を算出するための最小距離算出手段と、複数のグループのうち、最小距離算出手段により算出された距離の下限値が所定のしきい値以下のもののみを選択するためのグループ選択手段と、グループ選択手段により選択されたグループに含まれる用例の中で、単語列編集距離により定義される入力文との距離が最小の解を探索するための探索手段とを含む。

40

50

【 0 0 2 5 】

好ましくは、用例検索装置はさらに、探索手段により解が見出されたことに応答して、探索手段による探索と並行して、所定のしきい値を見出された解の距離で置換するための手段を含む。

【 0 0 2 6 】

さらに好ましくは、複数のグループの各々に含まれる用例の、第 1 の言語の単語列は、一つの単語グラフ形式で表現され、探索手段は、単語グラフの先頭ノードから最終ノードまでの可能な全経路について、当該経路に現れる単語列と入力単語列との照合を行なうことによって、入力単語列との間の単語列編集距離が最小となる経路を選択するための経路選択手段を含む。

10

【 0 0 2 7 】

経路選択手段は、各グループに含まれる単語グラフの各経路に現れる単語列と、入力文との間の照合の途中経過を状態とする問題状態集合の中から、A * アルゴリズムを用いて単語列編集距離を最小にする目標状態を探索するための手段を含んでもよい。

【 0 0 2 8 】

本発明の第 4 の局面に係る用例検索コンピュータプログラムは、コンピュータ読取可能な、第 1 及び第 2 の言語の対訳関係にある用例からなる用例コーパスを記憶するための第 1 の記憶装置と、それぞれコンピュータ読取可能な、第 1 及び第 2 の言語の対訳辞書並びに第 1 及び第 2 の言語のシソーラスを記憶するための第 2 の記憶装置と、対訳コーパスから検索された用例を用いて、第 1 の言語の入力文を第 2 の言語の文に翻訳する用例機械翻訳手段とを備えたコンピュータ上で実行されると、当該コンピュータを、第 1 の言語の入力文を受け、シソーラスを用いて用例コーパスから入力文と所定の関係を有する用例を検索する用例検索装置として動作させる、用例検索コンピュータプログラムである。用例コーパスは、各用例に含まれる第 1 の言語の単語列に含まれる内容語数及び機能語数に基づいて複数のグループに分割されている。当該用例検索装置は、複数のグループの各々について、予め定義された単語列編集距離を用い、当該グループに含まれる第 1 の言語の単語列と入力文との間で距離の下限値を算出するための最小距離算出手段と、複数のグループのうち、最小距離算出手段により算出された距離の下限値が所定のしきい値以下のもののみを選択するためのグループ選択手段と、グループ選択手段により選択されたグループに含まれる用例の中で、単語列編集距離により定義される入力文との距離が最小の解を探索するための探索手段とを含む。

20

30

【 0 0 2 9 】

好ましくは、用例検索装置はさらに、探索手段により解が見出されたことに応答して、探索手段による探索と並行して、所定のしきい値を見出された解の距離で置換するための手段を含む。

【 0 0 3 0 】

さらに好ましくは、複数のグループの各々に含まれる用例の、第 1 の言語の単語列は、一つの単語グラフ形式で表現され、探索手段は、単語グラフの先頭ノードから最終ノードまでの可能な全経路について、当該経路に現れる単語列と入力単語列との照合を行なうことによって、入力単語列との間の単語列編集距離が最小となる経路を選択するための経路選択手段を含む。

40

【 0 0 3 1 】

経路選択手段は、各グループに含まれる単語グラフの各経路に現れる単語列と、入力文との間の照合の途中経過を状態とする問題状態集合の中から、A * アルゴリズムを用いて単語列編集距離を最小にする目標状態を探索するための手段を含んでもよい。

【 発明を実施するための最良の形態 】

【 0 0 3 2 】

- 構成 -

図 1 に、本発明の一実施の形態に係る用例に基づく用例機械翻訳装置 30 のブロック図を示す。図 1 を参照して、この用例機械翻訳装置 30 は、言語資源としてコンピュータ読

50

取可能な対訳コーパス 42、第 1 の言語のシソーラス 44、第 2 の言語のシソーラス 45 及び対訳辞書 50 を含む。これらはいずれもハードディスク等の記憶装置に格納される。

【0033】

対訳コーパス 42 は、翻訳方向における原言語と目的言語との文の対の集合である。両言語の文は互いに対訳関係にある。対訳コーパス 42 中の文は単語に分割されそれぞれ品詞情報が付与されている。

【0034】

この用例機械翻訳装置 30 では、特許文献 1 に記載の方式と同様に、対訳コーパス 42 中の対訳関係にある文のペアを利用して翻訳を実行する。以下、この文のペアを「用例」と呼ぶ。

【0035】

対訳辞書 50 は、後述するように翻訳パターンの抽出と訳語の置換処理において使用される。

【0036】

シソーラスとしては、第 1 の言語（原言語）のシソーラス 44 と第 2 の言語（目的言語）のシソーラス 45 とが用意されている。シソーラスは、単語を単語間の意味の近さに基づいてツリー状の階層関係に配置したものである。第 1 の言語のシソーラス 44 は、用例検索及び翻訳パターン抽出処理において使用される。第 2 の言語の言語のシソーラス 45 は、翻訳パターンの抽出に用いられる。

【0037】

用例機械翻訳装置 30 はさらに、原言語の入力文 40 を受け、対訳コーパス 42 及びシソーラス 44 を用いて、対訳コーパス 42 の中で入力文 40 と最も類似した原言語の文を持つ用例を検索するための用例検索部 46 と、用例検索部 46 により検索された用例から、第 1 のシソーラス 44、第 2 のシソーラス 45 及び対訳辞書 50 を参照して翻訳パターンを抽出するための翻訳パターン抽出部 48 と、用例検索部 46 により検索された用例が複数ある場合、翻訳パターン抽出部 48 から出力される複数の翻訳パターンの中で所定のものを選択するための翻訳パターン選択部 52 と、翻訳パターン選択部 52 により選択された翻訳パターンの中で、変数に束縛された単語の訳語を対訳辞書 50 から引き、その訳語でもって目的言語パターンの変数を具体化する処理を行なって出力文 56 を生成するための訳語置換部 54 とを含む。

【0038】

以下簡単に用例検索部 46、翻訳パターン抽出部 48、翻訳パターン選択部 52 及び訳語置換部 54 について説明する。用例検索部 46 の詳細については、後にさらに詳述する。

【0039】

用例検索部 46 は、全用例の原言語文を走査する。入力文と用例原言語文の単語列間の距離を測り、最小距離の用例を選び出す。ただしこの最小距離が大きければ、検索された用例は翻訳処理に有用ではない。そのため距離にしきい値を設ける。しきい値以内の距離の用例が存在しなければ用例検索及び翻訳処理は失敗に終わる。

【0040】

単語列間の距離には意味距離の加味された単語列編集距離が使われる。この単語列編集距離 $dist$ は次の式で表される。

【0041】

【数 1】

$$dist = \frac{I + D + 2 \sum SEMDIST}{L_{input} + L_{example}}$$

ここで L_{input} は入力文の単語数、 $L_{example}$ は用例原言語文の単語数、 I は用例原言語文を入力文に変換するために必要な挿入単語数、 D は同じく必要な削除単語数、 $SEMDIST$

10

20

30

40

50

I S T は同じく必要な置換により置換される語の間の意味距離を示す。

【 0 0 4 2 】

この式に従って、挿入語と削除語の数、及び置換語の意味距離が足し合わされ、入力文と用例原言語文の長さの和でもって正規化して単語列編集距離 $d i s t$ が算出される。2 単語が同品詞の内容語である場合のみ置換の対象となる。この場合には 2 単語間の意味距離が単語列編集距離の計算に使われる。意味距離計算においては、二つの単語に関して第 1 の言語のシソーラス 4 4 中の概念階層における意味概念間の位置関係によって意味距離を計算する。意味距離は 0 ~ 1 までを値域とし、0 に近いほど 2 単語が意味的に類似していることを示す。

【 0 0 4 3 】

以下、日英翻訳における用例検索の例を示す。(1 j) は入力文、(2 j) は用例の原言語文とする。このうち入力文 (1 j) の「色」と用例の原言語文 (2 j) の「デザイン」の部分が両文の差分となる。

【 0 0 4 4 】

(1 j) 色 / が / 気 / に / 入り / ません

(2 j) デザイン / が / 気 / に / 入り / ません

ここで、「色」と「デザイン」とがシソーラス 4 4 上で完全に異なった語であるものとすると、単語間の意味距離は 1 となる。従ってこの 2 文間の単語列編集距離は $(0 + 0 + 2 * 1) / (6 + 6) = 0 . 1 6 7$ となる。

【 0 0 4 5 】

翻訳パターン抽出部 4 8 は、用例検索部 4 6 により検索された全ての用例に対し、原言語文中の、入力文と異なる箇所を変数で置換し、用例目的言語文中の対応する箇所に同じ変数を当てはめた翻訳パターンを生成する。両言語の文の間で対応をとる際は、変数となる単語のみ対象とし、全ての単語の対応をとる必要はない。つまり、変数部分以外の箇所は全体として対応していると仮定する。このため用例のほとんどの部分は変更されず、訳文の組み合わせ時に発生する誤りや不自然さの回避が期待される。

【 0 0 4 6 】

この原言語と目的言語の単語間の対応をとるには、様々な単語アライメント手法を適用できる。本実施の形態では、対訳辞書 5 0、第 1 の言語のシソーラス 4 4 及び第 2 の言語のシソーラス 4 5 に基づいて単語間の対応関係を判断している。

【 0 0 4 7 】

先の例の (2 j) に対応する目的言語文を (2 e) とする。

【 0 0 4 8 】

(2 j) デザイン / が / 気 / に / 入り / ません

(2 e) I d o n o t l i k e t h e d e s i g n .

このフェーズでは目的言語文 (2 e) 中で「デザイン」に対応する箇所が探し出され、「デザイン」と「design」との間の対応が取られる。この結果、以下に示すような原言語パターン (2 j p) と目的言語パターン (2 e p) とからなる翻訳パターンが作られる。入力文によるその変数束縛は (1 j b) となる。

【 0 0 4 9 】

(2 j p) X / が / 気 / に / 入り / ません

(2 e p) I d o n o t l i k e t h e X .

(1 j b) X = 「色」

翻訳パターン選択部 5 2 は、用例検索部 4 6 により複数の用例が検索され、その結果翻訳パターン抽出部 4 8 によって複数の翻訳パターンが抽出された際に、その中の一つを選択するための処理を行なうものである。翻訳パターン選択部 5 2 は、複数の翻訳パターンから一つを選択するために、(1) より多くの用例検索結果から同じ翻訳パターンが抽出された方を選ぶ、(2) 翻訳パターン中に現れる単語のコーパスでの出現頻度の合計が大きい方を選ぶ、というヒューリスティクスを使用して翻訳パターンを選択する。これらで一つの翻訳パターンを決定できない場合には、翻訳パターン選択部 5 2 は任意の一つの翻

10

20

30

40

50

訳パターンを選ぶ。

【0050】

訳語置換部54は、翻訳パターン選択部52により選択された翻訳パターンの変数に束縛された単語の訳を対訳辞書50から引き、その訳語でもって目的言語パターンの変数を具体化する。先の例に基づいて説明すると、目的言語側の変数束縛は以下の(1 e b)となり、訳文(1 e)が得られる。

【0051】

(1 e b) X = "color"

(1 e) I do not like the color.

- 用例検索部46の詳細 -

10

図1に示す用例機械翻訳装置30の各処理の中で翻訳実行時間の大きな割合を占めるのは、用例検索部46による用例検索である。用例の選択基準には、前述した単語列編集距離が使われる。

【0052】

用例検索処理は、用例の原言語文を候補文とし、入力文との単語列編集距離がしきい値以内で最小の候補文をすべて求める。単語列編集距離は2文間の関係で定義され、二つの単語列のDPマッチにより計算可能である。従って各候補文と入力文間のDPマッチを逐次的に繰り返すことで単語列編集距離が最小の候補文を求めることができる。

【0053】

しかし単純にこの方法を使おうとすれば、用例数に比例した処理時間がかかる。そのため大規模コーパスを利用したリアルタイムの翻訳処理をそうした方法で実現することは通常のコンピュータでは困難である。そこで本実施の形態では、以下に述べるような実装により用例検索を効率的に短時間で行なえるようにした。

20

【0054】

[候補文集合の分割]

図2を参照して、対訳コーパス42に含まれる候補文を、その内容語数と機能語数とを元に複数(M個)のグループ70-1、70-2、... 70-Mにグループ分けする。このようにグループ分けすることにより、入力文の内容語数と機能語数及び距離しきい値を用いて検索対象の候補文数を絞ることができる。具体的には以下のように用例検索を行なう。

30

【0055】

まず、機能語同士、内容語同士はすべて一致すると仮定した場合の、グループごとに可能な最小距離を求める。最小距離が距離しきい値の範囲内で小さいグループから順に、単語列編集距離が最初に定めたしきい値以内で最小の候補文を検索する。あるグループ中に解が見つかれば、その解の距離を新たなしきい値として検索対象のグループをさらに絞ることができる。グループ内での用例検索については後述するが、グループ内では全ての候補文の内容語数と機能語数が等しい、つまりは単語数も等しいという事実が用例検索の前提条件として利用されている。

【0056】

[単語グラフ]

内容語数と機能語数とを基準に分けられたグループ70-1, 70-2, ..., 70-Mの各々に対し、複数の候補文を一つの単語グラフ72-1, 72-2, ..., 72-Mにまとめる。すなわち一つのグループごとに一つの単語グラフが作成される。図3に単語グラフの例を示す。

40

【0057】

図3に示されるように、単語グラフ80は有向グラフであり、先頭ノード90から最終ノード108に至る可能な道筋がそれぞれ一つの候補文に対応する。図3に示す単語グラフ80は、先頭ノード90及び最終ノード108を含めて、全部で10個のノード90、92、94、96、98、100、102、104、106及び108を含む。ノードの間を結ぶリンクが単語に対応する。

50

【 0 0 5 8 】

例えば「全部売り切れました」という文はノード 9 0、9 2、9 6、1 0 2 及び 1 0 8 という道筋に対応する。ノード列だけで見れば「全部届きました」という単語列も同じ道筋に対応する。ただし先の文では「売り切れる」というリンクによりノード 9 2 及び 9 6 が接続されていたが、この文では「届く」というリンクでこれらノードが接続されている点異なる。このようにして、複数の文で共通な単語列をグラフ中では一つにまとめていることにより、1 グループに含まれる全ての候補文を一つの単語グラフ 8 0 で表すことができる。

【 0 0 5 9 】

単語グラフを利用することにより、グループ内の全候補文を同時並行的に調べながら、10
入力文との距離が最小の候補文を検索することができる。

【 0 0 6 0 】

[A * アルゴリズム]

二つの単語列を照合した結果を示す単語の一致、置換、挿入、削除の列を「照合列」と呼ぶこととする。グループ内において単語列編集距離を最小とする候補文又は解を検索することは、単語グラフの先頭ノードから最終ノードまでの可能な全経路について、各経路に現れる単語列と入力単語列との照合列の中から単語列編集距離を最小にするものを探索することである。本実施の形態では、この探索問題の解放に A * アルゴリズムを用いている。

【 0 0 6 1 】

一般に A * アルゴリズムでは、問題状態集合の中から最小コストの下限の推定値が最小のものが選ばれ、継続状態に展開される。ここで対象とする問題では、状態は、単語グラフの経路と入力文との照合の途中経過を意味する。また「継続状態に展開」とは、選ばれた状態から遷移可能な全ての状態を生成し考慮の対象とすることを意味する。20

【 0 0 6 2 】

[探索]

ここでは、ある単語グラフを用いた、単語列編集距離最小の経路を探索する処理について説明する。図 3 に示すように、ある単語グラフはノードとリンクとを含む。リンクは単語をラベルとして持ち、一つの始点ノードと一つの終点ノードとを結ぶ。例えば図 3 に示す「売り切れる」というリンクはノード 9 2 と 9 6 とを結んでいる。単語グラフ全体で一つの先頭ノードと一つの最終ノードとを持つ。図 3 ではこれらはノード 9 0 及び 1 0 8 に相当する。30

【 0 0 6 3 】

対象となる問題状態空間は以下で説明する状態、作用素、初期状態及び目標状態により構成される。

【 0 0 6 4 】

(1) 状態

状態は `paths` , `node` , `input` , `trans` という属性を持つ。各属性の内容は以下の通りである。

- ・ `paths` : その時点までの照合列のリスト 40
- ・ `node` : 単語グラフのノード。このノードまで照合が進んだことを示す。
- ・ `input` : 入力単語列のうち、まだ照合に使われていない部分
- ・ `trans` : 適用可能な作用素。

【 0 0 6 5 】

`paths` 内の各照合列の一致、置換、挿入、削除をそれぞれ (E : 単語)、(S : グラフ側単語、入力側単語)、(I : 入力側単語)、(D : グラフ側単語) の形式で表し、それぞれ E レコード、S レコード、I レコード及び D レコードと呼ぶ。

【 0 0 6 6 】

状態のコストは `paths` 内の任意の一つの照合列のコストである。 `paths` 内のどの照合列も等しいコストを持つ。照合列のコストは、それに含まれるレコードのコストの 50

和である。Eレコードのコストは0、Iレコードのコストは1、Dレコードのコストは1と定義する。Sレコードのコストは、置換される2単語間の意味距離を2倍した値であるが、意味距離が0の場合には0でなく小さな正の値を与える。これは、類似語関係と同一語関係とを区別するためである。この値がSレコードの最小コストとなる。

【0067】

状態に作用素を適用することにより継続状態が生成される。一般に一つの状態に複数の作用素を適用することが可能である。従って一つの状態からいくつかの継続状態が生成される。

【0068】

(2) 作用素

5種類の作用素、T作用素、E作用素、S作用素、I作用素及びD作用素を以下のように定義する。T作用素とI作用素は状態に適用されるが、E、S、Dの各作用素は状態、及びその状態のnodeを始点とするリンクの組に適用される。T作用素は実際に照合を進める作用素ではなく、trans属性とともにE、S、I、Dの各作用素の適用順序を制御する役目を持つ。なお以下の説明では、作用素が適用される状態をs、リンクをl、生成される継続状態をs'と表し、各作用素について適用条件とどのような継続状態が生成されるかを示す。

【0069】

・T作用素：

- 条件：s.transがE作用素又はS作用素である。

【0070】

- 生成：s'.trans = s.transがE作用素ならばS作用素とNILとから選択（説明は後述）、s.transがS作用素ならばNIL。s'の他の属性値はsと同じ。

【0071】

・E作用素：

- 条件：s.transがE作用素である。

【0072】

かつ、s.inputが空リストでない。

【0073】

かつ、lのラベルと、s.inputの先頭とが同一語である。

【0074】

- 生成：s'.paths = s.pathsの各要素にEレコードを追加した値
s'.node = lの終点

s'.input = s.inputから先頭を消去した値

s'.trans = E作用素とS作用素とNILとから選択（説明は後述）

・S作用素：

- 条件：s.transがS作用素である。

【0075】

かつ、s.inputが空リストでない。

【0076】

かつ、s.inputの先頭と、lのラベルとが同品詞の内容語であり、かつ同一語ではない。

【0077】

かつ、これら2単語の意味距離は1未満である。

【0078】

- 生成：s'.paths = s.pathsの各要素にSレコードを追加した値
s'.node = lの終点

s'.input = s.inputから先頭を消去した値

s'.trans = E作用素とS作用素とNILとから選択

10

20

30

40

50

・ I 作用素 :

- 条件 : $s.trans$ が NIL である。

【 0 0 7 9 】

かつ、 $s.input$ が空リストでない。

【 0 0 8 0 】

- 生成 : $s'.paths = s.paths$ の各要素に I レコードを追加した値

$s'.node = s.node$

$s'.input = s.input$ から先頭を消去した値

$s'.trans = E$ 作用素と S 作用素と NIL とから選択

・ D 作用素 :

- 条件 : $s.trans$ が NIL である。

【 0 0 8 1 】

かつ、 $s.paths$ に最新レコードが I レコードでない要素がある。

【 0 0 8 2 】

- 生成 : $s'.paths = s.paths$ から最新レコードが I レコードである要素を除き、残った要素に D レコードを追加した値

$s'.node = l$ の終点

$s'.input = s.input$

$s'.trans = E$ 作用素と S 作用素と NIL とから選択

上の記載において、「 S 作用素と NIL とから選択」とは、 s' に S 作用素を適用できる可能性があれば $s'.trans$ の値を S 作用素とし、可能性がなければ NIL とすることを意味する。本実施の形態では、 $s'.input$ の先頭が内容語であり、その語と同一の語を除く同品詞語をラベルとし $s'.node$ を始点とするリンクが存在する場合に S 作用素を適用できる可能性があると判断する。

【 0 0 8 3 】

また、「 E 作用素と S 作用素と NIL とから選択」とは、 $s'.input$ の先頭語をラベルとし、 $s'.node$ を始点とするリンクが存在すれば $s'.trans$ の値を E 作用素とし、そうでなければ S 作用素と NIL とから選択する。

【 0 0 8 4 】

D 作用素の 2 番目の条件は I レコードの後に D レコードが来るのを防いでいる。つまり、I レコードと D レコードとが連続する場合、D レコードが先に来るようにし、実質的に同じ削除と挿入が入れ替わっただけの異なる状態が現れる冗長性を排除する。

【 0 0 8 5 】

(3) 初期状態と目標状態

初期状態では、 $paths$ は空リストを要素とするリスト、 $node$ は先頭ノード、 $input$ は入力単語列全体、 $trans$ は E 作用素である。目標状態は、 $node$ が最終ノード、かつ $input$ が空リストであるような状態である。

【 0 0 8 6 】

[探索アルゴリズム]

上記の初期状態、作用素及び目標状態で表現される状態空間からコスト最小の目標状態を探索する。初期条件としてコスト上限値が与えられる。コスト上限値は入力文長と候補文長の和を距離しきい値に乗じた値である。

【 0 0 8 7 】

[評価関数]

状態空間探索時に使用する評価関数 f^* を次のように定義する。

【 0 0 8 8 】

$$f^*(s) = g(s) + h^*(s)$$

$g(s)$ は初期状態から状態 s に達するまでにかかったコストを示す。つまり $g(s)$ は先に定義した状態のコストであり、 $s.paths$ から計算できる。目標状態では $f^*(s) = g(s)$ となる。 $h^*(s)$ は状態 s から目標状態までにかかるコストの下限で

10

20

30

40

50

ある。

【0089】

一つの単語グラフを構成する全候補文の内容語数及び機能語数はそれぞれ同一である。従って状態 s において入力文側とグラフ側の未処理の内容語数及び機能語数がそれぞれ一意に決まる。それぞれの個数を C_{input} 、 C_{graph} 、 F_{input} 、 F_{graph} として、残り語数に基づく最小コスト $h'(s)$ を次のように計算する。

【0090】

$$h'(s) = |C_{input} - C_{graph}| + |F_{input} - F_{graph}|$$

さらに、 T 作用素の適用が先行する場合を含めて状態 s に最初に適用可能な E 、 S 、 I 及び D の各作用素について、それが適用されたと仮定したときの目標状態までにかかるコストの下限を次の値とする。

【0091】

・ E 作用素： $h'(s)$

・ S 作用素： $h'(s)$ に S レコードの最小コストを加えた値。

【0092】

・ I 作用素： $s.input$ の先頭が内容語の場合は、 $|C_{input} - 1| - C_{graph}| + |F_{input} - F_{graph}|$ に 1 を加えた値、機能語の場合は $|C_{input} - C_{graph}| + |(F_{input} - 1) - F_{graph}|$ に 1 を加えた値。

【0093】

・ D 作用素： $|C_{input} - (C_{graph} - 1)| + |F_{input} - F_{graph}|$ と $|C_{input} - C_{graph}| + |F_{input} - (F_{graph} - 1)|$ を求め、その小さいほうの値に 1 を加えた値。ただし、 $s.node$ を始点とするリンクのラベルが内容語のみであるか機能語のみであれば、対応する一方の値に 1 を加えた値。

【0094】

これらを使って $h^*(s)$ を以下のように計算する。

【0095】

(1) $s.trans$ が E 作用素のときは、 E 作用素が適用されたときのコストの下限

(2) $s.trans$ が S 作用素のときは、 S 作用素、 I 作用素又は D 作用素が適用されたときのコストの下限の最小値

(3) $s.trans$ が NIL のときは、 I 作用素又は D 作用素が適用されたときのコストの下限の最小値。

【0096】

[探索処理]

探索処理は以下のように行なわれる。以下に示す探索処理は、コンピュータ上で実行されるプログラムにより上記した用例検索を行なう際のプログラムの制御の流れを示すものである。なお以下の説明中において、 $OPEN$ は未展開状態を、 $CLOSED$ は展開済状態を、それぞれ保持するためのリストを示す。また「同じ状態」とは、 $paths$ を除く属性値が等しい状態を意味する。

【0097】

図 4 は、この処理全体のフローチャートである。まずステップ 120 で、グループごとに可能な最小距離を算出する。ステップ 122 で、可能な最小距離がしきい値以下のグループを選択する。ステップ 124 で、選択されたグループを最小距離の昇順でソートする。ステップ 126 で、選択されたグループの数をグループ数の値 N_{MAX} に代入する。

【0098】

ステップ 128 以下がグループごとの繰り返し処理である。まずステップ 128 で繰り返し制御変数 J に 0 を代入する。ステップ 130 で変数 J に 1 を加算する。ステップ 132 で J がグループ数 N_{MAX} を超えたか否かを判定する。超えていれば処理を終了する。超えていなければステップ 134 でグループ (J) の可能な最小距離がしきい値以下か否かを判定する。しきい値以下であればステップ 136 に進み、それ以外の場合には処理を終了する。

【0099】

ステップ136ではグループ(j)内で距離しきい値の解を探索する処理を行なう。この処理の詳細については図5を参照して説明する。

【0100】

ステップ138で解が存在したか否かについての判定が行なわれる。解が存在していれば制御はステップ140に進み、それ以外の場合にはステップ130に戻る。

【0101】

ステップ140では、求められた解の距離を新たなしきい値に代入し、ステップ130に戻る。

【0102】

図5に、ステップ136で行なわれる処理の詳細について示す。ステップ160で、コスト上限 C_{MAX} に、所与の値(入力文長と候補文長の和を距離しきい値に乘じた値)を代入する。

【0103】

ステップ162でOPENに初期状態のみを入れる。

【0104】

ステップ164でOPEN内にコストが C_{MAX} 以下の状態があるか否かについての判定を行なう。条件を満たす状態がなければ処理を終了する。条件を満たす状態があればステップ166に進む。

【0105】

ステップ166ではOPENから評価関数 f^* を最小にする状態 s を取り除き、CLOSEDに入れる。

【0106】

ステップ168で、状態 s が目標状態か否かを判定する。目標状態であればステップ174に進み、それ以外の場合にはステップ170に進む。

【0107】

ステップ174では状態 s を解の一つとし、続くステップ176でコスト上限 C_{MAX} を状態 s のコストで置換し、ステップ164に戻る。

【0108】

一方ステップ170では、状態 s の全ての継続状態を生成する。そしてステップ172で、各継続状態 s' について図6に示す処理を実行する。

【0109】

図6を参照して、ステップ190で $f^*(s')$ がコスト上限 C_{MAX} 以下か否かを判定する。条件が充足されていなければ処理を終了する。条件が充足されていればステップ192で、OPEN及びCLOSED中の同じ状態と比較し、条件により以下の処理を行なう。

【0110】

(a) 同じ状態がなければ、状態 s' をOPENに追加する(ステップ194)

(b) 状態 s' よりコストの大きい同じ状態がOPEN又はCLOSEDに既存であれば、この既存状態を消去し(ステップ196)、状態 s' をOPENに追加する(ステップ198)。

【0111】

(c) 状態 s' とコストの等しい同じ状態がCLOSEDに既存であれば、この既存状態を消去し(ステップ200)、状態 s' をOPENに追加する(ステップ202)。

【0112】

(d) コストの等しい同じ状態がOPENに既存であれば、この既存状態の $paths$ に $s'.paths$ をマージする(ステップ204)。

【0113】

以上の処理が終了したら図6の処理を終了し、図5のステップ164に戻る。

【0114】

10

20

30

40

50

[単語グラフの特徴の利用]

単語グラフの形状の特徴として、開始ノードを始点とするリンク数が他のノードを始点とするリンク数よりも圧倒的に大きくなる傾向がある。そのため `node` 属性に開始ノードを持つ状態に `D` 作用素が適用されると、多くの継続状態が作られることとなり計算時間が大きくなる。これは、照合列の先頭要素が `D` レコードとなる場合である。この展開数の増大を避けるため、単語グラフ中、先頭ノードから数段階の仮のリンクとノードを加える。先頭ノードを持つ状態から `D` 作用素によって第 1 の仮のノードを持つ状態へ遷移する。第 1 の仮のノードは、全候補文について 2 番目の語をラベルとするリンクの始点となり、通常の単語グラフのノードに合流する。第 1 の仮のノードにある状態は `E` 作用素又は `S` 作用素の適用により通常のノードの状態、`D` 作用素によって第 2 の仮のノードを持つ状態に遷移する。

10

【 0 1 1 5 】

何段階まで仮のノードを用意するかは、用例検索時に使われる可能性のある距離しきい値の最大値から計算できる。候補文の長さを L とすると、照合列の先頭に `D` レコードが d 個並ぶという条件で、候補文とその距離を最小にする入力文は、候補文から先頭 d 語を除いた文である。そのときの距離は $d / ((L - d) + L)$ である。この値が距離しきい値の最大値を越える場合は探索する必要はない。距離しきい値の最大値を $d / ((L - d) + L)$ から $d \leq L / (1 + \quad)$ が導かれる。この式を満たす d の最大の整数値が用意すべき仮のノードの段数である。

【 0 1 1 6 】

20

- 動作 -

以下、本実施の形態に係る用例機械翻訳装置 30 の動作について説明する。図 1 において、翻訳パターン抽出部 48、翻訳パターン選択部 52 及び訳語置換部 54 の動作は特許文献 1 において提案されている用例機械翻訳装置の動作と同様である。従ってここではそれらについての詳しい説明は繰り返さず、用例検索部 46 による用例検索の詳細について説明する。

【 0 1 1 7 】

用例検索部 46 の動作を説明するために、用例検索の実行例を示す。ここでは、図 3 に示す単語グラフ 80 から入力文「全部揃いました」の類似文を検索することにする。以下の説明では状態を $[\text{paths}, \text{node}, \text{input}, \text{trans}, f^* \text{関数値}]$ の形式で記述する。`node` 値には図 3 中でノードにつけた参照番号 (90 - 108) を用いる。を `S` レコードの最小コストとする。また「揃う」と「売り切れる」の意味距離を 1.0、「揃う」と「届く」との意味距離を 0.7 であると仮定する。

30

【 0 1 1 8 】

初期状態 s_0 は次のようになる。

【 0 1 1 9 】

$s_0 = [(()) , \text{ノード 90} , (\text{全部} , \text{揃う} , \text{ます} , \text{た}) , \text{E 作用素} , 0]$

状態 s_0 に適用可能な作用素は `E` 作用素と `T` 作用素とである。これらの作用素を適用して継続状態 s_1 と s_2 とが得られる。`OPEN` は $\{ s_1 , s_2 \}$ となる。

【 0 1 2 0 】

40

$s_1 = [(((\text{E} , \text{全部}))) , \text{ノード 92} , (\text{揃う} , \text{ます} , \text{た}) , \text{S 作用素} , \quad]$

$s_2 = [(()) , \text{ノード 90} , (\text{全部} , \text{揃う} , \text{ます} , \text{た}) , \text{NIL} , 2]$

ここで `OPEN` の中から、 f^* 関数値の小さな状態 s_1 が選ばれ展開される。状態 s_1 に適用可能な作用素は `S` 作用素と `T` 作用素とである。「売り切れる」と「届く」とをラベルとする二つのリンクに関して `S` 作用素の適用条件がテストされる。ここでは、1 未満の意味距離の得られる「届く」のリンクについてのみテストが成功する。`S` 作用素の適用により継続状態 s_3 が得られる。また `T` 作用素の適用により継続状態 s_4 が生成される。`OPEN` は $\{ s_2 , s_3 , s_4 \}$ となる。

【 0 1 2 1 】

$s_3 = [(((\text{E} , \text{全部}) , (\text{S} , \text{届く} , \text{揃う}))) , \text{ノード 96} , (\text{ます} , \text{た}) , \text{E}$

50

作用素， 1 . 4]

s 4 = [(((E , 全部))) , ノード 9 2 , (揃う , ます , た) , N I L , 2]

ここで O P E N の中から、f * 関数値が最小の状態 s 3 が選ばれ展開される。状態 s 3 に E 作用素を 2 回適用した状態 s 7 が解となる。

【 0 1 2 2 】

s 7 = [(((E , 全部) , (S , 届く , 揃う) , (E , ます) , (E , た))) , ノード 1 0 8 , () , N I L , 1 . 4]

この例の状態遷移の様子を図 7 に示す。図 7 に示されるように、T 作用素の適用による遷移を除き、解である状態 s 7 に向かって一直線に探索が進んでいる。

【 0 1 2 3 】

以上のようにこの実施の形態に係る用例機械翻訳装置 3 0 によれば、大規模なコーパスから短時間で入力文との編集距離が最小の用例を探索することができる。探索された用例から翻訳パターンを抽出し、翻訳パターン中の変数部分を訳語置換することにより、入力文 4 0 (図 1 参照) に対応する翻訳出力文 5 6 を得ることができる。

【 0 1 2 4 】

特に、本実施の形態は、(1) シソーラスを使わず、単語列編集距離の定義において意味距離を定数とする、(2) 内容語・機能語の分類を行なわない(全ての語が一方に属するとする。)、又は他の分類を使う(文字種で分ける場合など。)、並びに(3) 1 文字ずつを 1 単語として扱う、という場合に編集距離に基づいたストリング検索装置を提供し、特許文献 2 が挙げている各応用課題において、大量の候補の中から解を効率よく検索する手段となる。

【 0 1 2 5 】

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味および範囲内でのすべての変更を含む。

【 図面の簡単な説明 】

【 0 1 2 6 】

【 図 1 】 本発明の一実施の形態に係る用例機械翻訳装置 3 0 のブロック図である。

【 図 2 】 本実施の形態における対訳コーパス 4 2 のグループへの分割と単語グラフの作成とを模式的に示す図である。

【 図 3 】 単語グラフの例を示す図である。

【 図 4 】 本実施の形態に係る用例機械翻訳装置 3 0 の用例検索部 4 6 で実行される処理のプログラムフローチャートである。

【 図 5 】 用例グループ内で解を探索する処理のプログラムフローチャートである。

【 図 6 】 探索処理において継続状態ごとに行なわれる処理のプログラムフローチャートである。

【 図 7 】 本発明の一実施の形態に係る用例機械翻訳装置 3 0 の用例検索部 4 6 の動作例を説明するための図である。

【 符号の説明 】

【 0 1 2 7 】

4 0 入力文、 4 2 対訳コーパス、 4 4 第 1 の言語のシソーラス、 4 5 第 2 の言語のシソーラス、 4 6 用例検索部、 4 8 翻訳パターン抽出部、 5 0 対訳辞書、 5 2 翻訳パターン選択部、 5 4 訳語置換部、 5 6 出力文、 8 0 単語グラフ

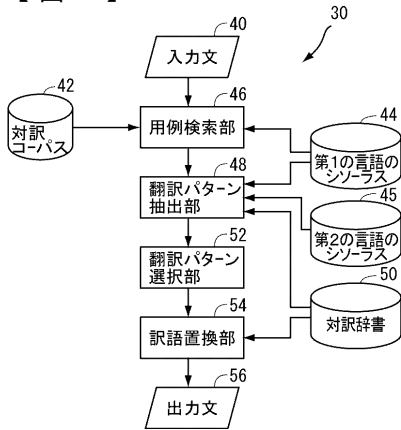
10

20

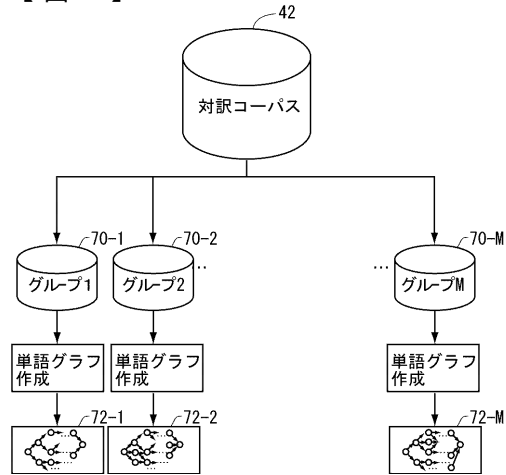
30

40

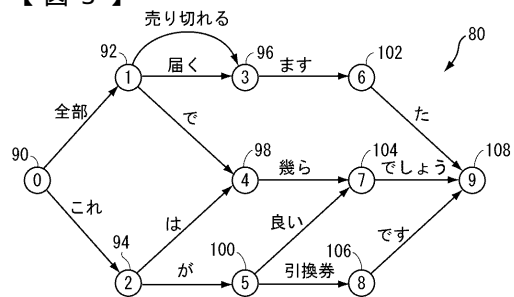
【図1】



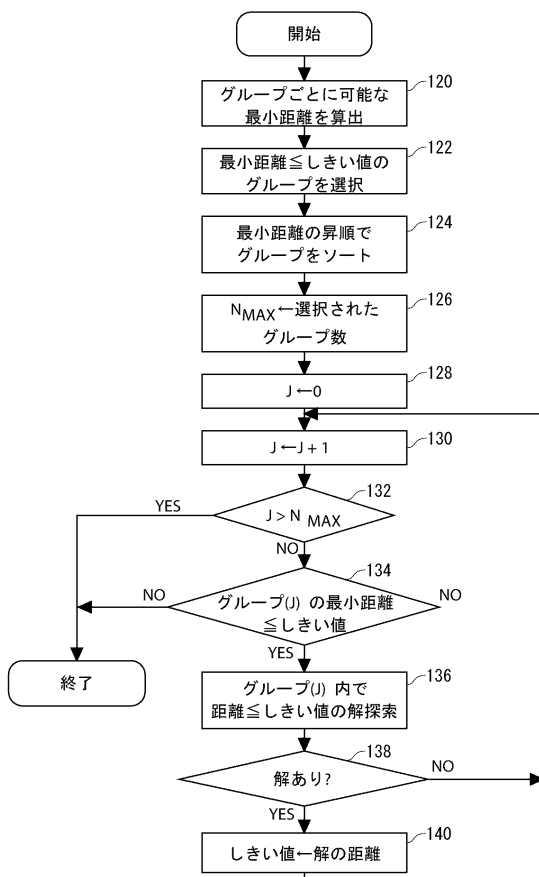
【図2】



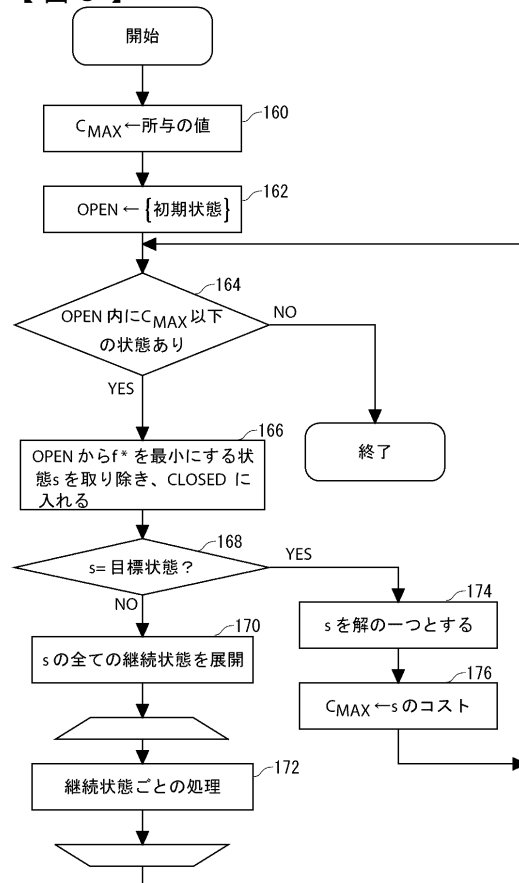
【図3】



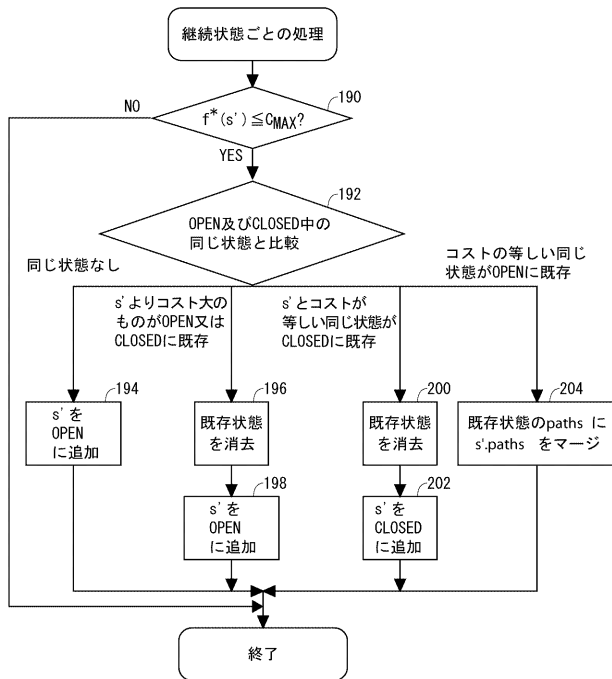
【図4】



【図5】



【 図 6 】



【 図 7 】

