

(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2004-362249

(P2004-362249A)

(43) 公開日 平成16年12月24日(2004.12.24)

(51) Int.Cl.⁷

G06F 17/28

F I

G06F 17/28

U

テーマコード (参考)

5B091

審査請求 有 請求項の数 17 O L (全 29 頁)

(21) 出願番号 特願2003-159662 (P2003-159662)

(22) 出願日 平成15年6月4日(2003.6.4)

(出願人による申告) 国等の委託研究の成果に係る特許出願(平成15年度通信・放送機構、研究テーマ「大規模コーパススペー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条の適用を受けるもの)

(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 今村 賢治

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 BA14 CA21 CC03 CC05 CD11
EA01

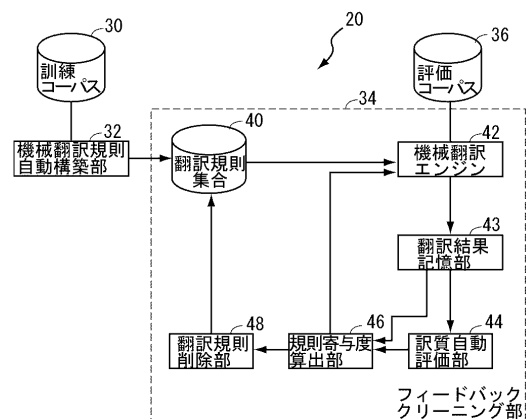
(54) 【発明の名称】 翻訳知識最適化装置、翻訳知識最適化のためのコンピュータプログラム、コンピュータ及び記憶媒体

(57) 【要約】

【課題】 対訳コーパスから自動獲得された翻訳規則をクリーニングしてより翻訳品質を向上させる事ができる翻訳知識最適化装置を提供する。

【解決手段】 翻訳知識最適化装置34は、翻訳知識を記憶する翻訳規則集合記憶部40と、評価コーパス36と、翻訳規則集合記憶部40に記憶された翻訳知識を利用して、評価コーパス36中の原言語の文を目的言語に翻訳する機械翻訳エンジン42と、機械翻訳エンジン42による翻訳結果の品質を、評価コーパス36を参照して自動的に評価する訳質自動評価部44と、訳質自動評価部44の出力する評価値が極値をとる様に、翻訳規則集合記憶部40内の翻訳知識の集合の最適化を行なう規則寄与度算出部46及び翻訳規則削除部48を含む。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

機械翻訳のための翻訳知識を最適化するための翻訳知識最適化装置であって、
翻訳知識の集合を記憶するための翻訳知識記憶手段と、
原言語と目的言語との対訳文を複数個含む、機械読取可能な対訳コーパスを記憶するための手段と、
前記翻訳知識記憶手段に記憶された前記翻訳知識の集合を利用して、前記対訳コーパス中の前記原言語の文を前記目的言語に機械翻訳するための機械翻訳手段と、
前記機械翻訳手段による翻訳結果の品質を、前記対訳コーパスを参照して自動的に評価して評価値を出力するための訳質自動評価手段と、
前記訳質自動評価手段の出力する評価値が極値をとる様に、前記翻訳知識の集合の最適化を行なうための最適化手段とを含む、翻訳知識最適化装置。

10

【請求項 2】

前記翻訳知識は、前記原言語の構文パターンから前記目的言語の構文パターンへの構文変換規則を含む、請求項 1 に記載の翻訳知識最適化装置。

【請求項 3】

前記最適化手段は、
前記翻訳知識の集合に含まれる翻訳知識の各々について、その規則寄与度を算出するための手段と、
前記規則寄与度が予め定める条件を満足する翻訳知識を、前記翻訳知識の集合から削除するための手段とを含む、請求項 1 に記載の翻訳知識最適化装置。

20

【請求項 4】

前記規則寄与度を算出するための手段は、
前記翻訳知識の集合の全体を用いて、前記機械翻訳手段による翻訳、及び当該翻訳の結果の前記訳質自動評価手段による訳質評価を行ない、初期評価値を得るための手段と、
前記翻訳知識の集合中の翻訳知識ごとに、前記翻訳知識の集合から当該翻訳知識を削除して得られる部分集合を用いて、前記機械翻訳による翻訳、及びその翻訳結果の前記訳質自動評価手段による訳質評価を行ない、削除後評価値を得るための手段と、
前記削除後評価値と前記初期評価値との差分を、前記ある翻訳知識の前記規則寄与度として算出するための手段とを含む、請求項 3 に記載の翻訳知識最適化装置。

30

【請求項 5】

前記最適化手段は、
前記翻訳知識の集合の全体を用いて、前記機械翻訳手段による翻訳、及び当該翻訳の結果の前記訳質自動評価手段による訳質評価を行ない、初期評価値を得るための手段と、
予め定められた方法に従って、前記翻訳知識の集合から複数の部分集合を作成するための手段と、
前記複数の部分集合の各々を用いて前記機械翻訳手段による翻訳、及びその翻訳結果の前記訳質自動評価手段による訳質評価を行ない、その評価値が前記初期評価値に対し所定の条件を満足するか否かを判定するための判定手段と、
前記判定するための手段により前記評価値が前記所定の条件を満足すると判定された部分集合の各々について、その補集合に属する翻訳知識を前記翻訳知識の集合から削除するための手段とを含む、請求項 1 に記載の翻訳知識最適化装置。

40

【請求項 6】

前記部分集合を作成するための手段は、前記翻訳知識の集合から予め定められる数の翻訳知識を除いて得られる部分集合を複数個作成するための手段を含む、請求項 5 に記載の翻訳知識最適化装置。

【請求項 7】

前記部分集合を複数個作成するための手段は、前記翻訳知識の集合から一つの翻訳知識を除いて得られる部分集合を複数個作成するための手段を含む、請求項 6 に記載の翻訳知識最適化装置。

50

【請求項 8】

前記部分集合を作成するための手段は、前記翻訳知識の集合から予め定められる数の翻訳知識を除いて得る事が可能な全ての部分集合を作成するための手段を含む、請求項 5 に記載の翻訳知識最適化装置。

【請求項 9】

前記機械翻訳手段は、原言語の文を機械翻訳する際に、前記翻訳知識の集合内のどの翻訳知識を使用したかについての情報を出力する機能を持ち、

前記翻訳知識最適化装置はさらに、前記初期評価値を得る際に翻訳された文ごとに、前記機械翻訳手段から出力される、翻訳の際に使用した翻訳規則を特定する情報を記憶するための手段を含み、

10

前記判定手段は、

前記記憶するための手段に記憶されている、前記翻訳規則を特定する情報を参照して、前記複数の部分集合の各々について、当該部分集合の補集合に含まれる翻訳規則を用いて翻訳された前記原言語の文の集合を特定するための手段と、

前記部分集合の各々を用いて、当該部分集合の補集合に含まれる翻訳規則を用いて翻訳された前記原言語の文の集合を前記機械翻訳手段により再び機械翻訳するための手段と、

前記部分集合の各々に対し、前記初期翻訳結果のうち、当該部分集合の補集合に含まれる翻訳規則を用いて翻訳された翻訳結果を、前記再び機械翻訳するための手段による翻訳結果で置換え、当該置換え後の初期翻訳結果に対して前記訳質自動評価手段による訳質評価を行なって、当該部分集合による翻訳結果の評価値を得るための手段と、

20

前記部分集合の各々に対し、当該部分集合による翻訳結果の評価値が前記初期評価値に対し前記所定の条件を満足しているか否かを判定するための手段とを含む、請求項 5 に記載の翻訳知識最適化装置。

【請求項 10】

前記判定するための手段は、前記部分集合の各々に対し、当該部分集合による翻訳結果の評価値が、前記初期評価値を上回っているか否かを判定するための手段を含む、請求項 9 に記載の翻訳知識最適化装置。

【請求項 11】

予め準備された、前記原言語と前記目的言語との対訳文からなる訓練コーパスから、各々が訓練サブコーパス及び評価サブコーパスを含む複数のサブコーパス対を作成するための手段と、

30

予め定められる翻訳規則の構築方式に従って、与えられる対訳コーパスから翻訳規則を自動的に構築するための翻訳知識自動構築手段と、

前記翻訳知識自動構築手段を用いて前記訓練コーパスから翻訳知識を自動構築し、基本翻訳知識として記憶するための基本翻訳知識記憶手段と、

前記複数のサブコーパス対の各々に対して、前記訓練サブコーパスから前記翻訳知識自動構築手段を用いて翻訳知識の集合を自動構築し、当該翻訳知識の集合に対し、前記評価サブコーパスを前記機械読取可能な対訳コーパスとして、前記翻訳知識記憶手段、前記機械読取可能な対訳コーパスを記憶するための手段、前記機械翻訳手段、前記訳質自動評価手段、及び前記最適化手段による最適化を行なうための手段と、

40

前記最適化を行なうための手段によって最適化された、前記複数のサブコーパス対の各々に対して得られる翻訳知識の集合を、一つの翻訳知識の集合に集約するための手段とをさらに含む、請求項 1 に記載の翻訳知識最適化装置。

【請求項 12】

前記集約するための手段は、

前記基本翻訳知識記憶手段に記憶された前記基本翻訳知識に含まれる翻訳知識の各々について、前記最適化手段により算出された差分を、前記複数のサブコーパス対の全てにわたって合計するための差分合計手段と、

前記差分合計手段により合計された差分が所定の条件を満足する翻訳知識を削除する様に前記基本翻訳知識記憶手段に記憶されている前記基本翻訳知識を更新するための手段とを

50

含む、請求項 1 1 に記載の翻訳知識最適化装置。

【請求項 1 3】

前記基本翻訳知識を更新するための手段は、前記差分合計手段により合計された差分が負となる翻訳知識を削除する様に前記基本翻訳知識記憶手段に記憶されている前記基本翻訳知識を更新するための手段を含む、請求項 1 2 に記載の翻訳知識最適化装置。

【請求項 1 4】

前記複数のサブコーパス対を作成するための手段は、

前記訓練コーパスを予め定める個数に実質的に等分して前記予め定める個数の評価サブコーパスを作成するための手段と、

前記予め定める個数の評価サブコーパスの各々に対して、前記訓練コーパスから当該評価サブコーパスを除いたコーパスを作成し、当該評価サブコーパスと対となる訓練サブコーパスを作成するための手段とを含む、請求項 1 1 に記載の翻訳知識最適化装置。 10

【請求項 1 5】

コンピュータにより実行されると、当該コンピュータを、請求項 1 から請求項 1 4 のいずれかに記載の翻訳知識最適化装置として動作させる、翻訳知識最適化のためのコンピュータプログラム。

【請求項 1 6】

請求項 1 5 に記載のコンピュータプログラムによりプログラムされたコンピュータ。

【請求項 1 7】

請求項 1 5 に記載のコンピュータプログラムを記録した、コンピュータ読取可能な記憶媒体。 20

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

この発明は翻訳規則等の翻訳知識を用いた機械翻訳装置のための翻訳知識の作成装置に関し、特に、訓練コーパスから自動構築した翻訳知識等の様に誤り又は冗長な情報を含む知識を取捨選択する事により、的確な翻訳知識の集合を自動的に作成するための装置及びそのためのコンピュータプログラムに関する。

【0002】

【従来の技術】

機械翻訳の一手法として、構文トランスファ方式のものが知られている。構文トランスファ方式は、予め原言語の単語又は句から目的言語への単語又は句などへのマッピング規則（変換規則）及び単語の対訳等を準備しておき、原言語の入力文を解析した後にこのマッピング規則と単語の対訳とを適用して目的言語の翻訳文を得ようとするものである。構文トランスファ方式の機械翻訳システムの構築において最も手間がかかる作業は、この翻訳規則及び単語の対訳の様な翻訳知識の作成である。

【0003】

元々は翻訳規則は人手で準備されるものであった。しかし、原言語と目的言語との対訳文の集合である対訳コーパスの充実に伴い、翻訳規則を対訳コーパスから自動的に獲得する手法が提案されている。翻訳規則を自動的に獲得できれば、機械翻訳システムの構築のための作業量が大いに削減される。 30 40

【0004】

翻訳規則を対訳コーパスから自動的に獲得する手法として複数のものが提案されている。しかし、その様に自動獲得した規則には以下の様な問題がある。

【0005】

例えば、従来の翻訳規則の自動構築方法は不完全であり、構築された翻訳規則にはどうしても誤りが含まれる事が避けられない。たとえば、後掲の非特許文献 1 では対訳コーパスから翻訳規則の基になる句の対応関係を自動抽出しているが、約 8 % の対応関係が誤っていると報告されている。誤りを含む規則が翻訳時に使用されると誤訳を生じる。

【0006】

また、通常は一つの原文でも複数通りの翻訳を行なう事が可能である。対訳コーパスにその様な対訳群が含まれていると、その多様性のために多数の冗長な規則が獲得される。その結果、互いに競合する複数の規則が獲得されてしまう。

【 0 0 0 7 】

例えば言換え表現が存在すると、それらの表現ごとに異なる翻訳規則が作成される。その結果、機械翻訳を行なう際のあいまい性が増大する。あいまい性が増大すると、適切な翻訳を生成する事が困難になる。すなわち、対訳コーパス中の言換え表現により、機械翻訳の精度が低下する。

【 0 0 0 8 】

また、対訳コーパス中に、文脈に依存する訳又は状況に依存する訳が含まれていると、過剰な省略を行なったり、湧き出し語を生じたりする翻訳規則が得られてしまう。そうした翻訳規則は誤訳の原因となる。

【 0 0 0 9 】

従来、こうした冗長 / 競合規則を処理するためのアプローチとして、以下の二つが提案されている。第1のアプローチは、翻訳時に適切な規則を選択する事によりあいまい性を解消する方法である。第2のアプローチは、翻訳規則を自動獲得した後に、後処理として競合規則を取捨選択し、翻訳規則をよりの確なものにするという方法である。本発明は、この第2のアプローチをとる。

【 0 0 1 0 】

この第2のアプローチによる競合規則の整理及び最適化（これを以後「翻訳規則のクリーニング」又は単に「クリーニング」と呼ぶ。）として知られているものに、後掲の非特許文献2において提案されているものと、非特許文献3において提案されているものとがある。

【 0 0 1 1 】

非特許文献2において提案されている手法は、自動獲得された翻訳規則のうち、同じパターンの出現する頻度が所定の値（例えば2）以上の規則のみを採用するという、規則の出現頻度に基づく手法である。非特許文献3において提案されている手法は、特に多数出現するパターンのみを処理対象とし、さらに二乗検定による仮説検定を行なう事によって統計的に信頼性が高い規則のみを抽出するという手法である。

【 0 0 1 2 】

【 非特許文献 1 】

イマムラ、K. (2001). 構文解析と融合した階層的句アライメント. 第6回自然言語処理パシフィックリムシンポジウム (NLP RS 2001) 予稿集 377 頁から 384 頁 (Hierarchical phrase alignment harmonized with parsing. In Proceedings of the 6th Natural Language Processing Pacific Rim Symposium (NLP RS 2001), pp. 374 - 384)

【 0 0 1 3 】

【 非特許文献 2 】

メネツェス、A.、リチャードソン、スティーブン D. (2001). バイリンガルコーパスからの変換マッピングの自動抽出のための最良優先アルゴリズム. MTサミットVIIにおける『用例ベース機械翻訳ワークショップ』予稿集、35 頁から 42 頁 (Menezes, A., Richardson, Stephen D. (2001) A best first alignment algorithm for automatic extraction of transfer mappings from bilingual corpora. In Proceedings of the 'Workshop on Example-Based Machine Translation' in MT Summit VII, pp. 35 - 42)

【 0 0 1 4 】

【 非特許文献 3 】

10

20

30

40

50

イマムラ、K. (2002). パターンベース機械翻訳のための句アライメントにより得られた翻訳知識の応用. 第9回機械翻訳における理論的及び方法論的問題に関する会議予稿集、74頁から84頁 (Imamura, K. (2002). Application of translation knowledge acquired by hierarchical phrase alignment for pattern-based MT. In Proceedings of the 9th Conference On Theoretical and Methodological Issues in Machine Translation (TMI-2002), pp. 74 - 84)

【0015】

10

【発明が解決しようとする課題】

上記した非特許文献2に記載の手法では、規則の数はクリーニング前の1/9程度になり、かつ翻訳品質の若干の向上が見られたという例が非特許文献3に報告されている。しかし、冗長規則が大幅に削減されたにもかかわらず、それに見合う様な翻訳品質の向上は得られていない。

【0016】

また、非特許文献3で提案された手法では、統計的に信頼できる規則として得られるものの数が、コーパスサイズに比べて少ない。そのため、十分な数の翻訳規則を得るためには超大規模コーパスを必要とする問題点がある。また統計的に信頼でき、かつ機械翻訳に十分な数の規則を作成可能な超大規模コーパスは現在存在しない。

20

【0017】

それゆえにこの発明の目的は、対訳コーパスから自動獲得された翻訳規則をクリーニングしてより翻訳品質を向上させる事ができる翻訳知識最適化装置及びそのためのコンピュータプログラムを提供する事である。

【0018】

この発明のほかの目的は、通常規模の対訳コーパスから自動獲得された翻訳規則をクリーニングして、より翻訳品質を向上させる事ができる翻訳知識最適化装置及びそのためのコンピュータプログラムを提供する事である。

【0019】

この発明のほかの目的は、通常規模の対訳コーパスから自動獲得された翻訳規則を比較的短時間にクリーニングして、より翻訳品質を向上させる事ができる翻訳知識最適化装置及びそのためのコンピュータプログラムを提供する事である。

30

【0020】

【課題を解決するための手段】

本発明の第1の局面に係る翻訳知識最適化装置は、機械翻訳のための翻訳知識を最適化するための翻訳知識最適化装置であって、翻訳知識の集合を記憶するための翻訳知識記憶手段と、原言語と目的言語との対訳文を複数個含む、機械読取可能な対訳コーパスを記憶するための手段と、翻訳知識記憶手段に記憶された翻訳知識の集合を利用して、対訳コーパス中の原言語の文を目的言語に機械翻訳するための機械翻訳手段と、機械翻訳手段による翻訳結果の品質を、対訳コーパスを参照して自動的に評価して評価値を出力するための訳質自動評価手段と、訳質自動評価手段の出力する評価値が極値をとる様に、翻訳知識の集合の最適化を行なうための最適化手段とを含む。

40

【0021】

好ましくは、翻訳知識は、原言語の構文パターンから目的言語の構文パターンへの構文変換規則を含む。

【0022】

さらに好ましくは、最適化手段は、翻訳知識の集合に含まれる翻訳知識の各々について、その規則寄与度を算出するための手段と、規則寄与度が予め定める条件を満足する翻訳知識を、翻訳知識の集合から削除するための手段とを含む。

【0023】

50

規則寄与度を算出するための手段は、翻訳知識の集合の全体を用いて、機械翻訳手段による翻訳、及びその翻訳の結果の訳質自動評価手段による訳質評価を行ない、初期評価値を得るための手段と、翻訳知識の集合から、ある翻訳知識を削除して得られる翻訳知識の集合を用いて、機械翻訳による翻訳、及びその翻訳の結果の訳質自動評価手段による訳質評価を行ない、削除後評価値を得るための手段と、削除後評価値と初期評価値との差分を、ある翻訳知識の規則寄与度として算出するための手段とを含んでもよい。

【0024】

さらに好ましくは、最適化手段は、翻訳知識の集合の全体を用いて、機械翻訳手段による翻訳、及びその翻訳結果の訳質自動評価手段による訳質評価を行ない、初期評価値を得るための手段と、予め定められた方法に従って、翻訳知識の集合から複数の部分集合を作成するための手段と、複数の部分集合の各々を用いて機械翻訳手段による翻訳、及びその翻訳の訳質自動評価手段による訳質の評価を行ない、その評価値が初期評価値に対して所定の条件を満足するか否かを判定するための判定手段と、判定するための手段により評価値が所定の条件を満足すると判定された部分集合の各々について、その補集合に属する翻訳知識を翻訳知識の集合から削除するための手段とを含む。

10

【0025】

部分集合を作成するための手段は、翻訳知識の集合から予め定められる数の翻訳知識を除いて得られる部分集合を複数個作成するための手段を含んでもよい。

【0026】

好ましくは、部分集合を複数個作成するための手段は、翻訳知識の集合から一つの翻訳知識を除いて得られる部分集合を複数個作成するための手段を含む。

20

【0027】

さらに好ましくは、部分集合を作成するための手段は、翻訳知識の集合から予め定められる数の翻訳知識を除いて得る事が可能な全ての部分集合を作成するための手段を含む。

【0028】

機械翻訳手段は、原言語の文を機械翻訳する際に、翻訳知識の集合内のどの翻訳知識を使用したかについての情報を出力する機能を持ち、翻訳知識最適化装置はさらに、初期評価値を得る際に翻訳された文ごとに、機械翻訳手段から出力される、翻訳の際に使用した翻訳規則を特定する情報を記憶するための手段を含み、判定手段は、記憶するための手段に記憶されている、翻訳規則を特定する情報を参照して、複数の部分集合の各々について、当該部分集合の補集合に含まれる翻訳規則を用いて翻訳された原言語の文の集合を特定するための手段と、部分集合の各々を用いて、当該部分集合の補集合に含まれる翻訳規則を用いて翻訳された原言語の文の集合を機械翻訳手段により再び機械翻訳するための手段と、部分集合の各々に対し、初期翻訳結果のうち、当該部分集合の補集合に含まれる翻訳規則を用いて翻訳された翻訳結果を、再び機械翻訳するための手段による翻訳結果で置換え、当該置換え後の初期翻訳結果に対して訳質自動評価手段による訳質評価を行なって、当該部分集合による翻訳結果の評価値を得るための手段と、部分集合の各々に対し、当該部分集合による翻訳結果の評価値と初期評価値とが所定の条件を満足しているか否かを判定するための手段とを含んでもよい。

30

【0029】

好ましくは、判定するための手段は、部分集合の各々に対し、当該部分集合による翻訳結果の評価値が、初期評価値を上回っているか否かを判定するための手段を含む。

40

【0030】

好ましくは、翻訳知識最適化装置はさらに、予め準備された、原言語と目的言語との対訳文からなる訓練コーパスから、各々が訓練サブコーパス及び評価サブコーパスを含む複数のサブコーパス対を作成するための手段と、予め定められる翻訳規則の構築方式に従って、与えられる対訳コーパスから翻訳規則を自動的に構築するための翻訳知識自動構築手段と、翻訳知識自動構築手段を用いて訓練コーパスから翻訳知識を自動構築し、基本翻訳知識として記憶するための基本翻訳知識記憶手段と、複数のサブコーパス対の各々に対して、訓練サブコーパスから翻訳知識自動構築手段を用いて翻訳知識の集合を自動構築し

50

、当該翻訳知識の集合に対し、評価サブコーパスを機械読取可能な対訳コーパスとして、翻訳知識記憶手段、機械読取可能な対訳コーパスを記憶するための手段、機械翻訳手段、訳質自動評価手段、及び最適化手段による最適化を行なうための手段と、最適化を行なうための手段によって最適化された、複数のサブコーパス対の各々に対して得られる翻訳知識の集合を、一つの翻訳知識の集合に集約するための手段とを含む。

【0031】

さらに好ましくは、集約するための手段は、基本翻訳知識記憶手段に記憶された基本翻訳知識に含まれる翻訳知識の各々について、最適化手段により算出された差分を、複数のサブコーパス対の全てにわたって合計するための差分合計手段と、差分合計手段により合計された差分が所定の条件を満足する翻訳知識を削除する様に基本翻訳知識記憶手段に記憶されている基本翻訳知識を更新するための手段とを含む。 10

【0032】

基本翻訳知識を更新するための手段は、差分合計手段により合計された差分が負となる翻訳知識を削除する様に基本翻訳知識記憶手段に記憶されている基本翻訳知識を更新するための手段を含んでもよい。

【0033】

好ましくは、複数のサブコーパス対を作成するための手段は、訓練コーパスを予め定める個数に実質的に等分して予め定める個数の評価サブコーパスを作成するための手段と、予め定める個数の評価サブコーパスの各々に対して、訓練コーパスから当該評価サブコーパスを除いたコーパスを作成し、当該評価サブコーパスと対となる訓練サブコーパスを作成するための手段とを含む。 20

【0034】

本発明の第2の局面に係るコンピュータプログラムは、コンピュータにより実行されると、当該コンピュータを、上記したいずれかの翻訳知識最適化装置として動作させるものである。

【0035】

本発明の第3の局面に係るコンピュータは、上記したコンピュータプログラムによりプログラムされたコンピュータプログラムによりプログラムされたコンピュータである。

【0036】

本発明の第4の局面に係る記憶媒体は、上記したコンピュータプログラムを記録した、コンピュータ読取可能な記憶媒体である。 30

【0037】

【発明の実施の形態】

以下、本発明の実施の形態について説明する。以下の説明では、同じ部品には同じ参照番号を付す。それらの機能も同一である。従って、それらについての詳細な説明は繰返さない。

【0038】

なお以下の説明では、第1及び第2の実施の形態を説明する。これらの実施の形態の基本的な考え方は以下の通りである。すなわち、自動構築された翻訳規則を用いて評価コーパス中の原言語の文を機械翻訳する。機械翻訳した結果に対し、非特許文献4に記載されている様な訳質の自動評価を行ない、自動評価値を得る。この自動評価値を向上させる様に翻訳規則の取捨選択を行なう事により、最適な翻訳規則の組合せ（最適な翻訳規則集合）を得る。 40

【0039】

以下の実施の形態では、最適な翻訳規則の組合せには山登り法を使用する。この際、組合せごとに得られる自動評価値を評価関数の出力とみなす。

【0040】

特に以下の実施の形態では、自動構築された翻訳規則集合に対し規則の削除だけを行なう事により、翻訳規則集合の最適化を行なう。この様に規則の削除に限定する事により、クリーニングのための処理が早くなるという効果がある。 50

【 0 0 4 1 】

また、以下の実施の形態では英語から日本語に翻訳する際の翻訳規則集合を最適化する場合について説明する。しかし、本発明はこうした言語の組合せに限定されるわけではなく、翻訳規則を適用する事により翻訳できる言語の組合せであればどのようなものに対しても適用できる。

【 0 0 4 2 】

[第 1 の実施の形態]

構成

図 1 は本発明の第 1 の実施の形態に係る翻訳規則抽出装置 2 0 の機能的構成を示すブロック図である。図 1 を参照して、翻訳規則抽出装置 2 0 は、原言語（英語）と目的言語（日本語）との対訳文を多数含む訓練コーパス 3 0 と、訓練コーパス 3 0 から機械翻訳規則を自動的に構築するための機械翻訳規則自動構築部 3 2 と、機械翻訳規則自動構築部 3 2 が構築した翻訳規則集合に対して、後述する様なフィードバッククリーニング処理を行なうためのフィードバッククリーニング部 3 4 と、フィードバッククリーニング部 3 4 がフィードバッククリーニングを行なう際に、訳質評価のために参照する評価コーパス 3 6 とを含む。評価コーパス 3 6 中の対訳文は、英語の原文と、原文を人が日本語に翻訳した結果（参照訳と呼ぶ。）とからなる。

10

【 0 0 4 3 】

フィードバッククリーニング部 3 4 は、機械翻訳規則自動構築部 3 2 により訓練コーパス 3 0 から自動的に構築された翻訳規則の集合を記憶するための翻訳規則集合記憶部 4 0 と、翻訳規則集合記憶部 4 0 に記憶された翻訳規則を用いて評価コーパス 3 6 中の全ての英語の原文を目的言語の文に翻訳するための機械翻訳エンジン 4 2 とを含む。機械翻訳エンジン 4 2 は構文トランスファ方式のものである。

20

【 0 0 4 4 】

フィードバッククリーニング部 3 4 はさらに、機械翻訳エンジン 4 2 による翻訳結果を、各文の翻訳の際に使用された翻訳規則を特定する情報とともに記憶するための翻訳結果記憶部 4 3 を含む。翻訳結果記憶部 3 5 はまた、翻訳結果とともに各文の翻訳の際に使用された翻訳規則を特定する情報も記憶する。

【 0 0 4 5 】

フィードバッククリーニング部 3 4 はこれに加えて、翻訳結果記憶部 4 3 に記憶されている日本語の文（翻訳文）の訳の品質（訳質）を、評価コーパス 3 6 を用いて自動的に評価するための訳質自動評価部 4 4 と、翻訳規則集合記憶部 4 0 に含まれる規則ごとに、その規則を削除した後の自動評価値を算出し、削除前の自動評価値との差分（この差分をその規則の「規則寄与度」と呼ぶ。）を算出するための規則寄与度算出部 4 6 とを含む。規則寄与度算出部 4 6 は、寄与度の算出の際には、訳質自動評価部 4 4 による評価値と、翻訳結果記憶部 4 3 に記憶されている、翻訳の際に使用された翻訳規則を特定する情報とを用いる。

30

【 0 0 4 6 】

フィードバッククリーニング部 3 4 はさらに、翻訳規則のうち、寄与度算出部 4 6 が算出した規則寄与度が所定の条件を充足した翻訳規則（本実施の形態では規則寄与度が負の値である翻訳規則）を翻訳規則集合記憶部 4 0 中の翻訳規則の集合から削除するための翻訳規則削除部 4 8 を含む。

40

【 0 0 4 7 】

本実施の形態では、機械翻訳規則自動構築部 3 2 による翻訳規則の自動構築には、前述の非特許文献 3 に記載された方法を使用する。

【 0 0 4 8 】

本実施の形態では、機械翻訳エンジン 4 2 は、構文トランスファ方式であって、後掲の参考文献 1 に記載されたものを使用する。機械翻訳エンジン 4 2 は、英語の構文構造を日本語の構文構造に変換する翻訳規則を使用する。機械翻訳エンジン 4 2 が使用する翻訳規則の例を図 2 に示す。この例では、1 つの規則は、構文規則と、原言語パターンと、目的言

50

語パターンと、用例とを含む。

【 0 0 4 9 】

構文カテゴリは、この規則が適用される英語構文ノードのカテゴリを表す。

【 0 0 5 0 】

原言語パターンは、この規則が適用される英語構文構造のパターンを示す。原言語パターンは、X , Yなどの非終端記号(変数)と、単語又はマーカなどの終端記号との列である。

【 0 0 5 1 】

目的言語パターンは、この規則が適用された場合に生成される日本語構文構造のパターンを示す。原言語パターンに対応する変数(X '、Y ' など)と、単語で表現された終端記号との列である。 10

【 0 0 5 2 】

用例は、訓練コーパス中に現れた変数の実例である。変数の数と一致する主辞単語の組である。本実施の形態における翻訳規則集合記憶部 4 0 中の各規則の用例は、訓練コーパス 3 0 中での出現例となる。

【 0 0 5 3 】

翻訳規則集合記憶部 4 0 が記憶する翻訳規則は機械翻訳エンジン 4 2 が使用する翻訳規則のフォーマットに従ったものである。

【 0 0 5 4 】

図 2 に示す規則のうち、たとえば規則番号 1 のものは、英語の " p r e s e n t a t t h e c o n f e r e n c e " という句に適用され、「会議(c o n f e r e n c e の訳)で発表する(p r e s e n t の訳)」という訳を生成する事を表している。 20

【 0 0 5 5 】

訳質自動評価部 4 4 は、後掲の参考文献 2 に記載の B L E U を用いる。B L E U の様な機械翻訳の自動評価法についてはこの他にもいくつか提案されている。これらは、機械翻訳システムの開発時、従来主観評価を行っていた部分を置換える事により、開発サイクルのスピードアップを狙ったものである。これらは全自動で行なわれるため、従来考えられていた様な開発支援ばかりではなく、本実施の形態の様に翻訳システムの自動チューニングにも利用できる。

【 0 0 5 6 】

本実施の形態で訳質自動評価に使用する B L E U は、評価コーパスの原文を機械翻訳エンジン 4 2 により機械翻訳した結果と、評価コーパス 3 6 中の参照訳との類似度を計算し、訳質をスコア(B L E U スコア)として出力する。類似度は、両者の N - g r a m 一致数で測定される。N は可変であるが、本実施の形態では 1 - g r a m から 4 - g r a m までを用いる。 30

【 0 0 5 7 】

ここで注意すべきは、B L E U スコアを本実施の形態の様に機械翻訳規則集合の評価に用いるためには、ある程度の大きさを持った文集合を用いる必要がある事である。B L E U スコアを 1 文ごとに算出する事も可能ではあるが、そのままでは主観評価とのずれが大きい。個々の類似度を翻訳結果集合に含まれる翻訳文の全体について計算し総和をとる事により、個々の誤差を相殺できる。 40

【 0 0 5 8 】

規則寄与度算出部 4 6 は、次の様にして規則ごとに規則寄与度を算出する。まず、機械翻訳エンジン 4 2 による評価コーパス 3 6 の原言語の全ての文の翻訳結果に対し訳質自動評価部 4 4 が算出したスコアを用い、基準となる自動評価値を得る。この値を削除前自動評価値と呼ぶ。この翻訳により、どの文の翻訳にどの規則が使用されたかという情報も得られる。

【 0 0 5 9 】

次に、翻訳規則集合記憶部 4 0 内の翻訳規則ごとに、翻訳規則集合記憶部 4 0 からその規則を削除して得られる部分集合を用いて評価コーパス 3 6 の原言語の全ての文の翻訳を行 50

なった場合のスコアを計算する。このスコアと削除前自動評価値の差分が規則寄与度である。本実施の形態では、削除後のスコアの計算を以下の考え方に従って行なう。なお、この例では、当然の事ながら、削除される一つの翻訳規則からなる集合と、その翻訳規則を削除する事により形成される部分集合とは互いに補集合の関係にある。

【0060】

基本的考え方に従って、翻訳規則集合記憶部40内の規則の組合せ（部分集合）ごとに評価コーパス36を全て翻訳する事も理論的には考えられる。しかしその場合には翻訳回数が非常に多くなる。よほど計算機資源に恵まれていないと合理的な時間内に結果を得る事ができない。そこで、以下の様にして計算量を少なくする。

【0061】

機械翻訳エンジン42による機械翻訳では、1文を翻訳すると、その翻訳に使用された翻訳規則を特定できる。その情報は翻訳結果記憶部43に記憶されている。逆にいうと、評価コーパス36全体を翻訳すると、各規則が使われる文を特定できる。

【0062】

ある規則を翻訳規則集合から削除して得られる部分集合を用いて機械翻訳エンジン42により翻訳を行なうとき、それによって変化する翻訳文は、そのある規則の削除前にその規則を使用して翻訳された翻訳文だけである。他の文は別の規則を使用して翻訳されたので、削除対象の規則が削除された後の翻訳規則集合を用いた翻訳を行なっても翻訳結果は変化しない。

【0063】

従って、翻訳規則集合からある規則を削除した場合、削除前にその規則を使用して翻訳した文のみを削除後の翻訳規則集合を用いて翻訳し、他の訳文とあわせて参照訳との類似度を求めれば削除後のBLEUSコアが得られる。全ての文を翻訳する必要はない。

【0064】

以上から、翻訳規則の削除のみを行なう事により、合理的な時間内に結果を得る事が可能になる。

【0065】

すなわち規則寄与度算出部46は、訳質自動評価部44による削除前自動評価値と、翻訳にどの規則が使用されたか（どの規則がどの文の翻訳に使用されたか）に関する、翻訳結果記憶部43に記憶されている情報を得る。規則ごとに、その規則を用いて翻訳された文を、その規則以外の規則を用いて再翻訳した場合の、訳文全体の自動評価値を算出する。この評価値と削除前自動評価値との差分（削除前自動評価値 - 削除後の評価値）を算出し、それをその規則の規則寄与度とする。規則寄与度算出部46はさらに、こうして算出された規則寄与度が負となる（つまり、削除する事により自動評価値が大きくなる）規則の規則番号を翻訳規則削除部48に与える機能を持つ。なお、規則寄与度算出部46では、その処理の収束を早めるため、削除される規則同士は互いに独立であると仮定し、1回の繰返しで、削除すべき規則を全て決定し削除している。

【0066】

より具体的には、規則寄与度算出部46は以下の様にして規則寄与度を算出する。翻訳規則集合のうち、機械翻訳エンジン42による翻訳の際に使用された翻訳規則の各々について、その規則を翻訳の際に使用した文の集合を求める。その文の集合が空集合でなければ、基の規則集合からその翻訳規則を取除いて得られる部分集合を用いて、その文の集合内の各文について機械翻訳エンジン42による翻訳を再度行なう。翻訳結果記憶部43に記憶された翻訳結果のうち、この翻訳規則を用いて翻訳が行なわれたものを、再翻訳したものと置換える。そして再度訳質自動評価部44によって訳質の自動評価を行なう。こうして得られた削除後の評価値と削除前自動評価値との差分がこの翻訳規則の規則寄与度となる。

【0067】

この処理を、翻訳規則集合記憶部40内の全ての翻訳規則に対して行ない、規則寄与度が負の規則を特定する事により、削除すべき翻訳規則が決定される。

10

20

30

40

50

【 0 0 6 8 】

翻訳規則削除部 4 8 は、翻訳規則集合記憶部 4 0 内の規則のうち、規則寄与度算出部 4 6 から与えられた情報に対応する翻訳規則を削除する機能を持つ。

【 0 0 6 9 】

動作

第 1 の実施の形態に係る翻訳規則抽出装置 2 0 は以下の様に動作する。訓練コーパス 3 0 及び評価コーパス 3 6 は予め準備されているものとする。機械翻訳規則自動構築部 3 2 は、訓練コーパス 3 0 内の各対訳文から翻訳規則を自動構築し、翻訳規則集合記憶部 4 0 に記憶させる。

【 0 0 7 0 】

機械翻訳エンジン 4 2 は、評価コーパス 3 6 に含まれる対訳文のうちの原文の全てを、翻訳規則集合記憶部 4 0 に記憶されている翻訳規則を用いて翻訳する。翻訳結果は、翻訳の際に使用された翻訳規則を特定する情報とともに翻訳結果記憶部 4 3 に記憶される。

【 0 0 7 1 】

訳質自動評価部 4 4 は、翻訳結果記憶部 4 3 に記憶されている翻訳文の訳質を、評価コーパス 3 6 に記憶されている参照訳を用いて自動的に BLEU スコアとして評価し、その結果を規則寄与度算出部 4 6 に与える。

【 0 0 7 2 】

規則寄与度算出部 4 6 は、訳質自動評価部 4 4 から与えられた BLEU スコアを削除前自動評価値とする。次に規則寄与度算出部 4 6 は、翻訳規則集合記憶部 4 0 内の各翻訳規則について、上記した方法に従って規則寄与度を算出する。そして、規則寄与度が負となる規則を特定し、その情報を翻訳規則削除部 4 8 に与える。

【 0 0 7 3 】

翻訳規則削除部 4 8 は、この情報に従って翻訳規則集合記憶部 4 0 に記憶されている翻訳規則集合内の規則を削除する。削除処理後の翻訳規則集合記憶部 4 0 に記憶されている翻訳規則集合は、クリーニングされ最適化されたものとなる。

【 0 0 7 4 】

具体例

翻訳例及び規則寄与度の算出の具体例を示す。なお、削除前自動評価値は 0 . 2 3 3 3 6 3 とする。

【 0 0 7 5 】

翻訳例 1

図 2 の規則 5 は、文脈依存訳から作成された誤り規則の例である。" the nearest subway station " と「最寄りの地下鉄」から作成された規則であり、原文の " station " の訳が日本語では省略されている。

【 0 0 7 6 】

英語 " Please tell me where the nearest rail road station is . " を翻訳すると、この規則 5 が適用されて、日本語「最寄りの鉄道はどこにありますか、教えていただけますか。」と翻訳される。

【 0 0 7 7 】

規則 5 を削除すると、この翻訳は「最寄りの鉄道の駅はどこにありますか、教えていただけますか」に変化する。削除後自動評価値は 0 . 2 3 3 5 4 9 となる。

【 0 0 7 8 】

従って、規則 5 の規則寄与度は $0 . 2 3 3 3 6 3 - 0 . 2 3 3 5 4 9 = - 0 . 0 0 0 1 8 6$ となる。従って規則 5 は削除される。削除の結果、" the nearest rail road station " は「最寄りの鉄道の駅」と正しく翻訳されるようになる。

【 0 0 7 9 】

翻訳例 2

図 2 の規則 6 は、翻訳規則自動構築誤りによって作成された誤った規則の例である。自動構築時、" rent two bicycles " を解析した結果、" rent two

10

20

30

40

50

” が動詞句、” b i c y c l e s ” が名詞句になった例である。正しくは、” r e n t ” が動詞句、” t w o b i c y c l e s ” が名詞句であるが、翻訳規則の自動構築の際にはこの種の誤りの発生を完全に防止する事はできない。

【 0 0 8 0 】

英語 “ I w a n t t o r e n t t w o r a c k e t s ” を翻訳すると、規則 6 が適用されて「ラケットを 2 借りたいのですが」と翻訳される。規則 6 を削除すると、この翻訳は「ラケットを 2 本借りたいのですが」に変化する。すると、規則 6 の削除後の自動評価値は 0 . 2 3 3 5 2 9 となる。規則 6 の規則寄与度は - 0 . 0 0 0 1 6 6 となり、規則 6 は削除される。

【 0 0 8 1 】

翻訳例 3

図 2 の規則 7 及び規則 8 は、言換え表現から作られた規則の例である。どちらも正しい規則であるが、互いに競合する規則である。

【 0 0 8 2 】

英語 “ P l e a s e c a s h t h i s t r a v e l e r ' s c h e c k . ” を翻訳する際には、規則 7 又は規則 8 のいずれかが適用される。今回は規則 7 が選ばれたものとする。翻訳結果は「このトラベラーズチェックを現金にしたいのですが」となる。

【 0 0 8 3 】

規則 7 を削除すると、この翻訳は「このトラベラーズチェックを現金にしてください」に変化する。すると削除後自動評価値は 0 . 2 3 3 5 8 5 となる。これは、評価コーパス 3 6 中に、規則 8 に一致する対訳文が、規則 7 に一致する対訳文よりも多く含まれている事を示す。

【 0 0 8 4 】

規則 7 の規則寄与度はこの場合 - 0 . 0 0 0 2 2 2 となる。その結果、規則 7 が削除され、評価コーパス 3 6 中により多く出現する表現に一致する翻訳が行なわれる様になる。

【 0 0 8 5 】

実施の形態 1 の効果

以上の第 1 の実施の形態の翻訳規則抽出装置 2 0 では、フィードバッククリーニング部 3 4 の機能により、対訳コーパスから自動構築された翻訳規則群を、訳質自動評価部を用いて自動的にクリーニングする事ができる。その結果、機械翻訳結果に悪影響を及ぼす翻訳規則が排除されるので、自動構築された翻訳規則を用いる翻訳システムの翻訳結果の品質が向上するという効果が得られる。現実に、クリーニング後の翻訳規則を用いて翻訳を行なった結果に対しては、未クリーニングの翻訳規則を用いた翻訳結果よりもよい評価が得られた。

【 0 0 8 6 】

コンピュータによる実現

以上述べた第 1 の実施の形態に係る翻訳規則抽出装置 2 0 は、コンピュータ及びその上で実行されるソフトウェアによっても実現される。図 3 に翻訳規則抽出装置 2 0 を構成するコンピュータの外観図を、図 4 にそのブロック図を、それぞれ示す。

【 0 0 8 7 】

図 3 を参照して、翻訳規則抽出装置 2 0 を構成するコンピュータシステムは、C D - R O M (C o m p a c t D i s c R e a d - O n l y M e m o r y) ドライブ 7 0 及び F D (F l e x i b l e D i s k) ドライブ 7 2 を有するコンピュータ 6 0 と、いずれもコンピュータ 6 0 に接続されたモニタ 6 2、キーボード 6 6、及びマウス 6 8 とを含む。

【 0 0 8 8 】

図 4 を参照して、コンピュータ 6 0 はさらに、C P U (中央演算処理装置 : C e n t r a l P r o c e s s i n g U n i t) 7 6 と、C P U 7 6 に接続されたバス 8 6 と、バス 8 6 を介して C P U 7 6 と相互に接続された R A M 7 8、R O M 8 0、及びハードディスク 7 4 とを含む。バス 8 6 には C D - R O M ドライブ 7 0 及び F D ドライブ 7 2 も接続

10

20

30

40

50

される。C D - R O M ドライブ 7 0 には C D - R O M 8 2 が、F D ドライブ 7 2 には F D 8 4 が、それぞれ装着され、C P U 7 6 等との間のデータの入出力を行なう事ができる。

【 0 0 8 9 】

図 3 及び図 4 に示すコンピュータは、以下に述べる様な制御構造を有するコンピュータプログラム（以下単に「プログラム」と呼ぶ。）を実行する事により、図 1 に示す翻訳規則抽出装置 2 0 として動作する。このプログラムは、たとえば C D - R O M 8 2 上にコンピュータ読取可能なデータとして記録されて流通する。この C D - R O M 8 2 を C D - R O M ドライブ 7 0 に装着し、プログラムを讀出してハードディスク 7 4 に記憶する事により、コンピュータ 6 0 はいつでもこのプログラムを実行する事ができる。なお、訓練コーパス 3 0、評価コーパス 3 6 などはハードディスク 7 4 に記憶しておく。C P U 7 6 はまた、必要なデータはハードディスク 7 4 から讀出して R A M 7 8 に格納する。

10

【 0 0 9 0 】

プログラムの実行時には、ハードディスク 7 4 に記憶されているプログラムを R O M 8 0 にロードする。C P U 7 6 は、図示しないプログラムカウンタにより示されるアドレスの命令を R O M 8 0 から讀出して実行する。C P U 7 6 は、実行結果を所定のアドレスに出力し、あわせて実行結果に従ってプログラムカウンタの内容を更新する。

【 0 0 9 1 】

こうした処理を繰返し行なう事により、最終的な翻訳規則の集合が得られる。得られた結果は、本実施の形態では最終的にハードディスク 7 4 に格納される。

【 0 0 9 2 】

なお、コンピュータ 6 0 の動作自体は周知であるので、ここではその詳細については繰返さない。

20

【 0 0 9 3 】

プログラムの制御構造

図 5 を参照して、フィードバッククリーニング部 3 4 を実現するプログラムは以下の制御構造を有する。まず、このプログラムは、起動されるとステップ 1 0 0 で削除規則集合 R r e m o v e を空集合とする。ステップ 1 0 2 で、機械翻訳エンジン 4 2 を用いて評価コーパス 3 6 の全ての原文を翻訳規則集合記憶部 4 0 の翻訳規則を参照して翻訳し、翻訳結果集合 D o c を得る。このとき同時に、翻訳するためにどの規則が使われたかを記録する。この記録に基づき、ある規則 r を用いて翻訳された原文集合を求める。この原文集合を、規則 r に対して S [r] とする。続いてステップ 1 0 4 で、この翻訳結果集合 D o c から、訳質自動評価部 4 4 を用いて初期（削除前）自動評価値 s c o r e を算出する。

30

【 0 0 9 4 】

続いて以下に述べるステップ 1 0 8 ~ ステップ 1 2 0 までの処理を、翻訳規則集合記憶部 4 0 内の全ての翻訳規則 r について繰返す。まずステップ 1 0 8 では、規則 r を用いた原文集合 S [r] が空集合か否かを判定する。空集合の場合にはこの規則 r に対しては何も行なわない。S [r] が空集合でない場合、制御はステップ 1 1 0 に進む。

【 0 0 9 5 】

ステップ 1 1 0 では、原文集合 S [r] に含まれる原文の全てを、翻訳規則集合から規則 r を取除いたものを用いて、機械翻訳エンジン 4 2 により翻訳する。その結果得られる訳文の集合を T [r] とする。続くステップ 1 1 2 で、ステップ 1 0 2 で求めた翻訳結果集合 D o c 中の、規則 r を用いて翻訳された文の集合を集合 T [r] で置換えた新たな翻訳結果集合 D o c [r] を求める。ステップ 1 1 4 で、この翻訳結果集合 D o c [r] に対する、訳質自動評価部 4 4 による自動評価値 s c o r e [r] を算出する。この自動評価値 s c o r e [r] が削除後自動評価値である。ステップ 1 1 6 で、初期自動評価値 s c o r e からこの削除後自動評価値 s c o r e [r] を減算し、その結果を規則寄与度 c o n t r i b [r] に代入する。

40

【 0 0 9 6 】

ステップ 1 1 8 では、規則寄与度 c o n t r i b [r] が負か否かを判定する。規則寄与度 c o n t r i b [r] が負であれば、制御はステップ 1 2 0 に進み、この規則 r を削除

50

規則集合 `R r e m o v e` に追加する。規則寄与度 `c o n t r i b [r]` が負でなければその規則については何もしない。

【 0 0 9 7 】

以上のステップ 1 0 8 ~ 1 2 0 の処理を全ての規則 `r` について繰返し行なった後、制御はステップ 1 2 4 に進む。ステップ 1 2 4 では、削除規則集合 `R r e m o v e` が空集合でないか判定する。削除規則集合 `R r e m o v e` が空集合であればこのプログラムの実行を終了する。削除規則集合 `R r e m o v e` が空集合でない場合には、ステップ 1 2 6 でこの削除規則集合 `R r e m o v e` に含まれる規則を翻訳規則集合記憶部 4 0 に含まれる翻訳規則集合から削除する。この後、制御は先頭のステップ 1 0 0 に戻り、ステップ 1 2 4 で削除規則集合 `R r e m o v e` が空集合であると判定されるまで、以上の処理を繰返す。

10

【 0 0 9 8 】

以上の様な制御構造を有するプログラムを図 3 及び図 4 に示すコンピュータ 6 0 で実行する事により、図 1 に示す第 1 の実施の形態の翻訳規則抽出装置 2 0 を実現する事ができる。

【 0 0 9 9 】

変形例

上記した第 1 の実施の形態では、翻訳規則の全てについてその規則寄与度を算出して削除するか否かを判定している。しかし、全ての翻訳規則についてこうした処理を行なう必要はなく、一部の規則のみに対して行なってもそれなりの効果が得られる。しかし、翻訳規則の全てについて規則寄与度を算出して削除するか否かを判定した方が、明らかに最終的に得られる翻訳規則に誤った規則又は冗長な規則が含まれる可能性が低くなる。従って、翻訳規則の全てについて上記した処理を行なう方が好ましい。

20

【 0 1 0 0 】

また上記した実施の形態では、一度に一つずつの翻訳規則についてその規則寄与度を算出している。この様にすると、翻訳規則の各々について削除すべきか否かを判定できるので、翻訳規則の最適化を目指す上では好ましい。しかし、この判定を翻訳規則の一つずつについて行なう事が必須というわけではない。原理的には、一度に複数の翻訳規則を削除した場合を想定してその寄与度を算出し、その結果に従ってそれら複数の翻訳規則をまとめて削除する事も可能であり、そうした方法によってもある程度は上記した実施の形態と同様の効果を奏すると考えられる。

30

【 0 1 0 1 】

また、削除すべきか否かを決定する翻訳規則の数は、上記した実施の形態では「 1 」に固定されている。この様に数を固定する事により、処理が簡単になるので、実際にはこうした形で本発明を実施する事が多いと思われる。しかしこの数も常に同じ数である必要はない。たとえば何らかの基準によってその都度決められる数の翻訳規則を処理対象として、その規則寄与度を算出する様にしてもよい。

【 0 1 0 2 】

本発明では、翻訳規則の集合の任意の部分集合（当初の翻訳規則内の翻訳規則の任意の組合せ）を取出し、どの部分集合を用いて機械翻訳を行なえば翻訳結果の訳質として最もよい評価値が得られるか、を確認し、その結果によって最終的な翻訳規則の集合を決定する、という考え方を基本的枠組みとしている。その基本的枠組みの中で、さらに計算機資源を節約しつつどの様にすれば効率的にある程度好ましい基本規則の集合が得られるか、についての一つの実施の形態が上記した第 1 の実施の形態である。この基本的枠組みの中で、第 1 の実施の形態とは細部で異なる実施の形態が他にもあり得る事、及びそうした実施の形態が上記した第 1 の実施の形態についての詳細な説明に基づいて容易に実施する事ができる事は、当業者であれば容易に理解できるであろう。

40

【 0 1 0 3 】

[第 2 の実施の形態]

概略

第 1 の実施の形態の装置によりクリーニングした翻訳規則集合を用いる事により、翻訳の

50

品質はかなり向上する。しかし、未だ向上の余地があると思われる。また、第1の実施の形態では、訓練コーパスとは別に評価コーパスを準備する必要がある。評価コーパスについては、原文に対する参照訳が必要となるため、できれば評価コーパスを特に準備する必要がないほうが望ましい。

【0104】

また、一般的には、訓練コーパスに比べ、評価コーパスはサイズが小さい場合が多い。そのため、たとえ大域最適解を発見する事ができても、評価コーパスではすべての規則をテストできず、クリーニング漏れが発生する。その様なクリーニング漏れの発生を防止できる事が望ましい。

【0105】

そこでこの第2の実施の形態の装置では、第1の実施の形態の装置で用いたフィードバッククリーニング部34によるクリーニング結果に対し、交差検定と同様な考え方を採用し、より最適解に近いものを得るためのクリーニングを行なう。本明細書では、こうしたクリーニングの仕方を「交差クリーニング」と呼ぶ。

【0106】

一般的にN分割交差検定とは、データをN個のサブデータにほぼ等分し、一つをあるモデルのパラメータ推定に用い、推定されたモデルの当てはまりのよさを残りのデータで評価する事をN個のサブデータの全てについて行なう、という方法である。この交差クリーニングにより、上記した様なクリーニング漏れを防止する事ができる。

【0107】

図6に、この実施の形態で行なわれる交差クリーニングの概要を示す。以下、この処理の概要を説明する。

【0108】

ステップ1. 訓練コーパス140をN個に分割する。

【0109】

ステップ2. 分割によって得られたN個のサブコーパスを評価サブコーパス162A、162B、...とする。元の訓練コーパス140から一つの評価サブコーパス(例えば評価サブコーパス162A)を除いたN-1個のサブコーパス(評価サブコーパス162Aの場合、評価サブコーパス162B、162C、...)を一つにまとめ、訓練サブコーパス160Aを作成する。評価サブコーパス162Aと訓練サブコーパス160Aとを対にする。

【0110】

同様に、各評価サブコーパス162B、162C、...に対し、訓練サブコーパス160B、160C、...を作成し、それらを元の評価サブコーパス162B、162C、...と対にする。

【0111】

以上の処理の結果、N個のサブコーパス対150A、150B、...が形成される。これらN個のサブコーパス対150A、150B、...に含まれる訓練サブコーパス160A、160B、...の各々から、第1の実施の形態と同様にして翻訳規則の自動構築151を行なう。その結果、N個の自動構築翻訳規則集合152A、152B、...が得られる。

【0112】

ステップ3. さらに、これら自動構築翻訳規則集合152A、152B、...に対し、それぞれ評価サブコーパス162A、162B、...を用いて、第1の実施の形態と同様のフィードバッククリーニング153を行なう。その結果、N個のクリーニング後規則集合154A、154B、...が得られる。

【0113】

ステップ4. 最後に、N個のクリーニング後規則集合154A、154B、...に対して機械翻訳規則集約処理156を行ない、最終的な交差クリーニング後翻訳規則集合158を作成する。

【0114】

10

20

30

40

50

通常の交差検定との相違点はステップ 4 である。本実施の形態では、規則毎に規則寄与度の総和を算出し、それが 0 以上である場合に最終翻訳規則集合にその規則を出力する。逆にいえば、規則寄与度の総和が 0 未満の規則は翻訳規則集合から削除する。

【0115】

構成

図 7 にこの第 2 の実施の形態の翻訳規則抽出装置 180 の機能的ブロック図を示す。図 7 を参照して、この翻訳規則抽出装置 180 は、訓練コーパス 140 と、訓練コーパス 140 から自動的に翻訳規則を構築するための機械翻訳規則自動構築部 198 と、機械翻訳規則自動構築部 198 により自動構築された翻訳規則の集合（これを「基本翻訳規則集合」と呼ぶ。）を記憶するための基本規則集合記憶部 196 とを含む。機械翻訳規則自動構築部 198 は第 1 の実施の形態で使用されている機械翻訳規則自動構築部 32 と全く同一の機能を持つ。

10

【0116】

翻訳規則抽出装置 180 はさらに、訓練コーパス 140 を N 個に分割し、その一つからなる評価サブコーパス 162 と、他の N - 1 個からなる一つの訓練サブコーパス 160 とに分ける機能を持つ訓練コーパス分割部 190 と、訓練サブコーパス 160 から翻訳規則を自動構築するための機械翻訳規則自動構築部 32 と、機械翻訳規則自動構築部 32 の出力する翻訳規則集合を評価サブコーパス 162 を用いて第 1 の実施の形態と同様にしてフィードバッククリーニングするためのフィードバッククリーニング部 34 とを含む。フィードバッククリーニング部 34 及びその各部の機能は、第 1 の実施の形態におけるフィードバッククリーニング部 34 及びその各部の機能と同じである。従ってそれらの詳細な説明はここでは繰返さない。

20

【0117】

翻訳規則抽出装置 180 はさらに、機械翻訳規則自動構築部 32 による翻訳規則の自動構築及びフィードバッククリーニング部 34 による翻訳規則のフィードバッククリーニングを N 回繰返して実行する様に、訓練コーパス分割部 190、機械翻訳規則自動構築部 32、及びフィードバッククリーニング部 34 を制御するための繰返制御部 192 を含む。繰返制御部 192 による繰返は、訓練コーパス分割部 190 により選択される評価サブコーパス 162 を一つずつ入替えながら行なわれる。

30

【0118】

翻訳規則抽出装置 180 はこれに加えて、フィードバッククリーニング部 34 の規則寄与度算出部 46 により算出された規則寄与度を規則ごと及び繰返しごとに記憶するための規則寄与度記憶部 202 と、機械翻訳規則自動構築部 32 及びフィードバッククリーニング部 34 により作成された N 個のフィードバッククリーニング済みの翻訳規則集合を集約し、最終的な一つの交差クリーニング後翻訳規則集合を基本規則集合記憶部 196 内に作成するための翻訳規則集約部 194 とを含む。翻訳規則集約部 194 は、規則寄与度記憶部 202 に記憶されている規則ごと及び繰返しごとの規則寄与度を用いて、基本規則集合記憶部 196 に記憶されている基本翻訳規則集合から不要な規則を削除する事により規則の集約を行なう。

40

【0119】

機械翻訳規則自動構築部 32 及びフィードバッククリーニング部 34 の機能はそれぞれ第 1 の実施の形態で説明したものと同一である。

【0120】

訓練コーパス分割部 190 は、訓練コーパス 140 を以下の様に繰返しごとに異なる形で分割する。まず、前述の様に訓練コーパス 140 は N 個のサブコーパスにほぼ等分に分割される。それらをそれぞれ第 1 のサブコーパス、第 2 のサブコーパス、... 第 N のサブコーパスと呼ぶ事にする。

【0121】

繰返しの第 1 回目では、訓練コーパス分割部 190 は第 1 のサブコーパスを評価サブコーパス 162 とし、第 2 のサブコーパスから第 N のサブコーパスまでをまとめて訓練サブコ

50

ーパス 160 とする。繰返しの第 2 回目では訓練コーパス分割部 190 は、第 2 のサブコーパスを評価サブコーパス 162 とし、第 1 のサブコーパス、及び第 3 のサブコーパスから第 N のサブコーパスまでをまとめて訓練サブコーパス 160 とする。繰返しの第 3 回目では訓練コーパス分割部 190 は、第 3 のサブコーパスを評価サブコーパス 162 とし、第 1 のサブコーパス、第 2 のサブコーパス、及び第 4 のサブコーパスから第 N のサブコーパスまでをまとめて訓練サブコーパス 160 とする。以下同様にして、繰返しの第 N 回目では訓練コーパス分割部 190 は、第 N のサブコーパスを評価サブコーパス 162 とし、第 1 のサブコーパスから第 N - 1 のサブコーパスまでをまとめて訓練サブコーパス 160 とする。

【0122】

10

以上が訓練コーパス分割部 190 の機能である。

【0123】

翻訳規則集約部 194 は、次の様にしてフィードバッククリーニング後の翻訳規則を集約する。機械翻訳規則自動構築部 198 により、訓練コーパス 140 の全体から基本翻訳規則集合が自動構築される。この基本翻訳規則集合は基本規則集合記憶部 196 に記憶される。

【0124】

次に、繰返制御部 192 による N 回のフィードバッククリーニングにより、訓練コーパス 140 の N 個の訓練サブコーパス 160 より N 個の翻訳規則集合が得られる。これらを第 1 の翻訳規則集合、第 2 の翻訳規則集合、... 第 N の翻訳規則集合と呼ぶ事とする。そして、これらの翻訳規則集合を作成する際に規則寄与度算出部 46 により計算された各規則の規則寄与度が規則寄与度記憶部 202 に繰返しごとに別々に記憶される。規則 r についての i 回目の繰返しの際に計算された規則寄与度を $contrib[i][r]$ と表す ($i = 1 \sim N$ 、 $r = 1 \sim$ 基本規則数)。

20

【0125】

翻訳規則集約部 194 は、全てのフィードバッククリーニングが終了すると、規則寄与度記憶部 202 を参照して、翻訳規則 r ごとに、規則寄与度記憶部 202 に記憶されている規則寄与度の総和 $contrib[r] = \sum_i contrib[i][r]$ を計算する。そして、総和 $contrib[r]$ が負であればその規則 r を基本規則集合記憶部 196 に記憶されている基本規則集合から削除する。この処理を全ての規則 r に対して実行する事により、基本規則集合記憶部 196 に記憶されている基本規則集合に対するクリーニングが行なわれ、最終的な交差フィードバッククリーニング後の翻訳規則集合が得られる。

30

【0126】

動作

この第 2 の実施の形態に係る翻訳規則抽出装置 180 は以下の様に動作する。訓練コーパス 140 は最初に準備されているものとする。また訓練コーパス 140 を N 個にほぼ等分する方法も予め決定されているものとする。まず機械翻訳規則自動構築部 198 が訓練コーパス 140 から翻訳規則を自動構築する。構築された翻訳規則集合 (基本規則集合) は基本規則集合記憶部 196 に記憶される。

【0127】

40

以下の繰返し処理は、繰返制御部 192 による制御の下で実行される。まず訓練コーパス分割部 190 は、訓練コーパス 140 から第 1 のサブコーパスを選び、それを評価サブコーパス 162 とする。訓練コーパス分割部 190 はさらに、残りの N - 1 個のサブコーパスをまとめて訓練サブコーパス 160 とする。機械翻訳規則自動構築部 32 は、訓練サブコーパス 160 から翻訳規則を自動構築する。構築された翻訳規則集合は翻訳規則集合記憶部 40 に記憶される。

【0128】

機械翻訳エンジン 42 は、翻訳規則集合記憶部 40 に記憶されている翻訳規則を用いて、評価サブコーパス 162 中の原文集合に対する翻訳を行なう。訳質自動評価部 44 は、機械翻訳エンジン 42 による翻訳結果の訳質を自動評価し、スコアとして規則寄与度算出部

50

46に与える。

【0129】

規則寄与度算出部46は、第1の実施の形態で説明した通り、翻訳規則集合記憶部40に記憶されている各規則について、規則寄与度を算出する。算出された規則寄与度は、規則寄与度記憶部202に規則ごと、繰返しごとに $contrib[i][r]$ として記憶される。

【0130】

上記した処理をN回繰返す事により、規則寄与度記憶部202には、規則寄与度 $contrib[i][r]$ ($1 \leq i \leq N$ 、 $1 \leq r \leq$ 基本翻訳規則数)が記憶される。

【0131】

翻訳規則集約部194は、基本規則集合記憶部196に記憶されている各規則について、前述した通り規則寄与度の総和 $contrib[r] = \sum_i contrib[i][r]$ を計算する。 $contrib[r]$ が負の場合、その規則は基本規則集合記憶部196内の基本規則集合から削除される。

【0132】

翻訳規則集約部194が、基本規則集合記憶部196に記憶されている全ての翻訳規則に対して以上の処理を実行する事により、最終的に基本規則集合記憶部196には、交差クリーニング後の基本規則集合が得られる。

【0133】

第2の実施の形態の効果

この第2の実施の形態の翻訳規則抽出装置180によって交差クリーニングした後の翻訳規則集合を用いて機械翻訳を行なったところ、第1の実施の形態により得られたものよりもさらによい結果が得られた。また、第1の実施の形態の翻訳規則抽出装置20では、訓練コーパスとは別に評価コーパスを準備する必要があった。それに対してこの第2の実施の形態の翻訳規則抽出装置180では、訓練コーパス140のみを使用し、それと別に評価コーパスを用意する必要はない。従って、翻訳規則のクリーニングが、限られた対訳コーパスを用いて行なえ、その結果得られた翻訳規則集合を用いて、精度の高い機械翻訳を行なう事が可能になる。

【0134】

コンピュータによる実現

この第2の実施の形態に係る翻訳規則抽出装置180も、図3及び図4に示すコンピュータと、その上で実行されるプログラムとにより実現可能である。図8に、この第2の実施の形態に係る翻訳規則抽出装置180を実現するためのプログラムの制御構造をフローチャート形式で示す。

【0135】

図8を参照して、このプログラムは、訓練コーパス140から基本規則集合を自動構築するステップ210と、訓練コーパス140を均等にN個のサブコーパスに分類するステップ212とを含む。これらN個のサブコーパスを $EC[i]$ ($1 \leq i \leq N$)とする。

【0136】

このプログラムはさらに、以下のステップ216からステップ220を、変数iを1からNまで1ずつ増加させながら繰返すステップを含む。まずステップ216では、訓練コーパス140からサブコーパス $EC[i]$ を取除き、訓練サブコーパス160を作成する。この訓練サブコーパスを $TC[i]$ とする。

【0137】

続いてステップ218で、訓練サブコーパス $TC[i]$ から翻訳規則集合 $R[i]$ を自動構築する。さらにステップ220で、サブコーパス $EC[i]$ を評価コーパスとみなして翻訳規則集合 $R[i]$ をフィードバッククリーニングする。このフィードバッククリーニング処理の内容は、図5に示した第1の実施の形態のものと同様である。ただしこの際、図5のステップ116で算出された規則寄与度 $contrib[r]$ を $contrib[i][r]$ として記憶しておく事に注意する必要がある。

10

20

30

40

50

【0138】

ステップ216からステップ220までの処理をN回繰返した後、今度は以下に説明するステップ226からステップ232の処理を、ステップ210で自動構築された基本規則集合内の全ての規則 r について繰返し行なう(1 r 基本規則集合内の規則数)。

【0139】

ステップ226では、翻訳規則集合 $R[i]$ (1 i N)から、規則 r の規則寄与度 $contrib[i][r]$ を取得する。具体的には、前述した通り図5のステップ116で記憶されていた規則寄与度を記憶領域から取出す。ステップ228で、基本規則 r の寄与度 $contrib[r] = \sum_i contrib[i][r]$ を算出する。

【0140】

続くステップ230では、ステップ228で算出された寄与度 $contrib[r]$ が負か否かを判定する。負であればステップ232でこの規則 r を基本規則集合から取除く。負でない場合には何もしない。

【0141】

以上のステップ226からステップ232までの処理を、基本規則集合内の全ての規則に対して行なう事により、最終的に交差フィードバッククリーニングが行なわれた翻訳規則が得られる事については前述した通りである。この交差クリーニングにより、第2の実施の形態の説明の冒頭で説明した様なクリーニング漏れを防止する事ができる。

【0142】

第2の実施の形態の変形例

上記した第2の実施の形態の装置では、機械翻訳規則自動構築部32とは別に機械翻訳規則自動構築部198を設けている。しかしこれらは必ずしも別個のものとする必要はない。同じ機械翻訳規則自動構築部を用いて、その入力及び出力の接続先を切替える様にしてもよい。

【0143】

また、上記した実施の形態の装置では訓練コーパス140をN個のサブコーパスにほぼ等分する事により、訓練サブコーパスと評価サブコーパスとを作成している。しかし本発明はその様な実施の形態に限定されるわけではない。例えば、訓練コーパス140を必ずしも等分する必要はない。実質的に大きさの異なったコーパスに分割し、後は上記した通りの処理を行なう様にしてもよい。ただしその場合には、翻訳規則集約部194で規則を集約する際の規則寄与度の総和計算において、コーパスの大きさに従った重みを各寄与度に乘じた後に加算する事が望ましい。

【0144】

共通の変形例

上記した二つの実施の形態では、機械翻訳エンジン42として参考文献1に記載されたものを使用している。しかし本発明はその様な実施の形態に限定されるわけではない。翻訳規則を用いた構文トランスファ方式の機械翻訳エンジンであればどの様なものを用いてもよい。

【0145】

さらに、上記した二つの実施の形態では、訳質自動評価部44による訳質の自動評価にBLEUを用いた。しかし訳質の自動評価にはBLEUのみが使用可能なわけではない。例えば、後掲の参考文献3又は参考文献4に記載のものを用いる事も可能である。

【0146】

自動評価値として、本実施の形態では評価コーパス内の訳文との類似度が高い場合に評価値が高くなるものを使用した。しかし自動評価値としてはその様なものには限定されず、類似度が高い場合に評価値が低くなる様なものでもよい。また、評価コーパス内の訳文との類似度が高くなるほど、特定の値に近くなる様な評価値を用いてもよい。

【0147】

なお、ソフトウェアの流通形態は上記した様に記憶媒体に固定された形には限定されない。たとえば、ネットワークを通じて接続された他のコンピュータからデータを受取る形で

10

20

30

40

50

流通する事もあり得る。また、ソフトウェアの一部が予めハードディスク５４中に格納されており、ソフトウェアの残りの部分をネットワーク経由でハードディスク５４に取込んで実行時に統合する様な形の流通形態もあり得る。

【０１４８】

一般的に、現代のプログラムはコンピュータのオペレーティングシステム（ＯＳ）によって提供される汎用の機能を利用し、それらを所望の目的に従って組織化した形態で実行する事により前記した所望の目的を達成する。従って、以下に述べる本実施の形態の各機能のうち、ＯＳ又はサードパーティが提供する汎用的な機能を含まず、それら汎用的な機能の実行順序の組合せだけを指定するプログラム（群）であっても、それらを利用して全体的として所望の目的を達成する制御構造を有するプログラム（群）である限り、それらが本発明の技術的範囲に含まれる事は明らかである。

10

【０１４９】

参考文献リスト

[参考文献１] 古瀬蔵、山本和英、及び山田節夫（１９９９）．構成素境界解析を用いた多言語話し言葉翻訳．自然言語処理、６（５）：６３－９１。

【０１５０】

[参考文献２] ペネニ、Ｋ．、ルーコス、Ｓ．、ウォード、Ｔ、及びツ－、Ｗ．－Ｊ．（２００２）．Ｂｌｅｕ：機械翻訳の自動評価方法．第４０回計算言語学学会第４０回年次大会予稿集、３１１頁から３１８頁（Ｐａｉｎｅｎｉ、Ｋ．、Ｒｏｕｋｏｓ、Ｓ．、Ｗａｒｄ、Ｔ．、and Zhu、Ｗ．－Ｊ．（２００２）．Ｂｌｅｕ：a method for automatic evaluation of machine translation．In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics（ACL）, pp．311 - 318）

20

【０１５１】

[参考文献３] ヤスダ、Ｋ．、スガヤ、Ｆ．、タケザワ、Ｔ．、ヤマモト、Ｓ．、及びヤナギダ、Ｍ．、（２００１）．パラレルコーパスから検索された翻訳候補を用いた翻訳品質の自動評価法、機械翻訳サミット予稿集ＶＩＩＩ、３７３頁から３７８頁（Ｙａｓｕｄａ、Ｋ．、Ｓｕｇａｙａ、Ｆ．、Ｔａｋｅｚａｗａ、Ｔ．、Ｙａｍａｍｏｔｏ、Ｓ．、and Yanagida、Ｍ．、（２００１）．An automatic evaluation method of translation quality using translation answer candidates queried from a parallel corpus．In Proceedings of Machine Translation Summit VＩＩＩ, pp．373 - 378）

30

【０１５２】

[参考文献４] アキバ、Ｙ．、イマムラ、Ｋ．、及びスミタ、Ｅ．、（２００１）（Ａｋｉｂａ、Ｙ．、Ｉｍａｍｕｒａ、Ｋ．、and Sumita、Ｅ．、（２００１）．複数編集距離を用いた機械翻訳の自動評価．機械翻訳サミット予稿集ＶＩＩＩ、１５頁から２０頁（Ｕｓｉｎｇ ｍｕｌｔｉｐｌｅ ｅｄｉｔ ｄｉｓｔａｎｃｅｓ ｔｏ ａｕｔｏｍａｔｉｃａｌｌｙ ｒａｎｋ ｍａｃｈｉｎｅ ｔｒａｎｓｌａｔｉｏｎ ｏｕｔｐｕｔ．In Proceedings of Machine Translation Summit VＩＩＩ, pp．15 - 20）

40

【０１５３】

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及び範囲内でのすべての変更を含む。

【図面の簡単な説明】

【図１】本発明の第１の実施の形態に係る翻訳規則抽出装置２０の機能的ブロック図であ

50

る。

【図 2】 翻訳規則の例を示す図である。

【図 3】 翻訳規則抽出装置 20 を実現するコンピュータの外観図である。

【図 4】 図 3 に示すコンピュータの回路構成を概略的に示す図である。

【図 5】 第 1 の実施の形態に係る翻訳規則抽出装置 20 をコンピュータで実現するためのプログラムの制御構造を示すフローチャートである。

【図 6】 本発明の第 2 の実施の形態における交差クリーニング法の概略を説明するための図である。

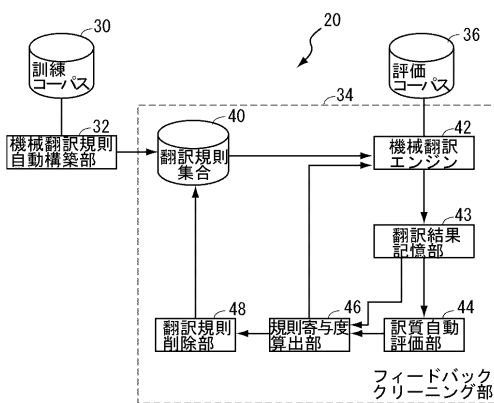
【図 7】 第 2 の実施の形態の翻訳規則抽出装置 180 の機能的ブロック図である。

【図 8】 翻訳規則抽出装置 180 を実現するためのプログラムの制御構造を示すフローチャートである。 10

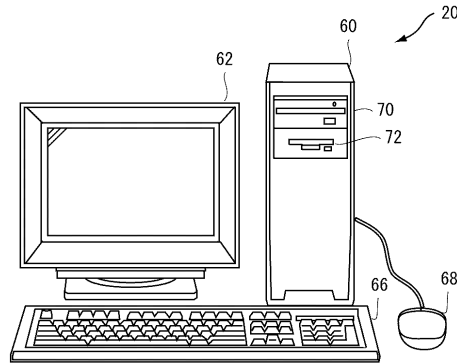
【符号の説明】

20, 180 翻訳規則抽出装置、30, 140 訓練コーパス、32, 198 機械翻訳規則自動構築部、34 フィードバッククリーニング部、36 評価コーパス、40 翻訳規則集合記憶部、42 機械翻訳エンジン、43 翻訳結果記憶部、44 訳質自動評価部、46 規則寄与度算出部、48 翻訳規則削除部、160 訓練サブコーパス、162 評価サブコーパス、190 訓練コーパス分割部、192 繰返制御部、194 翻訳規則集約部、196 基本規則集合記憶部、202 規則寄与度記憶部

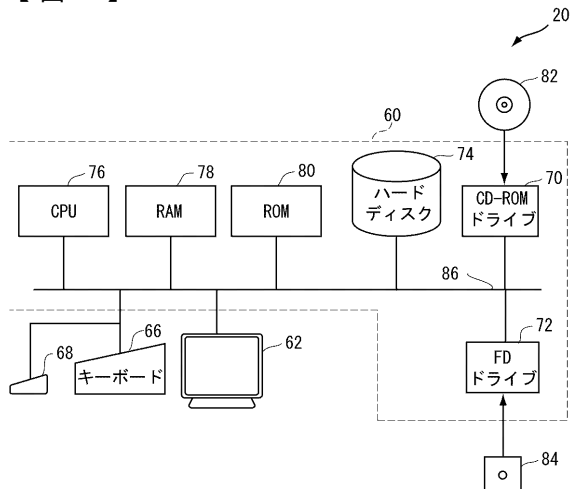
【図 1】



【図 3】



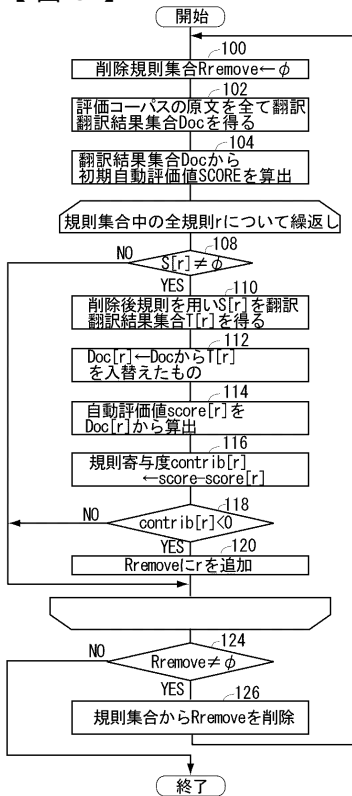
【図 4】



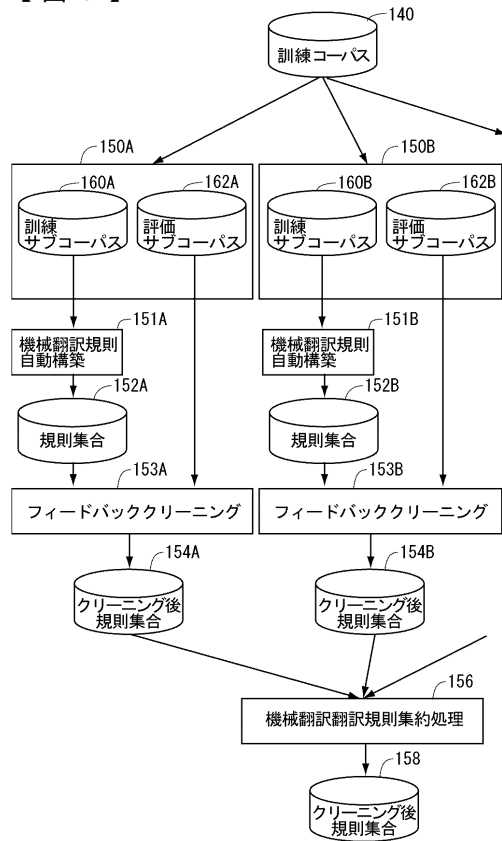
【図 2】

規則 番号	構文 カテゴリ	原言語 パターン	目的言語 パターン	用例
1	VP	Xvp at Ynp	⇒ Y' でX'	((present, conference), ...)
2	VP	Xvp at Ynp	⇒ Y' にX'	((stay, hotel), (arrive, p.m.), ...)
3	VP	Xvp at Ynp	⇒ Y' をX'	((look, it), ...)
4	NP	Xnp at Ynp	⇒ Y' のX'	((man, front desk), ...)
5	NP	the Xnmp Ynp station	⇒ X' Y'	((near, subway), ...)
6	VP	Xvp <num-noun> Ynp	⇒ Y' をX'	((rent, bicycle), ...)
7	S	please Xvp	⇒ X' たいのですが ((send), (reserve), ...)	
8	S	pelase Xvp	⇒ X' てください ((give), ...)	

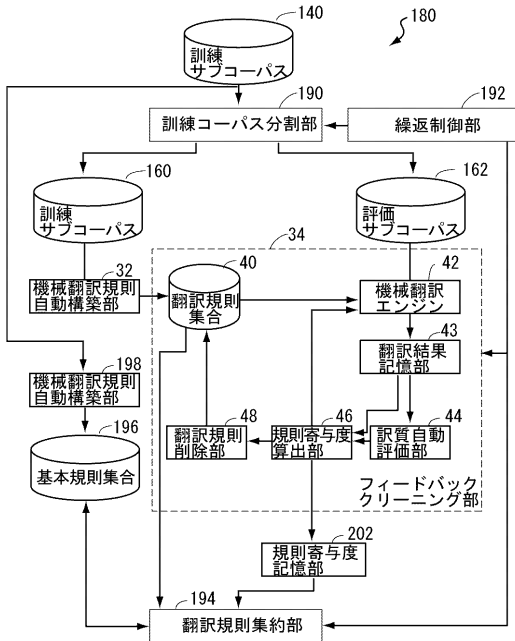
【図 5】



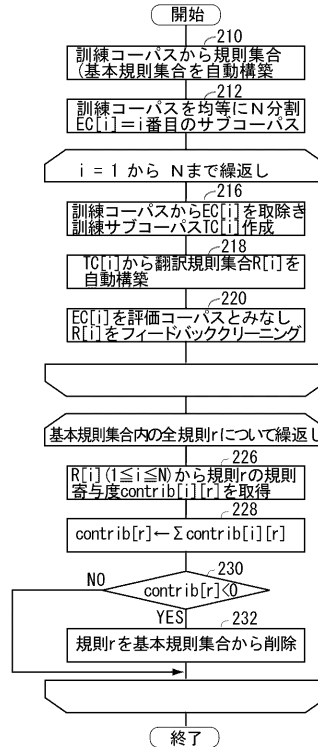
【図 6】



【図 7】



【図 8】



【手続補正書】

【提出日】平成16年6月24日(2004.6.24)

【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】0035

【補正方法】変更

【補正の内容】

【0035】

本発明の第3の局面に係るコンピュータは、上記したコンピュータプログラムによりプログラムされたコンピュータである。

【手続補正2】

【補正対象書類名】明細書

【補正対象項目名】0038

【補正方法】変更

【補正の内容】

【0038】

なお以下の説明では、第1及び第2の実施の形態を説明する。これらの実施の形態の基本的な考え方は以下の通りである。すなわち、自動構築された翻訳規則を用いて評価コーパス中の原言語の文を機械翻訳する。機械翻訳した結果に対し、後掲の参考文献2に記載されている様な訳質の自動評価を行ない、自動評価値を得る。この自動評価値を向上させる様に翻訳規則の取捨選択を行なう事により、最適な翻訳規則の組合せ(最適な翻訳規則集合)を得る。

【手続補正3】

【補正対象書類名】明細書

【補正対象項目名】0043

【補正方法】変更

【補正の内容】

【0043】

フィードバッククリーニング部34は、機械翻訳規則自動構築部32により訓練コーパス30から自動的に構築された翻訳規則の集合を記憶するための翻訳規則集合記憶部40と、翻訳規則集合記憶部40に記憶された翻訳規則を用いて評価コーパス36中の全ての英語の原文を目的言語の文に翻訳するための機械翻訳エンジン42とを含む。機械翻訳エンジン42は構文トランスファ方式のものである。

【手続補正4】

【補正対象書類名】明細書

【補正対象項目名】0044

【補正方法】変更

【補正の内容】

【0044】

フィードバッククリーニング部34はさらに、機械翻訳エンジン42による翻訳結果を、各文の翻訳の際に使用された翻訳規則を特定する情報とともに記憶するための翻訳結果記憶部43を含む。翻訳結果記憶部43はまた、翻訳結果ととともに各文の翻訳の際に使用された翻訳規則を特定する情報も記憶する。

【手続補正5】

【補正対象書類名】明細書

【補正対象項目名】0124

【補正方法】変更

【補正の内容】

【0124】

次に、繰返制御部192によるN回のフィードバッククリーニングにより、訓練コーパス

140のN個の訓練サブコーパス160よりN個の翻訳規則集合が得られる。これらを第1の翻訳規則集合、第2の翻訳規則集合、...第Nの翻訳規則集合と呼ぶ事とする。そして、これらの翻訳規則集合を作成する際に規則寄与度算出部46により計算された各規則の規則寄与度が規則寄与度記憶部202に繰返しごとに別々に記憶される。規則rについてのi回目の繰返しの際に計算された規則寄与度を $contrib[i][r]$ と表す($1 \leq i \leq N$ 、 $1 \leq r \leq$ 基本翻訳規則数)。

【手続補正6】

【補正対象書類名】明細書

【補正対象項目名】0125

【補正方法】変更

【補正の内容】

【0125】

翻訳規則集約部194は、全てのフィードバッククリーニングが終了すると、規則寄与度記憶部202を参照して、翻訳規則rごとに、規則寄与度記憶部202に記憶されている規則寄与度の総和 $contrib[r] = \sum_i contrib[i][r]$ を計算する。そして、総和 $contrib[r]$ が負であればその規則rを基本規則集合記憶部196に記憶されている基本規則集合から削除する。この処理を全ての規則rに対して実行する事により、基本規則集合記憶部196に記憶されている基本規則集合に対するクリーニングが行なわれ、最終的な交差フィードバッククリーニング後の翻訳規則集合が得られる。

【手続補正7】

【補正対象書類名】明細書

【補正対象項目名】0131

【補正方法】変更

【補正の内容】

【0131】

翻訳規則集約部194は、基本規則集合記憶部196に記憶されている各規則について、前述した通り規則寄与度の総和 $contrib[r] = \sum_i contrib[i][r]$ を計算する。 $contrib[r]$ が負の場合、その規則は基本規則集合記憶部196内の基本規則集合から削除される。

【手続補正8】

【補正対象書類名】明細書

【補正対象項目名】0139

【補正方法】変更

【補正の内容】

【0139】

ステップ226では、翻訳規則集合 $R[i]$ ($1 \leq i \leq N$)から、規則rの規則寄与度 $contrib[i][r]$ を取得する。具体的には、前述した通り図5のステップ116で記憶されていた規則寄与度を記憶領域から取出す。ステップ228で、基本規則rの寄与度 $contrib[r] = \sum_i contrib[i][r]$ を算出する。

【手続補正9】

【補正対象書類名】明細書

【補正対象項目名】0147

【補正方法】変更

【補正の内容】

【0147】

なお、ソフトウェアの流通形態は上記した様に記憶媒体に固定された形には限定されない。たとえば、ネットワークを通じて接続された他のコンピュータからデータを受取る形で流通する事もあり得る。また、ソフトウェアの一部が予めハードディスク54中に格納されており、ソフトウェアの残りの部分をネットワーク経由でハードディスク74に取込んで実行時に統合する様な形の流通形態もあり得る。

【手続補正 10】

【補正対象書類名】明細書

【補正対象項目名】0150

【補正方法】変更

【補正の内容】

【0150】

【参考文献 2】 パピネニ、K. , ルーコス、S. , ウォード、T. , 及びツー、W. - J. . (2002) . Bleu : 機械翻訳の自動評価方法 . 第40回計算言語学学会第40回年次大会予稿集、311頁から318頁 (Papineni , K. , Roukos , S. , Ward、T. , and Zhu , W. - J. . (2002) . Bleu : a method for automatic evaluation of machine translation . In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL) , pp . 311 - 318)

【手続補正 11】

【補正対象書類名】明細書

【補正対象項目名】0152

【補正方法】変更

【補正の内容】

【0152】

【参考文献 4】 アキバ、Y. , イマムラ、K. , 及びスミタ、E. , (2001) (Akiba , Y. , Imamura , K. , and Sumita , E. , (2001)) . 複数編集距離を用いた機械翻訳の自動評価 . 機械翻訳サミット予稿集VIII、15頁から20頁 (Using multiple edit distances to automatically rank machine translation output . In Proceedings of Machine Translation Summit VIII , pp . 15 - 20)

【手続補正 12】

【補正対象書類名】図面

【補正対象項目名】図2

【補正方法】変更

【補正の内容】

【図2】

規則 番号	構文 カテゴリ	原言語 パターン	目的言語 パターン	用例
1	VP	X _{VP} at Y _{NP}	⇒ Y' でX'	((present, conference), ...)
2	VP	X _{VP} at Y _{NP}	⇒ Y' にX'	((stay, hotel), (arrive, p.m.), ...)
3	VP	X _{VP} at Y _{NP}	⇒ Y' をX'	((look, it), ...)
4	NP	X _{NP} at Y _{NP}	⇒ Y' のX'	((man, front desk), ...)
5	NP	the X _{NMP} Y _{NP} station	⇒ X' Y'	((near, subway), ...)
6	VP	X _{VP} <num-noun> Y _{NP}	⇒ Y' をX'	((rent, bicycle), ...)
7	S	please X _{VP}	⇒ X' たいのですが	((send), (reserve), ...)
8	S	please X _{VP}	⇒ X' てください	((give), ...)

【手続補正 13】

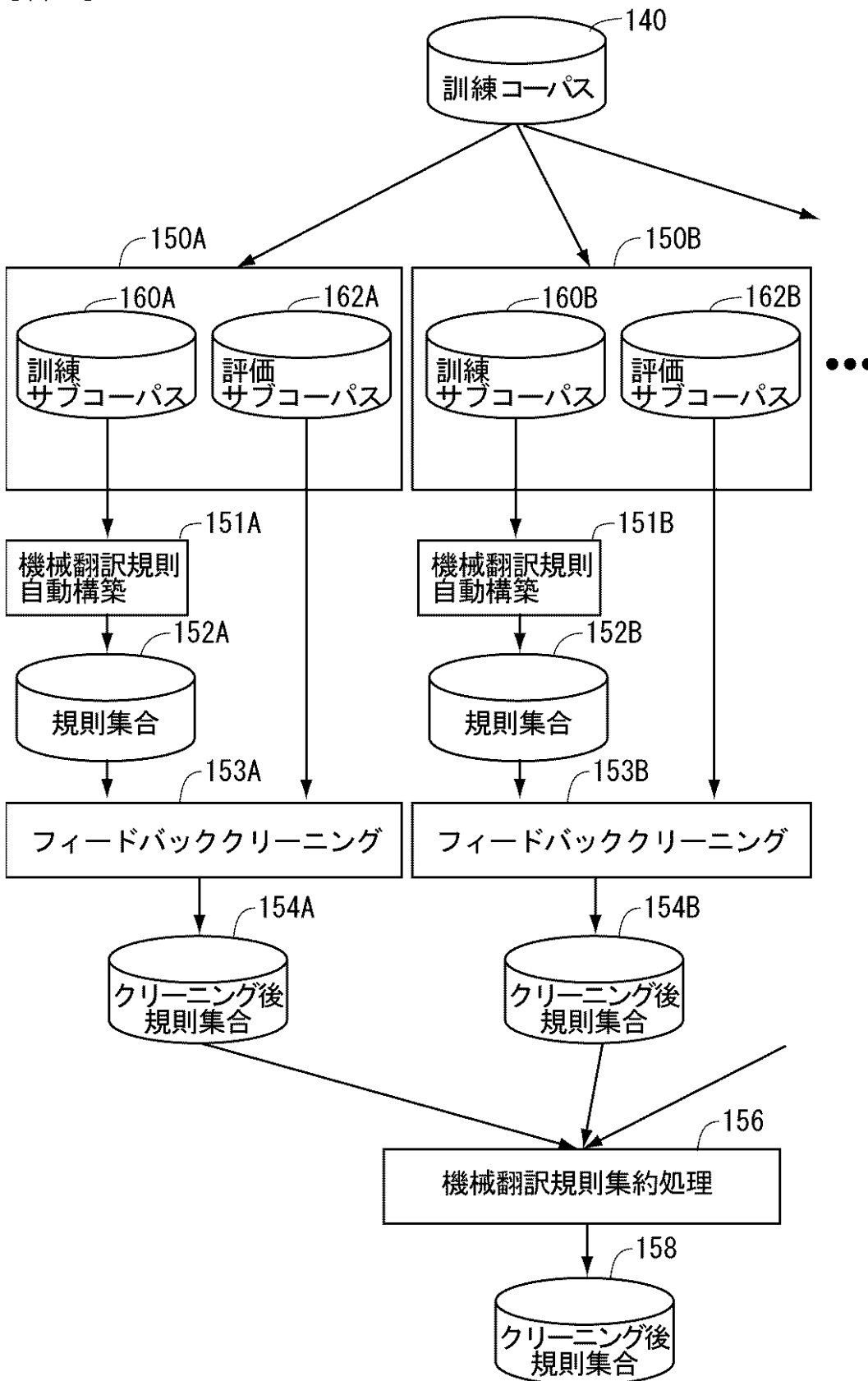
【補正対象書類名】図面

【補正対象項目名】図6

【補正方法】変更

【補正の内容】

【図 6】



【手続補正 1 4】

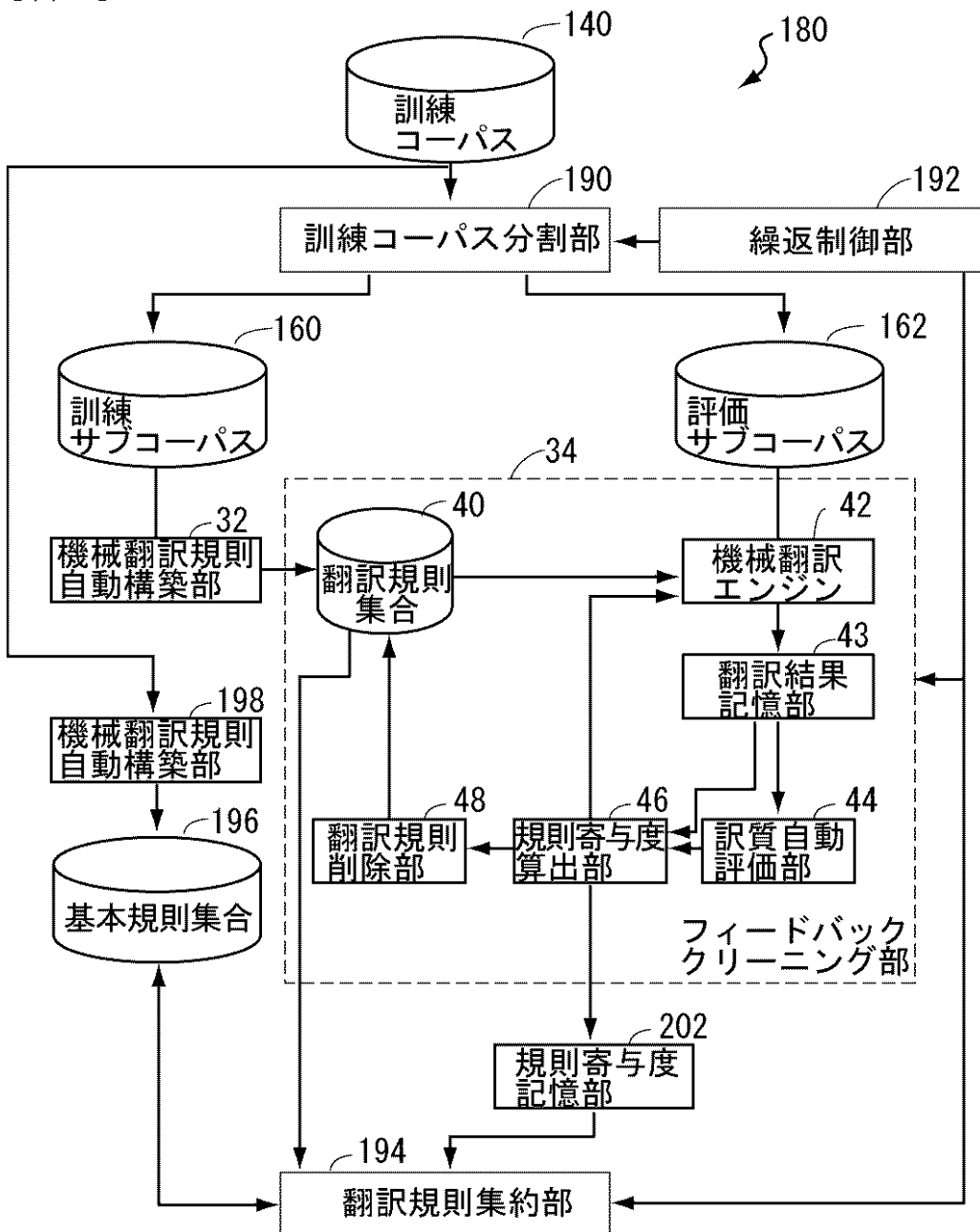
【補正対象書類名】図面

【補正対象項目名】図 7

【補正方法】変更

【補正の内容】

【図 7】



【手続補正 15】

【補正対象書類名】図面

【補正対象項目名】図 8

【補正方法】変更

【補正の内容】

【 図 8 】

