

(19)日本国特許庁 (J P)

(12)特許公報 (B 2)

(11)特許番号

特許第3388393号

(P 3 3 8 8 3 9 3)

(45)発行日 平成15年 3月17日(2003.3.17)

(24)登録日 平成15年 1月17日(2003.1.17)

(51)Int.Cl.⁷
G06F 17/27

識別記号

F I
G06F 17/27

M

請求項の数 2 (全 7 頁)

(21)出願番号 特願平11 - 238579

(22)出願日 平成11年 8月25日(1999.8.25)

(65)公開番号 特開2001 - 67357(P 2001 - 67357 A)

(43)公開日 平成13年 3月16日(2001.3.16)

審査請求日 平成11年 8月25日(1999.8.25)

特許法第30条第 1 項適用申請有り 村田真樹・馬青・内元清貴・井佐原均「用例ベースによるモダリティの日英翻訳」情報処理学会研究報告99 - N L - 130 - 16 , V o l . 99 , N o . 22 , p . 121 - p . 128(1999.03.05)

(73)特許権者 301022471
独立行政法人通信総合研究所
東京都小金井市貫井北町 4 - 2 - 1

(72)発明者 村田 真樹
兵庫県神戸市西区岩岡町岩岡558 - 2
郵政省通信総合研究所関西支所内

(72)発明者 内元 清貴
兵庫県神戸市西区岩岡町岩岡558 - 2
郵政省通信総合研究所関西支所内

(72)発明者 馬 青
兵庫県神戸市西区岩岡町岩岡558 - 2
郵政省通信総合研究所関西支所内

(74)代理人 100082669
弁理士 福田 賢三 (外 2 名)

審査官 和田 財太

最終頁に続く

(54)【発明の名称】データベースを利用したテンス、アスペクトあるいはモダリティに関する翻訳装置

1

(57)【特許請求の範囲】

【請求項 1】 第一の言語から第二の言語へのデータベースを利用した翻訳装置で、第一の言語に属する複数の用例と第二の言語に属する複数の用例からなり且つ個々の第一の言語に属する用例は、第二の言語に属する用例との間に少なくとも一つ以上の対応付があり、テンス、アスペクトあるいはモダリティ情報が付加されたことを特徴とする第一のデータベースを備え、

第一の言語に属する第一の用例と、第一のデータベースの第一の言語に属する第二の用例との間の文末からみて連続する共通の文字列の数をを用いて第一のデータベースの第一の言語に属する第二の用例との間の類似性を評価した値を導く手段を備え、

前記の手段は、(1) 該類似性を評価した値が高いほど類似性が高いとしてその類似性の高い順で、第一の言語

2

に属する第一の用例に対する第一の言語に属する第二の用例群を第一のデータベースから予め決められた数だけ選択するという第一の方法で選択し、(2) 第一の言語に属する用例の、第二の言語に属する用例への対応から、第二の言語に属する第一の用例群を第一のデータベースから選択し、(3) この選択された第二の言語に属する第一の用例群を代表するテンス、アスペクトあるいはモダリティについて、そのテンス、アスペクトあるいはモダリティを個々の用例のテンス、アスペクトあるいはモダリティの多数決で決定するという第二の方法で決定し、(4) この決定されたテンス、アスペクトあるいはモダリティを、第一の言語に属する第二の用例の翻訳のテンス、アスペクトあるいはモダリティとして用いる、という構成を備えることを特徴とする、データベースを利用したテンス、アスペクトあるいはモダリティに

10

関する翻訳装置。

【請求項 2】 第一の言語から第二の言語へのデータベースを利用した翻訳装置で、第一の言語に属する複数の用例と第二の言語に属する複数の用例からなり且つ個々の第一の言語に属する用例は、第二の言語に属する用例との間に少なくとも一つ以上の対応付があり、テンス、アスペクトあるいはモダリティ情報が付加されたことを特徴とする第一のデータベースを備え、

第一の言語に属する第一の用例と、形態素解析を行なって形態素を認識し、各形態素についてシソーラスの分類番号を付して、シソーラスの分類番号による構成に変換された該第一の言語に属する第一の用例と、同様に変換された第一のデータベース内の第一の言語に属する第二の用例との間の、文末からみて連続する共通の文字列の数をを用いることを特徴とする第一の言語に属する第二の用例との間の類似性を評価した値を導く手段を備え、前記の手段は、(1) 該類似性を評価した値が高いほど類似性が高いとしてその類似性の高い順で、第一の言語に属する第一の用例に対する第一の言語に属する第二の用例群を第一のデータベースから予め決められた数だけ選択するという第一の方法で選択し、(2) 第一の言語に属する用例の、第二の言語に属する用例への対応から、第二の言語に属する第一の用例群を第一のデータベースから選択し、(3) この選択された第二の言語に属する第一の用例群を代表するテンス、アスペクトあるいはモダリティについて、そのテンス、アスペクトあるいはモダリティを個々の用例のテンス、アスペクトあるいはモダリティの多数決で決定するという第二の方法で決定し、(4) この決定されたテンス、アスペクトあるいはモダリティを、第一の言語に属する第二の用例の翻訳のテンス、アスペクトあるいはモダリティとして用いる構成を備えることを特徴とする、データベースを利用したテンス、アスペクトあるいはモダリティに関する翻訳装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 この発明は、データベースを利用した翻訳装置に関するものである。特に、時制などが文末に表現される言語、例えば日本語、から他の言語、例えば英語、に翻訳する際に問題となる文末表現のテンス(時制)、アスペクト(相)あるいはモダリティ(様相)を翻訳する時に用いる。

【0002】

【従来の技術】 従来の技術の例を、データベースを用いた翻訳装置における翻訳方法の従来例を第一の従来技術に、用例を基にしたデータベースを用いた翻訳装置での翻訳方法の従来例を第二の従来技術に、用例間の類似性を見る方法として文末から数えた一致文字列の数をを用いた従来例を第三の従来技術として以下に示す。

【0003】 まず、データベースを用いた翻訳装置の翻

訳方法の第一の従来技術を図 2 のフローチャートに示す。図 2 のフローチャートでは、次の四段階の手続を示している。

【0004】 従来の文末表現の日英翻訳は、人手で作成した規則によってなされてきた。このため、まず、次の作業を行う必要があった。

【0005】 1) 解析以前に予め、人手による規則集の作成をする。例えば、連用形+動詞「いる」ならば、アスペクトが「進行相」となる。このような規則を、他の組み合わせに対しても作成し、規則集を作成する。また、テンス(時制)、アスペクト(相)あるいはモダリティ(様相)についても、その規則集を作成する。この様な人手により規則集を作成する場合は、規則の不備が残ってしまい、常にメンテナンスを続けて洗練化する必要がある。

【0006】 次に解析作業として、
2) 解析における手続 1 入力文の翻訳のための入力文の形態素解析や構文解析を行う。例えば「希望をいただいている。」が入力文の場合は、下記のような結果を得る。
希望 <名詞> を <助詞> いただいて <動詞> <連用形> いる <動詞>

【0007】 ここで、形態素解析部や構文解析部を変更すると、上記の規則集にも影響があり、適切な翻訳を維持するためには、上記の規則集にも変更すべき点が発生してしまう。

【0008】 3) 解析における手続 2 形態素解析や構文解析の結果と、規則を照合して、テンス(時制)、アスペクト(相)あるいはモダリティ(様相)を確定する。上記の場合、文末表現が、<連用形>+動詞「いる」の形になっているので、予め作成した規則により、「進行相」と確定される。

【0009】 続いて、次の様に文全体を構成するため、合成作業を行う。

4) 解析における手続 3 テンス、アスペクトあるいはモダリティの翻訳以外の部分は、従来の既によく知られた翻訳方法のどの方法を用いてもよく、それらのどれかを用いて翻訳し、テンス、アスペクトあるいはモダリティの翻訳は、上記の方法により翻訳し、これらを合成することにより、文全体の翻訳を完成する。

【0010】 ここで示した人手で作成した規則によって翻訳する方法では、規則のメンテナンスに多大な人的資源を投入する必要があるという欠点がある。

【0011】 次に用例を基にしたデータベースを用いた翻訳装置での翻訳方法の先行例である第二の従来技術を示す。

【0012】 本発明の方法でも、用例を集めたデータベースを利用しており、用例ベース手法に分類される。この用例を基にした手法を日英翻訳に利用した従来例としては、報告書(Eiichiro Sumita, Hitoshi Iida, and Hideo Kohya

ma、Translating examples: A new approach to machine translation, The third International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Language, No.3, (TMI, 1990), pp.203-212)がある。ただし、この技術は課題が「AのB」であって、「AのB」の日英翻訳では名詞Aと名詞Bの意味情報を複雑に組み合わせて利用していた。この研究と本発明とは以下の点で異なっている。

【0013】1) 請求項1の第一のデータベースに相当する用いるデータベースの構成、
2) 用例間の類似性を評価する方法、
3) ひとつの用例の中で注目する位置。

【0014】最後に、第三の従来技術として、用例間の類似性を見る方法として文末から数えた一致文字列の数をを用いた従来例を示す。

【0015】本発明でも、文末から数えた一致文字列の数を、類似性を評価した値(類似度)として用いているが、この方法自体には従来例がある。文末の省略表現の補完を行なう研究(村田真樹、長尾真、日本語文章における表層表現と用例を用いた動詞の省略の補完、言語処理学会誌、Vol.5, No.1, (1998))があり、文末一致文字列の文字数を類似度とする用例ベース利用を利用して行っている。

【0016】この従来例と本発明とは、以下の点において異なっている。

1) 従来例では、対象とする問題が省略表現の補完であり、本発明の様に、異なる言語間の翻訳ではない。
2) 請求項1の第一のデータベースに相当する用いるデータベースの構成が言語及び項目において相異している。

【0017】以上の従来技術は、本発明の方法とは部分的に一致しているに過ぎず、これらの先行技術を単に組み合わせても、本発明の翻訳方法を容易に思いつくものではない事は明らかである。

【0018】

【発明が解決しようとする課題】従来のデータベースを利用した翻訳装置における翻訳方法では、従来の文末表現の翻訳は、人手で作成した規則によってなされてきた。しかし、人手で作成した規則によって翻訳する方法では、翻訳精度を向上させるために行う規則のメンテナンスに多大な人的資源を投入する必要があるという欠点があった。

【0019】本発明は上記に鑑み提案されたもので、人手による規則集の作成をする必要が無くデータベースを利用した翻訳の知識が無くても翻訳精度の向上を図ることができるデータベースを利用した翻訳装置を提供する

ことを目的とする。

【0020】

【課題を解決するための手段】上記目的を達成するために本発明で用いる手段を、フローチャートで説明すると図1の様になる。また、以下に本発明による方法を簡潔に記述する。

【0021】1) 解析以前に予め、第一の言語の用例と、それに対応する第二の言語の用例を集めたデータベースを作成する。また、この際に第二の言語の用例のテンス、アスペクトあるいはモダリティの分類を付与する。この付与は、人手で行っても良いし、既によく知られている形態素および構文解析システムを補助として用いることも出来る。

【0022】2) 解析における手続1入力文の翻訳のための検索で、文末からの一致文字列が最も長い用例を、上記のデータベースから検索する。検索方法は、よく知られた2分検索の方法を使うことが出来る。

【0023】3) 解析における手続2翻訳の確定で、手続1から取り出した用例の英訳側の動詞部分のテンス、アスペクトあるいはモダリティの分類を、入力文のテンス、アスペクトあるいはモダリティと確定する。

【0024】4) 解析における手続3翻訳文の構成方法で、テンス、アスペクトあるいはモダリティの翻訳以外の部分は、従来の翻訳方法のどの方法を用いてもよく、それらのどれかを用いて翻訳し、テンス、アスペクトあるいはモダリティの翻訳は、上記の方法により翻訳し、これらを合成することにより、文全体の翻訳を完成する。

【0025】従って、上記目的を達成するために、請求項1に記載の発明は、第一の言語から第二の言語へのデータベースを利用したテンス、アスペクトあるいはモダリティに関する翻訳装置で、第一の言語に属する複数の用例と第二の言語に属する複数の用例からなり且つ個々の第一の言語に属する用例は、第二の言語に属する用例との間に少なくとも一つ以上の対応付があり、テンス、アスペクトあるいはモダリティ情報が付加されたことを特徴とする第一のデータベースを備え、第一の言語に属する第一の用例と、第一のデータベースの第一の言語に属する第二の用例との間の文末からみて連続する共通の文字列の数をを用いて第一のデータベースの第一の言語に属する第二の用例との間の類似性を評価した値を導く手段を備え、前記の手段は、(1)該類似性を評価した値が高いほど類似性が高いとしてその類似性の高い順で、第一の言語に属する第一の用例に対する第一の言語に属する第二の用例群を第一のデータベースから予め決められた数だけ選択するという第一の方法で選択し、(2)第一の言語に属する用例の、第二の言語に属する用例への対応から、第二の言語に属する第一の用例群を第一のデータベースから選択し、(3)この選択された第二の言語に属する第一の用例群を代表するテンス、アスペク

トあるいはモダリティについて、そのテンス、アスペクトあるいはモダリティを個々の用例のテンス、アスペクトあるいはモダリティの多数決で決定するという第二の方法で決定し、(4) この決定されたテンス、アスペクトあるいはモダリティを、第一の言語に属する第二の用例の翻訳のテンス、アスペクトあるいはモダリティとして用いる、という構成を備えることを特徴としており、用例を基にした方法でテンス、アスペクトあるいはモダリティを適切に翻訳するものを提案している。

【 0 0 2 6 】また、請求項 2 に記載の発明は、データベースを利用したテンス、アスペクトあるいはモダリティに関する翻訳装置であり、その類似性を評価する構成の特徴は、形態素解析を行なって形態素を認識し、各形態素についてシソーラスの分類番号を付して、シソーラスの分類番号による構成に変換された該第一の言語に属する第一の用例と、同様に変換された第一のデータベース

内の第一の言語に属する第二の用例との間に、文末からみて連続する共通の文字列の数をを用いることを特徴とする、第一の言語に属する第二の用例との間の類似性を評価した値を導く手段を備えている。

【 0 0 2 7 】

【発明の実施の形態】以下にこの発明の実施の形態を詳細に説明する。先ず第 1 の実施形態を、表 1 を用いて説明する。

【 0 0 2 8 】今「彼は私の知り合いだ」の時制を翻訳することを考える。このとき日英の翻訳対を大量に集めたデータベースに対して「彼は私の知り合いだ」と文末からの文字列一致が多い用例を上位から 1 0 個集めたものが表 1 のものだったとする。

【 0 0 2 9 】

【表 1】

| | | 日 本 語 文 | 分 類 | 英 語 文 |
|-----|-----|------------------|----------|---|
| 入力文 | | 彼は私の知り合いだ | 現在 | I am acquainted with him. |
| 番号 | 類似度 | 用 例 デ ー タ | | |
| 1 | 6 | 彼とは長年の知り合いだ | 現在 完了 | I have known him for a long time. |
| 2 | 6 | ふたりは長い間の知り合いだ | 現在 | The two are acquaintances of long standing. |
| 3 | 1 | 彼とは 1 0 年余の顔見知りだ | 現在 完了 | I have known him for over ten years. |
| 4 | 1 | 彼らは多年の知己だ | 現在 | They are friends of many years' standing. |
| 5 | 1 | 彼はこのクラブの恩人だ | 現在 | He is a benefactor of this club. |
| 6 | 1 | 彼は私の命の恩人だ | 現在 | I owe him my life. |
| 7 | 1 | 彼はかたい人だ | 現在 | He is reliable. |
| 8 | 1 | 彼はだれにも人当たりのいい人だ | 現在 | He is affable to everybody. |
| 9 | 1 | なんと男振りのいい人だ | 現在 | He is a gentlemanly person. |
| 9 | 1 | なんと男振りのいい人だ | 現在 | What a handsome looking man he is! |

【 0 0 3 0 】表の類似度は文末からの一致文字列の数を示している。また、ここで k 近傍法を用いる。k 近傍法とは 1 個の最も類似した用例を用いるかわりに、類似度の上位から順に取り出した k 個の用例の多数決により求める方法である。

【 0 0 3 1 】類似度が等しい用例がある場合は k の値に関わらず類似度が等しい用例はすべて用いて多数決を行

なう必要がある。さらに、ここでは処理の簡単のため、用例は多くても 1 0 個しか調べないこととする。

【 0 0 3 2 】また、上記の表 1 のうち、分類の欄は英語文の該当する動詞句より求めるものであるがこの部分はよく知られた処理プログラムを用いて自動で行なっても良いし、データベースを作成する際に人手であらかじめ分類を記入しておいてもよい。

【0033】まず、 $k = 1$ の場合を考える。このとき最も類似度の大きい 1 個の用例を用いて解析するわけだが、ここでは 1 番と 2 番が同じ類似度のため、1 と 2 番の用例を用いて解析を行なう。これで多数決を行なうと分類は「現在完了」が 1、「現在」が 1 と意見がわかれ、意見がわかれたときには先に上がった分類を解とすると決めておくと、解は先に上がった「現在完了」となり、不正解となる。

【0034】次に、 $k = 3$ の場合を考える。このとき最も類似度の大きい 3 個を選ぶわけだが、3 番の用例以降はすべて類似度が等しいので、10 個すべての用例を用いることになる。これで多数決を行なうと分類は「現在完了」が 2、「現在」が 8 と意見は分かれるが、数の大きい「現在」となり、これは正解の「現在」と一致し正解となる。

【0035】次に、 $k = 5, 7, 9$ の場合も同様に 10 個の用例すべてが用いられ解は「現在」となり、これも

正解となる。

【0036】この問題ではシステムは、 $k = 1$ のとき、誤った解を出力し、 $k = 3, 5, 7, 9$ のときに正しい解を出力するということになる。 k の値については装置を実際に作成する時に適切なものを選択するとよい。この方法による k の値は通常、多数決の都合上、奇数が望ましく、さらに 3 あるいは 5 で十分な場合が多い。データベースの用例が増えるに従って、より小さい k の値を用いる事ができる。

10 【0037】次に第 2 の実施形態を、表 2 を用いて説明する。今「彼は私の知り合いだ」の時制を翻訳することを考える。このとき日英の翻訳対を大量に集めたデータベースに対して「彼は私の知り合いだ」と文末からの文字列一致が多い用例を上位から 10 個集めたものが次の表 2 のものだったとする。

【0038】

【表 2】

| | | 日 本 語 文 | 分類 | 英 語 文 |
|-----|-----|-----------------|------|---|
| 入力文 | | 彼は私の知り合いだ | 現在 | I am acquainted with him. |
| 番号 | 類似度 | 用 例 デ ー タ | | |
| 1 | 25 | 彼とは長年の知り合いだ | 現在完了 | I have known him for a long time. |
| 2 | 24 | ふたりは長い間の知り合いだ | 現在 | The two are acquaintances of long standing. |
| 3 | 11 | 彼とは 10 年余の顔見知りだ | 現在完了 | I have known him for over ten years. |
| 4 | 11 | 彼らは多年の知己だ | 現在 | They are friends of many years' standing. |
| 5 | 10 | 彼はこのクラブの恩人だ | 現在 | He is a benefactor of this club. |
| 6 | 10 | 彼は私の命の恩人だ | 現在 | I owe him my life. |
| 7 | 10 | 彼はかたい人だ | 現在 | He is reliable. |
| 8 | 10 | 彼はだれにも人当たりのいい人だ | 現在 | He is affable to everybody. |
| 9 | 10 | なんと男振りのいい人だ | 現在 | He is a gentlemanly person. |
| 9 | 10 | なんと男振りのいい人だ | 現在 | What a handsome looking man he is! |

【0039】表 2 の類似度は、入力文の形態素解析を行なって形態素を認識し、各形態素についてシソーラスの分類番号を付して、シソーラスの分類番号による構成に変換された入力文を用意し、また日英の翻訳対を大量に集めたデータベースに対しても同様な変換を行ったものを用意し、これらの変換された後の文について、文末からみて連続する共通の文字列の数を示している。

【0040】解析は、第 1 の実施形態と同じく k 近傍法を用いることにする。

【0041】まず、 $k = 1$ の場合を考える。このとき最も類似度の大きい 1 番の用例だけを用いて解析を行なう。1 番の用例は分類が「現在完了」なので正解の分類「現在」と異なり、不正解となる。

50 【0042】次に、 $k = 3$ の場合を考える。このとき最

も類似度の大きい3個を選ぶわけだが、3番の用例と4番の用例の類似度が等しいので、4番の用例までの四つの用例を用いることになる。これで多数決を行なうと分類は「現在完了」が2、「現在」が2と意見がわかれ、解は先に上がった「現在完了」となり、これもまた不正解となる。

【0043】次に、k = 5の場合を考える。このとき最も類似度の大きい5個を選ぶわけだが、5番の用例以降はすべて類似度が等しいので、10個すべての用例を用いることになる。これで多数決を行なうと分類は「現在完了」が2、「現在」が8と意見はわかれるが、数の大きい「現在」となり、これは正解の「現在」と一致し正解となる。

【0044】次に、k = 7、9の場合も同様に10個の用例すべてが用いられ解は「現在」となり、これも正解となる。

【0045】この問題ではシステムは、k = 1、3のとき、誤った解を出力し、k = 5、7、9のときに正しい解を出力するということになる。kの値については装置を実際に作成する時に適切なものを選択するとよい。この方法によるkの値は通常、多数決の都合上、奇数が望ましく、さらに7あるいは9で十分な場合が多い。この場合も、データベースの用例が増えるに従って、より小さいkの値を用いる事ができる。

【0046】上記の実施形態に示されるように、本発明

の方法では、用例を集めたデータベースを整備して行くことによって翻訳精度の向上を図ることができ、従って、人手による規則集の作成をする必要が無くメンテナンスが容易であり、また、データベースを利用した翻訳の知識が無くても翻訳精度の向上を図ることができる。

【0047】

【発明の効果】この発明は上記した構成からなるので、以下に説明するような効果を奏することができる。

【0048】請求項1に記載の発明では、用例を基にした翻訳が可能となり、人手による規則集の作成をする必要が無くメンテナンスが容易であり、また、文末からみて、連続する共通の文字列の数であることとすることにより、簡単に類似性を評価することが出来、データベースを利用した翻訳の知識が無くても翻訳精度の向上を図ることができるようになった。

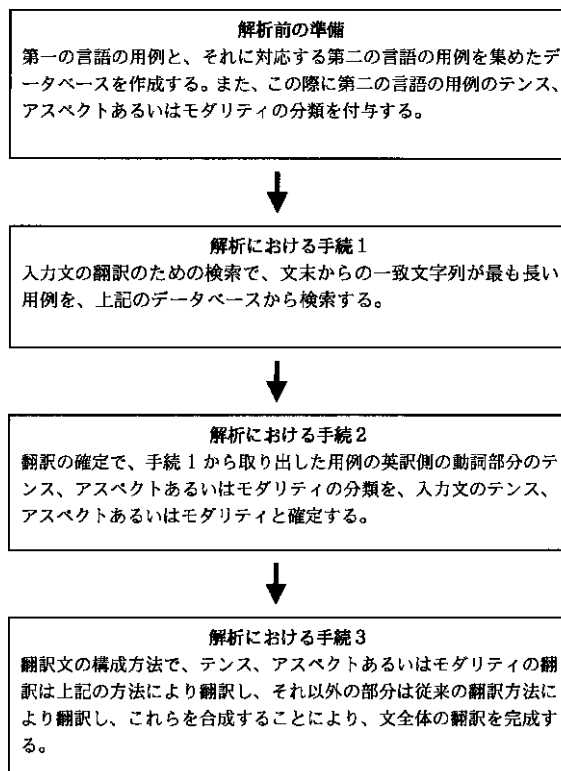
【0049】さらに、請求項2に記載の発明では、意味上の類似性を用いて類似性を評価することが出来るようになり、意味上からも適切な翻訳ができるようになった。

【図面の簡単な説明】

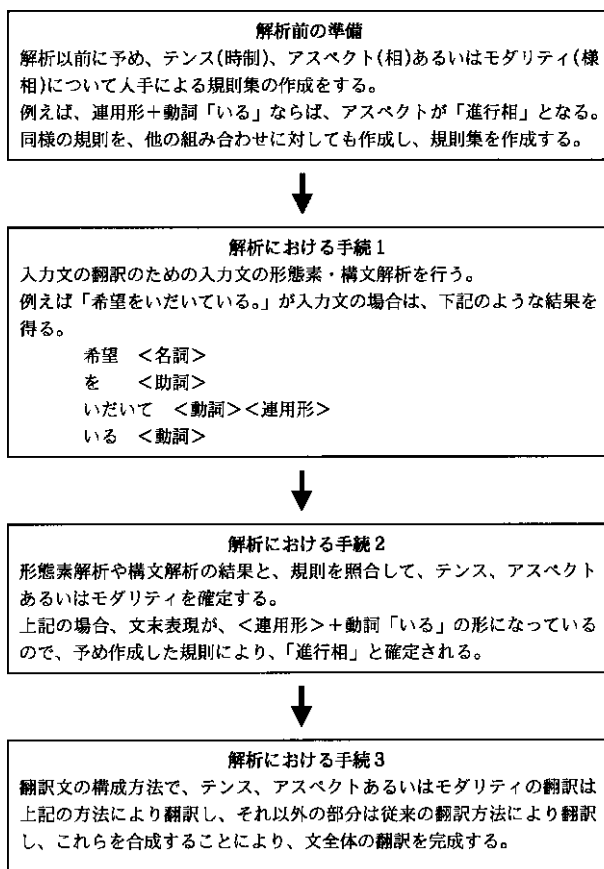
【図1】本発明における課題を解決するための手段を示すためのフローチャートである。

【図2】従来の技術における手段を示すためのフローチャートである。

【図1】



【 図 2 】



フロントページの続き

(72)発明者 井佐原 均
兵庫県神戸市西区岩岡町岩岡558 - 2
郵政省通信総合研究所関西支所内

(56)参考文献 特開 平 6 - 309352 (J P , A)

(58)調査した分野(Int.Cl.⁷, D B 名)
G06F 17/21 - 17/28
J I C S T ファイル (J O I S)