

非線形な表現構造に着目した日英文型パターン化

池原 悟* 阿部さつき* 徳久 雅人* 村上 仁一*

* 鳥取大学工学部 〒680-0945鳥取市湖山町南4.101 E-mail: {ikehara,tokuhisa,murakami}@ike.tottori-u.ac.jp
** NTTアドバンステクノロジー株式会社 〒210-0007 川崎市川崎区駅前本町12-1
川崎駅前タワービル satsuki@nlp.ntt-at.co.jp

あらまし 要素合成法を基本とした従来の機械翻訳方式の限界を突破する方法として、非線形な言語表現の構造を意味のまとまる単位にパターン化した文型パターン翻訳方式が期待される。本論文では、重文と複文を対象に、この方式の実現に必要な文型パターン辞書を試作した。具体的には、100万件の日英対訳コーパスから2つ又は3つの述部を持つ重文と複文合計15万件を抽出し、単語レベル(2.8万件)、句レベル(10.5万件)、節レベル(1.3万件)の3種類のグループからなる文型パターン辞書(合計24.6万件、異なり22.1万件収録)を作成した。各文型パターンは、いずれも形態素解析によって得られる文法情報を用いて記述することとし、対訳標本文に含まれる線形な表現要素を半自動的に変数化、関数化することなどにより作成したものである。従来、大規模な文型パターン辞書の開発は、文型パターン間の意味的排他性実現の困難性と膨大な開発コストが問題となるため、適用対象を限定するなど小規模な実現例しか見られなかった。しかし、今回の試作によって、ほぼすべての標本文(99%)が多くの線形要素(平均4~5カ所)を持つことが分かった。また、それらの要素を半自動的に関数化、変数化を行うことにより、文型パターンの開発コストは人手に頼る方法の約1/10に削減できた。これにより、実験的検討に必要な規模の文型パターン辞書を構築することができた。

キーワード: 機械翻訳, 文型パターン, 言語知識ベース, 要素合成法, 非線形要素, 汎化

Japanese to English Sentence Pattern Generations for Semantically Non-Linear Complex Sentences

Satoru Ikehara* Satsuki Abe** Masato Tokuhisa* Jin'ichi Murakami**

* Faculty of Engineering, Tottori University, Minami 4-101, Tottori City, Japan680-0945

**NTTAdvancedTechnology Corp., Ekimaehonchou, Kawasaki-ku, Kawasaki-shi, Japan210-0007 satsuki@nlp.ntt-at.co.jp

Abstract: In order to breakthrough the limitation of the conventional method based on Compositional Semantics, it is expected to realize a new translation method based on Sentence Patterns in which non-linear structures of linguistic expressions are represented as semantic units. This paper proposes the way to judge the linearity or non-linearity of linguistic expressions based on their definitions and how to generate sentence patterns from huge bilingual corpora. According to this method, three kinds of sentence patterns such as "word level", "phrase level" and "clause level" are generated in this order from Japanese to English corpus. In the experiments, 150,000 sentence pairs for complex and compound sentences are extracted from one million sentence pair corpora, and 128,000 patterns, 105,000 patterns and 13,000 patterns for each of three levels were generated from these sentence pairs. Due to the clarifications of decision process, the generation processes of the sentence patterns were mostly automated by using the results of morphological analysis and these 246,000 sentence patterns have been obtained in a year.

Key Words: Machine Translation, Sentence Pattern, Linguistic Knowledgebase, Compositional Semantics, Nonlinear Expression, Generalization

1. はしめに

従来、研究開発されてきた機械翻訳システムは、ほぼいずれもトランスファー方式を基本としている。この方式は、「原文の構文構造を目的言語の構造に変換する過程」と原文の各要素を翻訳する過程」を持ち、訳文は両者の結果を合成することによって得られる点に特徴がある^(1,2)。これは、構文構造と表現の意味を別々に変換するものであり、表現の構造と意味の関係が線形であることを前提とした要素合成方式が基本となっている。しかし、現実の言語表現には非線形なものも多く、表現が構成要素に分解される過程で全体の意味が次第に失われ、目的言語を生成する過程で復元できなくなることが問題であった⁽³⁾。

この問題を解決するには、「文構造とその意味を一体的に扱う仕組み」を実現することが重要である。文構造とその意味を一体的に扱う仕組みとしては、古くから「文型パターン翻訳」の方法が試みられてきた。文型パターン翻訳は「テンプレート翻訳」とも呼ばれている。パターンに適合する入力文に対して品質の良い訳文が得られることから、多くの商用システムでトランスファー方式と併用する形で実現されてきた。最近では「翻訳メモリ」とも併用される傾向にある。

しかし、これらの文型パターン翻訳で使用されている文型パターン数はいずれも少なく(200~300パターン程度)、特定の狭い分野の文書に適用される例が多い。これは、パターン作成のコストが大きいこと、また、パターン数を増やすとパターン間の意味的な相互作用が増加して翻訳精度が低下することによるためと考えられる。

これに対して、既に、構造と意味の関係を考慮した「多段階翻訳方式」が提案されている⁽⁴⁾。この方式は、原言語表現の構造を意味を失わないように目的言語に対応づける仕組みとして「結合値パターン」を使用している。パターンの意味的排他性の問題は、精密な意味属性体系を使用することで解決しており、単語レベルの翻訳において精度の良い訳文が得られている⁽⁵⁾。しかし、複文(埋め込み文を持つ文)、重文(接続のある文)の持つ非線形性が扱えないこと、また、原文に対して単一

の目的言語表現が対応づけられる仕組みであるため、文脈に応じた表現選択ができないことが問題として残されている。

これらの2つの問題を解決するため、最近、言語表現の意味類型化を基本とする「意味的等価変換方式」が提案された^(6,7)。この方式は非線形な言語表現の構造を意味的に類型化⁽⁷⁾した「意味類型パターン辞書」を使用するが、この辞書を構築するためにはやはり「大規模な「文型パターン対辞書」を作成する必要がある。

文型パターンは、言語表現の非線形な構造を取り出してパターン化したものであるが、大規模な対訳コーパスからこのような文型パターンを作成するには、与えられた言語表現のどの要素が線形要素であり、どの要素がそうでないかを判断する基準を明確にし、作業手順化することが重要である。

ところで、線形性と非線形性は、表現構造と意味の関係に対して定義されるものであるため、現実の言語表現に適用するには、与えられた各表現の意味の定義を必要とする。すでに、言語表現の意味については、関係意味論の立場から概念の二重性に着目した検討⁽⁹⁾が行われており、また、言語表現構造の線形性、非線形性の問題については、工学的立場からの検討⁽¹⁰⁾が行われている。

本検討では、その結果に基づいて日本語表現の意味を英語表現で表すこととし、大量の日英対訳例文から非線形な表現構造を取り出して文型パターンを作成するための作業基準と作業手順を定めた。また、それに従い、重文・複文を対象に構文レベルで24.6万件の日英文型パターンを作成したので、その結果について報告する。

以下、第2章では、機械翻訳の立場から、各文要素が文全体に対して線形であるか否かを判断するための原則を示す。第3章では、汎化作業の方針と単語レベル、句レベル、節レベルの3段階の文型パターン作成のための作業項目とその基準を述べる。最後に第4章では、重文と複文の対訳標本文から作成された文型パターンの内容を示す。

2. 文型パターン化の原則と方針

言語表現の線形性、非線形性を判定するための原則と非線形な表現構造を文型パターン化するための方針を示す。なお、以下では「文型パターン」を単に「文型」と称す。

2.1 言語表現における線形要素と非線形要素

文献(9)によれば、言語表現の線形性は以下の通り定義される。

定義1: 言語表現の線形要素

特定の概念(複合概念)を表現するための表現構造の要素のうち、他の要素(意味の異なる要素でも良い)に置き換えても表現構造全体の意味(複合概念¹⁾)が変わらないとき、その要素をその表現構造の「線形要素」と言う。

定義2: 表現構造の線形性と非線形性

線形要素のみから構成される表現構造を「線形な表現構造」と言い、1つ以上の非線形要素を有する表現構造を「非線形な表現構造」という。

本稿ではこの定義に従って対訳例文の線形要素を汎化し、文型を作成することとする。この定義を現実の言語表現に適用する方法と注意すべき点は以下の通りである。

(1) 英語による表現の意味の定義

通常、単語が単一概念を表現するのに対して、上記の定義は、「句、節、文等の表現は話者の認識の中で形成された複合概念の表現²⁾である」とを前提としている。

そこで、各言語は複合概念を表すための様々な形式を持っていることに着目し、与えられた日本語表現の意味(複合概念)を英語表現によって記述することとすると、定義1の「表現構造全体の意味が変わらない」とは、「対応する英語表現の構造が変わらないこと」と読み替えることができる。

その結果、与えられた日英対訳例文において、着目する日本文の文要素が線形であるか否かを判定するには、それを他の文要素に置き換えたとき、対応する英文全体の表現構造が変化するか否かを調べればよいことになる。

(2) 線形要素の重要な性質

文型化を考える上で、上記の定義は以下の意味を持つ。

<線形要素の制約条件>

第1は文型の線形要素の置き換え範囲(値域)の問題である。定義1は、他の要素に置き換えても表現構造全体の意味が変化しないしないうような要素を線形要素としているが、これは実際にどんな要素に置き換えても良いことを意味しない³⁾。元々日本語側で見て、意味をなさない表現になるような置き換えはできず、線形要素と言えども、置き換え可能な範囲には一定の制約がある。

<要素の選び方と全体の線形性>

第2は表現要素の線形性と表現全体の線形性の関係である。定義2によれば、すべての要素が線形な場合に限り表現は線形だとされている。これは、表現全体の線形性は、その要素分解の方法に依存して決まることを意味している。また、要素の単位を指定すれば、線形・非線形の区別は一意に決定できることから、汎化の程度に応じた文型が作成できる。

<文全体の非線形性と文要素自身の非線形性>

第3は、線形・非線形の区別は表現の部分と全体の関係と言うものであり、線形要素だと言ってもその要素自身が線形であることを意味しないことである。線形要素の内部構造は非線形であっても良い。このように線形・非線形の分類が再帰的な構造を持つことは、長文の構造が複数の非線形構造の組み合わせで解析できることを意味する点で大変重要な性質である。

2.2 文型化の原則

前節の定義に従って対訳コーパス中の例文に含まれる線形要素を抽出し、それを汎化することにより、日英文型対を作成

する。以下、日英対訳文から日英対訳文型を作成するための原則について述べる。

(1) 文型化の対象としない対訳文

現実には得られる対訳例文の品質は様々である。対応する英訳文の意味が単独で日本語文の意味に対応するものをパターン化の対象とし、前後の文脈から意識されているなど、与えられた日本語だけでは対応関係を持たないような対訳例は文型化の対象としない。

但し、文型化の目的が非線形な言語表現の意味を正しく翻訳することにあることから、意識された対訳例を文型化することは極めて重要である。そのような例文では、無理な汎化はせず対訳原文をそのまま文型としても良い。

(2) 文法レベルでの文型の記述

本稿では文型記述用の言語として、文献(1)で提案された「文型記述言語」を使用するが、文法レベルの情報を使用して文型を記述することとし、変数の変域に対する意味的な制約条件は付与しない。また、語順の変更や文型要素の移動可能指定の機能も使用しない⁴⁾。

これは、現段階では意味的制約条件付与の必要性和必要な意味の粒度などが不明なためである。意味レベルでの文型記述の必要性和その方法などについては、文法レベルで記述された文型の被覆率特性が明らかになった段階で検討する。

(3) 必須要素と任意要素の区分

対訳例文の要素を以下で示すような「必須要素」と「任意要素」に分類する。

<必須要素> : 日本語文型内にそれがないと対応する英語文型が定義できない要素を言う。

<任意要素> : 日本語文型内にそれがなくても英語文型は決定できる要素で、文型定義に使用するか否かによってさらに以下の2つに分類する。

- ・「原文任意要素」: 削除されても対応する英訳文は変化しない要素で、文型では陽に示されない。
- ・「パターン任意要素」: 陽に示しておかないと訳語や訳語挿入位置の決定が困難な要素で、文型の一部として陽に示される。

なお、「必須要素」と「任意要素」は、いずれも字面でも良いし後で述べるような変数や関数を含む表現でも良い。

(4) 変数化とその範囲

変数化する対象は、単語(自立語、複合語を含む)、句、節の3種類の表現でいずれも線形な文要素である。これに応じて英語表現中の対応する文要素も変数化する。

今回の文型試作の狙いは、総合的に被覆率の高い文型が得られるか、文型相互間の意味的独立性が確保できるかの2点についての指針を得ることである。相反するこれら2つの目標を調和的に実現するため対訳例文を段階的に汎化することによって文型を作成することとし、変数化された文要素の単位に応じて、文型を「単語レベル」、「句レベル」、「節レベル」の3種類のグループに分類する。以下変数化の基本原則を示す。

まず、変数化の対象となる文要素であるが、変数化する文要素はいずれも別途翻訳して英語文型に埋め込めるもので、それ自身が線形である必要はない。変数化判定の原則は以下の通りである。

英語側に対応する訳語(句と節を含む)を持つ場合

日本語と英語で変数化される要素は必ずしも文法的に同じ属性である必要はなく、品詞や活用形が異なっても良い。日英の対応する部分の品詞が異なる場合は、品詞変換の関数を使用して変数化する。

英語側に対応する訳語を持たない場合

英語側に対応する訳語を持たない要素は、前項の分類によって「原文任意要素」か「パターン任意要素」かのいずれか

¹⁾ 文献(6)によれば、言語表現は概念化された話者の認識を表し、単語は「単一概念」を表すのに対して複数単語からなる表現は「複合概念」を表す。即ち、概念化されていない認識は言語の直接的表現では表現されないため「言語表現の意味」と「言語表現の表す概念」を同義で使用する。

²⁾ 「言語表現の意味をいかなる記法で記述しても計算機から見れば単なる記号に過ぎないから、意味記述言語は、表現能力があり、相互矛盾のない体系であればよい。その点、自然言語は表現能力の高い言語であるの目的言語を使用して原言語の意味記述を行うこととするが、この方法は、機械翻訳システム構成上、便利な方法と考えられる。

³⁾ 例えば、「私は彼より背が高い。」の文において、「私」を「あなた」に置き換えても、この文が「2者比較」と言う「複合概念」を表す点での意味は変化しない。しかし、「私」を「川」や「月」などに置き換えると、表す複合概念が変化する以前に文としても成立しなくなる。

⁴⁾ 提案されている「文型パターン記述言語」では、変数の意味的制約条件の記述方法や文型要素の出現順序の可変性の指定方法等も規定されているが、本検討では、第1段階として文法的な属性のみの情報で記述された文型パターンを試作することを課題としている。

と判定する。

日本語と対応しない英語文型要素

逆に、英文中に日本語に対応づけられないような要素については、日英対訳文の意味的な関係を調べ、文脈なしにその対訳関係が成り立つなら、英語パターンの要素として残し、そうでない場合は削除する。

次に、変数化する要素の範囲であるが、入力文と日本語文型との照合時の結果、変数にバインド(代入)され、英語文型に持ち運ばれる範囲である。問題となるのは用言性の文要素であるが、汎用性の高い文型とするため、用言の変数化においては「語幹+活用形」の範囲を変数化し、時制、相、様相の情報を表す自動詞などについては、別途定められた関数を使用して記述する。従って通常、文型では用言の活用形は指定されないが、指定が必要な場合は、活用形指定関数を使用する。

また、日英文型間の変数の対応関係を明確にするため、まず、日本語文型内で使用する変数には通し番号を付与する。これに伴い、英語文型内の変数には、日本語文型内の意味的に対応する変数と同じ番号を付与する。日本語文型で使用された変数のすべてが英語文型で使用されている必要はない。なお、同一の文型要素の変数化では、同一の変数番号を使用するものとする。

(5)関数化などによる汎化

付属語(辞と辞相当語)要素のうち線形なものは、各種の関数および選択記号によって汎化するが、入力文と文型との照合の段階で意味的な曖昧性の発生しない方法(主として字面指定の関数または選択記号)で記述する。

2.3 文型化の個別の方針

文型化における個別の方針を示す。

(1)作業自動化の可能性の追求

膨大な対訳例文から文型を能率良く生成するため、対訳例文の形態素解析と構文解析情報を使用し、文型作成の半自動化を目指す。そのため、機械的な変数化と関数化ができるよう、文型作成に先立って元となる対訳標本文は形態素解析し、解析誤りは人手で修正しておくこととする。

(2)文型照合の容易性

従来の構文解析では解消困難な構文多義の問題を解決することも文型翻訳方式の目標の一つである。そこで、入力文と文型の照合では構文解析は使用せず、入力文の形態素解析結果を使用することとし、文型は入力文の形態素解析結果との照合において曖昧さの発生しない方法で記述する。これに伴い、形態素解析で文法的、意味的に解釈の確定しない表現要素については文法的、意味的な指定は行わず、字面もしくは字面の指定される関数を用いて文型を記述するものとする。

(3)変数書き換えによる文型の汎化

句レベルの文型化では、対訳標本文中の線形な句(既に単語変数化された表現を含んでも良い)を変数化するが、より汎化の範囲を拡大するため、単語レベルで変数化した単語変数についても句変数に置き換え可能なものを探して句変数化する。例えば、名詞を表す単語変数Nを名詞句変数NPに置き換えることができれば、より適用範囲の広い文型が生成できたことになる。

(4)離散記号による文型の汎化

原文任意要素の定義から分かるように、機械翻訳で使用される文型は、パターン要素のすべてが入力文に含まれていることが条件となるが、逆に入力文には、文型に定義されない要素を含むことも認められる。

そこで、文型がより多くの範囲の入力文に適合するようにするため、入力文と文型の照合を制御するための離散記号を使用し、可能な限りこの記号を使用した文型化を行う。

(5)字面レベルでの表記の揺れの吸収

日本語は表記の揺れが多く、このことが、文型との照合で大きな問題となる。この問題を解決するため、入力文と文型との照

合の段階で形態素解析プログラムの持つ標準表記認定機能を利用することを前提に、原則として標準表記を用いた文型記述を行うこととするが、同時に、格助詞、助動詞、副詞(いずれも相当語を含む)等について、可能な限り標準表記と異表記の関係をまとめた表を準備し、字面グループを指定する関数や選択記号を使用して使用可能な表記を指定する。

3. 文型化の具体的方法

前章で述べた原則と方針に従って、表1で示すような、単語レベル、句レベル、節レベルの文型を順に作成する。また、各レベルにおける作業項目一覧を表2に示す。以下、これらの作業の内容について述べる。

表1.各レベルにおける文型化の内容

文型	汎化の内容
単語レベル	(1)名詞、動詞、形容詞、副詞などの自立語の変数化、(2)線形で文型上不要な要素を任意化と文型の骨組みとなる要素の抽出、(3)字面要素についてのグループ化
句レベル	(1)適用範囲を品詞から句への拡大、(2)機能語の適用拡大(格助詞、格助詞相当語、等)、(3)英語句生成関数の適用
節レベル	(1)名詞節副詞節主節従属節の変数化による重文複文の基本構造のパターン化、(2)日本語節から英語句への変換関数の使用、(3)その他英語構造生成関数の使用

3.1 単語レベルの文型化の方法

すべての対訳例文に対して、単語レベルの文型を1パターンづつ作成する。文型化の作業は、「原文任意要素の削除」、「自立語の変数化」、「述部語尾表現の関数化」、「パターン任意要素の指定」、「表現要素のグループ化」、「各種加工」の6種類から構成される。以下、それぞれの具体的内容について述べる。

3.1.1 原文任意要素の削除

以下の2の条件を満たす要素を「原文任意要素」と判定し、原文から削除する。

対応する文要素が英文にあり、日英双方の文からそれらを取り去っても対訳関係が成り立つこと。

着目する文要素を単独で英訳することができ、英訳文に組み込めること。訳語選択、訳語位置決定の困難な要素は、対象外。

また、日本語と対応関係を持たない英文中の要素も原文から削除する。なお、の条件は満たすが、の条件を満たさないものが、「パターン任意要素」である。

3.1.2 自立語の変数化

語彙的な用法で使用された線形な自立語を以下の順序で変数化する。なお、機能的用法で使用された自立語は変数化しない。

まず、日本語表現内の自立語に対応する訳語が英訳文中の同一の品詞の語(複合的な表現でも良い)として存在し、対訳辞書などによってその対応関係が決定できるものを日英同一の変数(変数名と変数番号共に同一)に置き換える。

残された日本語表現内の自立語のうち、線形なものを変数化する。この場合、英語側に意味的に対応する自立語がなくても良い。また、英語側の自立語で、日本語側に対応する自立語を持たないものは変数化しない。

以下、自立語変数化の個別的方法を示す。

(1)名詞の変数化

語彙的な用法の名詞と複合名詞を名詞変数Nに置き換える。但し、数詞、時詞はNUM、TIMEを使用する。また、機軸語として使用された名詞の、もの、こと、人、とき、場合、原因、理由などは変数化しない。

*1 文型パターン照合では入力文の構文解析情報は使用しないが、文型パターンの作成では使用する。

*2 構文解析の精度は形態素解析に比べて十分とは言えないため、予め標本文すべての解析結果を手で修正するには多大なコストが必要となる。そこで、構文解析プログラムは、句変数化、節変数化などで構文情報が必要となる標本文に限って限定的に使用する。

*3 この記号(スラッシュ記号"/ ")は、文型パターン要素間(但し、文節境界)に挿入するもので、この記号がある位置は、原文任意要素が挿入されても良いことを意味する。

表2 文型化作業項目とその内容

#	汎化規則の分類	汎化規則の内容	備考		
1	原文任意要素の削除	原文任意要素を削除する			
2	単語 レベル の汎化	自立語 の 変数化	名詞一般	格要素の名詞の変数化、述語名詞の変数化、 数詞の変数化(数詞として識別する)、 複合名詞の単一名詞化	名詞機能語は対象外。底の抽象 名詞(とき, 原因, 理由, 等)も対象 外。
3			複合名詞	述語単独動詞の変数化	機能語として使用された和語動 詞は対象外
4			動詞一般	単独動詞連体形と連用形の変数化	
5			複合動詞	複合動詞の変数化	主動詞のみ変数化
6			形容詞・ 形容動詞	述語形容詞 形容動詞の変数化	
7				単独形容詞形容動詞連体形, 連用形の変数化	
8			副詞	文修飾, 用言修飾の副詞の変数化	変数化の可否に要注意
9			品詞変換関数 適用	品詞変換関数を使用する。(表記法の例) N(N2) N(V3),	活用形表記法も併用
10			述部語尾表現の関数化	英語アスペクト情報の関数化	(例) ing, ed, er, stなどの関数を 使用
11			英語様相情報の関数化		
12	任意 要素 の 指定	名詞修飾語	単独の動詞連体形, 形容詞, 形容動詞連体形, 連体詞の任意化	英語文型の決定に不要な線形要 素をバックス記号#n[]で囲む。任 意化する部分は, 変数を含んで 良い。	
13		用言修飾語	単独動詞連用形, 単独形容詞, 単独形容動詞連用形の任意化		
14		その他	副詞的用法の名詞(昨日, 今日)の任意化		
15			英語に訳出されるその他線形要素の任意化		
16		英語に訳出されない要素の任意化	日本語側を[]で囲む		
17	表現要素のグループ化	同種の格助詞相当語や副詞をグループ化	[A B] の記法による		
18	各種 加工	主語の補完	英語主語に相当する名詞が日本語側がない時, 日本語パターン に補う。例) [N1]英語側のパターンは[N1]heの形式で記述。	パターン適用時は, [N1]要素はな くても良い	
19		冠詞の削除	英語パターン内の変数された名詞の冠詞を削除	定まった冠詞は対象外	
20	句レベル の汎化	句の変数化	名詞句, 動詞句, 形容詞句, 形容動詞句, 副詞句の変数化	単純な構造を持つものが対象	
21		適用 範囲 拡大	現在形変換	現在形でも使用される文型を対象にpas()を削除する。	
22			丁寧表現の 標準化	丁寧表現をフラットな表現に変更する。 <注> 英語訳出されない丁寧表現が対象	接頭辞「お」, 「御」等を削除
23			NをNPへ拡張	名詞変数Nを名詞句変数NPに置換。それに伴い不要となったパ ターン任意要素を削除。	可否判断要注意
24			VをVPへ拡張	名詞変数Vを名詞句変数VPに置換。	可否判断に要注意
25		機能語の拡張	格助詞を同種の格助詞相当語に置き換え		
26			述部の語尾を同種語尾パターンに置き換え		
27	同じ意味の副詞等をパターンに追加		(A B) の記法による		
28	節レベル の汎化	節の変数化	名詞節, 連体節, 連用節, 引用節, 並列節の単一変数化		
29		節レベル関数の使用	日本語節を英語句に変換する関数の使用		
30			英語節構造を指定する関数の使用		

複合名詞の場合、英語文型に影響を与えるような接頭辞、接尾辞を複合名詞一部として名詞変数Nに含めるような変数化は行わない。そのような場合は、複合名詞を構成する要素の一部を字面で残し、「大手金融会社 大手N」などのように部分的に変数化しても良い。なお、NNのような変数化はパターン照合を難しくする割に効果がないので行わない。

(2) 動詞の変数化

語彙的用法の動詞の「語幹 + 活用語尾」を動詞変数Vに置き換える。機能動詞「ある, なる, いる, くる, する, 行く, 言う」などは、変数化しない。

V+Vの複合動詞は、アスペクトの意味を有することが多いので、全体を1変数化せず、「始める」, 「Vてみる」, などのように主動詞側を変数化し、アスペクト情報を担う動詞は変数化しない。また、N+Vのような複合動詞は、N, V共に変数化しても良いし、いずれか、片方を字面としても良い。なお、命令形と連用形など文型選択で活用語尾が手ごかりとなる文型については、活用形指示関数を使用する。

(3) 形容詞, 形容動詞の変数化

形容詞, 形容動詞の変数化は動詞と同様である。従って、形容動詞は「名詞+だ(助動詞)」とはしない。

(4) 副詞の変数化

名詞の場合と同じである。但し、変数化の基本原則を満たさない場合が多いと思われるので、注意が必要である。

(5) 品詞変換関数の適用

日本語側の文型要素が、英語文型で文法属性の異なる自立語に対応するものを対象に「変数関数」を使用した文型化を行う。対象となる表現部分は日本語側のN, Vがそれぞれ英語側でV, Nになるもので、英語側の要素を変数関数V(N), N(V)で表現する。

3.1.3 述部語尾表現の関数化

動詞, 形容詞, 形容動詞に接続する助動詞および助動詞と助詞の連鎖表現(連語)を「形式指定関数」を用いて書き換える。書き換える際の注意点は以下の通りである。

この段階では、線形な時制, 相, 様相の情報の任意化は行わない。

従って、日本語側, 英語側の文型において、これらの情報は、「形式指定関数」を使用するが、それが無いときは、関数を新設するか字面関数を使用する。

なお、対応する英語文型の箇所ですべて日本語文型で使用した関数を使用する必要はない。使用する関数の種類は英語側で独自に決定する。

3.1.4 パターン任意要素の指定

日英対訳文型を見比べ、下記の線形要素を「パターン任意要素」として指定する。

- 単独の動詞, 形容詞, 形容動詞の連体形と連体詞
- 単独の動詞, 形容詞, 形容動詞の連用形
- 英語に訳出されるその他の線形要素
- 英語に訳出されない要素

但し、英語文形パターンでこの要素が指定されたときは、対応する日本語文型に同一の番号を持つパターン任意要素が指定されていない。

3.1.5 表現要素のグループ化

助詞, 助詞相当語, または、副詞などの字面のうち、同一の意味で異なる表記を持つものを対象に、選択記号を用いて置き換え可能な表現として指定する。指定される要素は必須要素とパターン任意要素を共に含んでよい。

3.1.6 各種加工

前節までで得られた日英文型に対して下記の処理を行う。

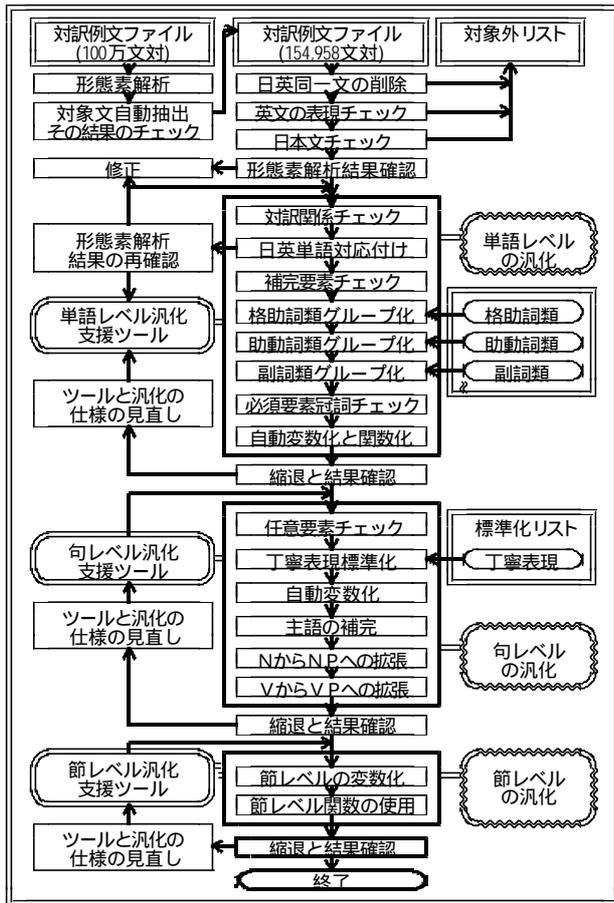


図1. 文型化の作業手順

(3) 関数化対象表現の抽出

関数化の対象となる字面を予め関数毎にグループ化し、グループ化された対応表と(1)の結果を参照して、関数化候補箇所を自動抽出する。

(4) その他の各種要素の記述支援

任意要素、補完要素の指定、冠詞と丁寧表現の汎化などについても可能な限り判断規則を設け、その候補となる要素を自動抽出する。

これらの方法で抽出された汎化候補に対して、線形要素であることが機械的に明確に判定できるものについては、機械的に変数化、関数化を行う。また、機械的な判断規則が得られないものは候補のまま残し、言語アナリストの判断にゆだねる。言語アナリストが判断に迷うものについては、設計者との協議により個別に扱いを決定する。

4. 重文と複文の文型化

述部を2つ又は3つ含む重文と複文(記述文)を対象に前章で述べた方法により文型化を行った。重文、複文を文型化の対象としたのは以下の理由による。すなわち、日本語表現のうち、単文の非線形構造については、既に日本語語彙大系⁽¹²⁾においてまとめられており、高品質の翻訳が可能となっているのに対して、重文複文の非線形構造については、類似の知識ベースがなく、訳文品質は依然として低いレベルにとどまっているからである。また、述部の数を2と3に制限したのは、重文、複文と言っても、現実の文では、4個以上の述部を持つ文全体が非線形であることは少なく、そのような文は、述部2または3の文型に分解して翻訳できる可能性が高いと考えられるためである。

4.1 対象例文と作業の状況

(1) 対象とする日英対訳例文

まず、辞書や日本語教材をはじめとする約30種類の対訳データファイルから作成した100万文の対訳コーパスから、重文、

複文15.5万件の対訳文を機械的に抽出した。その中には、会話文や文脈依存の訳文が含まれていたため、それらを人手で振り分け、12.9万件を標本文として3種類の文型を作成した。表3に抽出した対訳文と文型化の対象とした標本文の内訳を示す。

対象文に含まれる単語の種類と数を表4に示す。また、対訳標本文の平均文字数などを以下に示す。

日本語原文：平均文字数/文 = 23.3字(最大148字)
平均形態素数 = 12.9個/文(最大63個)
英文訳文：平均単語数 = 10.3語/文(最大59語)

表3 対訳標本文数と作成した文型数

文種別	説明	抽出した対訳文数	対象とした標本文数
文種別1	文接続1カ所を持つ文	72,018	57,235 (44%)
文種別2	文接続2カ所を持つ文	7,292	6,196 (5%)
文種別3	埋込み文1つを持つ文	54,931	46,907 (36%)
文種別4	埋込み文2つを持つ文	6,688	5,986 (5%)
文種別5	文接続と埋込み文各1つ	14,029	12,389 (10%)
-	合計	154,958	128,713 (100%)

表4 標本文中の単語の出現回数

#	種別	形態素数		#	種別	形態素数	
		延度数	異語数			延度数	異語数
1	名詞	417,886	56,861	8	接続詞	3,146	77
2	本動詞	223,178	10,324	9	感嘆詞	147	60
3	補助動詞	51,918	271	10	接頭辞	1,068	110
4	形容詞	31,681	915	11	接尾辞	1,749	336
5	形容動詞	19,587	2,562	12	助動詞	165,251	236
6	副詞	39,051	3,191	13	助詞	465,811	349
7	連体詞	32,585	731	14	記号	121,555	32
-	合計	-	-	-	合計	1,574,613	76,055

(3) 文型の適切性の検証

文型化作業の結果を確認するため、文型照合プログラムを用意し、作成された全文型に対する照合実験を行った。作成した文型の総数と不適切文型数の推移を表5に示す。

表5 不適切文型数の推移

種別	文型数	不適切パターン数	
		初期段階	現段階
単語レベル	128,071	58,194 (45.4%)	165 (0.13%)
句レベル	104,619	18,643 (17.8%)	1,810 (1.7%)
節レベル	13,031	4,638 (35.5%)	2,381 (18.3%)
合計	245,721	81,475	4,356 (1.8%)

この実験では、文型の作成に使用した標本文を入力文として文型辞書を検索し、適合する文型の中に、自分自身から作成された文型(自己パターン)と云うが含まれるかどうかを確認した。

その結果、最初の段階では、単語レベル、句レベル、節レベルで、それぞれ45%、18%、35%の文型に誤りがあることが判明したが、誤りの多くは、自動変数化プログラムに組み込まれた規則と文型記述仕様との不整合に起因するもので、人手修正の必要な文型の記述誤りは約5,000件であった。節レベルでは、まだ18%の不適合パターンを残しているが、この大半は、文型照合プログラムの機能不足により照合に失敗しているもので、文型記述の誤りは少ないと思われる。

以上の結果、文法レベルにおいて約24.6万件の文型をほぼ1年間(0.4人年)で作成することができた。人手作業に頼る場合に比べて、作業工数は、約1/10に減少したものと推定される。

4.2 作成された文型の例と数

作成した文型数の内訳を表6に示す。また、文型の例を表7に示す。

単語レベルのパターン化では、標本文数128,713文中、642文は変数化される部分がないため字面のみで文型として残され、

*1 日英対訳標本文から単語レベル、句レベル、節レベルの文型パターンの組を作成するのに30分かかると推定して、(30分/標本文)×15万標本文=450万分=37.5人年と見積られる。なお、いずれの場合も作業員としては、英語の素養のあるベテランの日本語アナリストの動員が必要である。

表 6 対訳標本文数と作成した文型数

文種別	説明	対訳原文データ		作成した文型数 () 重なり文型数			
		抽出文数	標本文数	単語レベル	句レベル	節レベル	合計
文種別 1	文接続 1カ所を持つ文	72,018	57,235 (44%)	56,883 (3,305)	45,241 (7,885)	6,343 (822)	108,467 (12,012)
文種別 2	文接続 2カ所を持つ文	7,292	6,196 (5%)	6,179 (99)	5,079 (127)	424 (7)	11,682 (233)
文種別 3	埋込み文 1つを持つ文	54,931	46,907 (36%)	46,684 (2,676)	38,654 (7,722)	3,807 (622)	89,145 (11,020)
文種別 4	埋込み文 2つを持つ文	6,688	5,986 (5%)	5,973 (84)	5,307 (223)	874 (63)	12,154 (370)
文種別 5	文接続, 埋込み文各 1 を持つ文	14,029	12,389 (10%)	12,352 (178)	10,338 (313)	1,583 (32)	24,273 (523)
- -	合計	154,958	128,713 (100%)	128,071 (6,342)	104,619(6,270)	13,031 (1,546)	245,721 (24,158)

表 7 .文型による非線形構造の記述例

区別	日本語文型又は日本語	英語文型又は英文
単語レベル	文型 それはN1にあるまじきN2.da.	Such N2 beunseemlyforN1.
	言語表現例 それは学生にあるまじき行為だ。	Such behavior isunseemlyforstudents.
句レベル	文型 N1は N2のN3/ことも#4[ある程度は]V5 .	(N1) can V5 whatN2V3#4[osomeextent].
	言語表現例 次郎の言うこともある程度わかる。	I can understand what Jiro saystosomeextent.
節レベル	文型 あれこれV1.temiru.taがN2がNP3.da.	All thingsV1.past, N2'sNP3.
	言語表現例 あれこれ考えてみたがそれがいちばんいい解決策だ。	All thingsconsidered, that'sthebestsolution.
節レベル	文型 N1は NP2もV3.naiとはNP4.gimu/ことだ。	It is AJP(VP4) that (N1)you shouldV3.notNP2.
	言語表現例 総理大臣の名前も知らないとはまことに哀れむべきことだ。	It is reallypitiablethatyoushouldnotknowthePrimeMinister'sname.
節レベル	文型 N1は /CL2とはV3.nai.kako.	(N1)I didV3.notCL2.past.
	適用例 彼があれほど英語が話せるとは思わなかった。	I didn'tknowhecouldspeakEnglishsowell.
節レベル	文型 CL1.teiru.nai.daとN2はVP3.kako.	N2VP3.past that CL1.not.
	言語表現例 彼女はもうぼくを愛していないのだとぼくは自分に言い聞かせた。	I convincedmyselfthatshedidnotlovemeanymore.

128,071件の文型が得られた。このうち6,342件(5%)は互いに同一のパターンとなったため、異なり文型数は121,729件である。

次に、句レベルのパターン化では、単語レベルで得られた文型のうちの約82%がさらに汎化され、104,619件の文型が得られた。このうち16,270件(16%)は、同一の文型となったため異なり文型数は、88,349件である。

これに対して、節レベルで作成された文型は13,031件で、単語レベルに比べて約1/10である。これは大半の対訳例文は非線形要素であり、汎化困難であること、すなわち、重文、複文構造の9割程度は、節の翻訳結果を組み合わせる従来の方法(要素合成法)では良い翻訳ができないことを示している。

なお、全文型のうち、同一のものは1,546件(11.9%)であった。記述レベルの違いによる文型の縮退の程度を見ると、単語レベル(6%)に比べて、句レベル(16%)は3倍以上大きい。これはほぼ予想されたところで、汎化するにつれて文型は縮退することを意味している。

以上、3段階の文型化で得られた文型の総数は、字面パターンを含め日本語文型24.6万件であった。そのうち同一の文型で縮退されるものは24,157件(9.8%)であるので、異なり文型の合計は221,564件である。

4.3 変数・関数の使用頻度

(1)変数の使用頻度

表 8に文型全体で使用された変数の種類と頻度を示す。

表 8 .変数の使用回数

変数種別	文型の種別(使用回数の割合)		
	単語レベル	句レベル	節レベル
N	303,319 (64.2%)	138,033 (43.4%)	10,135 (30.8%)
TIME	8,527 (1.8%)	5,187 (1.6%)	529 (1.6%)
NUM	6,036 (1.3%)	2,314 (0.7%)	189 (0.6%)
V	101,484 (21.5%)	48,036 (15.1%)	4,254 (12.9%)
REN	21,241 (4.5%)	2,158 (0.7%)	127 (0.4%)
ADV	11,491 (2.4%)	7,631 (2.4%)	603 (1.8%)
AJ	10,950 (2.3%)	6,193 (2.0%)	425 (1.3%)
AJV	9,473 (2.0%)	6,273 (2.0%)	434 (1.3%)
VP	---	58,908 (18.5%)	2,838 (8.6%)
NP	---	40,629 (12.8%)	1,985 (6.0%)
AJP	---	1,341 (0.4%)	78 (0.2%)
AJVP	---	935 (0.3%)	37 (0.1%)
ADVVP	---	117 (0.0%)	8 (0.0%)
CL	---	---	11,280 (34.3%)
合計	472,521 (100%)	317,755 (100%)	32,922 (100%)

表 4から、変数化の対象となり得る自立語は、名詞、本動詞、形容詞、形容動詞、副詞、連体詞を合わせて約76万語であったのに対して、表 8によれば単語レベルで変数化されたものは47万語であるから、自立語の62%が変数化されたことになる。文型当たりで見ると自立語5.9語中の3.7語が変数化され、非線形要素として字面のままとしたものは2.2語である。

句レベルでは1つの句変数の中に複数の自立語が縮退された場合を含め、59万語(78%)が変数化されている。また、標本文に含まれる節は28.2万件であったのに対して、それが節変数化されたものは、1.13万件(4%)にすぎない。

このことから、自立語、句、節のそれぞれ38%、22%、96%が非線形要素であったことになる。これを標本文数と対比すると、1文あたり、単語レベルでは平均2語あまり、句レベルでは平均1.3語程度の自立語が変数化されていない。標本文には節レベルではもちろん、自立語でも非線形なものがかかり多く存在していることになる。

(3)関数の使用頻度

文型全体で使用された関数の種類と頻度を表 9に示す。

表 8 .関数の使用回数

関数名	文型種別			関数名	文型種別		
	単語レベル	句レベル	節レベル		単語レベル	句レベル	節レベル
ta	33,155	33,138	5,614	teiku	399	427	57
teiru	9,737	11,090	2,126	temiru	391	399	60
reru	8,663	8,750	1,282	teoku	326	356	59
da	7,879	6,574	693	sugiru	326	317	44
nai	6,489	6,426	1,318	souda	291	311	63
teinei	4,304	4,815	1,029	utosuru	249	232	30
suitei	2,029	2,374	250	desu	215	311	28
you	1,880	1,799	259	kaishi	193	211	84
meirei	1,570	1,331	517	teshimau	189	298	75
tekudasai	1,126	994	444	tearu	189	205	64
tekuru	1,040	1,068	302	yotei	65	89	4
joutaihenka	930	1,417	358	rashii	53	72	23
tai	905	895	51	nisuru	46	198	16
tekureru	719	707	128	teyaru	45	61	7
teshimatta	639	526	58	teyoi	42	40	2
gimu	574	648	89	noda	38	38	23
dekiru	542	813	80	hougayoi	31	33	2
sase	521	694	71	katei	3	2	1
darou	505	534	184	kirezu	1	0	0
合計	86,295	88,193	15,495	合計	86,295	88,193	15,495

この表から、文型当たりの関数使用回数は単語レベルでは平均0.7回、句レベルでは0.95回、節レベルでは1.5回である。汎化のレベルが上がるにつれて関数の使用頻度が増えている点から、使用される変数の種類の違いだけでなく関数化の点でも汎化が進んでいることが分かる。

(3)線形な文型と非線形な文型

既に述べたように初期の段階で汎化不能として文型化の対象外としたものが624件あったが、実際に作成した文型のうち関数や変数が使用されず字面のままとなったものは単語レベルで302件存在した。いずれも線形要素を持たないと判定されたものである。対象標本12.9万件中、これらの割合が0.72%であることから、殆どの標本文(99%以上)は一つ以上の線形要素を持っていたこととなる。

また逆に、字面を含まない文型は、単語レベルで15件、句レベルで401件、節レベルで155件であった。このことから、いずれの場合も線形な文型は少ないが、特に、単語レベルで少ないことが分かる。

以上のように、文型全体では、殆どすべての重文・複文は要素合成法の適さない非線形な構造を持つものの、その内部には線形要素(平均4~5カ所)を持つものが多い。これは、多様な言語表現がカバーできるような文型が開発できる可能性を示しており、文型翻訳は、用例翻訳の限界を超える方法として期待される。

4.4 汎化による文型の同一化

(1)同一化する文型の割合

文型数と被覆率の関係を考えると、標本量が増大するにつれて同一化するパターンが増大することが期待される。しかし実際に12.9万件の標本文で得られた文型で同一化したものは単語レベル、句レベルで各々約5%、約16%でかなり少なく期待以下である。

このことから、入力文の任意要素を含むすべての要素が解釈できるような文型を網羅的に準備することは困難であることをが分かる。しかし、これは実用規模の文型辞書の開発の困難さを意味するものではない。機械翻訳では、必ずしも入力文の全要素が文型パターンで解釈される必要性はなく、入力文と必須要素の範囲で一致する文型が使用できる。機械翻訳への適用性効果を推定するには、今後、各文型の適合条件を考慮した被覆率評価試験によって文型数と被覆率の関係を明らかにする必要がある。

(2)文要素の任意化の効果

文型化では、各種の変数と関数のほか任意化などのための記号類が使用される。このうち、任意化の機能によって縮退する文型数を調べた。その結果によれば、前者では「パターン要素任意化」によって文型が縮退する割合は予想外に小さく0.2%にとどまった。この点については、任意化対象範囲とその定義方法について再検討する必要がある。そうである。

5.あとかき

言語表現の構造と意味の関係に関する線形性と非線形性の定義に基づき、日英対訳例文から、非線形な表現構造を取り出して文型化する方法を提案し、そのための具体的な作業手順を示した。また、その方法を重文、複文15万件の日英対訳例文に適用し、単語レベル、句レベル、節レベルの3種類の文型をそれぞれ12.8万件、10.5万件、1.3万件(合計24.6万件、異なり221.1万件)益試作した。

単語レベルと句レベルの文型では、標本とした対訳文に含まれる自立語(平均6語)のうち3~4語が線形要素として変数化されたこと、また、関数化される付属語が1件弱あったことから、文型当たりの線形要素は平均4~5カ所存在することなどが分かった。節レベルの文型では、変数化された節は4~5%に過ぎず、大半の節は非線形要素であった。

これらのことから、重文、複文の構造は、殆どが非線形で、節に分解して翻訳する要素合成法は適さないが、単語や句のレベルで見るとかなり多くの線形要素を含んでいることが分かる。

また、人手による文型の作成では、膨大な作業量が予想されたが、作業基準の明確化と各種のプログラムツールの準備などにより、作業の大半を自動化することができた。人手作業に頼

る場合に比べて作業量は約1/10に削減したと推定され、作業品質の均一化も達成することができた。

従来、汎用的で大規模な文型辞書の開発は困難と考えられ、対象分野を限定した用例翻訳などが試みられてきたが、上記のように標本文には多くの線形要素が存在すること、文型化作業の大半が機械化できることから見て、汎用的な文型辞書開発の展望が得られた。

ところで、今回作成した文型は、文法情報で記述されている。試行的な実験によれば、再現率はかなり高い値が得られる見込みであるものの、意味的な適合率はまだまだ低い値だと予想される。意味的な適合率を向上させるには、各種変数に対する意味的な制約条件の付与が必須と見られる。今回の試作によって多くの文型を得ることができ、種々の実験的検討が可能となった。今後は、(1)より高度な汎化手法とその効果の推定、(2)意味的制約条件付与の粒度とその効果の推定などについての実験的検討を進め、被覆率と意味的排他性ともに優れた文型を目指して改良を行っていく予定である。

謝辞

この研究は、科学技術振興機構(JST)の戦略的基礎研究事業(CREST)で行ったものである。ご議論を頂いた宮崎正弘氏(新潟大)、池田尚志氏(岐阜大)、新田義彦氏(日本大)、佐良木昌氏(長崎純心大)氏本和英氏(長岡技科大)、白井諭氏(元ATR)、横尾昭男氏(NTT)に感謝する。また、文型化作業を担当して頂いたNTTアドバンステクノロジー株式会社の皆様に感謝する。

参考文献

- (1) 長尾真:『自然言語処理』岩波書店(1996)
- (2) 長尾真,黒橋貞夫,佐藤,池原悟,中尾洋:『岩波講座「言語の科学」第9巻「言語情報処理」』岩波書店(1998)
- (3) 池原悟:『自然言語処理の基本問題への挑戦』人工知能学会誌,Vol.16, No.3, pp.522-430(2001)
- (4) 池原悟,宮崎正弘,白井諭,林良彦:『言語における話者の認識と多段翻訳方式』情報処理学会論文誌,Vol.28, No.12, pp.1269-1279(1987)
- (5) 金出地真人,徳久雅人,村上仁一,池原悟:『結合価文法による動詞と名詞の訳語選択能力の評価』情報処理学会,自然言語処理研究会,情処研報2003-NL-153-16, pp.119-124(2003)
- (6) 池原悟:『言語で表現される概念と翻訳の原理』電子情報通信学会,言語と思考研究会,信学技報,TL2003-25, pp.7-12(2003)
- (7) 池原悟,佐良木昌,宮崎正弘,池田尚志,新田義彦,白井諭,柴田勝征:『等価的類推思考の原理による機械翻訳方式』電子情報通信学会,言語と思考研究会,信学技報TL2002-34, pp.7-12(2002)
- (8) 衛藤純司,池原悟,池田尚志,佐良木昌,新田義彦,柴田勝征,宮崎正弘,白井諭:『意味類型構築のための接続表現の体系化について』情報処理学会,自然言語処理研究会,情処研報2003-NL-155, pp.31-38(2003)
- (9) 池原悟,宮崎正弘,佐良木昌,池田尚志,白井諭,村上仁一,徳久雅人:『機械翻訳のための日英文型パターン記述言語』電子情報通信学会,思考と言語研究会,信学技報TL2002-48, pp.1-6(2003)
- (10) 池原悟:『非線形な言語表現と文型パターンによる意味の記述』情報処理学会,自然言語処理研究会,情処研報2004-NL-159, pp.139-145(2004)
- (11) 池原悟,宮崎正弘,佐良木昌,池田尚志,白井諭,村上仁一,徳久雅人:『機械翻訳のための日英文型パターン記述言語』電子情報通信学会,思考と言語研究会,信学技報TL2002-48, pp.1-6(2003)
- (12) 池原悟,宮崎正弘,白井諭,横尾昭男,中岩浩巳,小倉健太郎,大山芳史,林良彦:『日本語語彙大系』第5巻「構文体系」岩波書店(1997)