

NEAC100による和文英訳

谷口道興

最近オフィス・オートメーション(OA)がブームになっている。このOAでは事務処理上のあらゆる作業をその対象としている。販売・仕入・経理、在庫管理、顧客管理、給与計算等。また外国文献の翻訳、外国向けビジネスレターの作成等も含まれる。前者についてはコンピュータ化が著しく進んでいるが、後者のコンピュータを使用している機械翻訳(二言語間あるいは多言語間)については、まだ個人の手作業に負うという感じがある。本稿では機械翻訳のための大型専用コンピュータではなく、小型の汎用コンピュータである長野大学NEACシステム100M80IIを使用して、日本語を英語に翻訳する方法論を述べる。使用言語はCOBOLである。

1. 和文英訳、機械翻訳のアルゴリズム

計算機(コンピュータ)は人間のように文章を読んで、その意味を理解するという水準まで達していない。計算機にそれを実行させるためには文章を符号化、数値化しなければならない。文章を符号化、数値化するために、文章を構成する各品詞をコード(code)によって表現することにする。

英語の品詞は名詞、代名詞、形容詞、副詞、動詞、前置詞、接続詞、間投詞の8品詞に分けられる。このうち動詞は動詞と助動詞の2つに分類する。これらの品詞に対して、名詞はN、代名詞はPRON、形容詞はA、副詞はADV、動詞はV、助動詞はAUX、前置詞はPREPとコード化しておく。さて、英文を構成する主構成要素は主語、目的語、補語である。この4つを主構成要素と呼ぶことにする。この主構成要素のうち、目的語は直接目的語と間接目的語の2つに分類する。この5つの構成要素がどのように配列されて英文が構成されるかを次の5つの文型で表わす。

文型1: 主語+動詞

文型2: 主語+動詞+補語

文型3: 主語+動詞+目的語

文型4: 主語+動詞+間接目的語+直接目的語

文型5: 主語+動詞+目的語+補語

この5つの文型の中で、すべての主構成要素4つを含んでいるものは文型5である。文型5の目的語を間接目的語と直接目的語の2つに分類したものを

文型6: 主語+動詞+間接目的語+直接目的語+補語とする。

この文型はすべての構成要素を含むもので、その構成要素がどのような順序で配列されるかを示している。そのため実際の文では、この中のいくつかは抜けることもある。

文型6について、それぞれの構成要素になり得る品詞を当てはめてみる。主語に対して名詞または代名詞が、動詞に対しては動詞が、間接目的語に対しては代名詞が、直接目的語に対しては名詞または代名詞が、補語に対しては名詞または代名詞または形容詞が当てはまる。すなわち名詞または代名詞+動詞+代名詞+名詞または代名詞+名詞または代名詞または形容詞となる。これが文型6を品詞で表わしたものである。これをコードで表わすと、NまたはPRON+V+PRON+NまたはPRON+NまたはPRONまたはAとなる。

さらに文の各構成要素には、それぞれ修飾語がつけられる。修飾語としては副詞、形容詞、前置詞、助動詞がある。これらの修飾語の機能を定義しておく、副詞は動詞や形容詞または他の副詞を修飾する語である。しかし、本稿の段階では動詞または形容詞を修飾する語と限定している。形容詞を修飾する場合は、形容詞の前に置かれるものであり、動詞を修飾する副詞は英文の並び方の

関係から文の最後にくる場合がある。形容詞は、名詞または代名詞の前に置かれてそれらを修飾する語。前置詞は、名詞または代名詞の前に、また名詞または代名詞が形容詞や副詞で修飾されている場合には、さらにその前に置かれる語。助動詞は、動詞の前に置かれてそれを修飾する語。

これらの修飾語を、定義に従って品詞で表わした文型6に当てはめると、

ADV + A + N または PRON + AUX + V +

PREP + ADV + A + PRON + PREP + ADV + A + N または PRON + N または PRON または A + ADV となる。これが品詞で表わした文型6の応用である。最後の N または PRON または A に対して修飾語が付いていないのは、補語には修飾語がほとんどないと思われるからである。

品詞で表わした文型6の応用をコードで表わすと図1となる。

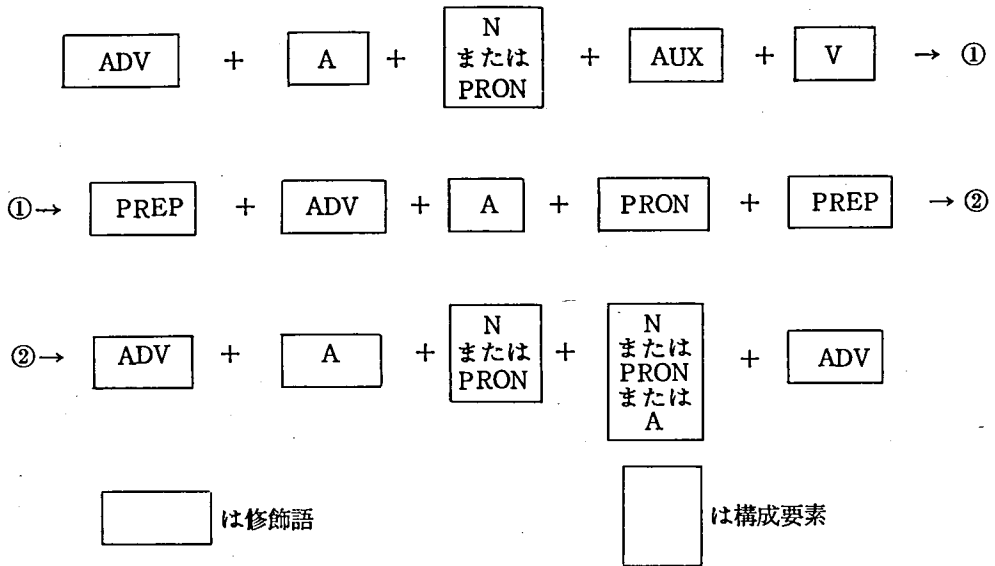


図1. 英語の文型

これに対して、日本語の品詞の分類は名詞、代名詞、形容詞、形容動詞、副詞、連体詞、動詞、接続詞、感動詞、助動詞、助詞となっている。また日本(語)文を構成する主な構成要素は、主語+述語+動詞である。このように日本語の品詞と構成要素の分類は、英語のそれとはかなり違っている。しかし機械翻訳を行うには、英語の場合と同じ品詞と構成要素を使用して、英語の場合と同じ方法で日本語の文型を作る。そのため形容動詞は形容詞に含め、助詞は英語の前置詞に相当するものとして考える。

英文の場合、4つの主構成要素をすべて含むものは文型5であった。日本語になった場合の主構成要素の並び方は、主語+目的語+補語+動詞である。目的語を間接・直接目的語に分けて、主語+

間接目的語+直接目的語+補語+動詞とする。これらの構成要素になり得る品詞は、英語の場合と同じで主語に対して名詞または代名詞、間接目的語に対して代名詞、直接目的語に対して名詞または代名詞、補語に対して名詞または代名詞または形容詞、動詞に対して動詞である。つまり名詞または代名詞・代名詞・名詞または代名詞・名詞または代名詞または形容詞・動詞となる。これらの品詞に修飾語が付くのであるが、日本語の場合には英語の場合と違う点がいくつかある。(例えば前置詞の位置が「学校へ」は「TO SCHOOL」と前置詞の位置が前後逆になる。また助動詞の位置も「行くだろう」が「WILL GO」というように前後逆である。さらに動詞を修飾する副詞は文の最後ではなく文中に入る。その他の修飾語については、英

語の場合とは同じである。)

これらのことからそれぞれの品詞に修飾語をつけてコード化したものが図2である。日本語を英(語)文に翻訳するためには、日本語の文型の配列

を文型6の応用の配列に並びかえればよい。しかし品詞で表わした文型を使用したのでは、コンピュータで実行はできない。コンピュータで実行させるためにはコード化した文型を使用することになる。

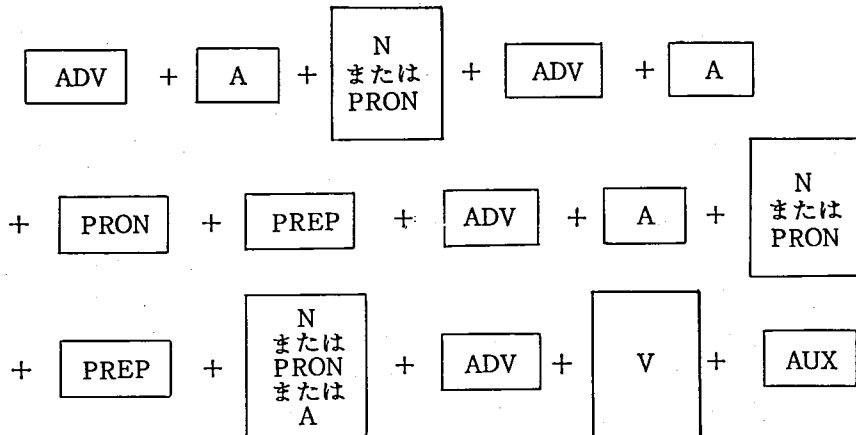


図2. 日本語の文型

2. ファイルについて

コンピュータに翻訳をやらせるための必要な要素の一つとして辞書がある。この辞書に相当するものがファイルである。ファイルはレコードの集合によって構成される。さらにレコードは項目の集合によって構成される。本稿のファイルは一つのレコードを3つの項目によって構成している。(例)「ガッコウ SCHOOL N」これは片仮名で表わした日本語の部分「ガッコウ」と、それに対応する英語の部分「SCHOOL」と「SCHOOL」という名詞をコード化した部分「N」の3つによって構成されていることになる。

上例において、3つの項目をそれぞれ JAPANESE, ENGLISH, M-CODE と呼ぶことにする。翻訳の不完全さを補うためにファイル内の JAPANESE を詳しく記述することにする。

ある一つの英単語についても、日本語の言い回しで違った言い方をすることがある。例えば HAPPEN という単語は、オコル、オキル、オコリマス、オキマスというように幾通りもの表現がある。したがって HAPPEN という一単語に対してレコードが4件分ほど必要ということになる。逆にそれが日本語についてもいえる。また「エキデ」

とか「クルマデ」とかのように場所を示すか方法を示すかの違いで、日本語の「で」という言葉が「AT」または「BY」と変わってくる。次に同音異義語について。日本語の「オキル」という語に対しては、WAKE UP, HAPPEN, RISE 等の英語が当てはまる。また、アクセントの判別ができないので「カエル」という日本語に対しては、FLOG, CHANGE, COMEBACK 等の英語が当てはまる。

このように、日本語の表現をできるだけ詳しくファイルに記述しておくことが必要で、言い換えればファイルを整備しておくことが重要となる。これらのファイルは必ずしも五十音順、またはアルファベット順に並んでいる必要はない。しかしプログラムの都合上、同音異義語だけは連続して並んでいなければならない。また、M-CODE の部分で名詞 (N) は単数形 (N-S) と複数形に、動詞も現在形 (V) と過去形 (V-P) に分けてある。これは翻訳のためのプログラムが高度になれば当然これらの識別が必要となるからである。しかしこの実験的なプログラムではその識別は除外し、次稿で検討することにしてはいる。

3. 入力部について

プログラムで使用する変数について。TABLE 1：ファイルから読み出された英単語の記憶場所。TABLE 2：ファイルから読み出されたコードの記憶場所。TABLE 3：コードの識別のために付けられたウェイトの記憶場所。TABLE 4：ファイルから読み出された日本語の記憶場所。TABLE 5：ファイルから読み出された英語を正しく並びかえたものの記憶場所。BAN：TABLE 1, 2, 3 の位置を示す。CNT：TABLE 4, 5 の位置を示す。ACP：キーボードから入力された仮名文字の記憶場所。WEIGHT：ウェイトの計算場所。ENGLISH 1とCODE 1：同意異義語とそのコードの記憶場所。ENGLISH 2とCODE 2：同音異義語とそのコードの記憶場所。W 1, W 2, W 3：単語の順序を入れかえるための記憶場所。

プログラムを実行させると、これらの変数がゼロまたは空白でクリアされる。次にBANに1が加算される。ついでカナ文字の入力が可能となる。このときの入力方法は一単語づつ区切って行なう。入力された単語はACPに格納される。次にファイルが1件づつ読まれてACPに格納されているカナ文字と一致するレコードを捜す。その際一致するレコードがファイルに無い場合は、画面上に“UNDISCOVERD IN MASTER-F”の文字が表示されるので、キーボードから必要な英単語と、その品詞を示すコードを入力することになる。これらはそれぞれTABLE 1とTABLE 2のBAN番目に格納される。またACPに格納されているカナ文字と一致するレコードがファイル内であれば、一致したレコード内のENGLISHとM-CODEを、それぞれENGLISH 1, CODE 1に格納する。つぎに同音異義語の有無を調べるために、ファイルをもう一件読む。そのレコード内のJAPANESEとACPが一致しなければ、ENGLISH 1とCODE 1に格納されている英語とコードが、TABLE 1とTABLE 2のBAN番目に格納される。またこのときACPに格納されていたカナ文字もTABLE 4のBAN番目に格納される。(以下、流れ図で説明)

4. 処理部について

このプログラムでは品詞の判別はコードによって行なっている。しかし一つの文中に名詞、形容詞等が複数個含まれているとNとN, AとAという具合に判別ができない。このような場合はN 1, N 2, A 1とA 2とすれば判別が可能である。コードの後に付けるこの数字をウェイトと呼ぶ。本節は、すべてのコードにウェイトをつけることから始まる。まずBANをゼロクリアする。BANに1を加算して、BANが1になる。コードが格納されているTABLE 2のBAN番目のコードは何であるかを調べる。もし主語となりうるNまたはPRONであったなら、WEIGHTに1を加算する。WEIGHTはゼロクリアされているので、 $0 + 1 = 1$ となる。NまたはPRONでなければWEIGHTへの加算は行なわない。再びBANに1が加算されBANが2となる。TABLE 2のBAN番目のコードは何であるかを調べる。もしNまたはPRONであれば、WEIGHTに1を加算する。WEIGHTをTABLE 3のBAN番目に格納する(以下、流れ図で説明)。各コードにウェイトを付け終ったならば、正しい英文になるように単語を並びかえる。まず主語が最初に置かれなければならない。主語となり得るものはNまたはPRONである。まずBANをゼロクリアする。BANに1を加算する。TABLE 2のBAN番目のコードは何であるかを調べる。NまたはPRONでなければBANに1を加算して、TABLE 2のBAN番目について調べる。もしNまたはPRONであれば今度はTABLE 3の、BAN番目のウェイトは何であるかを調べる。ウェイトが1であれば、そのNまたはPRONが主語であると判断する。そしてその主語はAまたはADVで修飾されているかを調べる。AはNまたはPRONの前に置かれてそれを修飾するのであるから、BANから1を減算してTABLE 2の(BAN-1)番目のコードを調べる。もしAでなければ、主語は修飾されていないものと考えて、主語となるNまたはPRONをTABLE 5の1番目に格納する。もしAであれば、そのAはさらにADVによって修飾されているかを調べる。ADVはAの前に置かれてそれを修飾するのであるから、BANから1を

減算して TABLE 2 の (BAN-1) 番目のコードを調べる。もし ADV でなければ、A は修飾されていないものと判断する。そして主語を修飾する A に対応する TABLE 1 の英語を TABLE 5 の 1 番目に、主語となる N または PRON に対応する TABLE 1 の英語を TABLE 5 の 2 番目に格納する。もし ADV であれば、そのコードに対応する TABLE 1 の英語を TABLE 5 の 1 番目に、A に対応する TABLE 1 の英語を TABLE 5 の 2 番目に、主語となる N または PRON に対応する TABLE 1 の英語を TABLE 5 の 3 番目に格納する。この 1, 2, 3……番目を示す変数が CNT である。CNT は 1 から始まり、TABLE 5 への格納があるごとに 1 ずつ加算される。

動詞の格納について。BAN をゼロクリアしてから (BAN+1) をくり返しなが、TABLE 2 の BAN 番目のコードを調べる。V は AUX を伴っているかを調べる。日本語の文型では、AUX は V の後に置かれている。そのため BAN に 1 を加算して TABLE 2 の (BAN+1) 番目のコードを調べる。AUX でなければ、TABLE 1 の BAN 番目の英語を、TABLE 5 の CNT 番目へ格納する。もし AUX であれば、V と AUX の位置を入れかえてから、TABLE 5 に格納する。その方法は、TABLE 1, 2, 3 の BAN 番目に格納されている英語、コード、ウェイトを、いったん W 1, W 2, W 3 に移し、TABLE 1, 2, 3 の (BAN+1) 番目に格納されている英語、コード、ウェイトを TABLE 1, 2, 3 の BAN 番目に格納する。そして W 1, W 2, W 3 に格納されていた英語、コード、ウェイトを TABLE 1, 2, 3 の (BAN+1) 番目に格納する。これで AUX が V の前に入れかわる。その後、TABLE 1 の BAN 番目の英語を TABLE 5 の CNT 番目に、TABLE 1 の (BAN+1) 番目の英語を TABLE 5 の (CNT+1) 番目に格納する。

目的語について。目的語になり得る品詞は、N または PRON である。(BAN+1) をくり返しなが、TABLE 2 の BAN 番目のコードを調べてゆく。もし N または PRON があれば、TABLE 3 の BAN 番目のウェイトを調べてみる。ウェイトが 1 の N または PRON は、主語として使用されているので、この場合はウェイトが 2 のものを捜

す。ウェイトが 2 の N または PRON が見つければ、それを目的語と判断する。この目的語は、A または ADV によって修飾されているかを調べる。A は目的語の前に置かれて、それを修飾するのであるから、BAN から 1 を減算して (BAN-1) 番目のコードを調べる。もし A でないならば、目的語は修飾されていないものと判断して、今度は PREP を伴っているかを調べる。日本語の文型では、前置詞は N または PRON の後に置かれている。そのため、さきほど 1 が減算された BAN に 1 を加算して元に戻し、さらに 1 を加算して (BAN+1) とする。TABLE 2 の (BAN+1) 番目のコードが PREP であるなら、PREP と N または PRON の位置を入れかえる。入れかえの方法は前述した方法による。そして TABLE 1 の BAN 番目の英語を、TABLE 5 の CNT 番目に、TABLE 1 の (BAN+1) 番目の英語を、TABLE 5 の (CNT+1) 番目に格納する。もし PREP を伴っていない場合は、目的語だけが格納される。もしこの目的語が A によって修飾されている場合は、さらに A が ADV によって修飾されているかを調べる。もし ADV によって修飾されていない場合は、A に対応する TABLE 1 の英語を、TABLE 5 の CNT 番目に、目的語を TABLE 5 の (CNT+1) 番目に格納する。もし ADV によって修飾されている場合は、そのコードに対応する TABLE 1 の英語を、TABLE 5 の CNT 番目に、A に対応する TABLE 1 の英語を TABLE 5 の (CNT+1) 番目に、目的語を TABLE 5 の (CNT+2) 番目に格納する。目的語は、間接目的語と直接目的語として、1 つの文中に 2 回含まれる場合がある。そのため、ウェイトが 3 の N または PRON について、再度この処理をくり返す。さらに以上の処理の中に含まれなかった ADV があるなら、それは V を修飾する ADV と判断して、そのコードに対応する TABLE 1 の英語を、TABLE 5 の文末に格納する。以上で英単語の並びかえは終了する。処理結果はプリンタで入力データ (日本語) と共に印字する。

5. 結論

このプログラムを実行させた結果、文型 1 から

文型5の文例の翻訳は成功した。従って本論で述べた方法論の範囲内の文であれば正しく翻訳ができる。逆に、予め用意されたファイルと構文に合致する文でなければ翻訳はできない。即ち、大きなファイルと複雑な文法情報が用意されなければこれらの問題は解決できない。

(謝辞) 本稿の作成にあたり、小林誠氏に多

大なる御協力をいただきましたことを感謝いたします。

(文献)

- 1) 山崎貞：新自修英作文典，研究社（1972）
- 2) 長尾真：着実に進展している機械翻訳，科学朝日（Oct. 1980）

〔処理部の流れ図〕

