

言語モデルにもとづいた
機械翻訳実験

通商産業省 電気試験所
西村 楚彦

1	概 要	2
2	文字列の処理	7
3	品詞列の処理	10
4	文字列の展席	16
5	例	19

1 概要

1.1 はじめに

筆者の作成した機械翻訳システムは、単に特定の言語を翻訳するためのだけのためのプログラムではない。一般的な記号の列を処理するためのプログラムである。

記号列の処理の算法は、いろいろの利用者が自分で記述しなければならない。その記述の形式は、メタ言語的な約束によって強く制限されている。逆にいえば、少数個のメタ言語的な約束だけにしかたがって、容易に対象言語の性質を記述することができると思われる。

ここで主として報告したいのは、いわゆる自然言語の翻訳を行なうために、メタ言語的な約束をとれほど設定しなければならないか、ということである。

1.2 構成

計算機の中には、辞書、文法表、処理プログラムの三者を入れておく。ここへ入力記号列を与えると、それにしたがって出力記号列が得られる。これが機械翻訳である。

辞書というのは、対象言語の性質をメタ言語的な約束にしたがって記述した規則の集合である。普通は、原語文字列、品詞、訳語文字列で一つの規則になる。

文法表というのは、対象言語の性質をメタ言語的な約束に
 ついて記述した規則の集合である。普通は、品詞の列と
 それで構成される句の品詞で一つの規則になる。

処理プログラムというのは、入力記号列を辞書や文法表の
 記述にしたがって、メタ言語的なしかたで解釈して変換し、
 出力記号列を作り出す算法である。

1.3 過程

翻訳は次のようにすすめられる。まず原語文字と呼ばれる
 記号の列、すなわち入力記号列 X を、品詞と呼ばれる記号の
 列 Y_0 に書き換える。記号列 Y は、書き換えが可能なあいだは
 重ねて何回も書き換えられ、記号列 Y_n となる。最後にこの記
 号列 Y_n が、訳語文字と呼ばれる記号の列、すなわち出力記号
 列 Z に展開される。

記号列の書き換えの規則は、あるきまつた値を有する記号
 列 R の「識別」と、それにたいする「処理」とからなってい
 る。

記号列 X を記号列 Y_0 に書き換える仕事は、文末品詞と呼ば
 れる特別の記号 y に行き当たったときに中断され、 Y の書き
 換えに移行する。

記号列 Y の書き換えは、そのなかに識別できる記号列 R が

もはやなくならぬときに中断され、記号列 Z の展開に移行する。

1.4 逐次処理

原文の要素の列 (X または Y) と文法規則の集合とを照合して文法的な解析処理をすすめてゆくのに、おおざっぱにいうと文走査型と逐次処理型とが区別できる。

文走査型では、文法規則の一つをとりあげて、それで文(始端、終端の限定された記号列)全体を端から端まで走査し、その規則があてはまる部分を処理する。これを、順序づけられた文法規則のすべてについてくりかえすものである。

逐次処理型では、文の端をとりあげて、文法規則の集合中にそれがあてはまる規則があるかどうかを走査する。文の端が処理できたら、(できなくても、)新しい端をとりあげてくりかえす。

筆者のシステムでは、まず文字列 X の処理を逐次処理で行ない、それが文の終端までゆくと、次に品詞列 Y の処理を逐次処理で行なう。どちらの場合にも照合は、逐次処理、見出しとの一致、最長一致、右方向の原則による。

文字列 X の処理は文の左から右へただ一回だけ行なわれる。このときの入力 は原文の文字列 X で、出力は語の品詞列 Y 。

である。ここで文のすかたは、文字の層から品詞（または語）の層にうつるのである。

品詞列 Y の処理は文の左から右へ反復して行なわれる。品詞列 Y に処理をくわえた結果はふたたび品詞列 Y になる。だからその品詞列 Y について重ねて文法規則との照合をとるのである。

文字列処理および品詞列処理の規則の集合はそれぞれ表の形になっている。その表の各項目はすべて、見出し（引数）と関数とからなっている。

1.5 要約

X : 原文の文字列

Y : 品詞の列

Z : 訳文の文字列

R : 識別された記号列。きまつた値を有し、長さは1以上

S : 処理（書き換え）される記号列。 R に含まれ、長さは0以上

T : 書き換える記号列

辞書の規則: R は X である。 S は常に R に等しいものとされる。 T には二種類ある。第一種の T は、長さ1の品詞列 Y と長さ0以上の訳語文字列 Z とからなっている。第二種の T は長さ1以上の原語文字列 X である。

文法表の規則: R は Y である。 S は R に含まれ、長さは0以上である。 T の長さは、0、1、または S の長さに等しい。 T の値は表面上は Y である。

普通品詞: 辞書中できまつた値の X にたいして与えられ、のちに訳語文字列 Z に展開される。

文末品詞: 普通の語の品詞とほとんど同じであるが、文字列処理段階を品詞列処理段階に切り替える働きを、重ねて有する。

空品詞: 辞書中で、きまつた値の X にたいして与えられ、書き換えた結果の T は空となる。

未知語品詞: 二つの識別された文字列ではさまれた、長さ1以上の、識別されない一つの文字列 X にたいして、固有の品詞が与えられる。のちに展開される

文字列 Z の値は X の値の先に1字の空白を付けたものに等しい。

文字列の書き換え： 辞書でこれが指定されると、識別された文字列 X が、品詞 Y ではなく、別の文字列 X' に書き換えられる(第二種のT)。

句： 長さ1以上の Y を、長さ1の Y で書き換える。これには屈折のしかたが幾通りがある。長さ1以上の Y を抹消することもできる。

文法語： 長さ0以上の Y を、長さ1の Y で書き換える。後者の Y は長さ0以上の Z を暗に含む。

2 文字列の処理

2.1 語の識別と有限状態文法

入力記号列 X の要素(単位)は文字である。文字の集合はホレリス文字と呼ばれる、英字、数字、特殊記号、空白を含む48種の文字からなる。それらのあいだには取扱い上、何の区別もない。如何なる類別も、対応もない。

原文の文字列 X は、文字列処理の規則にしたがって逐次処理される。文字列処理の規則の集合は辞書と呼ばれる。

辞書の見出しは原語文字列である。この文字列 R によって識別された文の構成要素 s を語という。辞書の定義はその語の品詞、およびその語を訳文に展開するための訳語文字列である。識別された文字列は語 y に変換され、もとの文字列 X は失われてしまう。

文字列 X を語 y に書き換える規則は、いわゆる有限状態文法にほぼしたがうものと考えられる。また本システムでは、文字、語、文という三種の単位しかみとめていないが、システムによっては、単語というような単位を設定し、それは、語幹およびそれにつながる 0 個以上の語尾からなるとすることもある。その場合の、語幹から語尾、語尾から語尾への連鎖の規則は、やはり有限状態文法にほぼしたがうとされている。

2. 2 文字列の書き換え

文字列 X は普通は語 y に書き換えられるが、そうではなくかたたび文字列 X に書き換えてよい。これがなされたときは、もとの文字列は失われ、新しい文字列が文の文字列の一部となる。だからここでは逐次処理はいわば足踏み状態になる。

書き換えた文字列 X について、重ねて何回でも書き換えが

おこつてもよいが、いつかは語 y に変換されなければならない。

文字列から文字列への書き換えはまったく任意でよい。したがつてこの規則を利用して解析あるいは生成される言語は有限状態言語をこえていてよい。一般の自然言語処理もこの規則によつて行なうこともできよう。

しかしその目的には具合のわるい制限もある。そのうちの四角をあげよう。

まず、展開規則との結びつきがわるいことである。つまり文字列の書き換えは真に書き換えであつて、書き換えられたもとの文字列の値は失なわれ、どのような形でも保存されることはない。書き換えの経過があつては要になつても利用はできず、文の解析につづいて生成をするとか、品詞情報を文字情報に変換するとかいった通常の翻訳には利用できない。

次に、かぎられた反復処理しかできないことがある。右方向の逐次処理で、ある文字列の書き換えが足踏みしてくりかえしているあいだはよいが、いったん処理が右へうつると、その左側の文字列は失なわれてしまい、あとからもう一度見直すことはできない。だから文字列の左端(処理端)は書き直す以外の方法では保存できない。

第三に、あるきまつた文字列が識別されたときには一意の書き換えしかできない。だからきまつた要素からいろいろな語を生成する、ランダムゼネレーションのような実験はできない。

第四に、文字列の識別、書き換えはその値を直接指定してしかできない。文字列中の一部が任意の値でよいとか、ある値を他の値で代入するとか、ある値を名前と呼ぶとかいうことができない。

筆者としては、この文字列の書き換えにそれほど大きな機能を予定してはいたわけではないので、うえのことは別に欠点ではない。つまり文字列の書き換えは、たとえば英和翻訳における、英語の不規則変化形の処理といったような、小さい仕事にだけつかうつもりになっている。

3 品詞列の処理

3.1 メタ言語

品詞列処理の規則は狭義の文法と呼ばれる。文の品詞列の左端の品詞列を文法表と照合する。もし照合がとれなければ

ばいひとつ右にずらした品詞列について照合する。

文の品詞列 γ の中にある品詞列 α , 文法表の見出しと一致したとき, それを識別された品詞列 R という。識別された品詞列 R 中の一つの部分列 β が T で書き換えられる。書き換えが済むと, β の位置の文の左端にもどって次の照合がくりかえされる。

識別された記号列 R 中の部分列 β を, 他の記号列 T で書き換える, この書き換えはまったく任意なものではなく, 記号列の長さや, 展開規則との結び付きによつて, 強く制限されている。

記号列 γ の中で記号列 R が識別されたときの処理の, メタ言語的な約束には七種類がある。そのうち六種類では記号列 β の書き換えがおこり, そのほかの一種類では書き換えはおこらない。

(1) 句構造の代入

β の長さは1以上であつて, T の長さは1である。 T は暗に β の情報をそのまま含み, それかのちに展開される。

$$a \ b \ \underline{c \ d \ e \ f} \ \rightarrow \ a \ b \ \underline{m \ f}$$

(2) 逆順の句構造の代入

β の長さは 2 以上であって, γ の長さは 1 である。表面上, γ に関しては (1) の規則と同じである。 γ は暗に記号列 β の情報を逆順に保存し, それらのちに展開される。

$$a \underline{b c d e f} \rightarrow a b \underline{m f}$$

(3) 語の代入

β の長さは 1 以上であって, γ の長さは 1 である。表面上は (1) の規則と同じである。 γ は記号列 β の情報を抹消し, かわりに固有の値 Z を暗に保有する。

$$a \underline{b c d e f} \rightarrow a b \underline{m f}$$

(4) 語の挿入

β の長さは 0 であって, γ の長さは 1 である。つまり, 識別された記号列 R 中の, 指定したある位置に, 一つの記号 y が割り込む。 y は固有の値 Z を暗に保有する。

$$a \underline{b c d e f} \rightarrow a b \underline{n c d e f}$$

(5) 抹消

S の長さは 1 以上であって、 T の長さは 0 である。つまり、識別された記号列 R 中の部分列 S を抹消する。抹消された値はどのようにしても、回復されることもないし、展開されることもない。

$$a \underline{b c d e f} \rightarrow c a b f$$

(6) 転置

S の長さは 2 以上である。 T の長さは S の長さに等しく、その値は、記号列 S の構成要素を逆順にならべたものに等しい。

$$a \underline{b c d e f} \rightarrow a \underline{f e d c b}$$

(7) 飛び越し

記号列 S の書き換えはおこなない。うちの (1) から (6) までの場合には、書き換えがすむと、記号列 Y と文法表との照合は、記号列 Y の左端にもどって再開される。

しかし、飛び越しの場合には、記号列 R が識別されると、次の照合は記号列 Y 中の S の位置から始られる。

いわゆる結合の強さの優先順位を処理する都合のために、

この指定を開発した。

3.2 評価

うゑに述べた七種類の規則を、メタ言語の立場から評価してみよう。

規則(1)と(2)とは、あきらかに句構造文法、それも見かけ上は、文脈検査型の句構造文法である。おそらく(3)代入と(7)飛び越しともまたそれに属するのであろう。

規則(6)転置は、いわゆる変形規則を満足させるために導入したものである。実際に利用してみると、転置がおこつたことをしるしづけるものが Υ の中に必要なことがあり、それで規則(4)挿入と関係づけて用いることが多い。

規則(5)抹消の意義はあきらかではないが、(3)代入、(4)挿入と関係があるらしい。

規則(4)挿入で気付いたことであるが、この機能はもつと拡大することができる。この規則で挿入することができるのは、固有の値 y 、 Z を有する語だけである。それをたとへば識別された記号列 R 中の任意の記号(または記号列)の値 Υ 、 Z を「複写」できるように拡張すると、このシステムの能力は格段に大きくなる。それによつてたとえば、記号微分、記号計算(括弧をはずす)、鏡像記号列の生成といったよ

うたことが可能になる。

85

現在のシステムでは、これらは一般的には可能ではない。

識別される記号列Rや書き換える記号列Tの指定は、きま
った値についてだけ可能である。それを記号列中の一部の記
号が任意の値でよいというようにすることも考えられ、SNOBOL
やCOMITでは普通の装置となっている。

その応用例としては、たとえば冠詞と名詞とのあいだに何
かがはさまっていて、全体として名詞句になるとか、接続詞
の両側に同じ品詞の句があったときには、全体としても同じ
その品詞の句になるとかがある。

集合の要素としての記号は、相互にまったく独立で、固有
の値を有するものとしてある。実際の言語事象では、大文字
と小文字とか、清音と濁音とかのあいだには、あきらかにあ
る種の関係がある。あるいは動詞の細分類というような例も
ある。記号について、集合の集合とか、部分集合とかいうと
らえ方を導入することも考えられる。

そのほかとくに英語に関係した話題として、一致(呼応)
の処理、多義語、多品詞語、多機能品詞などがあ

4 文字列の展開

品詞列 γ から文字列 Z への展開は、左から右へただ1回だけ行なわれる。訳文の文字列は語を構成する文字からなり、語と語とのあいだでは単純な連結だけかなされ、それ以外の文字の抹消、挿入、変更などは何もなされない。

未知語の直前に空白が1字入る以外には、語と語とのあいだにはどのような区切り文字も挿入されない。

語 y の文字列 Z は、語順 γ 中のその位置に展開されるだけで、そのほかの位置には影響しない。

訳語の語順は品詞処理における品詞の語順処理のすべてを、そしてそれだけを反映している。品詞処理のそれぞれが語順におよぼす影響はほとんど自明であるが、あらためて以下に述べる。

品詞 y には語の品詞と句の品詞とがある。語の品詞には一つの語の文字列 Z が直接に従属している。語の品詞および句の品詞の、1以上の長さの列に代入したのが句の品詞である。つまり句の品詞には句または語の品詞の列 γ が直接に従属している。

文字列の展開にあたって、各品詞の値は何の意義ももたない。そして文字列 Z 中に含まれることもない。すなわち、

語の文字列はどのような品詞の下に従属していてもまったく同じである。そして訳文は、語の文字列と未知語を区切る空白とだけからなる。

(1) 句構造の代入と展開語順

句構造の代入が行なわれると、その句の中の語順は固定され、それ以後変更することはできない。句の中の語順は正方向の展開か逆方向の展開かのいずれかだけを選択できる。だから長さ n の品詞列には $n!$ の順列通りの順序が可能であるが、それが一つの句構造として代入されるときには、二通りの語順のいずれかしか許されない。

(2) 逆順の句構造

句の中の語順を逆方向と指定したとき、その効果は、その句を直接に構成する品詞列についてだけ有効である。すなわち、句を構成する品詞がふたたび句の品詞であった場合に、語順の指定は下位の句同士のあいだの順序はさだめるが、下位の句の内部の語順を変更することはけつしてない。

(3) 語の代入

ある品詞列にたいして語を代入すると、その品詞列に従属す

る語の文字列はすべて抹消される。そしてそのかわりに、ただ一つの語の文字列がその位置を占める。抹消された語が回復することはけっしてない。

(4) 語の挿入

語を挿入するとその位置にその語の文字列が入る。前後の語には、文字位置以外の何も影響をおよぼさない。

(5) 抹消

ある品詞列を抹消すると、それに従属する語の文字列はすべて抹消される。抹消されたことをしるしづけるものは何も残らない。抹消された語はどのような形でも回復することはない。

(6) 転置

ある品詞列を転置すると、それに従属する語や句もまた転置される。転置はそれらの語や句のあいだの語順を逆にするが、それより下位の句の内部の語順を変更することはけっしてない。

転置がおこったことをしるしづけるものは何もないし、転置された品詞列の一部または全部がふたたび転置その他あら

ゆる処理をくわえられてかまわない。

(7) 飛び越し

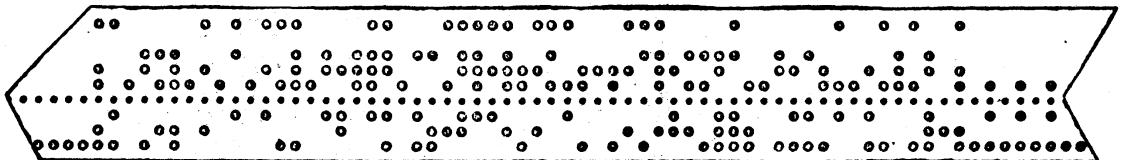
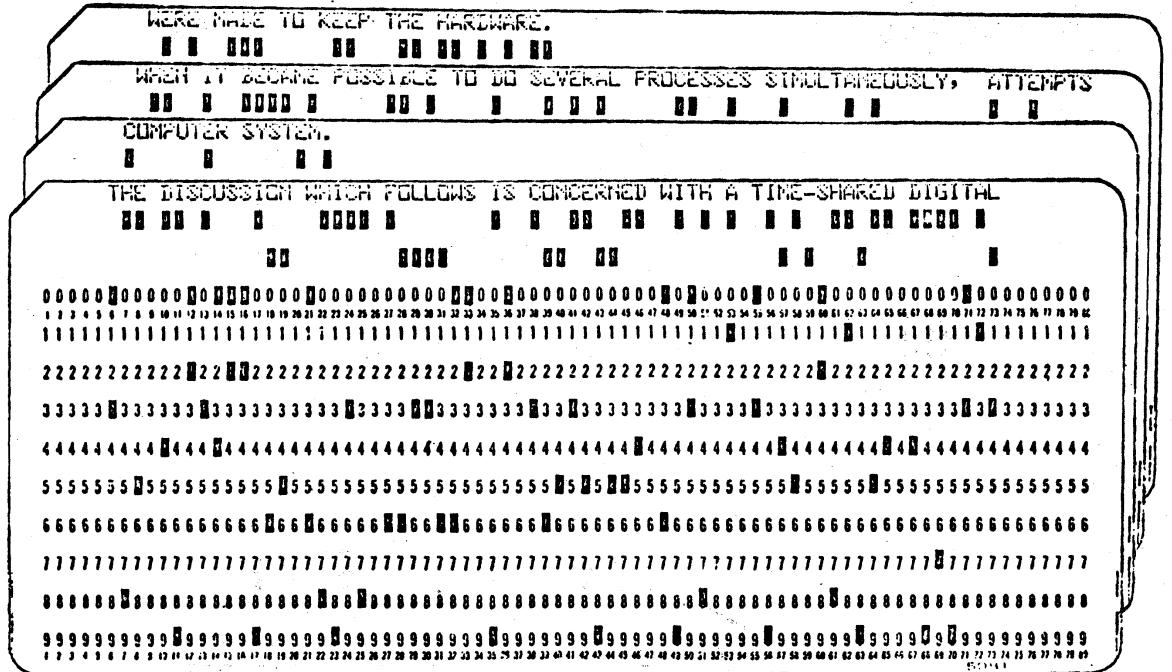
飛び越しは品詞列の照合に関係するだけで、文字列の展開には何も影響しない。

5 例

R				T	S
形	形			形	0 2 代
名				句	0 1 代
形	名			句	0 2 代
句	名			句	0 2 代
句	前	句		句	0 3 逆
句	前	句	名	句	2 2 代
句	前	形			2 Δ 飛
句	前	名		句	2 1 代

The discussion which follows is concerned with a time-shared digital computer system.

When it became possible to do several processes simultaneously (usually, but not necessarily, computation and input-output) attempts were made to keep the hardware (in particular, the central processor) as busy as possible.



後が続くところの議論は時分割計数型計算機組織に関係される。

同時に幾つかの過程をなすことが可能となった時、試みはハードウェアを保つために作られた。