

## 英日機械翻訳システムにおける解析手法について

中瀬 純夫

(株) カテナ・リソース研究所

本稿では、商用の英日翻訳システム *STAR* に用いられている解析メカニズムに関し、その骨子を紹介する。

この解析方式の特徴の1つは、ほかの解析処理に先立って、文脈自由型の統語解析だけが独立して行われるところにある。ここではおもに、表層統語構造に対する解析手法について、*WFS* (*Well-Formed Substring*) の概念を用いた定式化と各種の評価を行う。

そのほか、*STAR* の翻訳過程について、その概略を示しておく。とくに、英文中間構造の作成、語彙規則や構造規則などの適用、莫大な解析アンビグエィティの解消手法等について報告する。

## On Syntactic Analysis Technique in English-Japanese Machine Translation

*Sumio Nakase*

Catena-Resource Laboratories, Inc.

27, Ichiban-chou, Chiyoda-ku,

Tokyo, 102

Japan

This paper presents an outline of the analysis mechanism employed in *STAR*, an English-Japanese machine translation system.

What distinguishes this system from other similar systems is that it performs syntactic analysis with context-free rules separately from any other analysis process. Formalizations and evaluations of this analysis method are fully discussed using the concept of *WFS* (*Well-Formed Substring*).

A brief explanation about the internal modules in the machine translation process is shown. Also reported here are some subjects about a deeper analysis technique, such as the creation of an intermediate expression for a given English sentence, application of the lexical and structural rules, and some strategies for disambiguation.

## 1. はじめに

本稿では、英日翻訳システムSTAR ([1])の解析メカニズムに関する紹介を行う。ここでは、文脈自由型の統語解析について詳しく述べ、その延長上で、翻訳システムにまつわる諸々の要素が示唆される形をとった。

まず、2. でSTARの翻訳部分の構成要素について最低限の紹介を行った上で(図1にその概略を示す)、3. において、表層の統語解析について、その定式化、評価、そのほか詳細に議論する。4. では、言語処理の内容に関わる解析を行う部分について説明する。

言語処理としての扱い方、色々のレベルでのシステム機能に関わる議論、STARの開発、運用にともない累積している諸問題については、直接には今回は扱わない。

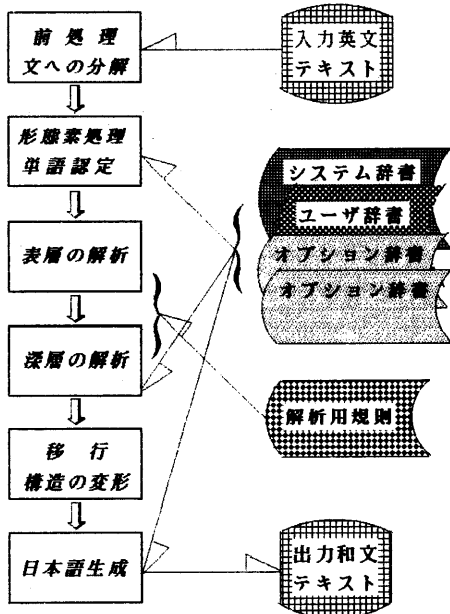


図1 STARにおける翻訳プロセス

## 2 STARにおける翻訳処理の概略

### 2.1 前処理

英文入力テキストを「文」の単位に分割する。それ以降は文単位に処理される。あわせて、最小限の書式情報、前編集記号などの解析を行う。

### 2.2 形態素解析

入力文に含まれる「単語」の認定を行い、それらに辞書情報、形態情報を付与して、単語間の前接/後接関係に基づくWFS構造(ホモグラフ構造)として出力する。

辞書による認定のほか、字句解析による慣用表現などの1語としての認定(LRパーザによる)、未知語に対する語尾推定などを用いた辞書情報の推定などを行う。

### 2.3 表層の解析

STARの解析処理では、表層の統語構造を文脈自由型の規則で捉え、そこで用いられた規則の随伴規則により、格文法的処理のしやすい依存構造を作成する。その構造の上で、辞書の記述や作成された構造そのものに基づいた種々の操作が施される。STARでは、これが2つのフェーズに分けられている。第1のフェーズは文脈自由型規則に基づいた解析だけを行い、全ての可能な解を出力する。この内容は、3. で詳しく述べられる。

### 2.4 深層の解析

表層の解析結果を受けて、解の選択、語彙規則や構造規則などの実行、中間構造の作成などを行う。4. において、その概要が示される。

### 2.5 移行

解析の結果である英文中間構造を評価して可能な範囲での日本語構文構造への変形を行う。述語に対する連用修飾成分の語順の入れ替え、デフォルトとしての日本語送り、活用の指示、否定範囲の変更等々の操作である。

### 2.6 生成

変形された中間構造をトラバースしながら、各単語の日本語辞書記述を参照して日本語文を出して行く。例えば、図7や図8の構造を見ると、自然に日本語が出力できそうなことが想像できる。このとき、各単語に対し、自分自身や共起する語に対する送りや活用のコントロールがなされる。

### 2.7 STARの辞書構成

STARの翻訳処理に用いられる辞書は、翻訳基本辞

```

=====
| - Dictionary Name : mydic
| - Edited by      : nakane
| - Date of Creation: Tue Nov  8 10:27:55 JST 1988
| - Purpose       :
=====
Meeting
--
S nota::EVENT: [ミーティング]
S nota::EVENT: [会議]
S nota::EVENT: [大会]
S nota::EVENT: [会議] [P(wish)N][との]
!report
--
S nota:: [報告]
S nota:: [報告] [P(ton,of)N][に関する]
--
S01::X[X[報告する]]
S01::X[X[通称]である] [SH+HUMAN+] [P(for)N+ORGANIZATION][の]
S01::X[X[報告する]] [P(ton)N][について]
S01::X[X[報告する]] [ON][と]
S01::X[X[入る]] [ON+HUMAN+][に] [P(ton)N][に] [P(for)N][のみ]
S01::X[X[出願する]] [ON+ORSELF+][に] [P(ton)N][へ]
S01::X[X[報告する]] [ON][と] [CV]S[と] | 分詞
S 10c::X[X[報告する]] [OC]S[と] | that 節
S 10c::X[X[報告する]] [ON][と] [CV]S[と] | と。下克上
S 10c::X[X[報告する]] [OV]S[と] | 動名詞

```

図2 ユーザ辞書の記述例

書とオプション辞書とに大別される。前者はシステム辞書がベースとなるが、そのほかに、システム辞書への追加・修正内容を「ユーザ辞書」として与えることができる(図2)。STARでは、ユーザ辞書を作成することにより、システム辞書の全ての内容が変更可能である。

オプション辞書は、図3の例にみるように、英日対訳をベースとした簡便な辞書形式であり、専門語辞書や私用辞書などにはこの形が多く用いられる。

1回の翻訳処理に際し、ユーザ辞書は高々1つ指定できるが、オプション辞書は0個以上任意個(実際はシステムの制約から10数個程度まで)用いることができる。オプション辞書の記述は、骨格の翻訳基本辞書記述の中に、優先された訳語として挿入される。基本辞書にない語については、基本辞書相当の記述が作成される。

これらの辞書データは、翻訳処理に先立って、辞書コンパイラを用いて内部形式に変換するものとする。特に、動詞の活用、名詞の複数などの屈折形は、コンパイルの時点で生成される。

relay (n,*) :	[DEVICE]	继电器
testing :	[ACTION]	テスト
individual (a) :		各々の
continuity :	[PROPERTY]	導通
energized (a) :		励磁
de-energized (a) :		非励磁
ground :		アース

図3 オプション辞書の記述例

### 3 WFSパーズング

#### 3.1 WFSについて

WFS (Well-Formed Substring) は、書換え規則による統語解析の過程で生まれる構文成分を特徴付けるものである。筆者の知る限り、S. Kuno が [2] において、プシュダウン・オートマトン相当のパーザの処理時間の爆発を抑えるためにこの用語を導入している。ここでWFSとは、句構造文法(書換え規則)に与えられたシンボルの一つと、対応する入力文の部分列(開始語位置と終了語位置)とだけにより特定されるものをいう。

たとえば、"I saw a man with a telescope." という文にたいする単純な文脈自由型の解析を例に取ってみよう。このとき、動詞句のWFSとして、"saw a man" という部分列に対応するものと、"saw a man with a telescope" という部分列に対応するものの2つのWFSが発生する。これにたいし、ボトム・アップ的に表現して、後者の「作られ方」は少なくとも2つある。"saw a man" という動詞句に "with a telescope" という前置詞句がついたものという解釈と、"a man with a telescope" という名詞句が動詞 "saw" についたものという解釈とである。しかし、WFSとしてはこの2つは区別されない。このことがWFSの本質的な性格であり、このことによって、文脈自由型文法に対しては、高々入力語数の $n^2$ の処理時間による解析が保証される。

もちろん、そこで用いる文法が非常に詳細なシンボル

を持ち、例えば、「with により修飾済みの動詞句」といったカテゴリーを明示的に持っていれば、それらは区別されて、3つのWFSを持つことになる。逆に、「作られ方」がそれ以降の書換え規則適用に反映するようなアルゴリズムでは、潜在的に莫大な数のシンボル集合を持った文法を用いているとも考えられる。本稿で「WFS」を正面に打ち出しているのは、こうした戦略は避けるべきであるという立場に立っているとも言える。しかし、これは更に吟味すべき基本的な論点であると考えられる。

なお、ここで述べた例は、いわゆる「フローティング・ストラクチャ」の問題とは別問題である。STARにおけるフローティング・ストラクチャの扱いとしては、例えば、動詞句と前置詞句がまとまるという規則を外し、文脈自由型規則の範囲内では常に目的語を前置詞句が修飾するという解析だけを行っておくといった方策に対応している。

#### 3.2 WFS構造

表層構造に対する我々の統語解析のアルゴリズムを定式化するために、「WFS構造」という概念を定式化しよう。有限束、ネットワーク・フローなど、いろいろの定式化ができるが、ここでは、M. Kay のチャート・パーズングの "inactive arc" のように、有向グラフ上のアークでWFSを表現してみる。その時、WFS構造は、ラベルのつけられたアークにより構成され、サイクル(方向を含めてのループ)を持たない有向グラフである。

WFS構造のノードは、アークの間の接点としての役割であり、抽象的なものでよいが、文脈自由型文法の場合は、「語位置」という呼び方で直感的な概念に一致するので、そちらを使うことも多い。

また、とくに、ソース・ノードとシンク・ノードという特別のノードを設定する。それらに対しては、「文頭」及び「文末」という呼び方もする。

WFS構造は1つの書換え規則を前提にする。それを次に文法Gとして示しておこう。なお、実際のインプリメントでは、明示的に与えられた文法をアルゴリズム中で拡大することが多いが、ここでは暗黙的に拡大されたものまで含めてGの中に記述されているものとする。

$$G = (S, R, s_0)$$

S: シンボルの集合(ここでは、ターミナル、ノンターミナルの区別はしないでおく)。

R: 書換え規則(一般的には、0型であつてもよい)

$s_0$ : 開始シンボル(文記号)。

G上のWFS構造 $\Omega$ は、次のようなものである。

$$\Omega = (G, P, p_1, p_2, W)$$

G: 上述の文法。

P: ノード集合(語位置の集合)。

$p_1$ : ソースノード(文頭位置)。

$p_2$ : シンク・ノード(文末位置)。

W: WFSの集合。次のようなラベル付きの有向アークである。

$$w \in W = w = (p_1, p_2, s)$$

$$: p_1, p_2 \in P, s \in S.$$

ただし、Wはサイクルを生まない。すなわち、 $p_2$ から $p_1$ へのパスはない。

本稿で述べるWFSパーズングは、初期WFS構造から出発し、書換え規則の右辺のシンボル列に対応するWFS構造上の任意のパスに対し、同じ始点ノードと終点ノードを持ち、その書換え規則左辺のシンボル列に対応するようなパスを追加することによるボトムアップ・アルゴリズムである。ここで、WFS（あるいは、WFS列）の同定を行うことがWFSパーズングの本質である。

### 3.3 AND-ORグラフによる解析履歴の構造

STARの表層解析のフェーズは、ほとんどWFSパーズングだけを行うという、論理的にはきわめて単純なものになっている。その出力は、WFSパーズングの解析処理の履歴になる。STARでは、その結果の範囲内で再びパーズング（探索）をすることにより、より深い解析やより深い構造の作成を行っている。これに対し、これらを表層の解析と同時に進行するようなインプリメントをしたとしても、以下に述べるような解析履歴の構造を意識しておくことは重要であろう。

1つの入力単語列（ターミナル・シンボル列）に対する文脈自由型パーズングの1つの解析結果（導出結果）は、通常、統語木として表現される。WFS構造は、このシンボル列を拡張した横の構造とすることができる。これに対し、WFSパーズングの経過と結果とは、統語木が重なりあった縦の構造であるということもできる。こうした縦の構造は、AND-ORグラフの形式で自然に表現される。

図4に、簡単な文脈自由型文法による解析結果のAND

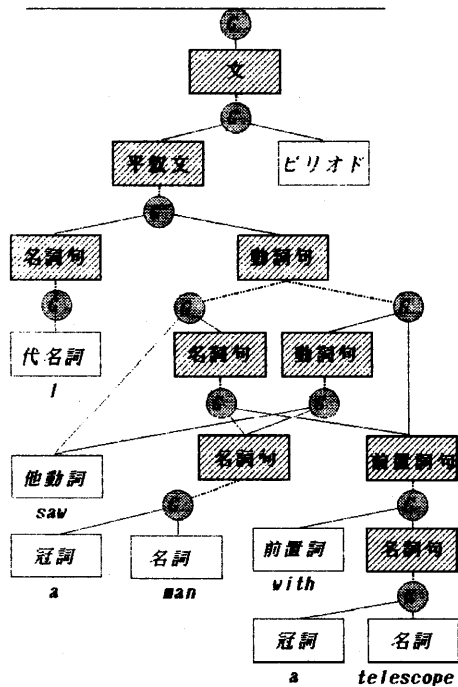


図4 解析結果のAND-ORグラフ

D-ORグラフを示す。但し、これは、解に貢献するWFSだけについての部分構造であって、それ以外は刈り取られていることに注意していただきたい。

ここで、WFS構造のアーキはORノードとして表現される。我々はこれをWノードと称している。Wノードの子ノードは、そのWノードに対応するWFSシンボルを右辺として適用された規則の適用に対応するところのANDノードである。これをGノードと呼ぶことにする。各ノードの子ノードへのリンクは順序づけられている。Gノードの子ノードの並びは、WFS構造での1つのパスに対応する。Wノードの子ノードは、後で述べる解析コストの小さい順に並べられている。

実際のデータ構造では、解のAND-ORグラフは、そこで用いられた規則と、子ノードに対応するWFSの並びとで表現されたGノード（生成情報）の集合として表現する。但し、ターミナル・シンボルについては、後続の処理で単語情報などを利用するために種々の情報をつけ加えておく。一般に、入力WFSに対してなんらかの「生成情報」を与えておくことにより、何段階かにわたりWFSパーズングを行うことができる。ターミナル・シンボルの辞書情報などは、このときの「生成情報」の特別な場合と考えることができる。

### 3.4 文脈自由型のWFSパーズング・アルゴリズム

ここで述べるアルゴリズムは、WFS構造を順次入力して、解析履歴のAND-ORグラフを出力するためのものである。入力のWFS構造には前節で述べたような「生成情報」が付加されているものとする。入力のWFSの条件は次の通りである。

(入力順序の条件)

2つのWFS、 $w_1$ と $w_2$ に対し、 $w_1$ の終端ノードから $w_2$ の開始ノードへのパスがあるなら、 $w_1$ は $w_2$ より前に入力されなければならない。

WFS構造に自然な半順序を導入すれば、「 $w_1$ の終了語位置が $w_2$ の開始語位置より小さいか等しいならば、 $w_1$ は $w_2$ より前に入力されなければならない」でもよい。

次に、アルゴリズム本体を示す。ここでは、入力されるかここで作成されるかするWFSを順に取り出して、1つずつ着目して行く。そのWFSを留めて、それより若いWFSよりなる部分木構造と文法規則集合との照合を行う。ここでは、対応する規則を順に取り出して、その1つずつと照合している。このとき、WFS側は対象WFSが、規則側は対象項が移り変わる。

一応、用いる規則を2型規則だけに限定しておく。0型、1型への拡張も容易であるが、AND-ORグラフの拡張が必要であるなど、若干の手直しが要求される。

- 1° WFSを1つ入力する。なければ、終了する。
- 2° 未着目のWFSを1つ取り出し、着目WFSとし、右辺最終項（最右の項）のシンボルがそのWFSのシンボルに一致する規則の集合を取り出す。未着目のWFSがなくなれば1°へ。
- 3° 規則集合から未着目の規則を1つ取り出し、その右辺最終項を対象項に、着目WFSを対象WFSにする。

未着目の規則がなくなれば $2^\circ$ へ。

- 4° 対象項が規則右辺第1項なら、規則適用は成功で、次のアクションをとったあと、 $6^\circ$ へ移る：

対象WFSの開始語位置と着目WFSの終了語位置とを開始語位置と終了語位置とするWFSで、規則の左辺シンボルをシンボルとするWFSを作る。ただし、同じWFSがあるなら作らない。

そのWFSの「生成情報」(Gノード)に今回の対象WFSの列と規則番号とを追加する(出力する)。

- 5° 対象WFSの開始語位置を終了語位置とするWFSで、そのシンボルが対象項の1つ前の項のシンボルに一致するものを1つとる。それがあれば、そのWFSを対象WFSに、その規則項を対象項にして $4^\circ$ へ移る。

- 6° 対象WFSが着目WFSなら $3^\circ$ へ移る。

そうでないなら、対象WFSと対象項とがそれぞれ1つ前のものに戻り、 $5^\circ$ へ移る。

結果として、文頭を開始語位置とし、文末を終了語位置として、開始記号をシンボルとするようなWFSがあれば、この文は受理されたことになる。このWFSの「生成」のGノードをルートにするAND-ORグラフが解の構造になる。

これがないとき、STARでは、入力文の「極大なWFS分割」を右辺とする仮定の規則を起動して、「救出」を行っている。

本節と同様のアルゴリズムはALICEシステム[3]のためのパーザ(0型規則までサポート)などにインプリメントされたが、我々がWFSの概念を全面的に採用したのはSTARの文脈自由型パーザが初めてである。

### 3.5 所要処理時間と記憶量の評価

まず、入力語数に対する処理時間と記憶量の形式的評価を見てみよう。ここでは、入力は線形のターミナル・シンボル列であると仮定する。こうしても、以下の評価では一般性を失わない。

入力語数を $n$ としたとき、WFSの個数のオーダは高々 $n^2$ である。実際、シンボルの種類を無視すれば、WFSの総数は明らかに $n(n+1)/2$ に抑えられる。シンボルの種類は高々これに対する乗数である。

文法が2型規則だけからなっており、その右辺項数が高々2であるならば、処理時間あるいはAND-ORグラフのサイズのオーダは高々 $n^3$ である。

これを見るために、1種類のシンボル「X」と、1つの規則「 $X \rightarrow X \cdot X$ 」だけからなる文法を考える。この単純化は最大値評価のオーダを変えない。

このとき、 $n-k$ 個ある長さ $k$ のWFSの「作り方」は $k-1$ 通りあるので、AND-ORグラフのGノードの個数は $n(n-1)(n+1)/6$ になる。これはまた、規則適用を試みる回数に一致する。

なお、このときの解の個数(アンビギュイティ)は、 ${}_{n-1}C_{2(n-1)}/n$ になる。これはリーフが $n$ である2分木の個数であり、 $n$ の指数関数のオーダになる。この形の数は、最終的なアンビギュイティには反映しなくとも、

バックトラック・メカニズムだけによる単純な解析における計算回数の評価してみると、この形のコンプレキシティが頻繁に発生することが判るであろう。

ここでは、規則右辺の項数を高々2と仮定したが、一般の場合につき前記のアルゴリズムを適用すると $n^3$ オーダの処理時間は保証されない。これを保証するには、規則を機械的に2項規則へ変換すればよい。それはまた、アルゴリズムにおいて吸収することもできる。それには、WFSのシンボルとして「完成途中」のシンボルを追加すればよい。これがチャート・パーズングの“active arc”などの考え方に一致することは明らかである。

しかし、実際に記述される規則で、右辺項数の多い規則が処理時間の爆発に寄与するという直感に反している。我々は、基本的に前節のアルゴリズムを踏襲し、理論的な上限を $n^3$ で抑えるという方策は採らなかった。

一方、指数関数的な計算回数の爆発は、実際の規則記述によっても発生する。また、統語規則だけによる詳細な解析の結果として莫大なアンビギュイティが発生することも確かめられている。実際、後で示す実験でも、2億近いアンビギュイティを持つ解が受理されている。

したがって、表層の解析と並行して、どこまで可能性のない解釈の刈り取りが可能であるのかということが論点になる。本稿の立場は、(1)そこで十分な保証のある刈り取り手法は確立されていない、(2)そして、前節に述べたような「盲目的な」アルゴリズムにより全ての可能な解釈の探索を行っても十分にリーズナブルな範囲での処理時間と記憶容量で納まるということで、(3)詳細な語彙情報、意味情報、別の枠組みでの文構造表現などを用いた解析は、独立のモジュールにおいて体系化して差し支えないであろうというものである。

最後に、STARのWFSパーザを用いて、英字新聞記事やエネルギー関連のアブストラクトの563文について処理を行った結果を図5に示す。ここには、解析に失敗した文や、語数5未満の文は含まれていない。そのほかには、データの検定などの分析は経ていない結果であることに留意していただきたい。本格的な分析は次の機会に譲る。

なお、ここで用いられているのは、文脈自由型規則として3000弱の規則である。ただし、その各規則は、次節で述べる「不定繰り返し」や「省略」の記述を含む。また、この実験はNEWS830(SONY)を用いているが、動作時の環境は最良のものではない。

この図を見ると、処理時間が入力語数の1.9乗程度と見積もれることが判る。しかし、中には40語前後で40秒以上を要している文も見受けられる。

ここでは、スペースの都合で処理時間の語数に対する関係だけを示したが、そのほかのいくつかの数値の語数との関係についても簡単に述べておこう。

入力のWFS(ホモグラフ)はほぼ入力語数に比例していたが(1.026乗)、これは直感的な予想に一致する。

解析中に作成されたWFSの数はおおそ入力語数の1.5乗強のオーダになっており、1000前後以下に集中している。この数は、WFSパーズングで最小限必要な記憶量に比例する。なお、ここで用いたSTARのバージョンでは、処理可能なWFSの数を8192以下に限っている。ちなみに、この時のGノードを含めたパーズング用の記憶量の合計は約300kバイトである。

解析中に作成されたGノードの個数は、語数の1.65乗程度とかなり少ないものであった。実際に成功する規則適用は、上限値に比べればはるかに少ない。

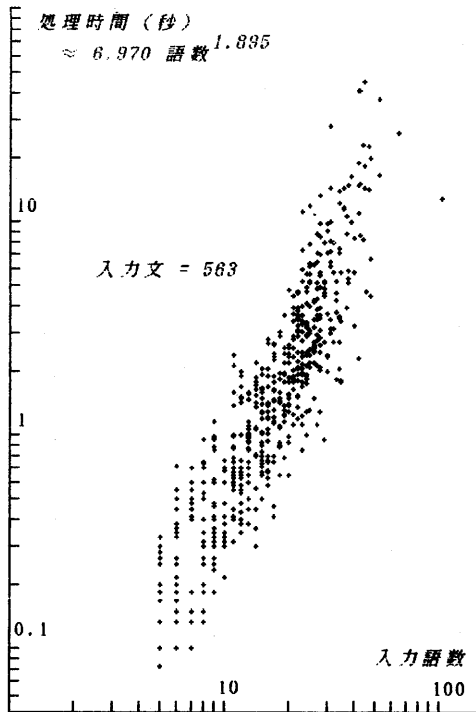


図5 WFSパージングの所要時間

解として取り出されたWFSの個数はちょうど語数の1.5乗であった。ほとんどは100前後以下に納まっているが、1000を越えるものも1つあった。

成功した解析のアンビギュイティを評価すると、語数の3.5乗となったが、これまでのものに比べらつきが非常に大きかった。65語の文で193,251,096ものアンビギュイティを持つものも見受けられた。この文は人名の並列と同格表現がほとんどを占めている文である。これに対し、最大の語数103語の文のアンビギュイティは1万程度である。これらのことは、アンビギュイティが組合せ的な爆発の性格を持ち、規則の書き方や入力文の性格へ依存する度合いが大きいことを示している。

なお、この65語の文では、599のWFSと817のGノードを出力し、103語の文では、432のWFSと556のGノードを出している。

解の中での最小コストの値を評価してみると、入力語数の0.95乗程度であった。これは、コスト関数の線形性から予想されることであるが、線形に密集した系列の左上方に「不自然さ」を持った解が分布している非常に興味深い分布であった。

### 3.6. 表層の解析におけるコスト関数(ウェイト)

アンビギュイティの解消のために、STARでは、ウェイト・メカニズムと呼ぶ機構を用いている。これは、解析に関わる種々のコストを計算し、最終的にそのコスト関数の値の最小なもの(あるいは最小からいくつか)を解として出力するものである。表層の解析においては、その基本値を与える。それは、後続の解析フェーズにお

いて、最適解探索のための推定コスト関数になる。

予め、辞書などに振り当てたウェイト値が計算の基礎となる。それらが、ターミナル・シンボルのWFSのコスト値である。そこから、次のように、WノードとGノードのコスト値が決定される。

Wノードwのコスト:

$$\phi(w) = \min_i \psi(g_i)$$

ただし、 $g_i$ はwの子ノード

Gノードgのコスト:

$$\psi(g) = c_i \phi(w_i) + c_a$$

ただし、 $w_i$ はgのi番目の子ノードで、 $c_a$ は使われた規則の基本ウェイト値、 $c_i$ は使われた規則の右辺i項目のウェイト係数の値。

ここで、STARの文脈自由型規則の記述では、次の例のような記述が行われる。

```
:c VP --> vt1 ( 2:PP ) NP
:b NP --> NP (+a cma NP)* (cma) and NP
```

はじめの例は、他動詞vt1と名詞句NPとから動詞句VPが構成されるが、任意的に前置詞句がそのあいだに挿入されるという規則である。このとき、規則の基本ウェイト値はcにより表現され、ウェイト係数は、PPだけが2で、そのほかは1であることを示している。

次は、名詞句の並列の規則で、単に"and"や"or"を挟んで名詞句が並列するほかに、接続詞の直前にコンマがおかれることや、あいだにコンマに分かれた名詞句が任意個挟まれて、いわゆるリスト構造を構成することもあることを示している。このときの基本ウェイト値は、bの表現する値に不定繰り返し起こった回数だけaの値を足したものである。

ここで、はじめの例のPPのウェイト係数は、間に挿入される前置詞句が長くなればなるほど不自然さが増すことを示し、2番目の例の括弧のウェイト値は、リスト構造が長くなるとコストが増すようなコントロールをしていることになる。

### 3.7. そのほかの付加機構

これまでの、WFSパージングとして、下からの情報の伝播や環境状態の利用は「ストックに」避けるべきことを述べてきた。しかし、STARにおいても、限られた範囲ではその原則からはずれた部分がある。

STARの規則記述では、終端記号は英字小文字3桁で、非終端記号は大文字2桁で記述することになっているが、それらの後に1桁の変数を指定することができる。この変数は、規則の左辺から右辺に合成・伝播し、その過程で規則の不受理を引き起こすことがある。現在のSTARでは、このメカニズムを主語と述語の間の数の一致におもに用いているが、動詞の時制・態などや、形容詞・副詞の級の扱いなどへの拡張が想定される。次の例を参照のこと。

```
NPx --> noux      (※1)
VPz --> vt1z NP   (※2)
NPw --> artz NPx  (※3)
SGw --> NPx VPz   (※4)
```

SE --> SGc prd (※5)  
 VPz --> vtvz SGy (※6)  
 NPx --> NPx and NPx (※7)

終端記号 nou の数は、左辺の NP へコピーされる (※1)。動詞類も同様である (※2)。※3 に対し、NP の「数」が「p」であると、冠詞「a(an)」の「数」は単数「s」であるので、組み込み関数によって「NIL」が戻り、この規則は適用されない。これに対し、NP の数が「s」や単複同形の「o」であると、結果は「s」になり、それは W に代入される。※4 の場合、主語が複数、述語が三単現なら、w は NIL になるが、主語が単数、述語が複数といったケースの場合は、w は「y」という限られた「数」になり全面的な拒否はされない。しかし、通常は「数」"y" は受理されない (※5)。ただし、「order」などの特別な動詞の目的節としてはそうした「数」が受理される (※6)。そのほか、名詞並列で1つでも複数のものがあればそれらをまとめたものも複数になるといった特別の変換が行われる (※7)。

こうした変数値の伝播は、シンボル種類の拡大を意味する。すなわち、変数値まで含めて WFS の同定が行われ、「数」の違うシンボルは別の WFS と見なされる。

次に、STAR では、文脈自由型の規則適用に反しているかのような動きをするものがある。例として、名詞の同格表現の次の規則を考えてみよう。

NP --> NP cma NP cma

実際の文章では、句切り記号(ピリオド、コンマなど)の直前にこの表現が現れたときには、規則右辺の最後のコンマは省略される。そこで、次のような文脈依存型の規則も必要になってくる。

NP prd --> NP cma NP prd  
 NP cma --> NP cma NP cma

実際には、これを文脈自由型の規則だけで表現することも可能であるが、それは莫大な規則数の増大を引き起こす。STAR では、直後にピリオドやコンマのあるとき、先の規則の最終項を省略した規則が適用可能であるという特別手続きをアルゴリズムに追加している。

そのほかにも、WFS パージングの「ストイシズム」に従わなくとも対応できるものはいくつか残されている。従属節などでの「スラッシュ要素」(「ミッシング・ストラクチャ」)、「フローティング・ストラクチャ」などがその代表であり、我々の検討事項になっている。

最後に、プログラミング的な特別処理として、次の2点だけに触れておく。一つは、アルゴリズムの説明でも触れた、解析失敗の場合の極大 WFS 分割をまとめる仮想規則の適応であり、もう一つは、G ノード・テーブルがあふれた場合の「ガベージ・コレクション」である。後者は、G ノードの選択肢の中で、最小コストの WFS 以外を棄却することで、一時的に G ノード個数を WFS 数にまで絞りこむ操作である。

我々はパージングの効率化を含む種々の実験を行ってきた。例えば、規則を予め探索して、環境による到達可能性を用いた WFS 抑制や規則適用の抑制といったメカニズムである。しかし、ある程度コントロールされた規則記述を用いているからでもあろうが、我々の試行では、それらの効果は高々2、3割の処理時間の短縮にしかな

らず、逆に、解析に失敗する文に対しては、小さい WFS しか作成されないことになる。その結果、解析失敗の時の翻訳結果の質が低下してしまう。ハードウェア性能の向上などを考慮するとき、この程度の効率化は余り有効でないと考ええる。

むしろ、我々の現在の関心は、より「深層」の解析との融合を「ストイックに」避けている WFS パージングのメカニズムが本質的に失っている要因をいかにして補って行くかにある。とくに、語彙レベルの個別情報を表層の解析で活用することは重要な検討課題であろう。

#### 4. より「深い」解析

##### 4.1 中間構造の作成

STAR の解析処理の中核で用いられるのは、文脈自由型文法に基づく統語構造ではなく、依存文法的な格支配構造である。また、日本語の生成に至るまで、原則として同じ形の構造が用いられる。この構造の原型は、WFS パージングに用いられた各々の規則の「随伴規則」により作られる。言い替えると、英語の解析規則は、表層の統語構造を把握するための条件部と、中間構造を作成する実行部との組になっているといってもよい。

図6に、簡単な文脈自由型規則と図で表現した随伴規則の組を示す。図7にこの規則の適用により作成された構造を示す。また、図8に埋め込み文のイントロデュサなどが下位の構造に移動する規則によった構造を示す。

STAR の中間構造は木構造をベースにしているので、共通の支配語を持つ要素の並列や形容詞節などでの明示的な支配構造表示はできない。

##### 4.2 語彙規則の適用

STAR の辞書記述では、各単語は品詞細分類 (SSC: Syntactic Sub-Category) 毎の記述の並びで構成されている。各 SSC のもとで、訳語の違いや、用法の違いなど、なお複数の語彙記述がなされる。その記述の、日本語に関わる部分(訳や日本語送り・活用などの情報)の違いを無視したものを「英語文型パターン」と呼ぶ。解析処理では作成中の中間構造を評価しながら、英語文型パターンの決定を試みる。ここでは、その適合性を数値で表現し、最も好ましいものを採用することになる。また、その結果は文全体の解析のコスト値にも反映する。

その典型的な例は、前置詞や副詞小詞との共起記述である。ここでは、条件に合致しない英語文型には大きいペナルティが課される。これにより、WFS パージングでの解の適・不適が逆転する事も少なくない。

さらに、着目している語に共起する語の属性値によって文型を選択することもある。STAR では、各単語の各個別記述に与えることのできる属性を内在素性と呼び、そこに与える記号をフィーチャ(またはフィーチャ・コード)と呼んでいる。一般的には、これは意味マーカなどと呼んでいるものに当たるが、現在の STAR では、フィーチャの値や体系は利用者に委ねられており、既存のものとの重複などによる混乱のない限り、どのような内容のフィーチャ体系を定義してもよい。

例えば、次のような記述をすることができる。

"take": [乗る] [SN<HUMAN>] [ON<VEHICLE>] {}  
 "catch": [ひく] [ON<KAZE>]  
 "kill": [自殺する] [ON<ONESELF>] {}

"look": {する}  
 [P(after)N<ANIMATE,HUMAN>] [の世話を]  
 "a": {一人の} [-N<HUMAN>]  
 "plane": :::VEHICLE: {飛行機}  
 "cold": :::KAZE: {風邪}  
 "himself": :::ONESELF: {}

フィーチャの制限を行うと、着目語の文型パターンが評価されるだけでなく、相手の語の文型パターンも限定されることになる。なお、STARの辞書では、英語の表層格を記述し、深層格は陽には現われない。

そのほか、「受動態の時にのみ優先的に利用すべき記述」などの指定が評価される。また、前置詞句に対する共起指定は、「仮に」別の構文要素にかかっているフローティング・ストラクチャに対しても試み、これが採用されるとその係先が決定する。

語彙規則による共起要素への関係は、英文中間構造の上に新しいリンクを作成する。例えば、上の"take"の例では、"take"から主語のノードと目的語のノードへアークが張られる。これは、母体の木構造とは独立の構造であり、このアークに沿って、着目語から対象語への日本語の活用、送りなどのコントロールが行われる。

#### 4.3 ディスアンビギュイティ

STARにおける解析のアンビギュイティの解消は、WFSパーズング結果のAND-ORグラフの上で最小コスト問題を解くことに帰着される。"analyser"と名付けられた、「より深い解析」を行うSTARのモジュールは、WFSパーズングの結果の推定コストが与えられたAND-ORグラフ上で、随伴規則、語彙規則、構造規則（解析上必要な範囲での構造の変形、情報の伝播、並列構造などの評価）を用いた中間構造の作成・評価を行いながら、最尤解析の探索を行うものである。このモジュールは、最終コスト値の最小のものから、指定された順位の解（英語中間構造）までを出力する。

現在のSTARでは、この探索処理を原則としてトップ・ダウン・パーズングの枠組みで、バックトラックを行いながら遂行する。このフェーズでだけは、探索回数爆発を避けるためのヒューリスティックスを導入しているが、希に多大な処理時間が要求されることがある。

#### 5. おわりに

ここで得られた諸々の成果は、STARの研究、開発、運用、実用化に色々の形で携わり、協力していただいた多くの方々に負うところが大きい。この場を借りて謝意を表したい。

#### [参考文献]

- [1] (株)カテナ・リソース研究所:「英日翻訳システムSTARシステム説明書」、1988
- [2] Kuno,Susumu: The Predictive Analyser and a Path Elimination Technique, Communication of the ACM, Vol. 8 No. 7, 1963
- [3] Nakai,Hiroshi et al.: "ALICE" A System for Analysing the Linguistic Structure of Japanese Sentences, COLING 80 資料, 1980

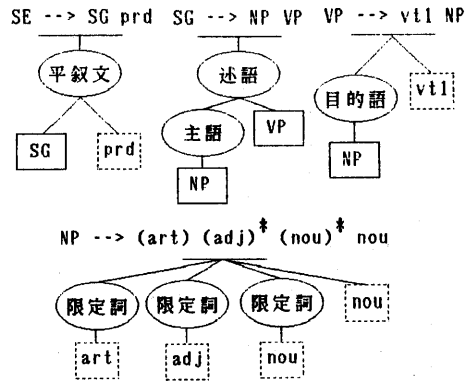


図6 構造作成規則

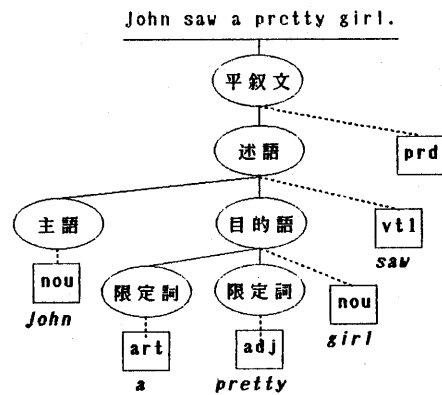


図7 随伴規則による中間構造

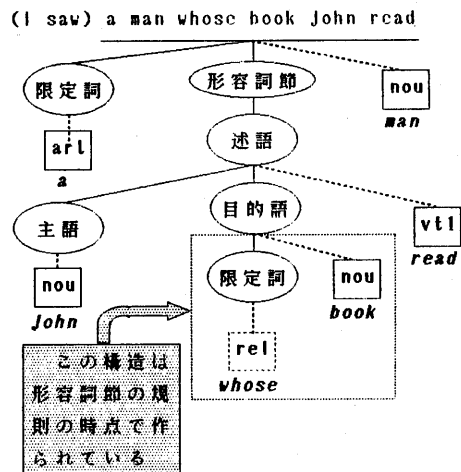


図8 繰り下げられた部分構造