

# 最長一致法と文節数最小法について

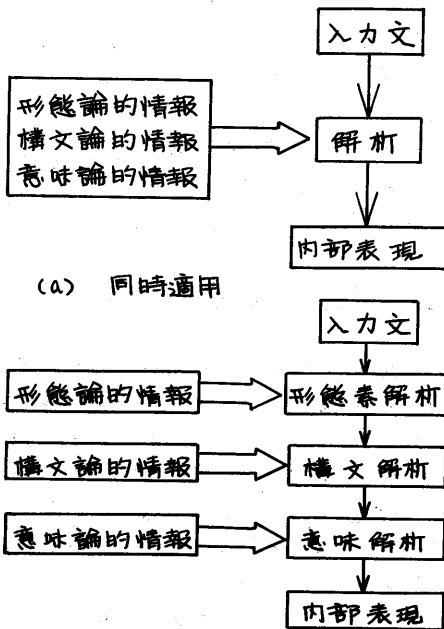
吉村賢治 日高 達 吉田 将  
(九州大学 工学部)

## 1. まえがき

電子計算機による自然言語理解の過程は、システムが持っている言語的情報の適用の仕方からいくつかの型に分類できる。図-1(a)はシステムが持っている形態論的情報、構文論的情報、意味論的情報を同時に使用する例であり、図-1(b)はこれらの情報を段階的に使用する例である。(a)は人間が行っている方法であるが、解析の途中で不十分な意味的情報を強く反映させることは、誤った解析に導く原因になることが考えられる。また、入力文が仮名又はローマ字で書かれた日本語文の場合、解の探索空間は非常に大きくなり、これに対して意味解析などの複雑な処理を行うことは効率の低下を

もたらすと考えられる。従って、システムの持つ意味論的情報が不十分な場合には、少なくとも形態素解析や構文解析などの純語論的解析と意味論的解析を分離して段階的な構成が望まれる。

以下、段階的な構成として図-1(b)の方式を用いた場合について考察する。このとき、各解析段階における解の探索空間の変化を図-2に示す。集合 $S_1$ は形態論的制約を満たす解の集合であり、集合 $S_2$ は構文論的制約を満たす解の集合である。一般に、日本語文の形態素解析に用いられる形態論的制約は、正規文法のクラスト属する規則で表現される。この規則は明確である反面正しい日本語文に対する制約が弱い。そのため、入力文として仮名又はローマ字表記の日本語文を用いた場合、集合 $S_1$ は多くの誤った解析を含む大きな集合となる。このような条件のもとで効率の良い解析を行うためには、集合 $S_1$ の尤度(もっともらしさの度合)



(b) 段階的適用

図-1 情報の適用例

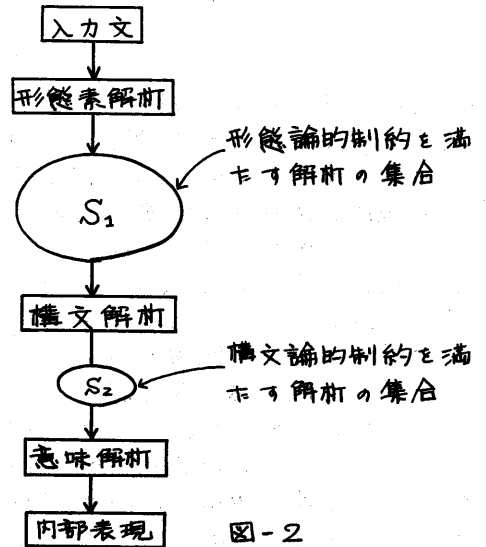


図-2 解の探索空間

が高い解析結果から構文解析を適用する必要がある。

本稿では、集合  $S_1$  の各要素に対してその尤度を評価するヒューリスティック情報である最長一致法と文節数最小法について考察し、両者の比較を行うために行った実験について報告する。

## 2. 最長一致法

従来、このヒューリスティック情報としては、入力文字列のある位置から始まる単語を認定する段階で複数の候補が存在する場合、入力文字列における単語の長さが長いものから優先する最長一致法が用いられている。この手法を用いた解析手順は、形態論的制約を正規文法としてとらえた場合、適用可能な書き換え規則に優先順位を付けて、その順位の高い書き換え規則から優先して適用するという手順になる。すなわち、 $X, Y_i$  を非終端記号、 $w_i$  を終端記号として  $w_i$  に対応する入力文字列の部分列の長さを  $|w_i|$  で表わすとき、ある時点で  $m$  個の書き換え規則、

$$P_1: X \rightarrow w_1 Y_1$$

$$P_2: X \rightarrow w_2 Y_2$$

⋮

$$P_m: X \rightarrow w_m Y_m$$

が適用可能ならば、 $|w_i|$  の値が最大となる書き換え規則  $P_i$  から優先して適用する手順となる。これは書き換え規則に優先順位を付けた文脈自由言語のバックトラック方式の構文解析アルゴリズムと同じ手順である<sup>(1)</sup>。このアルゴリズムの能率については文献(1)で議論されているが、その議論において文脈自由文法を正規文法で置き換えても正規文法の特長は活かすことができず、正規文法に対してと同様な結論を得る。従って、最長一致法を用いた場合の能率は入力文字列の長さ  $n$  に対して、最悪の場合に必要とするステップ数は

$O(n^2)$ 、メモリ数は  $O(n)$  となる。ここで  $C$  は  $n$  に無関係な定数である。

ここでは最長一致法は多くのシステムに用いられているが、二つの問題点が考えられる。第一の問題点は、長一致した単語を優先する根拠が明らかでないことである。最長一致法の歴史を調べると、海外の文献では英単語の語幹抽出アルゴリズム (stemming algorithm) において用いられている<sup>(2)</sup>。この場合は、接辞辞書の検索回数を減らすことが目的となっている。国内の文献では、文節単位に分ち書きされた入力文を対象とする仮名漢字変換において自立語の認定に用いられている<sup>(3)</sup>。ただし、この場合の最長一致法は尤度の評価としてではなく解析の順序付けのために用いられている。このように短い単位に対して各々の根拠を持って用いられてきた手法が、バテ書き日本語文に対して適用されていると考えられる。

第二の問題点は、単語の長さという局所的な対象を評価しているため、文全体に対する尤度の評価になっていないと考えられることである。

以上のような問題点はあるが、必要とするメモリ数が少ない、プログラミングが容易などの利点も持っている。

## 3. 文節数最小法

解析結果の尤度を解析結果の文節数を用いて評価する手法が文節数最小法である。すなわち、解析結果に複数の候補がある場合には文節数が少ないものから優先する手法である。このヒューリスティック情報を用いる場合には、表方式のアルゴリズムを使用することができ、その能率は最悪の場合に必要とするステップ数、メモリ数ともに入力文の長さ  $n$  に対して  $O(n^2)$  となる<sup>(4)</sup>。また、特別の場合として、文節数が最

小となる解析のみを求める場合の能率は、最悪の場合に必要とするステップ数、メモリ数ともに  $O(n)$  となる。

日本語文の統語規則は、文節を構成する単語の並びを規定する形態論的規則（文節構造規則と呼ぶ）と文を構成する文節間の係り受け関係を規定する構文論的規則（係り受け構造規則と呼ぶ）として与えられることができる<sup>(5)</sup>。従って、文節構造規則を満たす解析が正しい日本語文であるためには、少なくとも係り受け構造規則を満たさなければならぬ。係り受け構造規則は文節間の意味的呼応関係を規定するものであるから文節数が少ない解析結果の方が係り受け構造規則を満たす確率は高いと考えられる。これが文節数で尤度を評価する根拠である。このように考えると、最長一致法は解析結果の文節数を局所的な単語の長さで近似的に評価しているとみなすこともできる。

文節数最小法は文節数で尤度を評価しているため、文全体の尤度を評価できると考えられる。

#### 4. 実験

2., 3. で最長一致法と文節数最小法について簡単に考察したが、このようなヒューリスティック情報に基づいた手法の有効性を理論的な側面だけで与えられることは困難である。本章では、これら二つの手法を比較するために行った実験の概要を示し、5. でその結果を報告する。

##### 4.1 述語の定義

次節で実験に用いた最長一致法のアルゴリズムを示すが、本節では、その記述を簡明にするために幾かの述語を定義する。なお以下の議論では  $n$  長さの入力文字列、

$s = s(1) s(2) \dots s(n)$   
を固定する。

##### [定義1] 単語構造

単語  $w$  の綴り  $W$ （活用語の場合は終止形の綴り）、品詞  $H$ 、活用情報  $K$  からなる3項系列  $(W, H, K)$  を単語  $w$  の単語構造と呼ぶ。□

##### [定義2] 述語 $WS, J, E, C$

- i) 入力記号列  $s$  の部分列  $s(i+1)s(i+2)\dots s(j)$  に対して、 $w = s(i+1)s(i+2)\dots s(j)$  なる単語  $w$  が存在し、その単語構造が  $\alpha$  であることを  $WS(i, j, \alpha)$  で表わす。
- ii) 単語構造  $\alpha$  が自立語の単語構造であることを  $J(\alpha)$  で表わす。
- iii) 単語構造  $\alpha$  の単語が文節末の語になりうることを  $E(\alpha)$  で表わす。
- iv) 単語構造  $\alpha_1$  の単語  $w_1$  と単語構造  $\alpha_2$  の単語  $w_2$  の連接  $w_1 w_2$  が文節の部分列になり得ることを  $C(\alpha_1, \alpha_2)$  で表わす。□

これらの述語を用いて、一般に用いられる形態論的制約は次の文節構造規則として定義される。

##### [定義3] 文節構造規則

$n$  長さの入力記号列  $s$  に対して、  
 $i_0 (= 0) < i_1 < \dots < i_m (= n)$   
なる整数  $i_0, i_1, \dots, i_m$  と単語構造  $\alpha_1, \alpha_2, \dots, \alpha_m$  が存在して、次の i, ii, iii を満たすとき  $s$  は左文節を成すといひ、さらに iv を満たすとき文節を成すといふ。

- i.  $WS(i_{l-1}, i_l, \alpha_l)$  ( $l = 1, 2, \dots, m$ )
- ii.  $J(\alpha_1)$
- iii.  $C(\alpha_{l-1}, \alpha_l)$  ( $l = 2, 3, \dots, m$ )
- iv.  $E(\alpha_m)$  □

ただし、i, ii, iii, iv の条件は文節を成すための必要十分条件ではなく、必要条件にしかすぎない。

##### 4.2 最長一致法のアルゴリズム

最長一致法には単語最長一致法、文節最長一致法、二文節最長一致法などのいろいろな方式がある。そこで議論を明確にするため、実験に用いた最長一致法のアルゴリズムを示す。文節数

最小法を用いた表方式のアルゴリズムについては既に報告しているので本稿では割愛する(4)。

アルゴリズムの記述において、プッシュ・ダウン・スタックをSTACKとし、ANSLISTは整数 $i, j$ と単語構造 $\alpha$ からなる3項系列 $(i, j, \alpha)$ のリストとする。また、 $\Gamma(i), \Gamma_F(i), \Gamma_S(i)$ は同様な3項系列の集合であり、 $\Gamma_F(i), \Gamma_S(i)$ は次のように定義される。

[定義4] 集合 $\Gamma_F(i), \Gamma_S(i)$

$$\Gamma_F(i) = \{(i, j, \alpha) \mid i < j \leq n, WS(i, j, \alpha), \neg J(\alpha)\}$$

$$\Gamma_S(i) = \{(i, j, \alpha) \mid i < j \leq n, WS(i, j, \alpha), J(\alpha)\} \quad \square$$

集合 $\Gamma_F(i), \Gamma_S(i)$ を求めることは、各々付属語辞書検索ルーチン、自立語辞書検索ルーチンに入力記号列 $\alpha$ の部分列 $\alpha(i+1)\alpha(i+2)\dots\alpha(m)$ を与えて、その最左部分列としての付属語、自立語の位置及び単語構造を全て求めることに相当する。

[アルゴリズム]

初期状態 ANSLIST = NIL, STACK =  $\phi$   
 入力 記号列 $\alpha = \alpha(1)\alpha(2)\dots\alpha(n)$   
 出力 ANSLIST

- (1)  $\Gamma_S(0)$ を求める。
- (2)  $\Gamma_S(0)$ の各要素の優先順位を求める。優先順位が最高の要素 $(0, j, \alpha)$ について、

$$ANSLIST = ANSLIST \cdot (0, j, \alpha)$$

とし、残りの要素を優先順位の低いものから順にSTACKにプッシュする。

(以下ANSLISTの末尾を $(i, j, \alpha)$ とする)

- (3)  $j = n$ ならば、  
 $E(\alpha)$ ならばANSLISTを出力して停止する。  
 $\neg E(\alpha)$ ならば(9)へ。

- (4)  $\Gamma(j) = \text{中}$ とする。

- (5)  $\Gamma_F(j)$ を求め、 $\Gamma_F(j)$ の各要素 $(j, k, \beta)$

について $C(\alpha, \beta)$ ならば

$$\Gamma(j) = \Gamma(j) \cup \{(j, k, \beta)\}$$

とする。

- (6)  $E(\alpha)$ ならば $\Gamma_S(j)$ を求め、

$$\Gamma(j) = \Gamma(j) \cup \Gamma_S(j)$$

とする。

- (7)  $\Gamma(j) = \text{中}$ ならば(9)へ。

- (8)  $\Gamma(j)$ の各要素の優先順位を求める。優先順位が最高の要素 $(j, k, \beta)$ について、

$$ANSLIST = ANSLIST \cdot (j, k, \beta)$$

とし、その他の全ての要素を優先順位の低いものから順にSTACKにプッシュする。 (3)へもどる。

(バックトラックの手続き)

- (9) STACK = 中ならば、ANSLIST = NILにしてANSLISTを出力したのち停止する。

- (10) STACKの先頭をポップする。これを $(l, m, \beta)$ とする。

- (11) ANSLISTの末尾の項目の第2項目が $l$ に等しくなるまでANSLISTの末尾の項目を取り去る。

- (12) ANSLIST = ANSLIST  $\cdot (l, m, \beta)$ にして(3)へもどる。  $\square$

アルゴリズムにおいて記号'.'は連接の操作を示す。例えば、ANSLIST = (A B C)のとき、

$$ANSLIST = ANSLIST \cdot D$$

を実行した結果はANSLIST = (A B C D)となる。

アルゴリズムを実行した結果、ANSLIST = NILならば入力記号列は解析不能であり、その他の場合はANSLISTに入力記号列の形態素解析の結果が得られる。ここで示したアルゴリズムは一つの解析だけを求めるアルゴリズムであるが、複数個の解析が存在する場合は(3)でANSLISTを出力したのちに(9)から実行を再開することにより全ての解析を求めることができる。

ステップ(2), (8)における優先順位の評価は、

- i) 自立語よりも付属語を優先する。
- ii) 入力文における単語の長さが長いものを優先する。

iii) 活用する語を優先する。

とした。ここで評価の順序は i, ii, iii の順である。

このアルゴリズムは一般的なバックトラック方式のアルゴリズムであるから詳しい説明は省略し、ANSLISTの意味について述べる。

【定義5】 述語C

単語構造  $\alpha, \beta$  に対して  $C(\alpha, \beta)$  であることは、 $C(\alpha, \beta)$  であるか、又は  $E(\alpha)$  か  $J(\beta)$  であることと等価である。

【定義6】 述語LB

4変数述語LBを次のように再帰的に定義する。

i)  $WS(i, j, d)$  か  $J(\alpha)$  であることと  $LB(i, j, d)$  であることは等価である。

ii)  $LB(i_1, i_2, d)$  か  $WS(i_2, i_3, \beta)$  か  $C(\alpha, \beta)$  ならば  $LB(i_1, i_3, \beta)$  である。

iii) 述語LBは i, ii で定義されたものでなければならない。

整数  $i$  と単語構造  $\alpha$  が存在し  $LB(i, j, d)$  のとき入力記号列  $\alpha$  の部分列  $\alpha = \alpha(1)\alpha(2)\dots\alpha(j)$  は左文節列をなすといひ、さらに  $E(\alpha)$  であるような  $\alpha$  が存在するときは右文節列をなすという。

【定理1】

ANSLISTの末尾の項目が  $(i, j, d)$  ならば、 $LB(i, j, d)$  である。

証明は  $i$  に関する帰納法で容易にできるので省略する。

4.3 単語辞書

実験に用いた単語辞書は自立語辞書と付属語辞書からなる。自立語辞書には約8,300語の自立語が収録されており、付属語辞書には約300語の付属語、形式名詞、補助動詞が収録されている。本稿では単語辞書のデータ構造の詳細については述べないが、一回の検索ルーチンの実行で集合  $F(i)$ ,  $F_j(i)$  を能率良く求めることができるデータ

構造を用いている。

4.4 入力文

実験に用いた入力文としては、武者小路実篤の「人生論」から1,000文を用意した。入力文は仮名表記のバカ書きであり、一文の平均長さは44文字で、平均20文字毎に読点「、」で分かれ書きされている。

5. 実験結果

本章では4.で述べた最長一致法と文献(4)で述べた文節数最小法の比較を行った解析実験の結果について報告する。

5.1 文節数最小法の有効性

1,000文の入力文について実験した結果、960文は文節数が最小となる解析の中に正解が存在した。ここで正解とは原文と一致する解析である。ただし、原文との比較は品詞、活用情報の

原文	解析	件数
その人	其の日と	13
良くして	浴して	5
子が無い	漕がない	2
気がする	起臥する	2
〜の無い人	〜の内皮と	2
この蚊の	子のかの	1
その産む能力	園生無能力	1
死のそのもの	始祖のもの	1
見て嫌に	未更やに	1
もうあとに	盲啞とは	1
具合良く育って	具合よ養って	1
目匠者	名卑	1
本屋が良い本を	本屋がよ異本を	1
招き易い点で	招きや水天を	1
そう太いして	早速した	1
しかし金と	鹿レかぬと	1
どうかすると	同化すると	1
一フ一書く知はない	一フ一か基はない	1
どう言う人を	同意兩飛とを	1
甲の見た美が	甲のみ旅が	1
男が何か	男仮名とか	1

表-1

みを対象とし、綴りの区別はしないものとする。例えば、'さかい'に対しては'機械'、'機会'などの名詞が対応するが、これらの区別はしていない。文節数が最小となる解析に正解が得られなかった解析を表-1に示す。ここで'|'は文節の切り目を示す記号である。表-1において原文'その人'の正解が得られなかった原因は、自立語辞書に一単語として'其の日'が登録されていたことにある。このような複合語を自立語辞書に登録する場合には、文節数に相当する荷重の情報も付加する必要があると思わゆる。表-1からも明らかかなように実験した1,000文に対しては、正解は文節数が最小となる解析か、その次に文節数が少ない解析の中へすべて含まれている。このことは文節数最小法の有効性を十分に示している。

文節数が最小となる解析には、入力文字列の長さ10文字に対して平均2通りのあいまいな解析が存在する。

### 5.2 能率の比較

図-3に入力文字列の長さに対して、解析に必要としたメモリ量のグラフを示す。

文節数最小法に対しては、文節数が最小となる解析のみを求めるのに要するメモリ量が示してある。

現在の計算機技術においては、本稿で述べた形態素解析の処理時間はほとんど二次記憶の検索時間で費される。そこでステップ数の比較は自立語辞書の検索回数との比較で行った。入力文字

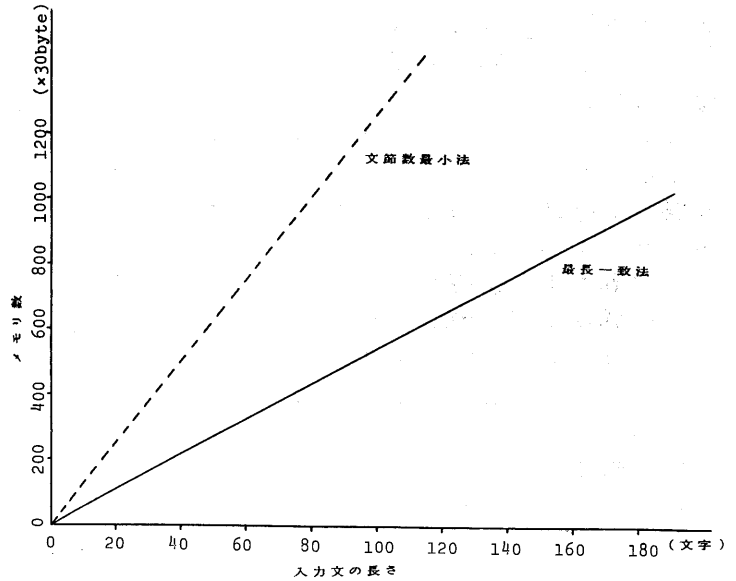


図-3 メモリ量の比較

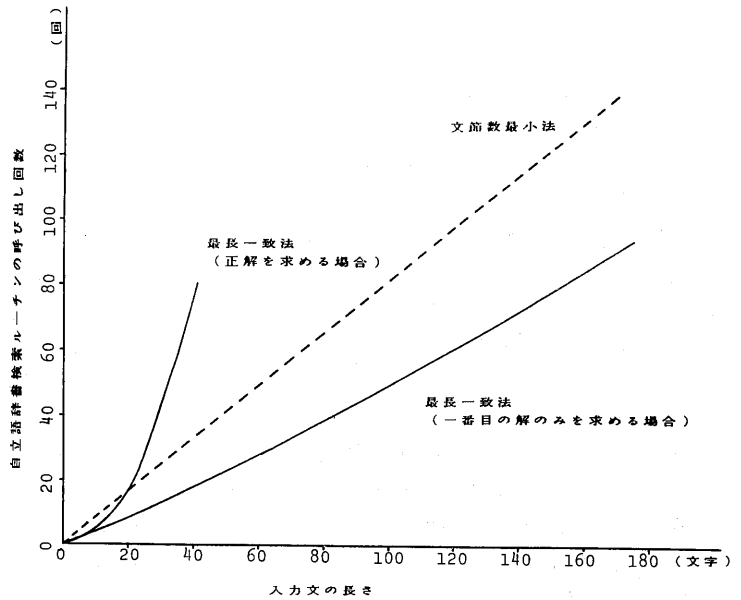


図-4 自立語辞書の検索回数との比較

列の長さに対して、解析中に行った自立語辞書の検索回数のグラフを図-4に示す。最長一致法については一番目の解を求めるのに要した検索回数と正解を求めるのに要した検索回数を示す。

### 5.3 解析の品質の比較

解析の品質の比較として、ここでは一番目に出力される解析の誤った部分

の入力文列における長さも比較する。誤った部分の入力文列における長さの入力文列長に対する割合の平均は、最長一致法で12.4%，文節数最小法では7.0%である。ただし、文節数最小法では第一番目の解析を選択するためには次のような評価をした。

- i) 文節数が少ないものを優先する。
- ii) 自立語の次は付属語を、付属語の次は自立語を優先する。
- iii) 活用のある語を優先する。

ここで評価の順序は i, ii, iii の順である。

#### 5.4 尤度の評価

2.で最長一致法は文全体に対する尤度の評価ができないことを述べたが、30個の入力文について、正解が一番目と二番目に出た出力文の尤度を比較したものを表-2に示す。この場合の文節数最小法による尤度の評価は5.3と同様である。

	最長一致法	文節数最小法
第一番目	73	124
第二番目	3	28

表-2 第一番目、第二番目に出た出力文の解析に正解が得られた文数

#### 6. あとがき

5.で示した実験結果より、第一番目の解を求めるのに要する処理時間、必要とするメモリ数の点では最長一致法が優れている。しかし、正解を求めるまでに要する処理時間や解析の品質、つまり尤度の評価の点では文節数最小法が優れている。これまで形態素解析の解析結果に対する尤度の評価には、ほとんどの場合最長一致法が用いられてきたが、文節数最小法も同様にそれ以上の能力を持っている。今後その利用法などに関する研究が必要であろう。

最後に、本実験に使用した自立語辞書の作成は初期の頃から着手し、以来

十数年間に亘り専念して来られた九州芸工大稲永結之講師、小西彬充即身、最長一致法の実験に関して多大な御協力をいただき、また日本ユニバックスの吉田正行氏をはじめ実験システムの実現のために多くの方々の協力をいただき、深く感謝の意を表す。

なお、本研究は昭和56~57年度21世紀文化学術財団学術奨励金および昭和56~57年度文部省科学研究費一般研究(B)によった。

#### 7. 参考文献

- (1) Aho, A.V., Ullman, J.D.: "The Theory of Parsing, Translation and Compiling, Volume 1: Parsing", Prentice Hall, 1972.
- (2) J.B. Lovins: "Development of Stemming Algorithm", MT, Vol. 11, 1968, 22-31.
- (3) 栗原, 黒崎: "仮名文の漢字混り文への変換について", 大工集, Vol. 39, 4, 1967, 659-664.
- (4) 吉村, 日高, 吉田: "表方式を用いた文節構造分析アルゴリズムとその効率について", 情報処理, 計算言語学研究 25-6, 1980.
- (5) Hitaka, T., Yoshida, S.: "A syntax parser based on the case dependency grammar and its efficiency", Proceedings of the COLING 80, 1980, 295-302.