

入力インターフェース用日本語サブセットの構成法とその応用

萩原 正也 徳田 雄洋
山梨大学 工学部 計算機科学科

本報告では、プログラミングシステムの立場から、日本語の文をUNIXシェルプログラムへ変換する翻訳系の作成に用いる人工的日本語サブセットの構成方法として、後置演算子文法を提案する。

後置演算子文法は、一種類の非終端記号と三種類の終端記号（演算子・変数・区切り記号）の関係を文法規則として記述したものである。この後置演算子文法は、その構文解析が常に一意にできることを特徴としている。さらに、単純後置翻訳法などを組み合わせることによって比較的簡単に翻訳系を作成することができる。

A method of constructing subsets of Japanese Language
using postfix operator grammars and its applications

Masaya Hagiwara and Takehiro Tokuda
Department of Computer Science, Faculty of Engineering,
Yamanashi University, 4-3-11 Takeda, Kofu 400, Japan

We propose a subclass of context-free grammars, which we call postfix operator grammars. Postfix operator grammars consist of productions of the following form, $S \rightarrow \text{var}$, or $S \rightarrow S a_1 \dots S a_n \text{op}$, where S is a nonterminal, var and op are terminals, and a_1, \dots, a_n are strings of terminals of length zero or more. Strings a_1, \dots, a_n work as delimiters when they are not empty. A terminal op belongs to a special class of terminals which work as postfix operators. Postfix operator grammars are unambiguous. Hence this grammar class is suitable for constructing subsets of Japanese language which can be used as input languages of programming systems such as a translator of Japanese sentences into command procedures.

1. はじめに

本論文の目的は、人工的日本語をawkやsedを用いたUNIXのシェルプログラムに翻訳する際
に使用する文法を簡易的に構成する方法を示すこと
である。

プログラミングシステムの分野では、構文解析手
法としてLL(k)系やLR(k)系やEarleyのアルゴ
リズムなどがあげられる[1]。LL(k)系やLR(k)系
の構文解析手法はそれぞれ文脈自由文法の部分集合
に対して効率がよい。これは、構文解析の準備とし
て構文解析表を作成するからである。しかし、文法
が与えられたときのLL(k)条件やLR(k)条件の判
定は、構文解析表を作成するまでは一般的には判明
しない。また、文法が大きくなると、それにともな
い構文解析表用の記憶領域は大きなものになる。さ
らに、文法を拡張するたびに構文解析表を生成し
直す必要がある。Earleyのアルゴリズムは一般的
な文脈自由文法に対して有効であり、曖昧さがあ
ってもなくても解析は可能である。しかし、文脈自
由文法の部分集合に対しての効率は、LL(k)系やLR
(k)系と比べてあまりよくない。

本論文で示す方法は、従来からプログラミングシ
ステムの分野で用いられてきた抽象木や演算子・被
演算子のモデルをもとにして、文脈自由文法の部分
集合であり、曖昧さがなく構文解析法が簡単な後置
演算子文法を提案し、それを日本語に対して導入す
ることにより、人工的日本語サブセットを簡易的に
構成するものである。この方法による人工的日本語
サブセットは、

- ・構文解析が単純な方法で、しかも一意に行
うことができる。
- ・自然な日本語表現に近いものが得られる。
- ・翻訳が比較的容易である。

といった特徴をもっている。

本論文の構成は次の通りである。第2章では、本
論文の動機について述べる。第3章では、一般的な
後置式を生成する文法について数学的証明も交えて
説明する。第4章では、第2章・第3章の理論を基
に、本論文で提案する後置演算子文法について述べ
る。第5章では、後置演算子文法の構文解析につい
て述べる。第6章では、後置演算子文法による人工
的日本語サブセットとその応用例として、

日本語によるCの型宣言を英語に翻訳する
日本語サブセット

UNIXのシェルプログラムへの翻訳を行う
簡単な日本語サブセット

を示す。第7章では、後置演算子文法の拡張につい
て述べる。第8章では、現状と今後について述べる。

2. 動機

プログラミングシステムのコンパイラ生成系や構
造エディタ生成系で取り扱っている基本モデルは、
図1のようなモデルである。

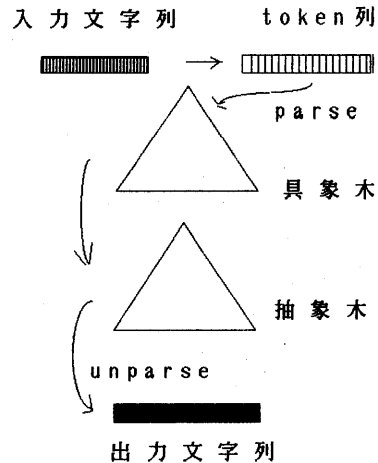


図1 基本モデル

例えば、カーネギーメロン大学のGandalfソフト
ウェア開発環境で使用されている構造エディタAL
OEでは、次のように抽象木を利用している[13]。

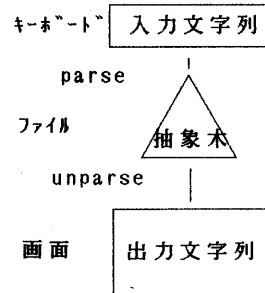


図2 ALOEにおける抽象木の利用

ALOEにおけるデータ構造は、単なる文字の集
まりではなく、内部表現として木構造になっている。
この木構造は、解析木から冗長な情報を消去した抽
象的な木になっている。この抽象木は、編集の対象
となるプログラミング言語に依存している。

しかし、ファイルの抽象木構造をそのまま端末で
みることはできず、抽象木に予約語やセミコロン・
かっこ・タブ・スペースなどを追加して端末上で見
やすい形にする操作を行う。これをアンバースと言
う。

UNIXは、AT&T ベル研究所の登録商標であ
る

このような抽象木の考え方は、プログラミングシステムの他の分野や自然言語処理の分野でも一般的である。例えば、ウィーン定義言語[19]や、変形生成文法[5]や格文法[6]、概念依存文法[14]などである。

もう一つの基本モデルの例として、中置式から後置式を翻訳するプロセスがある。そのプロセスは次のようになる。

《中置式から後置式を生成するプロセス》

G) 中置式の文法

$E \rightarrow E+T \quad E \rightarrow T$
 $T \rightarrow T * F \quad T \rightarrow F$
 $F \rightarrow i$

U) 後置式を生成するアンバース手続き

unparse(x) ... xを親ノードとする二分木のアンバースを行う手続き
 x : 具象木のノードを表す
 output(c) ... 文字 c を出力する手続き
 unparse(+) → unparse(左の子); unparse(右の子);
 output('+');
 unparse(*) → unparse(左の子); unparse(右の子);
 output('*');
 unparse(i) → output(i に付いているラベル);

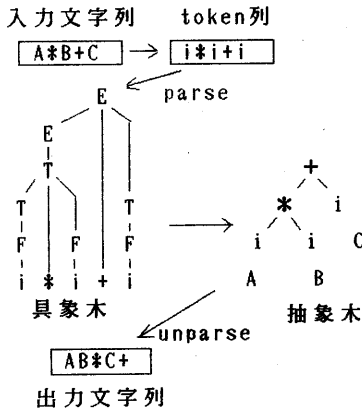


図3 unparse手続きの例

算術式はまず token列に変換され、中置式の文法による構文解析で具象木を生成する。具象木から不必要な情報を消去し抽象木を生成する。抽象木に対し後置式を生成するアンバース手続きを行うことによって後置式が生成される。

日本語の構文構造や意味構造も、プログラミングシステムで扱うように演算子・被演算子関係の抽象木にあてはめて取り扱い、抽象木にアンバース手続きを用いることによって翻訳を行いたい、というのが本論文の動機である。。

日本語の簡単な抽象木をつくるために、日本語の語句のかかり受けの関係を演算子と被演算子の関係に考えると、次のようになる。

1. 前置演算子:

例 形容詞と名詞 演算子 形容詞
 被演算子 名詞
 赤い 本

2. 中置演算子:

例 助詞と 演算子 と
 被演算子 名詞
 本 と ノート

3. 後置演算子:

例 目的語と動詞 演算子 動詞
 被演算子 目的語
 学校へ 行く

4. 0項演算子:

例 名詞 演算子 名詞
 被演算子 なし
 学校

これらの4種類の演算子の組合せは、一般には曖昧な文を生成する。

1) 中置演算子の存在は単独で曖昧となる。

例 文法 $S \rightarrow S + S$
 $S \rightarrow a$
 文 $a + a + a$

2) 前置演算子と後置演算子の共存は曖昧となる

例 文法 $S \rightarrow f S$
 $S \rightarrow S g$
 $S \rightarrow a$
 文 $f a g$

3) 同一の前置演算子 (または後置演算子) が異なる個数の被演算子を持つ場合は曖昧となる

例 文法 $S \rightarrow S -$ (負の符号)
 $S \rightarrow S S -$ (減算)
 $S \rightarrow a$
 文 $a a - -$

以上の考察から、一意に構文解析できる文法を得るためには、後置演算子と0項演算子の組合せ、または、前置演算子と0項演算子の組合せでなければならないことになる。

日本語 AとBの積とCの和

後置式 $AB * C +$

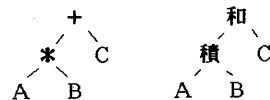


図4 後置式と日本語の類似

一方、図4に示すように日本語の表現と後置式の表現には、強い類似性がある。この点に注目し、後置演算子と0項演算子の組合せを用いることにより、一意的に構文解析できる人工的日本語文法を得ることができると考えられる。

3. 後置式を生成する文法

構文解析が一意的にできる人工的日本語文法を得るために、まず後置式を生成する文法について述べる。

3.1 後置式文法

【定義】 後置式の定義

後置式を次のように定義する。

- (1) 変数 1 個の式は後置式である。
- (2) 後置式を 1 個以上 k 個 ($1 \leq k$) 被演算子として並べ、最後に k 項演算子を 1 個書いた式は後置式である。(演算子と変数は異なるものである。)
- (3) (1), (2) で得られた式だけを後置式と呼ぶ。

【定義】 後置式文法の定義

次の形をした文脈自由文法を、後置式文法という。但し、文法記号で、 S は特定の非終端記号であり、 var, op は、終端記号である。

非終端記号 S
 終端記号 $VAR \cup OP$
 の集合 VAR : 変数の集合
 OP : 演算子の集合
 $(VAR \cap OP = \emptyset)$
 文法規則 $S \rightarrow var$
 $S \rightarrow S S \dots S op$
 $var \in VAR$
 $op \in OP$

このとき、右辺のそれぞれの S を被演算子と呼び、 S の個数は 1 以上である。 op は演算子である。

例1 解析が一意的な後置式文法

非終端記号 S
 終端記号 $VAR \cup OP$
 の集合 $VAR = \{var\}$
 $OP = \{+, *\}$
 文法規則 $S \rightarrow var$
 $S \rightarrow SS +$
 $S \rightarrow SS *$
 生成文の例 $var \ var \ var \ + \ *$

例2 解析が曖昧な後置式文法

非終端記号 S
 終端記号 $VAR \cup OP$
 の集合 $VAR = \{var\}$
 $OP = \{-\}$

文法規則 $S \rightarrow var$
 $S \rightarrow S -$ (負の符号)
 $S \rightarrow SS -$ (減算)
 生成文の例 $var \ var \ - \ -$

このとき例1は曖昧な文法ではないが、例2は曖昧な文法である。

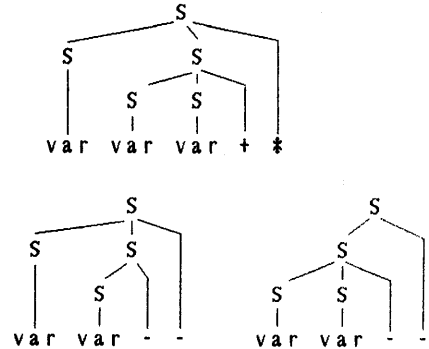


図5 例1と例2の解析木

後置式文法のうち、特に文法規則中の右辺の右端の op がすべて相異なる演算子であるならば、この場合の文法を限定後置式文法と呼ぶことにする。限定後置式文法では、文法規則の右辺に含まれる被演算子の個数が演算子記号に対して一意に決まるので、この被演算子の個数を、演算子の項数と呼ぶことにする。

3.2 限定後置式文法の無曖昧性

(定理1) 限定後置式文法で生成されるような文は、2通りの異なる構文解析木を持つことはない。つまり限定後置式文法は、曖昧ではない。

《定理の証明》

まず次の補助定理1を証明する

(補助定理1) 限定後置式文法で使用する終端記号 x に対して、関数 $g(x)$ を次のように定義する。
 $g(var) = +1$
 $g(op) = -(j-1)$ 但し j は演算子 op の項数である。
 $(1 \leq j)$

また、限定後置式文法から生成された終端記号列 $X_1 X_2 \dots X_n$ に対して、関数 $f(i) (1 \leq i \leq n)$ を、次のように定義する。

$$f(i) = g(x_1) + \dots + g(x_i)$$

このとき次のことがいえる。

$$f(i) \geq 1 \quad (1 \leq i \leq n-1)$$

$$f(n) = 1 \quad (1)$$

つまり1つの後置式の関数 f の値は1となり、1つの後置式の長さ $n-1$ 以下の前部分式の関数 f の値は1以上となる。

《補助定理1の証明》

限定後置式文法で生成された文に含まれている演算子 t に関する数学的帰納法を用いて証明する。

$t=0$ のとき、限定演算子文法で生成される文は var であり、

$$f(1)=g(var)=1$$

よって (1) 式が成立する。

$t < N$ のとき (1) 式が成立したとする。

$t = N$ のとき、次のことがいえる。

限定後置式文法から生成された終端記号列 $X_1 X_2 \dots X_p \dots X_i \dots X_s$ 、($1 \leq p, 1 \leq i, i \leq s$) について、終端記号 X_s は ℓ 個の後置式をもつ演算子であるとする

($1 \leq \ell$)。 ℓ 個の後置式に含まれている演算子の個数は最大 $N-1$ ($< N$) であるので、帰納法の仮定より各後置式の関数 f の値は1である。終端記号 X_i の左側で最も近いところにある演算子を op' ($= X_p$) そこまでにある後置式の個数を L 個とすると、(op' が存在しない場合は $L=0$ 、 $op' = X_i$ のときは $L=\ell$)

$$f(i) = g(x_1) + \dots + g(x_p) + g(x_{p+1}) + \dots + g(x_i) \\ = 1 * L + 1 * (i-p) \geq 1$$

$$f(s-1) = 1 * \ell \geq 1$$

$$f(s) = f(s-1) + g(x_s) \\ = \ell - (\ell - 1) \\ = 1$$

で、成立する。 補助定理1 証明終了

定理1の証明は、この補助定理1を用いて背理法で証明する。

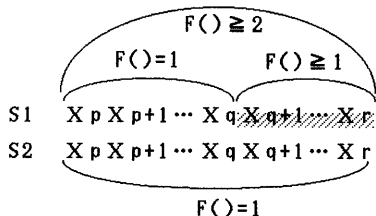
限定後置式文法で生成された終端記号列を $X_1 X_2 \dots X_n$ とし、 X_n が ℓ 個の後置式をもつ演算子であるとする。このとき終端記号列 $X_1 X_2 \dots X_n$ が ℓ 個の被演算子に一意に区切ることができないと仮定する。つまり、左端が一致して右端が一致しない被演算子

$$S1 = X_p X_{p+1} \dots X_q$$

$$S2 = X_p X_{p+1} \dots X_q X_{q+1} \dots X_r$$

が存在すると仮定する。

$S2$ に注目すると、 $S2$ は演算子 X_r に関する後置式だから関数 $f(r)$ の値は1である。



一方 $S1$ に注目すると、 $S1$ は演算子 X_q に関する後置式だから関数 $f(q)$ の値は1になる。このとき $X_{p+1} \dots X_q$ は $S1$ の次の後置式の前部分列にあたるので関数 f の値は1以上となる。よって $S1$ に注目すると $X_p \dots X_r$ の関数 $f(r)$ の値は2以上となる。

同じ終端記号列 $X_1 X_2 \dots X_n$ に対して一意な解を持つはずの関数 f が2つの異なる値を持ってしまう。これは矛盾である。

よって、左端が一致して右端が一致しない被演算子は存在しない。被演算子の区切りは一意である、つまり2通りの構文解析木を持つことはないことがわかる。 証明終了

このように、限定後置式文法は曖昧でなく構文解析が一意に行えることがわかる。次の第4章では、限定後置式文法の拡張を考える。

4. 文法の構成方法

4.1 後置演算子文法

本論文で示す文法の構成方法は、第3章で述べた限定後置式文法を・文法規則中に区切り記号を加える・演算子は、被演算子の個数を一意に決めることによって2つ以上の文法規則に含まれてもよい。の2点について拡張したものである。

構成方法はつぎのとおりである。

【構成方法】

条件1) 文法規則の形式は、次の形式のいずれかである。ただし、以下の文法規則のスキームにおいて、非終端記号は1つの特定の非終端記号 S であり、終端記号は、 var_i 、 op_i である。また、 a_{i1} 、 a_{i2} 、 \dots 、 a_{ik} は、終端記号列である。($1 \leq i, j, k$)

非終端記号 S

終端記号 $VAR UDEL UOP$

の集合 VAR : 変数の集合

DEL : 区切り記号の集合

OP : 演算子の集合

$$(VAR \cap DEL \cap OP = \emptyset)$$

文法規則 $S \rightarrow var_i$

$$S \rightarrow S a_{i1} S a_{i2} \dots S a_{ik} op_i$$

$$var_i \in VAR$$

$$a_{i1}, a_{i2}, \dots, a_{ik} \in DEL *$$

$$op_i \in OP$$

すなわち、被演算子は a_{i1} 、 a_{i2} 、 \dots 、 a_{ik} という終端記号列で区切られ、被演算子と区切り記号列の並びの最後には演算子 op_i が存在する。

条件2) 各終端記号 op_i は、その被演算子の個数が一意的に定まる。すなわち、同一の終端記号 op_i に対し、その被演算子の個数はどの文法規則でも同一である。

このとき3つの集合VAR, DEL, OPは互いに素であるならば、区切り記号の列 $a_{i1}, a_{i2}, \dots, a_{ik}$ に特に制約はない。文法全体ですべて同一の終端記号でもよいし、また空列でもよい。

この文法構成法で作成された文法を後置演算子文法という

4.2 後置演算子文法の無曖昧性

(定理2) 後置演算子文法で生成されるような文は、2通りの構文解析木を持つことはない。つまり後置演算子文法は、曖昧ではない。

《定理2の証明》

次の補助定理2を証明する

(補助定理2) 後置演算子文法で使用される終端記号 x に対して、関数 $g(x)$ を次のように定義する。

$$\begin{aligned} g(x) &= +1 & x \in \text{VAR} \\ g(x) &= 0 & x \in \text{DEL} \\ g(x) &= -(j-1) & x \in \text{OP} \end{aligned}$$

但し j は演算子 op の項数である。
($1 \leq j$)

また、後置演算子文法で生成された終端記号列 $X_1 X_2 \dots X_n$ に対して、関数 $f(i)$ ($1 \leq i \leq n$) を、次のように定義する。

$$\begin{aligned} f(i) &= g(x_1) + \dots + g(x_i) \\ f(i) &\geq 1 \quad (1 \leq i \leq n-1) \\ f(n) &= 1 \end{aligned} \quad (1)$$

つまり1つの後置式の関数 f の値は1となり、1つの後置式の長さ $n-1$ 以下の前部分式の関数 f の値は1以上となる。

《補助定理2の証明》

後置演算子文法で生成された文に含まれている演算子 t に関する数学的帰納法を用いて証明する。

$t=0$ のとき、文法 G で生成される終端記号列は var であり、

$$f(1) = g(\text{var}) = 1$$

よって (1) 式が成立する。

$t < N$ のとき (1) 式が成立したとする。

$t = N$ のとき、次のことがいえる。

終端記号列 $X_1 X_2 \dots X_p \dots X_i \dots X_s$ ($1 \leq s, 1 \leq p, i \leq s$) について、終端記号 X_s は l 個の後置式をもつ演算子であるとする ($1 \leq l$)。 l 個の後置式に含まれている演算子の個数は最大 $N-1$ ($< N$) であるので、帰納法の仮定より各後置式の関数 f の値は1である。終端記号 X_i の左側で最も近いところにある演算子を op' ($= X_p$) そこまである後置式

の個数を L 個とすると、(op' が存在しない場合は $L=0$ 、 $op' = X_i$ のときは $L=l$)

$$\begin{aligned} f(i) &= g(x_1) + \dots + g(x_o) + g(x_{o+1}) + \dots + g(x_i) \\ &= 1 \cdot l + 1 \cdot (i-p) \geq 1 \\ f(s-1) &= 1 \cdot l \geq 1 \\ f(s) &= f(s-1) + g(x_s) \\ &= l - (l-1) \\ &= 1 \end{aligned}$$

で、成立する。

補助定理2 証明終了

定理2の証明は、後置式文法の定理1の証明と同様であるので省略する。 定理2 証明終了

《文法例》 図4の日本語の抽象木に対する具象木を生成する文法

非終端記号	S
終端記号	VAR U DEL U OP
の集合	VAR = {A, B, C}
	DEL = {と, の}
	OP = {和, 積}
文法規則	S \rightarrow S と S の 和
	S \rightarrow S と S の 積
	S \rightarrow A
	S \rightarrow B
	S \rightarrow C
例文	A と B の 積 と C の 和

5. 構文解析方法

第4章で定義した後置演算子文法の構文解析方法は次の通りである。

【構文解析方法】

環境

入力文字列 $X_1 \dots X_n$ ($1 \leq n$)
 構文解析用スタック P
 文法規則照会用配列 M

解析手順

(1) 初期化

スタックP、配列Mを空にする。
 文字 X_1 をスタックPにプッシュする。

(2) 解析

スタックPの先頭文字 Y が

① $Y \in \text{VAR}$: $S \rightarrow Y$ で還元する。

② $Y \in \text{DEL}$: 何もしない

③ $Y \in \text{OP}$:

1) 演算子 Y が持つ被演算子の項数 m ($1 \leq m$) を調べ、スタックPの先頭から m 番目の S までの内容を配列Mにコピーする。

2) 配列Mの内容と右辺の右端の演算子が Y であるような文法規則のパターンマッチングを行う。

マッチするような文法規則があれば $S \rightarrow S a_{i1} \dots S a_{ik} Y$ で還元を行う。なければ、異常終了である。

④それ以外の場合：異常終了である。

(3)判定

入力文字列がまだ残っていれば、先頭文字を1つスタックPにプッシュし、再びステップ(2)へ行く。

入力文字列がなくなった場合、スタックPの内容がSだけならば正常終了である。それ以外であれば異常終了である。

6. 文法例と翻訳例

後置演算子文法を用いて構成された日本語サブセットの2つの文法例を示す。文法構成の際、演算子はキーワードとなる動詞や名詞など、区切り記号は助詞などが相当している。

G1. 言語Cの型宣言の日本語表現の文法例

非終端記号 S
 終端記号 VARUDELUOP
 の集合 VAR = {整数, 文字}
 DEL = {へ, の, を, に, 要素, 持つ, 返す}
 OP = {ポインタ, 配列, 関数}

文法規則
 $S \rightarrow S \text{ へ の ポインタ}$
 $S \rightarrow S \text{ を 要素 に 持つ 配列}$
 $S \rightarrow S \text{ を 返す 関数}$
 $S \rightarrow \text{整数}$
 $S \rightarrow \text{文字}$

この文法から生成される文の例には、次のようなものがある。

文字 へ の ポインタ を 要素 に 持つ 配列
 整数 へ の ポインタ を 返す 関数

G2. UNIXコマンドawk用日本語文法

非終端記号 S
 終端記号 VARUDELUOP
 の集合 VAR = {1, 2, myfile, yourfile}
 DEL = {を, と, の}
 OP = {出力する, 和, 積, 番目, フィールド}

文法規則
 $S \rightarrow S \text{ を 出力する}$
 $S \rightarrow S \text{ と } S \text{ の 和}$
 $S \rightarrow S \text{ と } S \text{ の 積}$
 $S \rightarrow S \text{ の } S \text{ の フィールド}$
 $S \rightarrow S \text{ 番目}$
 $S \rightarrow 1$

$S \rightarrow 2$
 $S \rightarrow \text{myfile}$
 $S \rightarrow \text{yourfile}$

この文法から生成される文の例には、次のようなものがある。

myfile の 1 番目の フィールド と yourfile の 2 番目の フィールド の 積 を 出力する。

myfile の 1 番目の フィールド と yourfile の 1 番目の フィールド の 和 と yourfile の 2 番目の フィールド の 積 を 出力する。

2つの文法例を翻訳する例を示す。T1は、日本語の型宣言を英語の型宣言になおしている。T2は、日本語のawkコマンド動作指定をUNIXコマンドに翻訳している。翻訳の定義は、属性文法[9]を用いている。

T1. G1.を用いたCの型宣言の日本語から英語への翻訳

$S0 \rightarrow S1 \text{ へ の ポインタ}$
 $\text{code}(s0) = \text{"a pointer to"} + \text{code}(s1)$
 $S0 \rightarrow S1 \text{ を 要素 に 持つ 配列}$
 $\text{code}(s0) = \text{"an array of"} + \text{code}(s1)$
 $S0 \rightarrow S1 \text{ を 返す関数}$
 $\text{code}(s0) = \text{"a function returning"} + \text{code}(s1)$
 $S \rightarrow \text{整数}$
 $\text{code}(s) = \text{"integer"}$
 $S \rightarrow \text{文字}$
 $\text{code}(s) = \text{"character"}$

翻訳例

文字 へ の ポインタ を 要素 に 持つ 配列
 $\rightarrow \text{an array of a pointer to character}$
 整数 へ の ポインタ を 返す 関数
 $\rightarrow \text{a function returning a pointer to integer}$

T2. G2.を用いた日本語からシェルプログラムへの翻訳 ただしcatsは、横方向に空白をいれてファイルを結合するコマンドである。

$S1 \rightarrow S2 \text{ を 出力する}$
 $\text{code}(s1) = \text{code}(s2) + \text{"\npr place}(s2) \n"}$
 $S1 \rightarrow S2 \text{ と } S3 \text{ の 和}$
 $\text{place}(s1) = \text{newname}()$
 $\text{code}(s1) = \text{code}(s2) + \text{code}(s3) + \text{"cats place}(S2) \text{ place}(s3) \n"}$
 $\text{awk '\{print }$1+$2}' > \text{place}(s1) \n"$

```

S1 → S2 の S3 の 積
place(s1) = newname()
code(s1) = code(s2) + code(s3) +
"cats place(S2) place(s3) |
awk '{print $1*$2}' > place(s1) ¥n"
S1 → S2 の S3 の フィールド
place(s1) = newname()
code(s1) = code(s2) + "cat place(s2) |
awk '{print part(s3)}' > place(s1)
¥n"
S1 → S2 番目
part(s1) = "$num(s2)"
S → 1
num(s) = "1"
S → 2
num(s) = "2"
S → myfile
place(s) = "myfile"
code(s) = ""
S → yourfile
place(s) = "yourfile"
code(s) = ""

```

翻訳例

myfile の 1 番目のフィールドと yourfile の 2 番目のフィールドの積を出力する。

```

→ cat myfile | awk '{print $1}' > temp1
cat yourfile | awk '{print $2}' > temp2
cats temp1 temp2 |
awk '{print $1*$2}' > temp3
lpr temp3

```

myfile の 1 番目のフィールドと yourfile の 1 番目のフィールドの和と yourfile の 2 番目のフィールドの積を出力する。

```

→ cat myfile | awk '{print $1}' > temp1
cat yourfile | awk '{print $1}' > temp2
cats temp1 temp2 |
awk '{print $1+$2}' > temp3
cat yourfile | awk '{print $2}' > temp4
cats temp4 temp5 |
awk '{print $1*$2}' > temp5
lpr temp5

```

UNIXシェルプログラムへの翻訳では、一つのコマンドの終了ごとにその結果を一時的な仮ファイルを生成している。これを、UNIXのパイプ(|)を使うことによって処理を行うことは、今後の課題である。

7. 議論

後置演算子文法は、UNIXシェルプログラム翻訳の日本語インターフェース用の日本語サブセットを構成するに際して有効である。後置演算子文法に対しては、さらに次の点に関しての拡張を考えている。

1. 被演算子の置換

後置演算子文法は、日本語で語順が比較的自由であると同様に、被演算子の順番に対して比較的自由である。被演算子の置換に対しては、構文解析方法を次のように行う。

【構文解析方法】

(2)解析

スタックPの先頭文字Yが

- ① $Y \in VAR : S \rightarrow Y$ で還元する。
- ② $Y \in DEL :$ 何もしない
- ③ $Y \in OP :$

1)演算子Yが持つ項数mを調べ、スタックPの先頭からm番目のSまでの内容を配列Mにコピーする。

2)配列Mの内容が、右辺の右端の演算子がYであるような文法規則の置換列になっていないかを調べる。置換列になっている文法規則があれば $S \rightarrow S \dots S \ a_i \ Y$ で還元を行う。なければ、異常終了である。

④それ以外の場合：異常終了である。

2. 被演算子の省略

後置演算子文法は、日本語で語の省略ができるのと同様に、被演算子を省略することも部分的には可能である。例えば、文法構成方法の条件に次の条件が加わる。

条件3)省略できる部分は前部分列だけとする。

条件4)各文法規則の区切り記号列の集合は、互いに素である。また、意味的に自分の文法規則の右辺の中に自分を埋め込むことはない。

このとき構文解析方法は次のように行う。

【構文解析方法】

(2)解析

スタックPの先頭文字Yが

- ① $Y \in VAR : S \rightarrow Y$ で還元する。
- ② $Y \in DEL :$ 何もしない
- ③ $Y \in OP :$

1)演算子Yが持つ項数mを調べ、スタックPの先頭からm番目のSまでの内容を配列Mにコピーする。

2)配列Mの内容と右辺の右端の演算子がYであるような文法規則の後部分列とパターンマッチ

ングを行う。マッチするような文法であれば $S \rightarrow S a_{i1} \dots S a_{ik} Y$ で還元を行う。そうでなければ被演算子と区切り記号の組を1組スタックからポップして、2)を行う。

④それ以外の場合：異常終了である。

8. おわりに

現在、本論文の後置演算子文法をもとにシェルプログラムへの翻訳系に使用する人工的日本語サブセットを構成中である。今後は、付録にあるようなユーザ登録手続きや、その他採点集計プログラムのようなコマンドプログラムを翻訳可能な人工的日本語文法を構成していく。

《参考文献》

- [1] Aho, A.V. and Ullman, J.D. : The Theory of Parsing, Translation, and Compiling, Vol.1 and 2, Prentice-Hall (1972 and 1973).
- [2] Aho, A.V. and Ullman, J.D. : Principles of Compiler Design, Addison-Wesley (1977).
- [3] Aho, A.V. and Sethi, Ullman, J.D.: Compilers principles, Techniques, and Tools, Addison-Wesley Publishing Company (1986).
- [4] Aho, A.V., Kernighan, B.W. and Weinberger, P.J. : Awk, -a pattern scanning and processing language, Unix User's Reference Manual, 4.3 BSD University of California Berkeley
- [5] Chomsky, N.: Aspects of the Theory of Syntax, MIT press (1957).
- [6] Fillmore, C.: The Case for Case, in Bach, E. and Harms, R. (eds) Universals in linguistic Theory, Holt, Reinhart and Winston (1968).
- [7] 淵一博 監修 古川康一・溝口文雄 共編 自然言語の基礎理論, 共立出版, (1986)
- [8] Knuth, D.E. : On the translation of Languages from Left to Right, Inf.Cont., Vol.8, No.6, pp.607-639 (1965).
- [9] Knuth, D.E. : Semantics of context-free languages, Math. Systems Theory, Vol.2, No.2 pp.127-145 (1968).
- [10] Lewis, P.M.II and Stearns, R.E. : Syntax directed transduction, J.ACM, Vol.15, No.3, pp.465-488 (1968).
- [11] Montague, R: The Proper Treatment of Quantification in Ordinary English, Approaches to Natural Language. D.Reidel Pub. (1973).
- [12] 長尾 真: 言語工学、昭晃堂 (1983)
- [13] Notkin, D. : The GANDALF Project, The Journal of Systems and Software, Vol.5, no.2, pp.91-106 (1985).
- [14] Shank, R.C., Abelson, R.P. : Scripts, Plans and Knowledge, IJCAI 4 (1975)
- [15] 田村直良、高倉伸、片山卓也: 自然言語処理を目的とした属性文法評価システム、コンピュータソフトウェア、Vol.3, No.3 (1986).
- [16] Teitelbaum, T. and Reps, T. : The Cornell Program Synthesizer, A syntax-directed Programming Environment, Comm, ACM, Vol.24, No.9, September, pp.563-573 (1981).
- [17] Teitelbaum, T. Reps, T. and Horwitz, S. : The Why and Wherefore of the Cornell Program Synthesizer, SIGPLAN Notices, Vol.16, No.6, pp8-16, The Proceedings of the ACM SIGPLAN/SIGOA, Symposium on Text Manipulation (1981).
- [18] Tomita, M. : An Efficient Context-Free Parsing Algorithm for Natural Languages and Its Applications, Computer Science Department Carnegie-Mellon University (1985).
- [19] Wegner, P.: The Vienna definition language Computing Surveys 4:1,5-63.
- [20] Winograd, T.: Language as a Cognitive Process, Addison-Wesley (1983).
- [21] Woods, W. : Augmented transition network grammar for natural language analysis, CACM, Vol.13, pp.589-602 (1970).

付録 ユーザー登録を行うシェルプログラム

ユーザー登録を行うメインプログラム

```
num="0"
count="1"
cp /etc/passwd p
cp /etc/group g
echo "How many people do you add?"
read num
echo "login_name UID group_name GID home_directory real_name"
while test $count -le $num
do
    adduser
    count= expr $count + 1
done
cp g /etc/group
cp p /etc/passwd
vipw
cat err
chmod +x own
chmod +x grp
own
grp
```

実際に登録ファイルに操作を行い登録を行う
プログラム

```
D="/usr/local/lib"
:
: echo "login_name UID group_name GID home_directory real_name"
read nnn1 nnn2 nnn3 nnn4 nnn5 nnn0
:
echo "$nnn1::$nnn2:$nnn4:$nnn0:$nnn5:/bin/csh"
> tmp
cat tmp
:
ttt= `
cat tmp p l
sed -e "s/:::/none:/g" l
sed -e "s/:// /g" l
awk '
    BEGIN { c=0
            d=0 }
    {if(NR==1) {a1=$1
               a3=$3}}
    {if(NR!=1 && $1==a1) c=1}
    {if(NR!=1 && $3==a3) d=2}
    END{printf "%1d%1n",c+d}
, ~
case $ttt in
0) ;;
1) echo "Use another user name;"$nnn1":"$nnn2"
:"$nnn3":"$nnn4 >> err
    exit ;;
2) echo "Use another user number;"$nnn1":"$nnn
2":"$nnn3":"$nnn4 >> err
    exit ;;
3) echo "Use another user name and user number
;"$nnn1":"$nnn2":"$nnn3":"$nnn4 >> err
    exit ;;
*) echo $ttt >> err
    exit ;;
esac
:
echo "$nnn3::$nnn4" > tmp1
cat tmp1
nnn= `
cat tmp1 g l
sed -e "s/:::/none:/g" l
sed -e "s/:// /g" l
awk '
    BEGIN {c=0
            d=0
            e=0}
    {if(NR==1) {a1=$1
               a3=$3}}
    {if(NR!=1 && $1==a1 && $3==a3) c=1}
    {if(NR!=1 && $1==a1 && $3!=a3) d=2}
```

```
{if(NR!=1 && $1!=a1 && $3==a3) e=4}
END{printf "%1d%1n",c+d+e}
,
case $nnn in
0) echo "$nnn3::$nnn4:$nnn1" > err ;;
1) ;;
2) echo "Incorrect group number;"$nnn1":"$nnn2
:"$nnn3":"$nnn4 >> err
    exit ;;
4) echo "Incorrect group name;"$nnn1":"$nnn2":
"$nnn3":"$nnn4 >> err
    exit ;;
*) echo $nnn >> err
    exit ;;
esac
:
if test -d $nnn5
then echo "Use another directory name($nnn5);"
"$nnn1":"$nnn2":"$nnn3":"$nnn4 >> err
    exit
fi
:
if mkdir $nnn5
then
    chmod 700 $nnn5
    echo $nnn1 >> name
    echo source $D/stdlogin > $nnn5/.login
    echo source $D/stdcshrc > $nnn5/.cshrc
    chmod 700 $nnn5/.login $nnn5/.cshrc
    echo chown -R $nnn1 $nnn5 >> own
    echo chgrp -R $nnn3 $nnn5 >> grp
    cat p tmp > pl
    cp pl p
    sed -e "/^~$nnn3:/s/^$/$nnn1/" g > gl
    cp gl g
    exit
fi
echo "cannot mkdir $nnn5;"$nnn1":"$nnn2":"$nnn
3":"$nnn4 >> err
```