

## 既存ソフトウェアへの日本語テキスト処理機能の追加

## - SNOBOL 4 処理系の場合 -

牛島 和夫, 吉田 和幸, 黒坂 輝彦\*

(九州大学工学部)

## 1. はじめに

近年, ハードウェア技術の進歩により, 計算機により日本語文字(漢字, ひらがな, カタカナ, ローマ字および日本語の文章に使用する特殊文字をまとめてこう呼ぶことにする)を扱うことが可能になってきた。最近国内で発表される計算機はそのほとんどが「漢字が扱えます」とうたっている。しかし, 実際に日本語文字を扱うプログラムを書こうとすると, 計算機の細部まで知っている必要があったり, 面倒な手順を要求されたりする場合が多い。

筆者らの利用する九州大学大型計算機センターでは1980年から富士通の提供するJEF (Japanese processing Extended Feature) [1]と云う日本語情報処理機能を導入している。JEFの利用形態として最もよく行われているものの一つに, 日本語ワードプロセッサとしての利用がある。しかしもっと高度な使い方(たとえば[2]にある文書作成支援ツールのようなこと)をしたければ, 自分でプログラムを組む必要が出てくる。このためにはCOBOL, FORTRAN, PL/Iの三つの言語でJEFの機能を使えるようになっていて, ところが, これらの言語は非常に「手続き的」であったり, 固定したフォーマットのファイルしか扱えなかったりして, 日本語文字処理を伴うプログラムを書くのは, かなりむずかしい。

SNOBOL 4 [3]は挿入, 削除を伴ったパターンマッチを一文で書け, その組み合わせにより文字列に関するかなり複雑な操作も簡単に記述できる。従ってテキストファイル処理などに際して適当なツールがない場合その場でプログラムを作ってみるというラピッドプロトタイピングに適した言語である[4]。我々の研究室ではFORTRANソーステキスト中からサブルーチンを切り出すといった, 汎用のテキストエディタの能力を超えた文字列処理のためのツールとしてSNOBOL 4を便利に使用している。

そこでこのSNOBOL 4に日本語文字処理機能を付加することで, 日本語文字, EBCDIC文字が混在しているテキストの高度な処理手続きを簡単に記述できるようにすることを考えた。以下, その実現

の過程と既存ソフトウェアの一部変更という形でこれを可能にした理由とについて述べる。第2章と第3章はその準備である。第2章ではJEFが提供する日本語文字処理環境の概要について述べる。第3章では対象となった既存のSNOBOL 4処理系の概要を示す。第4章でこの処理系に追加した日本語テキスト処理機能の設計とその実現, 第5章で処理速度の改善について述べる。第6章では改善の効果を実行例で示す。第7章では既存ソフトウェアを手直しすることによって日本語テキスト処理機能が比較的容易に実現できた理由を考察する。

なお, 日本語SNOBOL 4は九州大学大型計算機センターのFACOM M-200 OSIV/F4のもとで1983年3月から一般の利用に供している[5]。

## 2. 日本語文字処理環境

OSIV/F4が提供する日本語文字処理のソフトウェアとそれが用いるハードウェアについて簡単に述べる。

(1) 文字コード 従来から使われてきたEBCDICコード(1バイト)に加えて漢字等の日本語文字を表わすJEFコード(2バイト)が使われる。この2種類の文字コードの区別にはシフトコードを用いる。JEFコードは上位バイト, 下位バイトに分けて考えることができる。上位, 下位バイトはそれぞれ16進で41~FE, A1~FEの値をとる。

(2) データセット JEFが扱うデータセットには2種類ある。一つは従来のエディタが扱うデータセットと同様の形式を持つもので日本語文字は右筆というエディタで編集する。右筆ではEBCDIC文字も当然編集の対象となっている。他の一つは文章処理専用のデータセットで従来のエディタや右筆では扱えない構造になっている。このデータセットの編集にはFDMS和文エディタという専用エディタが使われる。

(3) 漢字端末 右筆と和文エディタが動作する端末はフルスクリーン型の端末で, そこでは漢字等の日本語文字はEBCDIC文字の2倍の幅で表示される。

\*現在エプソン株式会社

- (4) 漢字プリンタ ドットプリンタとレーザビームプリンタの2種類ある。ドットプリンタでは日本語文字はEBCDIC文字の2倍の幅で印刷され、レーザビームプリンタでは日本語文字の幅はEBCDIC文字の2倍と1.2倍の2種類の幅を選択できる。
- (5) プログラミング言語による日本語処理 JEFではFORTRAN, COBOL, PL/Iで日本語文字を扱う手段を提供している。これらの言語は手続き的であるため文字列処理を行うプログラムは複雑になりすぎる。また、JEFが提供する機能では、日本語文字とEBCDIC文字との混在は固定的なものしか許されていない。そのため文章ファイルのような日本語文字とEBCDIC文字が自由に混在しているデータの処理は一層複雑になる。そのため計算機を専門としない一般ユーザが文章ファイルを扱うプログラムをこれらの言語を用いて作成することはほとんど不可能であるといってもよい。JEFが提供する環境だけでは手軽に日本語文章を扱うプログラムを作成できるとはいえない。以下では日本語文字とEBCDIC文字が任意に混在するテキストを日本語テキストと称する。

### 3. SNOBOL 4 原処理系の概要

日本語テキスト処理機能の追加作業を説明する準備として、本章では既存のSNOBOL 4 処理系の構成を簡単に述べる。原著者の許可を得て入手したSNOBOL 4 処理系は、SIL (SNOBOL 4 Implementation Language) という抽象言語で書かれたもの (OS 360用) である [6]。

このSNOBOL 4 処理系はトランスレータ、インタプリタ、記憶管理、システムインタフェースという四つの部分から構成される [6]。トランスレータはSNOBOL 4 プログラムをプレフィックスコードに変換し、それをインタプリタが解釈実行する。記憶管理とシステムインタフェースはトランスレータとインタプリタの両方から使われる。記憶管理は領域の割り当て、データの記憶、記憶域の再構成等を行う。システムインタフェースは入出力やプログラム割込み等を受け持つ。

このSNOBOL 4 処理系を記述しているSILという抽象言語は図1のような階層構成により実現されている [6]。最も外側 (図1 (a)) がSILの世界で131個の命令を持つ。その内側 (図1 (b)) が131個の命令をアセンブラのマクロ機能を用いて定義し

たSILマクロの層である。その内側に入出力等OSに依頼する部分を実行時ルーチンとしてまとめた層がある (図1 (c))。これらのルーチンはアセンブリ言語で記述され、SILマクロから呼び出される。実際の入出力はFORTRANの実行時ルーチンに処理を依頼し、異常終了時処理、時間管理等のOSに依頼する部分は直接OSのマクロ命令を用いているものもある。これらが最も内側の層を構成する (図1 (d))。

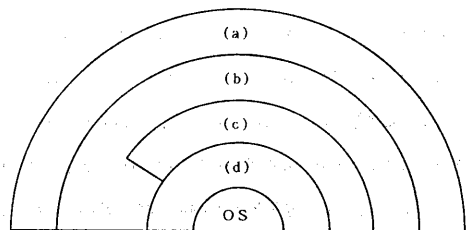
### 4. 日本語テキスト処理機能の追加

#### 4.1 日本語テキスト処理機能の設計

SNOBOL 4への機能追加にあたり次の目標を設定した。

- (1) 日本語文字とEBCDIC文字を同等に扱えるようにすること。シフトコードのような印刷されない文字をSNOBOL 4 利用者が意識しなくてもすむようにする。
- (2) 従来のエディタや右筆で作成したデータセットと同様にFDM S和文エディタで作成したデータセットも扱えること。このためには、可変長ファイルを自由に読み書きできる必要がある (和文エディタが作成するデータセットは可変長ファイルである)。原処理系ではFORTRAN実行時ルーチンを使用しているために入力データは固定長ファイルに入っているものしか扱えない。
- (3) 読み易いプログラム、わかりやすいプログラムを書くための一つ的手段としてプログラム中の変数、ラベル等に日本語文字が使用できるようにする。
- (4) TSS環境で使用するのに便利な機能を加える。

このうち、(1)、(2)が重要な目標である。これらを実現するた



(a) SIL言語

(b) SILマクロ定義

(c) 実行時ルーチン

(d) OSのマクロ命令、FORTRAN実行時ルーチン

図1. SIL処理系の階層構成

めには、

(a) 文字列処理のアルゴリズムの実現はかなり面倒なので SNOBOL 4 原処理系のアルゴリズムをなるべくそのまま使えるようにする。そのためには、日本語文字 1 字を EBCDIC 文字 1 字と同等に扱えるようにする。

(b) 入出力機能を強化し、可変長ファイルを扱えるようにするほか、和文エディタのファイルではその形式に合わせて読み書きできるようにする。

の 2 点を解決しなければならぬ。

この目標を実現するため、次の変更設計を行った。

(1) 文字の内部表現 SNOBOL 4 処理系の文字列処理のアルゴリズムをなるべくそのまま使えるようにするためには、文字コードはすべて同じビット長である必要がある。SNOBOL 4 処理系内部で日本語文字と EBCDIC 文字を同等に扱えるように EBCDIC 文字の左側 (上位バイト側) に 1 バイトのパディングバイト (X'00') を加え、形式的に 2 バイトで表現することにする。たとえば A, B, C を EBCDIC 文字とすると文字列

AB漢字C

の 16 進数表現はデータセット内で

C1C228B4C1BBFA29C3 (9 バイト、下線部はシフトコード)  
処理系内部で

00C100C2B4C1BBFA00C3 (10 バイト)

のようになる。データセット内と処理系内とで表現の差ができるので、SNOBOL 4 実行時ルーチン (図 1 (c) に属する) の入出力部でその変換を行う。また、1 文字が 2 バイトになるため SIL マクロ定義 (131 個) のうち文字列を扱うもの (20 個) を書き替える。そのうち 9 個は SIL マクロ定義内の変更だけでよく、他の 11 個は実行時ルーチン中にある下請けルーチンの変更が必要である。たとえば、前者の例ではデータ領域を確保する SIL マクロ BUFFER の定義

```
MACRO
&LOC BUFFER &N
&LOC DC (&N)X'40'
MEND
```

を次のように変更する。

```
MACRO
&LOC BUFFER &N
&LOC DC (&N)X'0040'
MEND
```

ここで MACRO, MEND はマクロ定義の開始, 終了を, BUFFER はマクロ名を表す。& のついた名前はマクロの仮引数である。DC はアセンブリ言語のデータ定義命令でここでは 16 進定数を定義している。

さらに文字コードを 8 ビットから 16 ビットに変更するのでコンパイル時および実行時に使用する字句解析のための表もそれに伴って変更する必要がある。コンパイル時に使用する表は 24 個あり、それらは変数名, ラベル名, 文字列等をソースプログラムから切り出すために使用される。表の内容は変更されない。実行時に使用する表は上の 24 個のほかパターンマッチ処理のために実行中に書き替えながら使用する動的な表が 1 つある。表に関係する SIL 命令は 3 個ある。そのうち 2 個はマクロ定義内の変更が必要であり、1 個はその下請けルーチンの変更が必要である。変数名, ラベル名等として日本語文字も使用できるようにするためにもこれらの表を変更する必要がある。

(2) 入出力 SNOBOL 4 原処理系では FORTRAN 実行時ルーチンの入出力機能を借りて実際の入出力を行っているが、FORTRAN 実行時ルーチンの仕様は公開されたものではないのでその機能の細部はわからない。したがって和文エディタで作成したファイル等の可変長ファイルからデータを FORTRAN 実行時ルーチンを用いて入力したとき、そのデータの長さを知る方法がない等の問題が生じる。そこで SNOBOL 4 実行時ルーチン自身 (図 1 (c)) が直接 OS のデータ管理マクロを使用して入出力を行うように書き改めることにした。これによって図 1 (d) の層から FORTRAN 実行時ルーチンはなくなる。

(3) TSS 環境での使い勝手 TSS 環境での使い勝手をよくするために次の機能を追加する [5]。

(a) SNOBOL 4 プログラム中で TSS コマンドからの実行時パラメータを受け取ることができる。

(b) SNOBOL 4 プログラム中から TSS コマンドの実行ができる。

このうち (a) は SNOBOL 4 処理系の初期設定時に TSS コマンドからのパラメータを SNOBOL 4 プログラムから読めるように処理系内部にコピーすることで実現する。(b) の実現には TSS コマンドを実行する外部関数を作成し、SNOBOL 4 の外部関数呼び出し機能を用いて SNOBOL 4 プログラムと結合する。そのため SNOBOL 4 処理系内部の変更はない。

入出力ルーチンを書き改めたことにより SNOBOL 4 の組み込み関数の一つである OUTPUT の仕様を次のように変更し、ほとんど使われず、実現が面倒な BACKSPACE 関数を廃止する必要があるので、

OUTPUT (出力変数, 論理機番, 制御文字)

従来は OUTPUT 関数の第 3 パラメータには FORTRAN の FORMAT 記述子を書き、改ページ、重ね打ち等を自由に行っていたが、入出力を実行時ルーチン中で直接行うようにするので、FORMAT 記述子を自由に使えなくなるためである [5]。OUTPUT 関数の第 3 パラメータは省略可能なパラメータである。SIL で書かれたもの以外の SNOBOL 4 処理系 (SPITBOL 等) の OUTPUT 関数の第 3 パラメータには FORMAT 記述子を使用しないので、互

換性を考えて作った、既存の SNOBOL 4 プログラム (たとえば [7]) にはこの仕様の変更は影響しない。

#### 4. 2 プログラム例

図 2 に日本語 SNOBOL 4 プログラムの例を実行回数情報 (図の左側の 3 列の数字が右側の各文の実行数、成功回数、失敗回数 [8]) 付きで示す。変数名とラベル名に日本語文字を用いた。この例は入力日本語テキスト中に現われる漢字列を抜きだしてその出現回数を計数するプログラムである。このプログラムでは '亜' が文字コードの上で漢字以外の文字との境界にある漢字として漢字の判定に利用している。このプログラムの入力テキストと出力結果の一部を図 3 に示す。

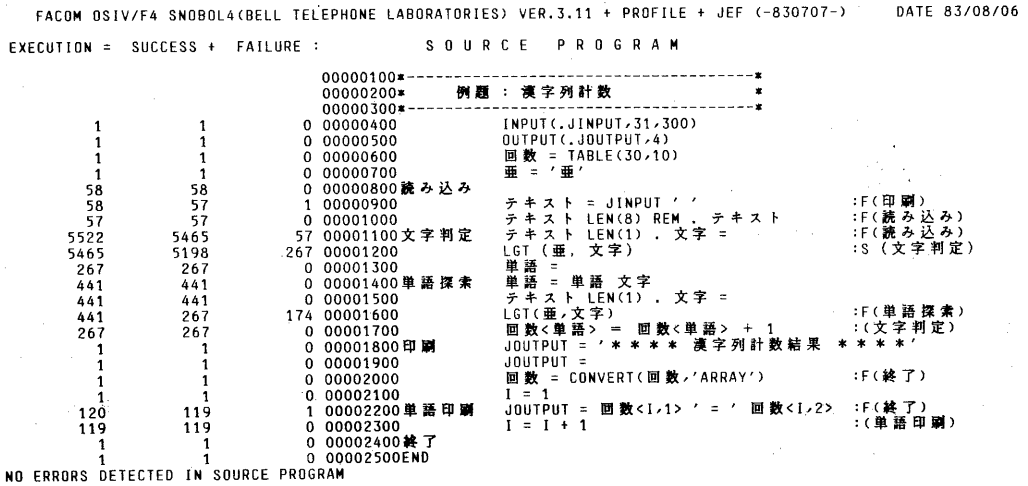


図 2. プログラム例 (漢字列計数プログラム)

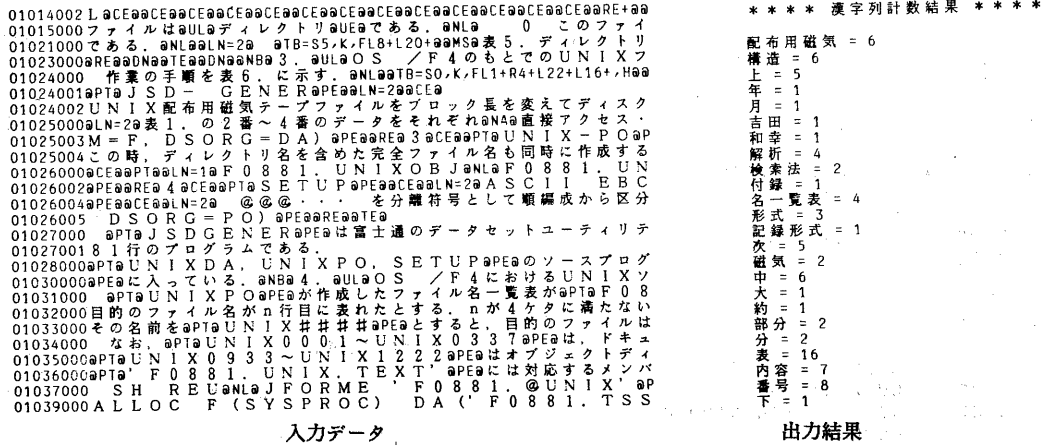


図 3. 漢字列計数プログラムの入力データと出力結果

## 5. 処理速度の改善

日本語 SNOBOL 4 処理系でプログラムを実行するときプログラム中でパターン関数 BREAK, SPAN, ANY, NOTANY を多用していると実行速度が極端に遅くなるのがわかってきた。これらの関数は上述の動的な表を書き替えて使用し、次の文字列とパターンマッチする。

BREAK : 引数に現れない文字から成るなるべく長い文字列。

SPAN : 引数に現れる文字から成るなるべく長い文字列。

ANY : 引数に現れる文字のうちのどれか 1 文字。

NOTANY : 引数に現れない文字のうちのどれか 1 文字。

これらの関数は表を検索する S I L 命令 S T R E A M を用いている。個々の計算機が実現しやすいアルゴリズムを採用できるように表の検索ルーチンは実行時ルーチン中 (図 1 (c)) にあり、S T R E A M 命令はそのルーチンを呼ぶように定義されている。

以下処理速度を遅くしている原因を究明するために表の構造と機能について調べる。

### 5. 1 表の構造と日本語文字への適用

SNOBOL 4 原処理系 (OS 360用) が字句解析に用いる表の構造は機械命令 TRT (翻訳テスト命令) に大きく依存している。TRT 命令は、第 1 オペランドの文字列を先頭から 1 バイトずつ第 2 オペランドで指す表を用いて変換して行き、変換結果が 0 でなくなったところまたは文字列が尽きたところでやめる。したがって第 2 オペランドで指定される表は 256 個のエントリを持つ。たとえば先頭番地 SNABTB と指された表で空白文字に対応する値 (文字に対応するコマンド) は SNABTB+X'40' に入れておく。SNOBOL 4 原処理系では動的な表は図 4 のように、256 バイトの項目部とそれに続く 3 個の機械命令からなる命令部で構成される。図 4 で stop や

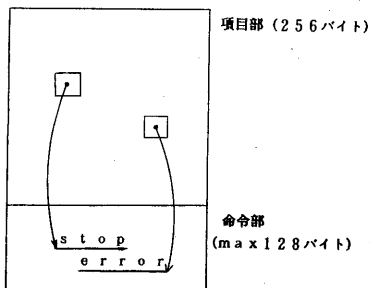


図 4. SNOBOL 4 原処理系の字句解析表

error がコマンドである。

この方法を日本語文字にそのままあてはめると、文字コードが 2 バイトなので一つの表が 64 K バイト (256 \* 256) になり、表全体では 1.6 M バイト (=64 K バイト \* 25) にもなってしまう。そこで、表の圧縮を試みた。TRT 命令は 1 バイトを単位として処理するため使用できない。

字句解析ルーチンがコンパイル時に実際に区別すべき字種は演算記号、カッコ等 20 種程度であるので、表には文字コードと対応するコマンドを並べて記入し、先頭から順に比較していくというアルゴリズムを採用して表の圧縮を行った。ただし、日本語文字、英字、数字は文字種としてはそれぞれをひとまとめにして考えてよいが、その数が多く、表中に全文字を入れることは大変なので表中の文字コードの項には日本語文字、英字、数字を表わすマークを書いておき、表検索ルーチン中でそのマークをみて当該字種の上限および下限と比較することで文字種の判定を行う。

静的な表はこれではよいが、このような構成では実行中に表を書き替えることは難しい。まず動作するものをつくるという方針に従って、動的な表の管理は SNOBOL 4 原処理系のアルゴリズムを (TRT 命令は使えないが) そのまま用いることにし、動的な表は 64 K バイトの領域をとりあえず持つことにした。この表を実行中に書き替える際は、まず表全体をクリアした後に必要なところに値を書き込むようになっていた。このためにこの表を使う ANY, NOTANY, BREAK, SPAN の実行には非常に時間がかかるのであった。

### 5. 2 表の構造の変更

実行速度が遅くなった原因は、64 K バイトの表をクリアするところにあると予想されるので、表自身を小さくする必要がある。表検索ルーチン中で文字コードを上位、下位バイトに分けて考え、表もそれに合わせて上位バイトで検索する主表とそれの各項からポインタで指される下位バイトのための副表に分けた (図 5 参照)。JEF コード (およびそのもととなっている JIS 漢字コード) は上位バイト、下位バイトそれぞれだけみると連続しているが 16 ビット全体として見ると不連続である。このため 16 ビットをひとまとめにして表を引くよりも上位、下位バイトに分け、上位、下位の順で表を引くほうがより JEF コード (JIS 漢字コード) にあった方法であろう。

この方法では、表の書き替えは上位バイトによって主表を引き副表を探し下位バイトによってその副表を書き替える。表中に JEF コー

ドの文字割当がない部分や下位バイトで区別する必要がない部分の副表は作らずにすませるので、64Kバイトあった表を大幅に圧縮できる。主表に書き込む値は原処理系と同様なコマンドか副表へのポインタ(X'80'~X'FC')かである。主表にcontinue等のコマンドが直接書き込まれている場合には下位バイトの値に関係なく(副表は作らずに)コマンドの指示通りに文字列の解析を続ける。副表にはコマンドが書き込まれる。

この方法を採用して作った動的な表と先に別の方法で圧縮していた静的な表との大きさの上での差はほとんど解消できる。2種類の違う構造の表を管理することは保守上から好ましくないので静的な表もすべて動的な表と同じ構造にした。

上記の変更に伴う処理系の書き替えについて触れる。4.1節(1)文字の内部表現の項でも述べた、表に関係する3つのSIL命令は次のものである。これらは、次の動作をする。

- CLERTB: 表のクリア
- PLUGTB: 表の文字に対応する位置に値を入れる。
- STREAM: 表を用いて文字列を解析する。

このうちSTREAMマクロだけが上で述べてきた表検索ルーチンを呼ぶ。そのルーチンのインタフェース仕様を変更しないのでSTRE

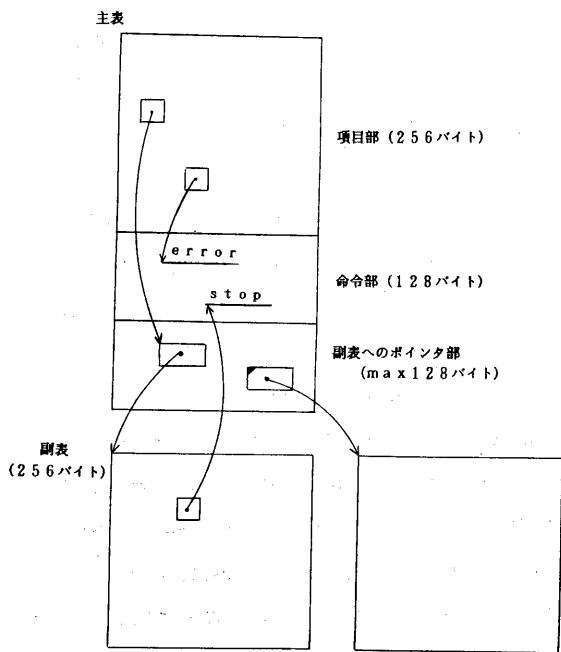


図5. 日本語SNOBOL 4 処理系 (改良版) の字句解析表

AMマクロも変更の必要がない。

CLERTBは主表256バイトのクリアする命令と副表へのポインタ領域のクリアする命令とをマクロ内で展開する。主表256バイトをクリアする部分は原版とはほぼ同じ形にもどった。

PLUGTBは従来の単一構造の表では文字コードに対応する位置に値を書き込むだけであったのでマクロ内で展開していたが、表を階層構造にしたために、値の書き込みが多少複雑になり、副表の管理も行うようになった。そのためその本体を実行ルーチン内に移しマクロはそのサブルーチンを呼び出すように変更した。

## 6. SNOBOL 4 処理系の実行速度の比較

SNOBOL 4 処理系への日本語テキスト処理機能の追加、字句解析表の変更による処理系の実行効率の変化を調べるために6個のSNOBOL 4 プログラムについてその実行時間の比較を行った。ここで比較した処理系は実行回数計数機能追加版SNOBOL 4 (S4Dと略す、原処理系に実行回数計数機能を追加したもの[8])、日本語SNOBOL 4 (S4J, JAPANESE)、表構造を変更したSNOBOL 4 (S4J (改良版))の3つの処理系である。コンパイル時間、実行時間の測定には各処理系の統計情報出力機能(SUMMARY)を用いた。さらにそれぞれの処理系で実行が正しく行われたことを確認する意味で実行回数計数機能を動作させている。なお、これらの測定はFACOM M-200 OSIV/F4のもとで行った。コンパイル時間、実行時間の測定結果を表1に示す。

STINGYは入力データを編集して出力するプログラムである。このプログラムのS4Jでの実行時間はS4Dでの時間に対して1割程度の増加で収まっている。S4Jがすべてのプログラムに対してこの程度の実行時間の増加で実行できるならばよかったのであるが、実は前章で述べたようにそうではなかった。LABELはFORTRANプログラムの文番号をつけなおすプログラム、POLYNOMIALは多項式の加算を記号的に実行するプログラムである。この2つのプログラムではSPAN関数、BREAK関数を多数使用している。そのためS4Jでの実行時間はS4Dでの実行時間と比べて極端に大きくなっている。SPAN関数、BREAK関数の実行回数が多いPOLYNOMIALプログラムの方がその差は著しい。

SPAN関数等が実行速度を遅くする原因であることを確かめるために作成したプログラムがSPANTESTプログラムである。この

プログラムはSPAN関数等を用いたパターンマッチを繰り返して実行するプログラムである。SPANTESTプログラムではS4DとS4Jとでの実行時間が3.5倍も違う。

S4J (改良版) での実行時間はLABELプログラム、POLYNOMIALプログラム、SPANTESTプログラムもSTINGYプログラムと同様にS4Dでの実行時間とほぼ同じである。なお、STINGY、LABELでS4Dでの実行時間よりもS4J (改良版) での実行時間が短くなっているのは入出力ルーチンをアセンブリ言語で書き直した(4.1節参照)ためであると考えられる。また、すべてのプログラムでコンパイル時間が短くなっているのも同じ理由と考えられる。

日本語テキスト処理機能を用いたプログラムとして日本語の文章中の英字列の相互参照表を作成するXREFAプログラムと漢字列計数プログラム(図2に示したもの)とのコンパイル時間、実行時間を示す。これらのプログラムはS4Dでは処理できない。漢字列計数プログラムの作成時にはSPAN関数等が実行時間に大きく影響することがわかっていたのでその作成に当たってそれらのパターン関数を使用しないようにしている。そのためS4JとS4J (改良版) との間で実行時間の大きな差は認められない。なお、漢字列計数プログラムではSPAN関数を用いることによって2196ミリ秒かかっていた実行時間を618ミリ秒に短縮できた。

表1. SNOBOL 4プログラムの各処理系における実行時間の比較

上: コンパイル時間 (ミリ秒)  
下: 実行時間 (ミリ秒)

処理系 プログラム名	S4D	S4J	S4J (改良版)
STINGY (39行)	65 2795	111 3059	54 2349
LABEL (192行)	218 15402	374 22264	142 13024
POLYNOMIAL (135行)	154 669	237 9781	97 703
SPANTEST (9行)	44 400	49 14344	29 416
XREFA (262行)	-- --	390 11343	142 4049
漢字列計数 (25行)	-- --	70 2399	34 2196

## 7. 議論

### 7.1 処理系の設計と日本語テキスト処理機能の追加

日本語機能の追加ではSNOBOL 4原処理系のSILで書かれた部分(図1(a)に属する)は変更の必要がなかった。SILで書かれたSNOBOL 4原処理系は10年以上も前に設計実現されたにもかかわらず、日本語機能の追加作業をこのように比較的容易に遂行できたのは、ドキュメント[6]がよく整備されていたこと、SILシステムの設計が明快でよく機能分割されていたことが大きな理由と考えられる。以下に要点をまとめる。

- (1) SILシステムの構成 3章にまとめたように階層的に構成されており、機械またはOSに独立な部分と依存する部分が明確に分離されている。たとえば4.2節で述べた字句解析を行うための表(その構造は機械、文字コードに依存する。)ではそれを操作するSIL命令は3個とそれらのマクロの下請けサブルーチンだけが表の構造に依存する。このようにSNOBOL 4処理系の作成に当たってデータ抽象化の概念が取り入れられているため一部の修正がプログラムの広域に影響を及ぼすことがなかった。
- (2) SIL言語の定義 SIL言語の文法に関するドキュメント[6]によりSIL言語の修得が容易であった。そのため文字の内部表現の変更によるSILマクロの書き替えもSIL言語の外部仕様に合わせてよくなかば機械的に行うことができた。
- (3) 実行時ルーチンの仕様 1つの実行時ルーチンは1つのSILマクロに対応しており、その外部仕様がSILマクロ定義やそのドキュメント[6]によりはっきりわかる。さらに実行時ルーチン間の副作用がないため書き替えが必要なルーチンは容易に書き替えられた。

### 7.2 入出力の設計

SNOBOL 4原処理系では入出力はFORTRAN実行時ルーチンを用いて行っていたが、今回、可変長ファイルに対して読み書きできるようにするために入出力部を全面的に書き替え、FORTRAN実行時ルーチンをSNOBOL 4処理系の中から排除した。これはFORTRAN実行時ルーチンの利用に関する情報がほとんどなく、むしろOSのデータ管理マクロのマニュアル等の方が完備しているので入出力を行うプログラムを書き易いためである。

このことに関してSNOBOL 4処理系の原著者であるGriswoldは[9]の中で「SNOBOL 4の設計時には移し換えを考慮

してどの計算機にもあるFORTRAN実行時ルーチンをSNOBOL 4の入出力部として採用した。FORTRANはよく知られた言語であるのでSNOBOL 4を移し換えようとする程度に計算機を知っている人ならばだれでもFORTRAN実行時ルーチンを手軽に使用できると考えたためである。しかし現時点で考えてみると、FORTRAN実行時ルーチンはだれでも使えるというものではない。むしろ、SNOBOL 4用の仕様の入出力ルーチンを設計した方がよかったのではないか。」と述べている。今回の日本語機能追加ではこのことがはっきりと現われた。

### 7. 3 日本語SNOBOL 4の移換性

SNOBOL 4原処理系は移し換えを考慮して設計されている。日本語テキスト処理機能の追加により移換性が損なわれていないかどうかを調べるのが今後の課題である。互換機といわれているFACOM, HITAC, IBMの計算機ではユーザ側から見たアーキテクチャやOSは同じであるがそのうえで作成されたソフトウェアには同じような機能を全く別の仕様で実現しているものもある。日本語処理機能もその一つであり、日本語テキスト保存ファイルの形式、日本語文字コード、日本語文字とEBCDIC文字を区別するシフトコード等に違いがあるため、日本語処理機能には互換性がない。日本語SNOBOL 4処理系がこれらのことに容易に対処できることが移し換えの容易さの点で要求されることである。

### 8. むすび

SILで書かれたSNOBOL 4処理系に日本語テキスト処理機能を追加し、さらにその効率改善を行ってきた。SILの設計は元来、多くの計算機システムの上でSNOBOL 4処理系を実働させることが第一義であったものと思われる。しかしその設計の故に、機能拡張を容易に行うことができたことを強調しておく。

本研究の一部は文部省科学研究費補助金（一般研究課題番号57460211）による。

### 参考文献

- [1] 武富, 高木, 川崎他: 日本語情報システムJEFの使用法. 九州大学大型計算機センター広報, Vol. 13, No. 4, pp. 406-468, 1980.
- [2] L.Cherry: Computer Aids for Writers, ACM SIGPLAN NOTICES, Vol.16, No.6, pp.61-67, 1981.
- [3] R.E.Griswold, J.F.Poage, I.P.Polonsky: The SNOBOL4 Programming Language, Prentice-Hall, Inc., 1968.
- [4] 木村: SOFTWARE TOOLとは何か, 情報処理, Vol. 20, No. 8, pp. 673-680, 1979.
- [5] 牛島, 黒坂他: JEF日本語処理機能を追加したSNOBOL 4について, 九大大型計算機センター広報, Vol. 16, No. 2, pp. 155-179, 1983.
- [6] R.E.Griswold: The Macro Implementation of SNOBOL4, W.E.Freeman & Co., 1972.
- [7] J.F.Gimpel: Algorithms in SNOBOL4, John Wiley & Sons, 1976.
- [8] 吉田, 牛島: 実行回数計数機能を追加したSNOBOL 4処理系とその移し換えについて, 情報処理学会論文誌探録決定, 1983.
- [9] R.E.Griswold: A History of the SNOBOL Programming Languages, ACM SIGPLAN NOTICES, Vol. 13, No.8, pp.275-308, 1978.