

## F O R T R A Nにおける日本語化の現状と今後の方向

黒田 幸明

NTTソフトウェア研究所

F O R T R A Nは、日本で最初に日本語機能の標準化が検討された言語であり、現在は、この結果を次期 F O R T R A Nの I S O規格 (F O R T R A N 8 X) に提案中である。

ここでは、F O R T R A Nの日本語機能に関する次の項目について説明する。

- (1) F O R T R A Nで日本語データを処理するには、新しい日本語文字型の導入が必須であること
- (2) エスケープシーケンスの隠ぺい等日本語機能として実現すべき項目
- (3) 実現方式
- (4) F O R T R A N 8 Xに提案している言語仕様の概要
- (5) ソースプログラムの日本語化等今後の検討課題

How to implement the national character handling facility on F O R T R A N and its proposal to F O R T R A N 8 X

Kohmei KURODA

NTT Software Laboratories

1-9-1 Kohnan Minato-ku Tokyo 108 Japan

Standardization of the national character handling facility was studied on F O R T R A N for the first time. This result is now proposed to F O R T R A N 8 X of I S O.

This paper describes as follows.

- (1) Necessity of national character data type
- (2) Requirements for the national character handling facility
- (3) Present implementation
- (4) Summary of language specifications proposed to F O R T R A N 8 X
- (5) Future tasks

## 1. はじめに

FORTRANは、日本で最初に日本語機能の標準化が検討された言語である<sup>1)</sup>。現在は、この結果を次期FORTRANのISO規格(FORTRAN8X)<sup>2)</sup>に提案中である。

ここでは、FORTRANの日本語機能の現状について紹介するとともに、今後の検討課題について述べる。

## 2. 日本語文字の特徴と問題点

### 2.1 特徴

- (1) 字数が10000字以上あるので、ASCII文字のように1バイトでは表現できない。日本では、普通2バイトで表現している。
- (2) 字数が多いので、文字のコード値による大小比較は、あまり意味がない。電話帳などは、特殊なソートを行っている。
- (3) 日本語文字のうち、漢字は1文字毎に意味をもっているため、変数名に漢字を使うとプログラムが読みやすくなる。

### 2.2 問題点

従来の文字型データで日本語文字を処理する場合、次の問題がある。

- (1) 日本語文字と単バイト文字を区別するために、エスケープシーケンスが必要になる(図1)。

見た目の文字数と計算機の内部表現に差が生じるので、レコードの設計などで、エスケープシーケンスの長さを考慮しなければならず、プログラミングがむずかしくなる。

漢字端末の表示 (14文字)	ISO標準FORTRAN
内部表現 (20バイト)	49 53 4F 1B 24 42 49 38 30 60 I S O エスケープ 標 準 シーケンス
	1B 28 4A 46 4F 52 54 52 41 4E Iスケープ° F O R T R A N シーケンス

図1 日本語文字と単バイト文字の区別の方法

- (2) 分離符と同じコード値をもつ日本語文字を定数に書くと翻訳エラーになる(図2)。

ソースプログラム '規則'

内部表現 27 1B 24 42 35 2C 42 27 1B 28 4A 27  
' Iスケープ° 規 則 Iスケープ° '  
シーケンス シーケンス

コンパイルの解釈 '\$B5,B'(J'

図2 文字型データには書けない日本語文字の例

## 3. 要求条件

### 3.1 プログラムの機能

(1) エスケープシーケンスが隠べいされること。具体的には、次の内容を実現する必要がある。

- ① 入力時のエスケープシーケンスの自動削除及び出力時のエスケープシーケンスの自動挿入
- ② 部分列処理などでの日本語文字の数が、バイト数でなく文字数であること

(2) すべての日本語文字を書ける日本語文字型が導入されること。

### 3.2 ソースプログラム

(1) 日本語文字の変数名が使えること。

(2) 日本語文字を含んだ注釈が使えること。

なお、(1)については、FORTRAN8Xの提案に含まれていないので、7. 今後の検討課題で説明する。

(2)については、既存の仕様のままで実現できるので、説明を省略する。

## 4. 実現方式

国内のほとんどのメーカーでは、FORTRANの日本語機能の標準化検討が開始される以前に、日本語機能を実現済みであった。その実現方式は、細かい仕様の差異はあるものの、大筋は同一の実現方式である。ここでは、その実現方式について述べる。

### 4.1 処理方式

次の処理を行うことにより、エスケープシーケンスをプログラムに入れられない方式とし、3.1の内容を実現している。

- (1) 日本語文字のみからなる日本語文字型を導入する。この情報を用いて、
- (2) 定数中のエスケープシーケンスは、翻訳時に自動削除する。
- (3) エスケープシーケンスは、入出力時に自動削除挿

入する。

#### 4. 2 日本語文字型の制約

4. 1 で示した方式は、3. 1 の要求条件を実現するために定めたものであり、従来の文字型データの機能を包含したものではない。

日本語文字型データは、制御コードを含まない日本語文字だけの文字列処理に用いる場合に有効である。

エディタなどのように制御コードを含んだ文字列処理を行う場合は、従来の文字型データを使う必要がある。

#### 4. 3 プログラム例

図3に、日本語機能がある場合とない場合のプログラム例を示す。図3(b)の例は、日本がISOに提案中の言語仕様案に基づいて書いたものである。

```
PROGRAM EXAM1
CHARACTER A*14,B*8
A='国内標準'
C      (Aの内容) 1B 24 42 国内標準 1B 28 4A
WRITE(2,100) A
B='際'
C      (Bの内容) 1B 24 42 際 1B 28 4A
A(6:7)=B(4:5)
WRITE(2,100) A
C      「国際標準」が出力される
100 FORMAT(1X,A14)
STOP
END
```

(a) 日本語機能がない場合

```
PROGRAM EXAM2
NCHARACTER A*4
A=NC'国内標準'
C      (Aの内容) 国内標準
WRITE(2,100) A
A(2:2)=NC'際'
WRITE(2,100) A
C      「国際標準」が出力される
100 FORMAT(1X,N4)
STOP
END
```

(b) 日本語機能がある場合

図3 プログラム例

#### 5. 言語仕様

日本がISOに提案中の言語仕様案の概要を次に示す。

##### 5. 1 ソースプログラム

(1) 日本語文字を含むときのソースプログラム1行に書ける文字数は、処理系依存である。

##### 5. 2 文字の大小順序と記憶単位

(1) 文字の大小順序は規定しない。

(2) 文字データと日本語文字データ間の記憶単位の関係は、規定しない。

(3) 日本語文字型データの記憶単位の共有は、同じ型のときのみ可能である。

##### 5. 3 定数と型定義

(1) 日本語文字の定数は、次のように書く。

NC'日本語' 又は NC"日本語"

(2) 例えば3文字の日本語文字型の変数Aは、次のように書く。

NCHARACTER A\*3

##### 5. 4 部分列

(1) 日本語文字の数え方は、文字数で行う。

##### 5. 5 演算

(1) 代入、結合(/)及び比較(.EQ.と.NE.)ができる。

##### 5. 6 入出力

(1) 書式付き記録に、単バイト文字と日本語文字が混在してもよいが、書式付き記録の長さは、単バイト文字に換算して計算する。換算方法は処理系依存である。

(2) 内部ファイルでは、日本語文字を扱えない。

(3) 日本語文字データの編集用に、N[w]形編集記述子を追加する。

(4) エスケープシーケンスの削除挿入の指定方法は規定していない。

##### 5. 7 組込関数

(1) 単バイト文字用と同等の組込関数を用意する。

#### 6. FORTRAN日本語機能の検討体制及び経緯 検討体制及び経緯を図4に示す。

現在は、FORTRAN8Xに対するISO/TC97/SC22でのDIP投票及びANSI/X3でのパブリックレビューの中で、日本の提案が検討されている。

#### 7. 今後の検討課題

##### 7. 1 型定義と定数の書き方

日本語文字型の定義方法として、次に示す2案があり、

ISO及びANSIでは、案1の支持者の方が多い。

(案1) 複数の多バイト文字集合を扱えるようにするために、文字集合の番号を指定する方法<sup>3)</sup>

例: CHARACTER (KIND=n) A\*3 (n=1,2,3...)

(案2) 5.3-(2)で示した方法<sup>4)</sup>

例: NCHARACTER A\*3

また、定数の書き方は、5.3-(1)で示した

NC'~'に対して、全言語N'~'に統一しようというのが国内の動向である。

実際のニーズ及びCOBOL等他の言語の動向を考慮すると、型定義は案2、定数はN'~'とするのが適切であると思われる。

### 7.2 ソースプログラムの日本語化

ソースプログラムの日本語化については、変数名に日本語を使えるようにすることに加えて、プログラムを入力しやすくするために、プログラム全体を2バイト文字で入力したいというニーズがある。

前者については、単バイト文字と日本語文字で同形の文字の扱いを規定すれば実現できる。

後者については、定数の分離符だけは単バイト文字でないと従来プログラムとの互換がとれなくなるという問題がある。

例えば、定数の中に日本語文字の'を書いた従来のプログラムを新しい処理系に入力すると、定数内の'が定数の分離符とみなされて翻訳エラーとなる。

これを解決するには、コード系を2バイトに統一するより方法がない。

### 7.3 処理系依存の項目の扱い

5.で示した言語仕様案には、処理系依存の項目が多い。これは、コード系の影響が言語仕様にはじみ出るからである。

これを解決するにも、コード系を2バイトに統一するより方法がない。

### 8. おわりに

4.2及び7.で示したように、現在提案中の日本語機能は、従来の文字型データの機能を完全に包含したものではない。根本的な解決をするには、計算機のコード系を2バイトに統一するより方法がないことを認識しておく必要がある。しかし、現実にはこのような対処は極めて困難であり、利用者のニーズ及び計算機のメモリとファイル効率を考えると、現在の提案が最も妥当な内容であるといえる。

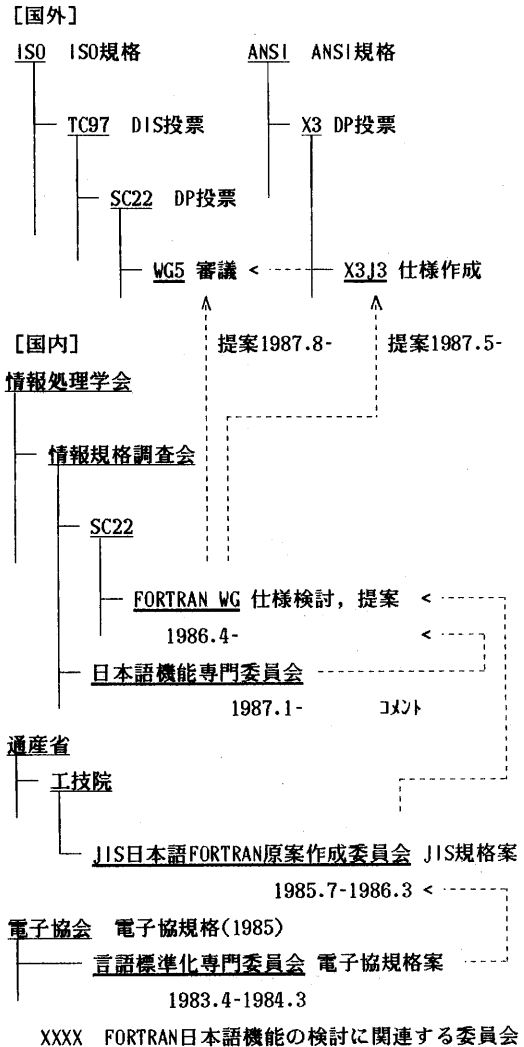


図4 検討体制及び経緯

#### 謝辞

本資料作成に当たり、貴重なご意見をいただいた国内SC22/FORTRAN WGの皆様には深謝いたします。

#### 参考文献

- 1) 電子協: 日本語FORTRAN JEIDA-42-1985, (1985).
- 2) ANSI/X3J3: FORTRAN8X S8.104, (1987).
- 3) 国内SC22/FORTRAN WG: Japan's Proposal to FORTRAN 8X, SC22/WG5リバプール会議資料, (1987).
- 4) 国内SC22/FORTRAN WG: Japanese Proposal to FORTRAN 8X--NCHARACTER type for National Character Handling--, ANSI/X3J3 106会議資料, (1987).