

プログラム言語で使う文字コード

和田 英一

東京大学工学部計数工学科

プログラム言語で書くプログラムのうち計算機入力用には金物表現が重要であり、それには使う文字セット、文字コードを規定する必要がある。プログラムの中でも注釈や文字列定数は比較的自由に文字セットが拡張できる。この報告では、これまでの主要なプログラム言語で使った文字セットとコード(FortranのEBCDIC, APL文字コードなど)について概観する。さらに、現在主として使われている文字セットのASCIIコードの構造を述べる。ASCIIコードを基本としたいくつかの8ビットコード系が整備されつつある。また日本語プログラミングの関係ではかつてはカナによるプログラム言語があった。データとしては漢字などの2バイトコードが一般的に使えるようになったが、Shift JISはなるべく使わないようにしたい。

Character codes used in programming languages

Eiti Wada

Department of Mathematical Engineering, University of Tokyo
Bunkyo, Tokyo 113 JAPAN

For programs in hardware representation to be input into computers, graphic character sets (and their codes) must be defined. Even when the syntax is defined strictly, it is relatively easy and often permitted to expand the character sets for literal constants and comments to include kanji and hiragana characters.

Character sets used in historical languages are surveyed. Then structure and related issues on ASCII code is discussed. Several 8 bit standard character sets are now being prepared which are likely used in programs. Structure of Shift JIS code is finally discussed and it is recommended not to use the Shift JIS code in programming environment.

はじめに

プログラム言語で記述したプログラムには人に見せるためのものと計算機に入れるためのものがあるとAlgol 60ではいっていた。つまり規準言語 (reference language)、出版言語 (publication language)、金物言語 (hardware language) があることになっていた [1]。規準言語は言語を定義するためのものという感じである。Algol 60の委員会はもっぱら規準言語にばかりかかわっており、金物言語はインプリメンター任せだったため、いろいろな流儀の金物言語が作られ、それもあり便利でないものだったので結局Algol 60のコミュニティは育たなかった。他の成功したプログラム言語は「はじめに金物言語ありき」という利用上のことを考えた発想で出発している。

金物言語というとそれを表現するための図形文字セット (とそのコード) が重要であり、現在は ASCIIコード全盛ともいう感じであるが、その他にもプログラム言語で使っている図形文字セットはいろいろあるので、この報告ではそれをサーベイしてみたい。なお、図形文字というのは復帰改行のような制御文字に対する言葉だがプログラム言語の場面では、まず図形文字が中心話題なので以下では単に文字セット、文字コードということにする。

プログラム言語における文字の使用区分

プログラムは文字の並んだものという構文規則もありうる筈だが実用上はもっと詳細に定義がしてある。それをみるとプログラム言語での文字の使われ方は文字の使われ方のうるさきについて次の 4通りに分類できよう。

1. 演算子 (+ := など)、綴りの記号(if begin など)、区切り記号(; など)
2. 識別名 (i instream foo など)
3. 文字列定数 ("argument type error")
4. 注釈 (comment standard name table initialization)

このうち下のものはどインプリメンターが文字セットを自由拡張できる。プログラム言語の文法書にそう記述してあるものもある。

テレタイプコード

最初の電子計算機は大学や研究所みたいなところで作られたから入出力にはフレクソライタ (電動タイプライタ) やテレタイプを使うことが多かった。国際のテレタイプは 5 単位で英字の大文字、数字、特殊記号だけの文字セットである (ITA 2 International Telegraph Alphabet No.2)。これは欧文モールスコードと同じ文字セットである。ケンブリッジ大学の Edsac はテレタイプを使っていたが、10進 2進のコード変換を簡単にするため、コードは国際テレタイプのそれとは同じではない。テレタイプのコードは紙テープにパンチしたとき、くずをなるべく出さないためか出現頻度の高い文字には 1 の少ないコードをわりあっていた。モールスコードが短かいのと同様である。従って QWER...に 1234... が対応していると数字は 2進法のパターンとはきれいには対応せず、Edsac ではそのためコードを変えてしまった。東大の TAC は Edsac をまねして作ったからコードも同じだったかと思う。

高橋研の PC-1 は国際テレタイプと基本的には同じで 6 単位のコードのものを使っていた。

EBCDIC Fortran 48文字セット

TAC が完成して少し経った頃 Fortran が登場した。これは IBM で開発したので IBM カードの文字を使う。これは図 1 に示すような 48 文字セットであるが (,), =, - は事務用のセットでは □, %, #, & になっていた。このうち Fortran が使っていたのは 47 文字で 8-4 パンチの - は何のためにあったのかよくわからない。

Fortran とほぼ同様に古い歴史の Lisp 1.5 は IBM 704 用に作られたので文字セットは Fortran と同じである。なにしろ記号アトムは英数字だしその他に () . しかいらないのだからこんなにシンタックスシュガーを取り入れない言語も珍しい。しかし Lisp 1.5 のマニュアルには次の記号アトムには対応する文字が値として結合していると書いてある。

```
dollar $ slash / lpar ( rpar ) comma , period .
pluss + dash -(11punch) star * blank eqsign =
```

- に 11 パンチと注がついているのは 8-4 パンチにもうひとつ使わない - があったことを示している。

Cobol もほぼ同じ文字セットであった。

日本語プログラムとカナ文字

我が国でも Fortran, Algol, Cobol の処理系が次第に作られるようになってくるとカナ文字の使えるものができるた。つまりカナ文字 Fortran とカナコボルである。面白いことにこのふたつはカナ文字の使える場所が違っていた。カナ文字 Fortran は綴りの記号のところをカナ文字化した。例えば Read をヨメのように書く。一方カナコボルは識別名にカナを許した。それにはコンパイラの字句読み込みルーチンを多少直し、名前表にカナ文字列も入るようにするだけでよい。Fortran の方はもともと行の途中の空白は無視するという構文規則だから、構文規則は特殊文字だけに頼らざるを得なかったのだが、文の識別にカナ文字という強力な助っ人がきたのでコンパイラ作りは容易になったものと察せられる [2] [3]。

カナ文字による日本語プログラミングはその後しばらく見なかつたけれども水谷氏の小朱唇も綴りの記号をカナ文字にするプログラム言語である。識別名はローマ字しか使えない [4]。

APL

APL が現れたときはその異質な文字セットと右から左という演算順序は人々を驚かせたものだ。これは特殊な端末

からしか使えない、つまりAlgol 60の金物表現の二の舞かと思われたが、IBMはAPL用のタイプエレメントまで用意したのでAPLは愛好者の間では使われるようになった。この文字セットのコード表は図2の通りで图形文字セットとして登録してある。

数学者は次々と記号、記法を発明し思考の節約をはかるものだが、既存の文字セットを後生大事にしてbegin endifなどやらず新しい文字セットと端末を作る態度にはある意味で感心させられる。こういうプログラム言語は恐らくAPLだけであろう。(BackusのFPもイリノイでインプリメントしたがどういう文字セットだったか忘れた。)

ASCIIコード

この業界で標準のような顔をしているのがASCIIコードである。ASCIIはもともとAmerican National Standard Code for Information Interchangeだから、ASCIIコードというのはコードがだぶって正しくないという人もいるがASCIIコードといわないとするわりが悪いのでそう呼ぶことにする(図3の左)。

これはISO 646という7ビットの標準で決めた基本コード表というのがあるのだがその仲間である。このコード表には10ヶ所ほど各国で自由に使っておられる場所(各國場所)がある。また2/3, 2/4はそれぞれナンバーサインかボンド、カレンシーサインかドルを入れることとなっている。もし各国で特に何も入れないときはそこに適当な記号を入れたIRV(International Reference Version)というコード表があり、この記号を使うことになっている。ANSIはIRVの2/4の場所を\$にしたもの自国標準としたこれがASCIIである。(よく見ると7/14も少し違う。)

IBMはパソコン以外はEBCDICで悠然としているが、世の中にはASCIIの端末が多く出回り、そのためASCIIに準拠したプログラム言語が増えてきている。Pascal, C, Ada, Common Lispなどがそうである。しかしこれらの言語もASCIIがなければプログラムが書けないというのではない。例えばPascalはFortranの48文字程度でプログラムが書けるよう、ある程度の文字の代用を決めている。[,],↑は(...),@を使ってよい。Adaもそうで|は#は:のように決めている。Adaのこの手当は各國場所を避けたかったからである。

ISOのコードについてはよく知らないのだが少し前のASCIIは5/14, 5/15が↑, ←になっていた。(今は^と_である) Smalltalkは値を返すのに↑を、代入に←を使っているからその仕様は←, ↑の時代に決まったものであろう。キーボードやfileoutフォーマットのファイルから^や_を読み込むと画面には←や↑ができるから一安心である。

プログラム言語で使うことを考えると^や_より↑や←の方が適しているような気がする。10年ほど前に当時院生だった近山君が8080用に「小さいLisp」を作った。これは←がcar, >がcdrのように特殊文字1字が基本関数名になっていた。↑や←はそれぞれrplaca, rplacdに使うと丁度よいという話もあったがそれはそうしなかったと思う。

ASCIIコードは今やRIVを駆逐せんばかりの勢いである。2/4にカレンシーサインの入ったIRVの端末など実在しないし誰も使おうとしないからである。

ISO 8859 8ビットコード

↑が^になったのは^を発音区別符号に使うためと思われる。IRVのうち"はダイアレスに、"はセリダに、"はアクチュートアクセントに、"はグレープアクセントにも使えるようにできている。そして^はサーカムフレックス兼用とした。これらの発音区別符号を使うためにはパックスペースが必要である。また欧州にはまだ変な字もありこれでは欧州各国で通信するには不便だというので、テキスト通信用のコードの標準をISOとCCITTで決めた。ISO 6937というのがそれである。ISO 6937は7ビットの使い方もあるが一応は8ビットコードで、左半分にIRVを置き、右にはノンスペーシングの発音区別符号やその他の文字を入れ、320くらいの字が使えるようになっている。しかしこれは使ってみるとテキスト通信にはなるほど便利だが文字によって1バイトったり2バイトだったりで処理には不便だということになった。欧州の一部ならもっと少ない字数でも役に立つというので西欧州用の8ビット1バイトコード(191字(左半分に95字、右半分に96字))というコード表を作った。これがISO 8859のパート1である。この頃から2/0のスペースも图形文字の字数に入れたりするようになった。この左半分はISOのコード表なのにASCIIになっている。同様にして8859の別のパートも作られているがパート1が8ビットシングルバイトの標準となるであろう。これがすぐプログラム言語用になるかどうかは予測がむずかしいが13/7と15/7に×と÷が入っているのでこれは使えると思う(図3)。

漢字コード

JIS X 0208(昔の呼び方ではJIS C 6226)が制定されたときは入出力はどうすればよいのだろうかと思ったが、その後ドットプリンタやビットマップディスプレイの普及のおかげで漢字処理は急に身近なものになってきた。

junetの通信も漢字かのものが圧倒的に多くなってきた。当然プログラム言語の中でも漢字を使いたいという希望ができて、もちろん当面は文字列定数と注釈が主だがいろいろ工夫が行なわれてきた。

その第一はシフトJISであった。これは8ビットの環境で1. JIS X 0201のローマ字またはASCII 2. JIS X 0201の片仮名、及び3. JIS X 0208の漢字の3種類をシフトコードなしに使おうとするものである。8ビットのコード表をみると列0, 1にはC0という制御文字が入っている。列2から7にはローマ字が入っている。列8, 9はC1の制御文字だがJIS X 0201を見るところは未定義である。片仮名は列10, 11, 12, 13に収まり(10/0は左へ入れると空白のところだから、これも使わず未定義)列14, 15も未定義である。そこで第1バイトが未定義なら漢字コードとし、漢字コードなら第2バイトはC0制御文字のところを除いて使うというふうにすればJIS X 0208の文字は収まるというのが骨子である。

これはMS-DOSで採用されたので、しばらくは一世を風靡したがいろいろ問題があつて見直されている。特にunixで

は \ をエスケープにしているのに、ここに第 2 バイトが現れて変なことを起こす。また C1 を使うのに不便だし、大体 JIS の表といちいち変換（つまりシフト）しなければならない。というわけで JIS X 0201 の片仮名（半角カナということにあり）を犠牲にし、右半分にゆったりと漢字を入れる EUC の考えが登場した。半角カナは濁点、半濁点の字は電報のように 2 字分使うので全角カナの方がずっと自然であるが、半角カナを使いたいという希望も後をたたず。

EUC ではこれを G2 に指示し SS2 で使うことにしている。漢字にはまだ、78年版と83年版の問題、補助漢字符号の問題などあるが、筆者が情報規格調査会の日本語機能委員会で提案しているのは左半分は ISO 8859 のように ASCII にする；右は83年版のJIS漢字にする（それも通常は78年版にある文字だけを使う）というものである。ASCII になると ¥ がなくなり \ が復活する。¥ は漢字の方を使うのである。

漢字コードの 3 区には全角の英字や数字がある。これでプログラムの本体を書いてもよいかとか、全角の foo と半角（ASCII）の foo は同じ識別名なのかとか話題は尽きないが、これは ASCII でも FOO と foo をどう見るかと同様な問題と思われる。

IBM PC 256 文字セット

筆者はその辺の事情に疎いが ECMA では最近 128 文字セットの検討をはじめた。それは IBM PC と System/2 がマルチ angular な 256 文字セットを使っているからという理由だそうだ（図 4）。これまでの規格の图形文字集合は 96 文字が単位だったから、ここでハタと困っているわけだが、一説には 32 文字コードを用意し、これで例えば列 8, 9 に呼び出すという案もあるらしい。まだこれはどういう方向に発展するかわからない。

128 文字セットといえば、Metafont の本にある Stanford や TEX のコードを思い出す。

プログラム言語で使う文字セットは落ち着くにはまだまだ時間がかかるであろう。

参考文献

- [1] 米田信夫、野下浩平：Algol 60 講義 共立出版 1979
- [2] 近藤頌子、原田賢一、土居範久：カナ文字 Fortran 情報処理 Vol.5, No.4 (July 1964) pp.222-223
- [3] 西村恕彦：日本語と Cobol 情報処理 Vol.8, No.3 (May 1967) pp.157-160
- [4] 水谷静夫：小朱唇言語仕様 東京女子大学日本文学科 1986

	12	11	0
1	+	-	0
2	A	J	/
3	B	K	S
4	C	L	T
5	D	M	U
6	E	N	V
7	F	O	W
8	G	P	X
9	H	Q	Y
8-3	=	.	\$
8-4	-)	*

図 1

	12	11	0
0	0	0	0
1	0	0	1
2	0	1	0
3	0	1	1
4	1	0	0
5	1	0	1
6	1	1	0
7	1	1	1
8	2	0	0
9	2	0	1
10	2	1	0
11	2	1	1
12	3	0	0
13	3	0	1
14	3	1	0
15	3	1	1

図 2

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	sp	0	ä	P	~	p	~	WSP	°	À	Đ	à	đ			
1	!	1	A	Q	a	q	!	i	±	Á	ñ	á	ñ			
2	"	2	B	R	b	r	~	ç	z	À	ò	à	ò			
3	#	3	C	S	c	s	£	3	À	ó	á	ó				
4	\$	4	D	T	d	t	~	¤	ó	À	ô	ô				
5	%	5	E	U	e	u	¥	¤	¤	Å	ø	ø				
6	&	6	F	V	f	v	!	!	!	€	ø	æ	ø			
7	'	7	G	W	g	w	S	-	ç	x	ç	÷				
8	(8	H	X	h	x	"	É	ø	é	ø					
9)	9	I	Y	i	y	©	!É	ü	é	ü					
10	*	:	J	Z	j	z	!	!É	ú	é	ú					
11	+	;	K	L	k	l	«	»	E	ö	ö					
12	,	<	L	\	l	l	-	¼	í	ü	ü					
13	-	=	M	J	m	j	½	í	ý	í	ý					
14	.	>	N	~	n	~	®	¾	í	p	í	p				
15	/	?	O	-	o	-	-	ë	í	ß	í	ÿ				

図 3

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
1	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
2	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
3	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
4	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
5	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
6	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
7	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
8	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
9	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
10	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
11	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
12	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
13	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
14	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ
15	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ	ñ

図 4