

## ISO10646と国際化された文字コードについて

太田 昌孝

東京工業大学総合情報処理センター

全世界の文字を含むコード体系としてISOで現在策定中の規格である10646は、各国漢字の混在、方向性の異なる言語の混在などへの対応はおざなりである。そこで、国際化に必要な各種付加情報を含めたエンコーディングであるICODEと、それを既存のバイトストリーム中で利用するための可変長エンコード技法であるIUTFを設計した。

### 1. はじめに

世界各国で使われている総ての文字に共通に使える文字コードを作ろうという動きの結果、ISOは、ISO10646 - Universal Coded Character Set (UCS) の制定作業を行っている。そのドラフト第一版のDIS10646 [1] は、もっぱらヨーロッパ系の国々からの、「同様の試みであるUnicode [2] と統一されるべきだ」という理由により、採用が否決された。その統一方法が検討された結果生まれたドラフト第二版であるDIS10646-1. 2は、第一版の面影をほとんどとどめず、ヨーロッパ言語に限った利用には有効であるが、国際化文字コードとしてはごく単純な用途でも使い物にならないものとなった。しかし、日本の反対にもかかわらず採用が票決され、現在IS化にむけての事務手続き中である。

DIS10646-1. 2の大きな欠点の一つは、Han Unificationである。Hanとは漢字のことであるが、漢字

を利用している主な国である中国(含台湾)、日本、韓国で使っている、共通の祖先をもつ漢字を同じものとみなして同じ文字コードを割り当ててしまおうというものである。これは、16ビットの文字コードとして米国計算機メーカーを中心に開発された文字コードであるUnicodeで採用された考えである。数万はある漢字を含めた文字を16ビットで表すためには当然必要となる処置ではあるが、このような文字コードは、当然のことながら、アジア圏では使い物にならない。日本語と中国語では、対応する漢字とはいえかなり形が違うものも多いにもかかわらず、それらには同じ文字コードが割り振られている。例えば、「草」という字の草冠は、中国漢字では、「十」が二つ横に並んだ形になるが、同一視されている。また、たとえ日中で同じ形と見なせる文字も、印刷の場合は異なる字形で表現するのが日中の言語が入り交じった場合の通常の表記法である。漢字仮名混じり文を表記する日本漢字と、漢文のための中国漢字では、当然のごとく字体が違う。が、文字コードを統一されてしまえば、この区別もできない。

また、もっぱらインドでさまざまな言語の表記に利用されているデバナガリ文字は、複数の母音と子音が組みあわさって一つの文字を形作る。ところが、この組みあわせの規則は、言語により異なる。それにもかかわらず、DIS10646-1.2では個々の母音と子音にのみコードが振られているため、言語毎に組みあわせ方法を変更することは不可能である。

こういった欠陥はヨーロッパ系の言語を使う限りはまったく問題とならないが、アジア圏では、なんらかの言語を識別する情報を付加しないことには、まともな表示もできないわけである。

DIS10646-1.2自体は16ビットまたは31ビットの固定長のコードである。32ビット目は常に0であり、ユーザー用のフラグビットとして利用できる。現在のところ下位16ビット以外のビットは総て0であり、16ビットコードであるUnicodeと全く同じものとなっている。各国の文字の割り当て状況を図1に示す。図は、下位8ビットが左から右に横軸、上位8ビットが上から下に縦軸となっており、左端の16進数が上位8ビットである。図でわかるとおり、ASCII文字(ISO-646-IRV)は、上位9ビットが0の位置に割り当てられており、7ビットのASCIIコードとまったく同じコードとなっている。図の、4E~9Fの部分が、問題のHan Unificationが行われている部分であるが、領域としては2万文字程度しか収納できず、5万を軽く越える漢字総てが表現できないのも深刻な問題である。もはや、Unificationの理由づけであった、「総ての文字を16ビットでエンコードする」という目的は、破綻しているわけである。

さて、DIS10646-1.2を既存のシステムの上で利用するため、ASCIIコードと共存可能な可変長コードUTFが提案されている。DIS10646-1.2にあ

るUTFを以後OUTF(=Original UTF)と呼ぶ。

OUTFは、利用に際して各種の不便があったため、X/Openが現在提案しているUTF(以後XUTFと呼ぶ)では、後に述べる各種の「改良」が施されている。付録Aのように、XUTFでは、ASCII文字は1バイトでそのまま表現可能である。いっぽう、文字コード2047までの文字は、2バイトで表現され、それ以上の文字は3バイトとなる。図1をみてのとおり、2バイトで表現可能なのは、ヨーロッパ、中近東の言語であり、アジア系の言語は冷遇されている。もちろん、ひらがな、カタカナも3バイトとなる。

というわけで、世界各国で使える共通の文字コードを作ろうという動きは、日本人の目からみて、うまくいったとはいいがたい。とはいえ、捕鯨問題同様、欧米人に道理は通じなかった。

## 2. ICODE

今後は、XUTFに従った製品が、外国メーカーから現実のものとして発売されることが予想される。

そこで、対応策の一つとして、国際化を考慮しUnicodeに言語情報等を付加し、Unicode上位互換なコードを設計した。このコードを、ICODE(=Internationalized CODE)と名付ける。

同時に、Unicodeと上位互換なエンコーディングとして、IUTF(=Internationalized UTF)を設計した。

上位互換ということで、ICODEは、下位16ビットは、まったくUnicodeと同じである必要がある。そこで、ICODEの下位16ビットは、Unicodeをそのまま、すなわちDIS10646-1.2の

Row-octet

00	ISO-646 IRV		Latin-1 Supplement	
01	Extended Latin-A		Extended Latin-B	
02	Extended Latin-B	IPA Extensions	Spacing Modifier Letters	
03	Combining Diacritical Marks		Greek	
04	Cyrillic			
05	Armenian		Hebrew	
06	Arabic			
09	Devanagari		Bengali	
0A	Gurmukhi		Gujarati	
0B	Oriya		Tamil	
0C	Telugu		Kannada	
0D	Malayalam			
0E	Thai		Lao	
10	Tibetan		Georgian	
1E	Additional Extended Latin			
1F	Greek Extensions			
20	General Punctuation	Super-/Subscripts	Currency Symbols	Comb. Diacritical Marks for Symbols
21	Letterlike Symbols	Number Forms	Arrows	
22	Mathematical Operators			
23	Miscellaneous Technical			
24	Control Pictures	O.C.R.	Enclosed Alphanumerics	
25	Box Drawing	Block Elements	Geometric Shapes	
26	Miscellaneous Dingbats			
27	Dingbats			
30	CJK Symbols And Punctuation		Hiragana	Katakana
31	Bopomofo	Hangul Jamo	CJK Miscellaneous	Combining Hangul Jamo
32	Enclosed CJK Letters and Months			
33	CJK Compatibility Words and Hours		CJK Compatibility Abbreviations and Days	
34	Hangul			
3D				
3E	Supplementary Hangul			
45				
46	Old Hangul			
4D				
4E	CJK Unified Ideographs			
9F				
A0				
DF				
E0	Private Use Area			
F7				
F9	CJK Compatibility Ideographs			
FA				
FB	Alphabetic Presentation Forms			
FC	Arabic Presentation Forms-A			
FD				
FE	CJK Compatibility Forms	Small Form Variants	Arabic Presentation Forms-B	
FF	Halfwidth And Fullwidth Forms			Specials

図1 DIS 10646-1.2の文字の割り当て(DISより引用)

下位16ビットをそのまま使うこととする。また、Unicodeでは総ての漢字は表現できないため、文字種の拡張のためのビットを1ビットその上に付加し、17ビットのコードとする。どのみちDIS10646-1.2は、17ビット程度への拡張は必要であろう。が、これがICODEの17ビット目と同じになるかどうかは、ここでは規定しない。ISOの今後の動き次第である。

1節で述べたように、言語識別情報等は処理コード中に必要である。言語識別情報としては、とりあえず漢字には3ビット必要となる。漢字の場合、

- 0 ANY
- 1 中国漢字
- 2 台湾漢字
- 3 日本漢字
- 4 韓国漢字
- 5 ベトナム漢字

のようにエンコードする。これを17ビットコードに付加して、全体で20ビットのコードを得る。

インド等には数百の言語があるといわれるが、文字種字体は漢字と比べて少ないので、その識別には下位17ビットの部分を利用できる。

さて、アラビア語等の言語は、他の多くの言語と異なり、横書きであるがもっぱら右から左に表記する。このような性質を方向性というが、方向性の異なる言語を混在させて表記するためには、方向性識別情報を導入する必要がある。一般には方向性は入れ子構造をもつ。つまり、ある方向性の言語内に、異なる方向性を持つ言語が埋め込まれ、さらにその中に方向性が異なるものが埋め込まれることがある。このような入れ子構造は文章の意味的構造を反映したものであるため、文字コードレベルで単純に対処できる問題ではない。そこで、方向性が異なる文字を混在させるた

めの最低限の処置として、本来の方向性とは逆転していることを示す方向性ビットを1ビット用意し：

- 0 本来の方向性
- 1 本来の方向性の逆

とする。そして、文字は表示の順に並べ、右から左に書きたい行内の文字は総て右から左の方向性、左から右に書きたい行内の文字は総て左から右の方向性となるように、方向性ビットで調節することとする。これだけの情報があれば、印刷などの最低限の用には足りる。

この結果全体で21ビットのコードが得られる(図2)。これは、IUTFで4バイトで表現できる。

場合によっては、さらにこの上に、DIS10646-1.2の上位14(15)ビットや、言語識別情報の上位ビットを付加し、ユーザ用フラグビット等の領域も考えて処理コードとすることになるが、これは今後の検討課題とし、当面は上位ビットは0とする。ただし、32ビット目にはDIS10646-1.2の32ビット目である、ユーザ用のフラグビットをそのまま置く。

結果のコードは、UCS4が31ビット全部を利用すると全体として32ビットを越える可能性もあるが、UCS2にすら空きがある現状からみて、実際には越えることはないであろう。万一越えても、既に64ビットプロセッサが実用化している現状では、深刻な問題とはならない。

### 3. IUTF

IUTFでは、2バイトで表現可能なコード数を極力増やし、ここに、言語識別情報を含めてエンコードされたよく使われる漢字等を置く。エンコード方針は4節で述べる。

XUTFの設計目標は、

- 1) UNIXファイルシステムとの整合性
- 2) 既存プログラムとの整合性
- 3) 処理コードとの間の変換の容易さ
- 4) 第一バイトでコード長が判明すること
- 5) バイト数が長くなりすぎないこと
- 6) ランダムアクセス後、文字区切りを効率良く発見できること

である。1)は、UNIXファイルシステムとの整合性を考慮し、ASCII文字スラッシュのコードは、スラッシュに対してのみしか現れないようにするということである。

IUTFの設計にあたっては、エンコード結果の同一性だけでなくこれらの目標にも配慮し、XUTF向けに書かれたアプリケーションプログラムが、ライブラリの入れ替えだけでそのままIUTFでも動作可能となるようにした。

IUTFでは、1)、5)、6)の性質は、そのまま保存されている。

2)の意味するところは、アスキー以外の文字のエンコーディングには、MSBの立ったものを使うということである。しかし、UTFは処理コードではないので、この要求はXUTF向けに書かれたプログラムに対しては無意味である。IUTFはこの性質はもたない。

3)の性質は、変換のためのライブラリプログラムの複雑さの点では犠牲になったが、速度的に遜色のないようにできる。XUTFにも含まれた表現の場合、処理速度はほぼ同等であるし、日本語等では、入出力量が少ないぶんかえて速くなるであろう。

4)の性質は、バイトの先読みを必要としないためには重要であるが、

4. 1) 第一バイトによりコード長が1バイトであるか2バイト以上であるか判明すること

4. 2) 第二バイトが存在する場合、そのバイトにより、コード長が判明すること

のように弱めても、やはり先読みは必要とらないので、そのようにした(3バイト目まで見る必要がある)。

IUTFでは、もっぱら2)の性質を犠牲にすることにより、XUTFに比べてより多くのエンコード領域が利用可能となった(詳細は付録B)。この領域は、IZONE (Internationalization Zone) と呼び、短縮したエンコードと、国際化されたアプリケーションに必要な、言語識別情報等を含めたエンコードを行う。

XUTFに現在表れる表現は、総てIUTFでも同じ意味を持つ。すなわち、0xff ff以下のICODEは、IUTFで読んでもXUTFと同じコードとして読まれる。

#### 4. IZONE

IZONEには、8256個の2バイト表現(以後、IZONE2と呼ぶ)と、35K個ほどの3バイト表現(以後、IZONE3と呼ぶ)が存在し、効率的なエンコーディングに利用できる。

このIZONEへの実際の文字の割り当ては今後の検討課題であるが、基本的な方針は以下のようなものであるべきである。

IZONE2は、極めて限られた大きさしかないので、国際協調の元に使う必要がある。日本語用には、2965文字もあるJIS X0208の第一水準漢字を全部入れることはあきらめ、かたかな、ひらがな、教育漢字、ごく少数の記号のみを入れよう。

頻出する中国漢字や、頻出するハングル(ハングルは既にXUTFにおいて、3バイトにエンコードされている)も、それぞれ数百~2000個程度エンコードされるべきである。

頻出するタイ文字やデバナガリ文字も、ハングル同様にprecombined characterとして、数百~2000程度エンコードされるべきである。

IZONE3には、その他のprecombined character等をエンコードする。

4バイト以上の領域は、ICODEをそのままエンコードすることに利用する。

## 5. 最後に

DIS10646-1. 2は、これまで挙げた以外にも種々問題をかかえているが、それらの問題にもなんとか対処できると思われる。

### 謝辞

本稿を用意するにあたって、筆者が所属する「ISO 10646の利用技術に関する研究会」の皆様にご多様な有益な御助言をいただきました。また、(株)シグマシステムには、研究会の運営のお世話になっています。ここに感謝の意を表明いたします。

### 付録A . XUTFの実際

1バイトで表現される値を次のように分類する。

C0:0<sup>32</sup>~127  
A :33<sup>126</sup>  
Tx:128<sup>191</sup>  
T1:192<sup>223</sup>  
T2:224<sup>239</sup>  
T3:240<sup>247</sup>  
T4:248<sup>251</sup>  
T5:252<sup>253</sup>  
Ty:254<sup>255</sup>(未使用)

この時、XUTFの文字のエンコードとして、

バイト列	文字コード
C0	0 <sup>32</sup> ,127
A	33 <sup>126</sup>
T1 Tx	128 <sup>2047</sup>
T2 Tx Tx	2048 <sup>2<sup>16</sup>-1</sup>
T3 Tx Tx Tx	2 <sup>16</sup> <sup>2<sup>21</sup>-1</sup>
T4 Tx Tx Tx Tx	2 <sup>21</sup> <sup>2<sup>26</sup>-1</sup>
T5 Tx Tx Tx Tx Tx	2 <sup>26</sup> <sup>2<sup>31</sup>-1</sup>

という組合せのみを許せば、バイト列中の任意の地点から、バイト列をせいぜい6バイト逆にたどることにより、文字の境界を確定することができる。

Unicodeのエンコードには、

T2 Tx Tx

までの3バイトで十分である。

### 付録B . IUTFの実際

1バイトで表現される値を次のように分類する。

C0:0<sup>32</sup>,127  
A :33<sup>126</sup>  
A':33<sup>46</sup>,48<sup>126</sup>  
C1:128<sup>159</sup>  
Tx:128<sup>191</sup>(C1を含む)  
T1:192<sup>223</sup>  
T2:224<sup>239</sup>(=S2+S3+S4+S6+S7)  
S2:224<sup>229</sup>  
S3:230<sup>235</sup>  
S4:236<sup>237</sup>  
S6:238  
S7:239  
U1:240<sup>255</sup>

この時、IUTFの文字のエンコードとして、

バイト列	文字コード
C0	0~32,127
A	33~126
T1 Tx	128~4095
T2 Tx Tx	4096~65535

[3] The Unicode Standard, Version 1.0, Vol. 1, The Unicode Consortium, Addison-Wesley, Oct. 1991

まではXUTFのものをそのまま利用し、さらに以下の表現を許す。

バイト列	文字数
C1	32
T1 A'	2976
T2 A'	1488
U1 A'	1488
U1 Tx	1024
T1 T2	512
T1 U1	512
U1 T2	256
S2 Tx A'	35712
S3 Tx A' Tx	>2 <sup>21</sup>
S4 Tx A' Tx Tx	>2 <sup>25</sup>
S6 Tx A' Tx Tx Tx Tx	>2 <sup>36</sup>
S7 Tx A' Tx Tx Tx Tx Tx Tx	>2 <sup>42</sup>

その結果、ASCIIのC1制御文字がそのまま利用できる他に、新たに8256個の2バイト表現と35712個の3バイト表現がえられることになる。

この時、XUTFと同様、バイト列中の任意の地点からバイト列をせいぜい8バイト逆にとどることにより、文字の境界を確定することができる。

#### 参考文献：

- [1] DIS 10646 Universal Multiple-Octet Coded Character Set (UCS), ISO/IEC JTC1 SC2, 1990
- [2] DIS 10646-1.2 Universal Multiple-Octet Coded Character Set (UCS), ISO/IEC JTC1 SC2, Oct. 1991

ビット位置	>= 21	20	19~17	16	15~0
内容	0	方向性	言語識別情報	拡張	Unicode

図2 I CODEのビット割り当て