

日本語文字コードについて

文字コードについて

文字コードを説明する前にコンピューターの中でどのように文字を表現するのかを簡単に説明します。よく耳にするビットとか、バイトとか、ワードといった単位のことです。

ビットとは何か？

コンピューターはデジタルです。

白か黒かの2種類のうちどちらかで、グレーはありません。

白か黒か？ → 0か1かに置き換えてみます。

一つのスイッチと考えれば、OFFの時0で通電(ON)しているときに1となります。

このスイッチ1ヶが1ビットとなります。

では、3ヶのスイッチ(ビット)があれば、幾通りの組み合わせが出来るでしょうか？

0 0 0

0 0 1

0 1 0

0 1 1

1 0 0

1 0 1

1 1 0

1 1 1

8種類の組み合わせが出来ますね。4ヶのスイッチでは16通りの組み合わせができます。8ヶでは256通りになります。この組み合わせで文字を表現するには、英数+記号は128あれば十分で予備を計算して8ヶのスイッチ(ビット)で足りる事になります。

バイトとは何か？

1文字を表現できるビットの集合をバイトと呼んでいます。8ビットですね。

文字コードって何？

コンピューターは英語圏で生まれ育ってきたので、それで良かったのですが、日本語などには対応できませんでした。

JISコード

そこで、考え出されたのがJISコードです。JISコードは2バイトを1文字に割り当てました。

JISコードでは7ビット X 2バイトで表現できる利点がありましたが、どこから漢字コードでどこが英数コードなのか判別できません、そこで、漢字コードの入口と出口に特殊なコード(エスケープシーケンス)を入れました。

SHIFT JISコード

このJISコードの欠点を克服するために考え出されたのがShiftJISコードです。JISコードの欠点を克服しエスケープシーケンスを使わずに漢字コードを表現できるようになりました、しかしこのSJISコードは8ビット X 2バイトを使わなくてはならないです。

そして

1バイト目が0x00 ~ 0x80、あるいは0xA0 ~ 0xDFに入っていたらそれをsingle-byteの文字としてそのまま表示する。0x20 ~ 0x7FはASCIIコードに準じ、0xA0 ~ 0xDFにはいわゆる半角カナ文字が割りあてられている。

1バイト目が0x81 ~ 0x9F、あるいは0xE0 ~ 0xFFに入っていたらそれは漢

字の 1 バイト目とみなし、次の 1 バイトと合わせて漢字を表示する。

このような複雑な処理をして漢字として表現するのですが、何か一つ取りこぼしがあるとそれ以後は文字化けしてしまうという事もあります。

また、文字コードが第 2 バイトの漢字コードが重なる所もあり表示が「侮ヲ」などによく化したものを見かけます。いわゆる 5 C 問題とも言われています

今までのパソコン WIN や MAC (OSX 以前) ではこの ShiftJIS コードが広く使われていました。

しかし、世界的な規模で広がった UNIX では EUC という漢字コードが使われ、日本用に EUC-JP というコードが考え出されました。

EUC-JP コード

漢字エスケープシーケンスを使わずに制御文字、ASCII コード文字と混在できますが、SHIFT-JIS 方式と同じく 8 ビット幅でないと表現できない、文字化けの危険性があるなどの欠点があります。表示アルゴリズムは SHIFT-JIS と似ており次のようになっています。

1 バイト目が 0x00 ~ 0x7F だったら、それを single-byte 文字としてそのまま表示する。

1 バイト目が 0x8E だったら、そのあとに続く文字を半角カナとして表示する (このときの 2 バイト目は SHIFT-JIS における半角カナ文字を表すコードと同じものが使われる)。1 バイト目が 0xA1 ~ 0xFE だったら、それは漢字の 1 バイト目とみなし、次の 1 バイトと合わせて漢字を表示します。

インターネットの普及で文字コードもグローバル化してきつつあります。

多国言語の使用が求められるようになってきた。それに伴い、漢字でも外字や拡張文字、異字形などを使う場面、要求も強くなってきた。

UTF-8

今までは漢字 1 文字を 2 バイトで表現したが、UTF-8 では最長 4 バイトで表現するようになる。最初の 1 バイト目は英数の ASC コードをそのまま表現できるので互換性の問題は発生しません。しかし容量が 1.5~2 倍に増えました。

EUC-JP の割り当て

1 バイト

ASCII の全て (実装系により JIS X0201/Windows-31J の当該エリアの場合あり)

2 バイト

JIS X0208 の非漢字の一部

3 バイト

JIS X0201 の 8 ビット文字 (半角カタカナ)

JIS X0208 の漢字エリアの全て

JIS X0212 の漢字エリアの全て

JIS X0213 の第 3・4 水準漢字の一部

Windows-31J の拡張文字エリア全て

4 バイト

Unicode の BMP 面以外全て

JIS X0213 の第 3・4 水準漢字の一部

5~6 バイト

Unicode の範囲外 (どんな文字が登録されるかという計画も無い)

このように JIS X0213 の漢字コードも含んでいるので、外字や拡張漢字なども表現できるようになってきました。MACOSX になってからは、この UTF8 が標準で使われるようになりました。

DTP 作業の中での文字コード

実際の DTP 作業の中で文字コードについてどのくらい認識していれば良いのでしょうか？知らなくても、知っていても組版は可能です。文字コードはそれを使っているソフトが吸収してくれるからです。しかし文字化けの修正とか文字コードで漢字を拾う場合などはある程度の知識が必要でしょう。

最近の動きについて

日本語文字セットとフォントの新しい環境

WindowVista の登場で文字環境も変わろうとしています。

同じ Windows でも XP と Vista では同じ文字でも違うように表示印刷されるようです。

飴 → 飴 樽 → 樽
溢 → 溢 晦 → 晦
鯖 → 鯖 葛 → 葛

ちなみに Mac では葛のように入りますが異形字を使ってヒラギノで表示すれば葛のようにも印刷できます。この文章も Mac を使って作成しました。詳しくは下記サイトが参考になります

http://www.jagat.or.jp/story_memo_view.asp?StoryID=9744