U Positivo - UTFPR - PUCPr - 11/02/2019 - 12:35:19.7 Prof Dr P Kantek (pkantek@gmail.com) Unicode e UTF-8 VIVO480a, V: 1.20 1

# Como trocar informações?

A informática convive - desde seu início - com uma dificuldade e tanto: como padronizar a utilização de códigos binários. O objetivo é nobre: permitir a troca de objetos binários e a sua correta interpretação por parte de todos os parceiros em uma comunicação qualquer. A necessidade é nova, a informática só apareceu na face da terra na década de 50 do século passado. No início, cada fabricante usava seus códigos. A IBM, por exemplo, construiu toda a sua família de equipamentos usando o código EBCDIC. Isto evoluiu para o código ASCII (início de 1960, comite ANSI-X3.2) que acabou sendo adotado como padrão. Ele privilegiou o idioma inglês, e padronizou seus caracteres nas primeiras 128 posições do código ASCII. Aos demais idiomas do planeta ficou disponível a segunda parte do código ASCII (mais 128 posições). É aqui que o código BRASCII (NBR-9614:1986 e NBR-9611:1991 da ABNT) colocou os caracteres acentuados do português. Assim como o Brasil, muitos outros países fizeram o mesmo. Isto apresentou pelo menos 2 problemas:

- Um arquivo composto no Brasil, quando impresso numa impressora fabricada na França, aparecia como uma série de caracteres malucos.
- Idiomas com muitos caracteres (katakana, hiragana, hangul, tagalog,...) simplesmente não tinham como ser representados no código ASCII.

#### Unicode

Em meados da década de 80, a academia e a indústria começaram a discutir propostas de resolver os 2 problemas acima. Diversos padrões foram sugeridos e eles começaram a ser usados, na base do "cada um por si". Uma primeira proposta sugeria usar 16 bits (2 bytes) o que permitiria  $2^{16}=65536$  caracteres. Parecia muito, mas vale citar Peter Norton, que disse há mais de 30 anos: na informática, não importa quanto você tem. Não é o suficiente.

Em 1991, criou-se o Consórcio Unicode e em outubro de 1991 o primeiro volume do padrão foi publicado. Em 1996, introduziu-se um mecanismo de codificação e com ele, a barreira dos 2 bytes foi ultrapassada. Agora, o espaço de endereçamento do Unicode ultrapassou o milhão de caracteres, o que permitiu incluir idiomas históricos (hieróglifos egípcios) ou mesmo idiomas muito pouco usados (Linear B por exemplo). O padrão atual consagra 1.114.112 códigos de ponto (codepoints), correspondendo ao intervalo hexadecimal de 00.00.00a  $10.\mathrm{FF}.\mathrm{FF}.$  Note que nem todo codepoint indica um caracter, já que alguns deles são reservados a outras finalidades.

Há alguns codepoints privados, que não têm significação no Unicode. Eles ficam para uma negociação entre remetente e destinatário e podem ser usados para qualquer coisa (são eles: Área privada: U+E000..U+F8FF com 6,400 caracteres; Área suplementar A: U+F0000..U+FFFFD com 65,534 caracteres e Área suplementar B: U+100000..U+10FFFD com 65,534 caracteres).

O Consórcio Unicode trabalha junto com a ISO (International Organization for Standardization) já que o padrão ISO/IEC 10646 é o próprio código Unicode. Um conceito importante é o de script. Trata-se de um conjunto de letras, sinais e regras de escrita que agregam partes do Unicode e que são usados em conjunto de *writing systems*. Na versão 7 do Unicode, são 123 os scripts, e há uma lista de candidatos a serem incluídos.

# UTF=Unicode Transformation Formats

Diversos mecanismos foram sugeridos para implementar o código Unicode, com a preocupação sempre presente de compatibilidade reversa. Foram propostos os UTF-1 (nunca chegou a ser muito usado, pois tinha problemas de performance), UTF-16 (usava 1 ou 2 palavras de 16 bits) , UTF-32 (cada caracter usava exatos 32 bits), mas todas elas acabaram migrando para a UTF-8 que se firmou como padrão de fato, principalmente porque muito cedo foi adotado como padrão pelo consórcio IMC (Internet Mail Consortium) e pelo W3C. O Google reporta que desde 2008, a maioria dos conteúdos HTML trocados na web usa a codificação UTF-8. O esquema todo foi escrito durante um jantar em 2 de setembro de 1992 por Ken Thompson e Rob Pike. Suas grandes qualidades são:

- $\bullet~$  Compatibilidade reversa, já que caracteres da primeira parte do ASCII não sofrem modificação (1º bit=0). Desta maneira qualquer arquivo ASCII que não use a segunda parte já é um arquivo UTF-8.
- multi-bytes. Estes contém um primeiro byte heading e um ou mais bytes de continuação. O heading tem 2 ou mais bits 1 seguidos por um zero enquanto os bytes de continuação sempre começam por 10.
- Auto-sincronização. Bytes únicos, header e continuação não compartilham valores, indicando que o início do caracteres está a no máximo 3 bytes para trás.
- Indicação do comprimento: O número de uns na alta ordem do byte heading indica a quantidade de bytes do caractere, sem ser necessário examinar os seus bytes.

Veia-se um resumo do exposto

veja-se um resumo do exposio							
bits	início	fim	tamanho	byte1	byte2	byte3	byte4
7	U+0000	U+007F	1	0xxxxxxx			
11	U+0080	U+07FF	2	110xxxxx	10xxxxxx		
16	U+0800	U+FFFF	3	1110xxxx	10xxxxxx	10xxxxxx	
21	U+10000	U+1FFFFF	4	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

Veja-se agora alguns exemplos de caracteres unicode de 1, 2, 3 e 4 bytes

(	Character	Binary code point			
\$	U+0024	0100100			
¢	U+00A2	000 10100010			
€	U+20AC	00100000 10101100			
報	U+24B62	00010 01001011 01100010			

Binary UTF-8	Hexadecimal UTF-8
00100100	24
11000010 10100010	C2 A2
11100010 10000010 10101100	E2 82 AC
11110000 10100100 10101101 10100010	FØ A4 AD A2

Uma tabela de conversão binário hexadecimal para ajudar:

hex	bin	hex	bin	hex	bin	hex	bin
0	0000	4	0100	8	1000	C	1100
1	0001	5	0101	9	1001	D	1101
2	0010	6	0110	A	1010	E	1110
3	0011	7	0111	В	1011	F	1111

#### Exemplos

1. Seja o seguinte conjunto de 5 caracteres UNICODE codificados em UTF-8.

#### F2B1B88DF4BAA8A7E192A9EF92A5F4A584BF

Seus codepoints em hexadecimal

### OB1EOD 13AA27 14A9 F4A5 12513F

E em decimal

# 728589 1288743 5289 62629 1200447

2. Seja o seguinte conjunto de 5 caracteres UNICODE codificados em

## 39F6B2B6A503E68FAFF59A84B1

Seus codepoints em hexadecimal

## 39 1B2DA5 03 63EF 15A131

E em decimal

# 57 1781157 3 25583 1417521

# Para você fazer

A seguir um conjunto de 5 caracteres UNICODE codificados em UTF-8. Você deve examiná-los e descobrir quais os 5 codepoints primeiro escritos em hexadecimal e depois em decimal (é só uma simples conversão).

## F6849DABC697CC8BDFA3F3BEBBA5

1 em hexa	2 em hexa	3 em hexa	4 em hexa	5 em hexa
1 1	2 1	9 1	4 1	,
1 em dec	2 em dec	3 em dec	4 em dec	5 em dec



