

UPOTREBA ODGOVARAJUĆEG KARAKTER SETA

Zašto postoji toliko različitih karakter setova za svaki od svetskih jezika, kako ih računari koriste i razumeju?

Osnovni jezik računara (kaže se i maternji) je engleski. Zaista, najveći zamajac celokupne softverske industrije nalazi se i dalje u Sjedinjenim Američkim Državama. Kako se računari koriste širom sveta potrebno je prilagođavanje potrebama korisnika sa različitih jezičkih područja.

Istorijski, sve je počelo kada se ukazala potreba da se u memoriji računara predstave slova. Morala se usvojiti neka konvencija prema kome bi se slova označila brojevima. Iako je jedan bajt dovoljan za prikazivanje 256 karaktera u početku je to izgledalo previše.

Jedan od karakter setova je **ASCII** (American Standard Code for Information Interchang). Ovaj set je 7-bitni ($2^7 = 128$ karaktera) i u sebe uključuje interpunkcijske oznake, cifre i slova engleskog alfabeta. ASCII set dovoljan je za prikaz informacija na engleskom jeziku, međutim nedovoljan je za prikaz većine drugih jezika jer mnogi jezici sadrže takozvane specijalne karaktere, npr. apostrofe, dvotačke, akcente... Dakle, bilo je potrebno proširenje karakter seta kako bi opslužio sve jezike, u cilju obezbeđenja razmene informacija između kompjuterskih sistema. Srećom prostora za proširenje je bilo dovoljno, jer je 128 mesta bilo prazno. Tako je nastao standardni MS DOS set karaktera ili CP 437 (kodna strana 437). Sa pojavom MS DOS-a 5.0 *Microsoft* se ozbiljnije bavi sistemskom podrškom za nacionalne setove karaktera. Definisane su 8-bitne nove kodne stranice. Našoj latinici odgovara kodna strana 852, a ćirilici CP 855. Međutim, nije vođeno računa o činjenici da se na nekim prostorima koriste podjednako oba pisma.

Sa lansiranjem *Windows*-a, 3.1. 1992. godine, ključna novina je pojava *True Type* fontova u čiju je osnovu ugrađena višejezička podrška. Definisane su nove kodne strane. Izostavljeni su grafički simboli, ali su ubačena slova karakteristična za gotovo sve zapadnoevropske i skandinavske jezike. Međutim, pri osmišljavanju kodnih stranica *Microsoft* se nije do kraja pridržavao standarda koji je već ranije usvojila ISO. Pojavom *Windows*-a 95 *True Type* fontovi su prošireni da podržavaju kodne stranice 1250, 1251, 1252, 1253, 1254 i 1257, tj. sadrže 652 karaktera – takozvani WGL4 (*Windows Glyph List 4*). Od pojave *Windows*-a NT može se koristiti Unicode standard.

ISO-8859 karakter set dizajniran je sredinom 80-tih od European Computer Manufacture's Association (ECMA) i potvrđen od International Standards Organisation (ISO). Ovo je osmorbitni karakter set – dakle sadrži 256 karaktera ($2^8 = 256$). ISO 8859 karakter set je serija od 10 i više standardizovanih višejezičkih grafičkih karakter setova. Najvažniji su:

1. Latin1 (Zapadna Evropa)
2. Latin2 (Istočna Evropa)
3. Latin3 (Južna Evropa)
4. Latin4 (Severna Evropa)
5. Cyrillic
6. Arabic
7. Greek
8. Hebrew
9. Latin5
10. Latin6

„I’m really terrified to see how difficult it can be for a non-latin1 person to be able to print in his/her own mother tongue!” — Akim Demaille, 1998.

(Bio sam zaista zastrašen da vidim kako je teško „ne-latin1“ osobi da može da štampa na svom maternjem jeziku!)

ISO-8859-1 (Latin1)

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	ı	Œ	ƒ	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
B0	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
C0	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î
D0	Ð	Ñ	Ò	Ó	Ô	Õ	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
E0	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î
F0	ä	ñ	ö	ó	ô	õ	÷	ø	ù	ú	û	ü	ý	þ	ÿ

Karakter set „Latin1“, poznat i kao zapadni karakter set (oznake ISO 8859-1) zadovoljava široke potrebe kao što je i sam Internet standard. Latin1 pokriva većinu zapadnoevropskih jezika : francuski, španski, portugalski, italijanski, albanski, nemački, norveški, finski, škotski, irski, engleski...

ISO-8859-2 (Latin2)

A0	À	Á	Â	Ã	Ä	Å	Š	Ś	Š	Ť	Ž	-	Ž	Ž
B0	à	á	â	ã	ä	å	š	ś	š	ť	ž	-	ž	ž
C0	Ř	Ā	Ĉ	Ď	Ě	Ĝ	Ĥ	Ĉ	Ě	Ĝ	Ĥ	Ĥ	Ĥ	Ĥ
D0	Ð	Ñ	Ń	Ō	Ű	Ų	×	Ŕ	Ū	Ŵ	Ŷ	Ŷ	Ŷ	Ŷ
E0	ř	ā	ĉ	ď	ě	ĝ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ
F0	đ	ñ	ń	ő	ű	ų	÷	ŕ	ű	ŵ	ŷ	ŷ	ŷ	ŷ

Latin2 pokriva centralnu i istočnu Evropu, odn. češki, mađarski, poljski, rumunski, hrvatski, slovački i slovenački jezik.

ISO-8859-3 (Latin3)

A0	Ĥ	Ĥ	Ĥ	Ĥ	Ĥ	Ĥ	Ĥ	Ĥ	Ĥ	Ĥ	Ĥ	-	Ĥ	Ĥ
B0	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ
C0	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā
D0	Ñ	Ō	Ű	Ų	Ŵ	Ŷ	×	Ŕ	Ū	Ŵ	Ŷ	Ŷ	Ŷ	Ŷ
E0	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā
F0	ñ	ő	ű	ų	ŵ	ŷ	÷	ŕ	ű	ŵ	ŷ	ŷ	ŷ	ŷ

Latin3 je popularan kod autora esperanta, a koristio se za turski jezik pre uvođenja Latin5.

ISO-8859-4 (Latin4)

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	Ā	Ķ	Ķ	Ķ	Ķ	Ķ	Š	Š	Š	Ē	Ģ	Ŧ	-	Ž	-
B0	ā	ķ	ķ	ķ	ķ	ķ	š	š	š	ē	ģ	ŧ	ņ	ž	ņ
C0	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā
D0	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā	Ā
E0	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā
F0	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā	ā

Latin4 predstavlja slova za estonski jezik.

ISO-8859-5 (Cyrillic)

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	Ё	Ђ	Ѓ	Є	Ѕ	І	Ї	Ј	Љ	Њ	Ћ	Ќ	-	Ў	Ў
B0	ё	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	н	о	п
C0	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю
D0	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о
E0	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю
F0	ё	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	н	о	п

Ova ćirilčna slova koriste se za bugarski, beloruski, makedonski, ruski i srpski jezik.

ISO-8859-6 (Arabic)

A0				A4								AC	AD		
				ﺥ								ﺀ	-		
												ﺀ			ﺀ
	C1	C2	C3	C4	C5	C6	C7	C8	C9	CA	CB	CC	CD	CE	CF
D0	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ
E0	-	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ
F0	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ	ﺀ

Ovo je arapski alfabet. Nažalost ne sadrži četiri eksta karaktera za persijski, ni osam eksta za pakistanski jezik.

Pojavio se 1992. godine. Predstavlja reorganizaciju Latin4 karakter seta. Naziva se još i nordijski.

Po pravilu se ASCII karakteri pojavljuju kao podskup većine drugih karakter setova. Na primer, prvih 128 karaktera određenog seta predstavlja sve ASCII karaktere. Ovo je neophodno zato što HTML tagovi u web stranicama moraju ostati isti, bez obzira na to kojim jezikom se prikazuje stranica.

U slučaju kreiranja sajta koji će imati verziju i na srpskom i na engleskom jeziku, može se koristiti ISO 8859-1 karakter set, koji će zadovoljiti obe verzije.

Karakter setovi su potrebni prvenstveno radi konverzije kodne oznake u karaktere (character encoding). Browser-u mora biti naznačeno koja vrsta kodiranja se koristi, da bi svaku kodnu oznaku preveo u odgovarajući karakter. Upravo ovome i služi vrednost karakter seta unutar navedenog <META> taga u okviru HTML-a.

Windows ANSI je osmorbitni karakter set dat kao set karaktera iznad i ispod 128 karaktera. Donji deo tabele je identičan ASCII karakter setu.

kodna strana iznad 128	1250 Istočna Evropa	1251 ćirilica	1252 zapadno evropski	1253 grčki ANSI	1254 turski	itd.
Ispod 128	ASCII	ASCII	ASCII	ASCII	ASCII	itd.

Unicode karakter set

Unicode standard određuje Unicode konzorcijum (Unicode Consortium - www.unicode.org). Isti standard obuhvaćen je i ISO specifikacijom, i to standardom ISO/IEC 10646.

*Unicode provides a unique number for every character,
no matter what the platform,
no matter what the program,
no matter what the language.*

(Sa Unicode web sajta)

*(Unicode dodeljuje jedinstven broj za svaki karakter
bez obzira na platformu
bez obzira na program
bez obzira na jezik)*

Unicode predstavlja najvažniju ekstenziju ASCII karakter seta. Unicode karakter set podržava karaktere gotovo svih govornih jezika, kao i matematičke i druge simbole. Pored engleskog i ostalih zapadnoevropskih jezika, podržava i japanski i kineski, kao i našu ćirilicu. Podrška svih subjekata Unicode standardu predstavlja najvažniji korak u internacionalizaciji Interneta.

Sledeća tabela prikazuje U+0400 ćirilični blok Unicode tabele. Prati redosled ISO-8859-5.

	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
10	11	12	13	14	15	16	17	18	19	1A	1B	1C	1D	1E	1F
20	21	22	23	24	25	26	27	28	29	2A	2B	2C	2D	2E	2F
30	31	32	33	34	35	36	37	38	39	3A	3B	3C	3D	3E	3F
40	41	42	43	44	45	46	47	48	49	4A	4B	4C	4D	4E	4F
50	51	52	53	54	55	56	57	58	59	5A	5B	5C	5D	5E	5F
60	61	62	63	64	65	66	67	68	69	6A	6B	6C	6D	6E	6F
70	71	72	73	74	75	76	77	78	79	7A	7B	7C	7D	7E	7F
80	81	82													
90	91	92	93	94	95	96	97	98	99	9A	9B	9C	9D	9E	9F
A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
B0	B1	B2	B3	B4	B5	B6	B7	B8	B9	BA	BB	BC	BD	BE	BF
C0	C1	C2	C3	C4	C5	C6	C7	C8	C9	CA	CB	CC	CD	CE	CF
D0	D1	D2	D3	D4	D5	D6	D7	D8	D9	DA	DB	DC	DD	DE	DF
E0	E1	E2	E3	E4	E5	E6	E7	E8	E9	EA	EB	EC	ED	EE	EF
F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FE	FF

Unicode fontovi:

Lucida Sans Unicode, od Bigelow & Holmes, isporučivan je sa nekim was verzijama. Sadrži 1775 simbola.

Times New Roman, Arial, Courier i Impact. Ovi fontovi sadrže WGL4 karakter set, koji podržavaju najveći deo latin-alfabeta, grčki i ćirilicu. Svaki sadrži 653 simbola, posebno Times New Roman sadrži malo iznenađenje sa 654-tim znakom.

Bitstream Cyberbit je realizovan kao TrueType font za Windows 95 i NT. Sadrži preko 8500 karaktera. Podržava većinu evropskih jezika uključujući grčki, ruski i turski. Takođe podržava japanski, korejski i kineski jezik.

TrueType GX font sadrži dosta karaktera, ali ne obavezno i Unicode tabelu.

Unicode je 16-bitni karakter set, što znači da svaki karakter sadrži dva bajta (DBCS - double byte character set). Na ovaj način moguće je mapirati $2^{16} = 65536$ karaktera. Karakteri koji se koriste u XML dokumentima pripadaju Unicode karakter setu. Kako svaki Unicode karakter zauzima dvostruko više mesta od njegovog ekvivalenta Latin1 seta, i sam XML bi trebalo da je dva puta veći od normalnih tekstualnih fajlova. Ipak, u najvećem broju slučajeva nije potrebno svih 16 bita po karakteru, pa XML dokumente možemo formirati i sa 8-bitnim karakter setom. Pri tome, XML procesor mora da prepozna i *UTF-8* (osmobarbitni) set i *UTF-16* (šesnaestobarbitni) set.

Rezime

Kako je osnovni jezik računara engleski i kako se računari koriste širom sveta, ukazala se potreba prilagođavanja rada na računarima korisnicima sa različitim govornih područja. Morala se usvojiti neka konvencija prema kojoj bi se slova označila brojevima. Iako je jedan bajt dovoljan za prikazivanje 256 karaktera, u početku je to izgledalo previše.

Jedan od karakter setova je ASCII. Ovaj set je 7-bitni. ASCII set dovoljan je za prikaz informacija na engleskom jeziku, međutim nedovoljan je za prikaz većine drugih jezika, jer mnogi jezici sadrže takozvane specijalne karaktere. Dakle, potrebno je proširenje karakter seta kako bi opslužio sve jezike, obezbeđujući razmenu informacija između kompjuterskih sistema. Dalji razvoj dovodi nas do osmобitnih karakter setova ISO-8859, Windows ANSI i, najзад, do Unicode šesnaestobitnog karakter seta. Po pravilu se ASCII karakteri pojavljuju kao podskup većine drugih karakter setova.

Unicode predstavlja najважнiju екстензију ASCII karakter seta. Unicode karakter set подржава karaktere готово свих govornih jezika, kao i matematičke i druge simbole. Podrška svih subjekata Unicode standardu predstavlja najважнiji korak u internacionalizaciji Interneta.

Оливера Михаилович

Употреба подходящего набора знаков

Резюме

Англијски језик јавља се основним језиком при пољзованију копмјутерима в целом мире и поштоу возникла потреба сделать работу на компјутере для пољзователей, говоращих различные языки, более удобной. Должно было усвоить какую-то конвенцию, согласно которой могли бы буквы обозначить числами. Хотя один байт является достаточным для представления 256 знаков, в начале это выглядело слишком много.

Одним из наборов знаков является ASCII. Этот набор состоит из 7 битов. ASCII набор достаточен для показания информации на английском языке, но недостаточен для показания большинства других языков, так как многие языки содержат, так называемые, специальные знаки. Значит, надо расширить набор знаков, чтобы он мог служить всем языкам, обеспечивая обмен информацией между компьютерными системами. Дальнейшее развитие приводит нас к особым наборам знаков ISO-8859, Windows ANSI и наконец до шестнадцатититного набора знаков Unicode. Как правило ASCII знаки являются подгруппой большинства других наборов знаков.

Unicode представляет собой самое важное расширение ASCII набора знаков. Unicode набор знаков поддерживает знаки почти всех языков, как и математические и другие символы. Поддержка всех субъектов в Unicode стандарте представляет важнейший шаг в интернационализации Интернета.

Olivera Mihailović

Usage of Aproprate Character Set

Summary

Since the fact that English language has become the elementary language of/for the Computer Technology and as the usage of Personal Computers is worldwide now, it became apparent that there is a need for adaptation of the data processing on computers, especially for Users from different language zone. A Convention had to be adopted to replace letters with numbers. Although one byte is sufficient for the representation of 256 characters, at the beginning it seems too much.

One of the character set is ASCII (American Standard Code for Information Interchange). This set employs a 7-bit character code. ASCII set is adequate for displaying the information in English language, but not adequate enough to display the majority of the world's languages, because most of the other languages contain the so-called special characters. Therefore, it is necessary to expand the character of the set, so it can be suitable for all languages, providing the exchange of information between computer systems and interfaces. The further development lead us to ISO-8859, a character set of 8-bit, Windows ANSI and, at the end, to UNICODE, a character set of 16-bit. According to rules, the ASCII characters appear as subset of the major other character sets.

UNICODE represents the most important extension of the ASCII character set. UNICODE character set supports almost all spoken languages, as well as mathematics' and other symbols. The support of the UNICODE standards, for all subjects, represents the most important step in the internationalization of the Internet.