

ГАПОУ МО «Мурманский колледж экономики и информационных технологий»

Развитие технологии одирования символьной информации в современных условиях

Девиз: «Таблиц кодировок много не бывает»

Автор: Багмет Е.Н.

г. Мурманск,
2016

□ Цель исследования:

Указать основные причины, приведшие к расширению используемых повсеместно кодовых таблиц до стандарта Unicode

□ Объекты исследования:

Кодовые страницы ASCII, Unicode

UNICODE

Обмен информацией

В процессе обмена информации по схемам:

□ **человек – человек;**

□ **человек – компьютер;**

□ **компьютер – компьютер**

представление **информации** происходит
в различных формах

Суть набор символов

Набор символов рассматривается, как таблица, задающая кодировку конечного множества символов алфавита (обычно элементов текста: букв, цифр, знаков препинания)

Символы в компьютере обычно кодируются одним или несколькими байтами и объединяются в **кодировки**

Код ASCII

ASCII Code Chart

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

Δ	b	d	L	z	f	n	Λ	M	x	λ	Σ	{		}	~	DEL
e	`	a	p	c	q	ε	ı	ä	μ	ı	ı	K	ı	ı	ı	ı

Latin-1

	.0	.1	.2	.3	.4	.5	.6	.7	.8	.9	.A	.B	.C	.D	.E	.F
8.	PAD 80	HOP 81	BPH 82	NBH 83	IND 84	NEL 85	SSA 86	ESA 87	HTS 88	HTJ 89	VTS 8A	PLD 8B	PLU 8C	RI 8D	SS2 8E	SS3 8F
9.	DCS 90	PU1 91	PU2 92	STS 93	CCH 94	MW 95	SPA 96	EPA 97	SOS 98	SGCI 99	SCI 9A	CSI 9B	ST 9C	OSC 9D	PM 9E	APC 9F
A.		ì	¢	£	¤	¥	¦	§	¨	©	ª	«	¬	-	®	—
B.	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
C.	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
D.	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
E.	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
F.	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

КОИ-8 – русская кодовая таблица

80	81	82	83	84	85	86	87	88	89	8A	8B	8C	8D	8E	8F
											»	«			
90	91	92	93	94	95	96	97	98	99	9A	9B	9C	9D	9E	9F
			ё												
A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
№			Ё												
B0	B1	B2	B3	B4	B5	B6	B7	B8	B9	BA	BB	BC	BD	BE	BF
Ю	а	б	ц	д	е	ф	г	х	и	й	к	л	м	н	о
C0	C1	C2	C3	C4	C5	C6	C7	C8	C9	CA	CB	CC	CD	CE	CF
П	я	р	с	т	у	ж	в	ь	ы	з	ш	э	щ	ч	ть
D0	D1	D2	D3	D4	D5	D6	D7	D8	D9	DA	DB	DC	DD	DE	DF
Ю	А	Б	Ц	Д	Е	Ф	Г	Х	И	Й	К	Л	М	Н	О
E0	E1	E2	E3	E4	E5	E6	E7	E8	E9	EA	EB	EC	ED	EE	EF
П	Я	Р	С	Т	У	Ж	В	Ь	Ы	З	Ш	Э	Щ	Ч	ТЬ
F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FE	FF

E0	E1	E5	E3	E4	E2	E6	E1	E8	E9	E7	E8	EC	ED	EE	EE
П	Я	Ь	С	Д	У	Ж	В	Р	Ы	З	Ш	Э	Щ	Ч	Р
E0	E1	E5	E3	E4	E2	E6	E1	E8	E9	E7	E8	EC	ED	EE	EE
Ю	У	Р	Ц	Д	Е	Ф	Л	Х	Н	Н	К	И	И	Н	О



Небольшой размер ASCII (256 символов), а также аналогичных таблиц не позволял закодировать большое количество символов национальных алфавитов

В настоящее время распространение получил новый стандарт Unicode, который отводит на каждый символ не 1 байт, а 2 байта

Unicode

Юникод, как система для линейного представления текста

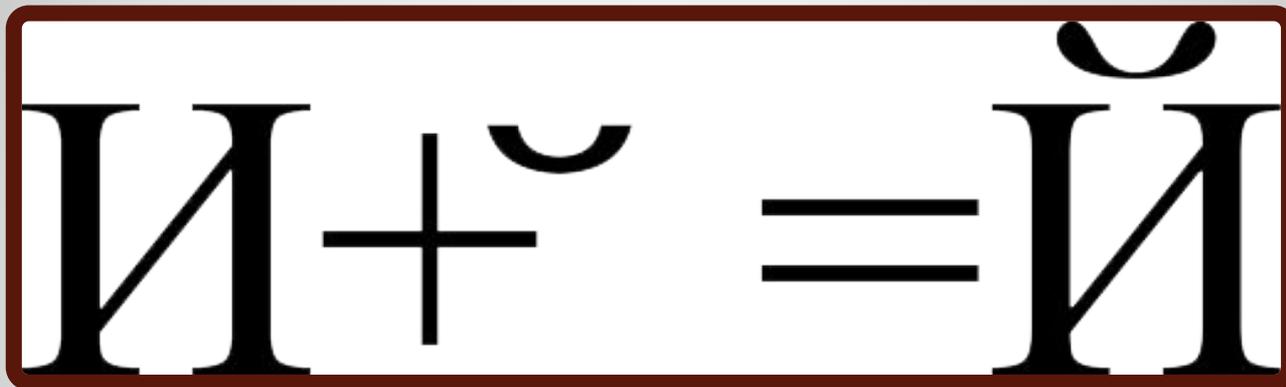
Символы, имеющие дополнительные над- или подстрочные элементы, могут быть представлены в виде построенной по определённым правилам последовательности кодов (составной вариант, *composite character*) или в виде единого символа (монолитный вариант, *precomposed character*).

На 2014г. считается, что все буквы крупных письменностей в Юникод внесены, и если символ доступен в составном варианте, дублировать его в монолитном виде не нужно.

Общее количество символов в Unicode равно 65536 ($2^{16}=1024*64$). Из них отдано:

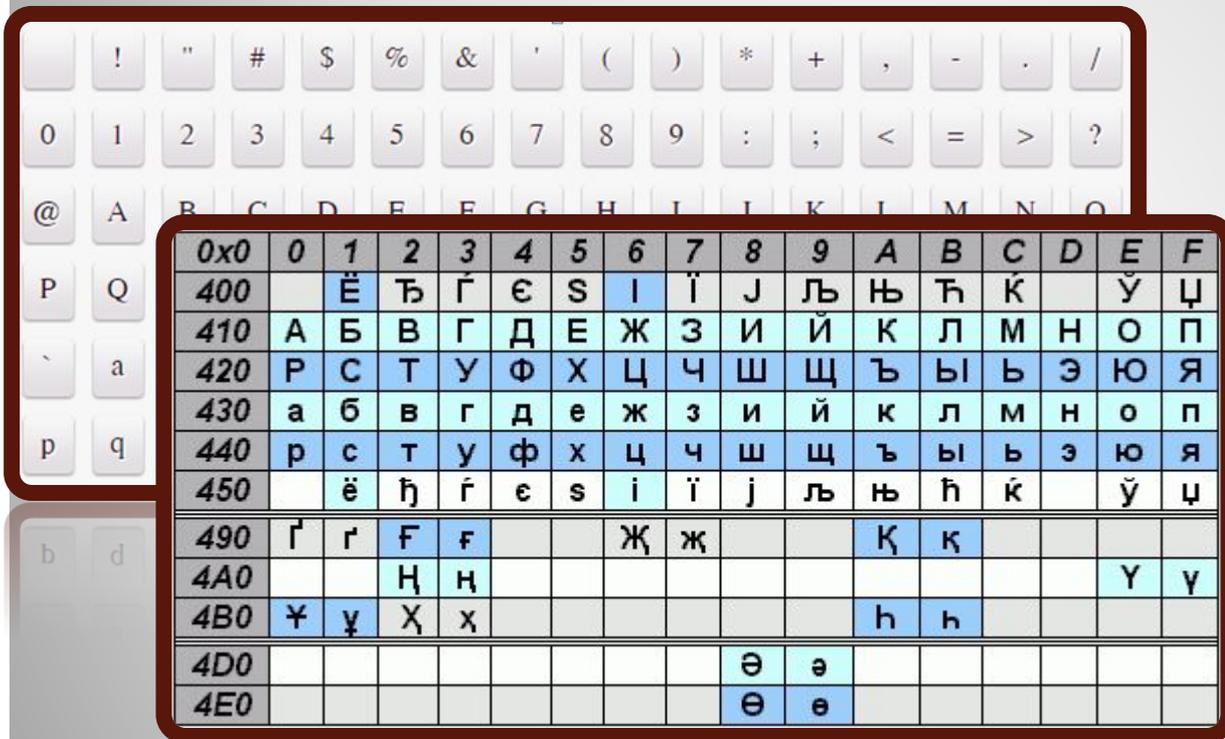
- 336 – на латынь;
- 256 – на русский язык;
- по 128 – на некоторые редкие национальные алфавиты

Отдельное кодирование диакритических
символов для экономии места



Представление символа «Й» в виде базового символа «И» и модифицирующего символа « ˇ »

Кодировки русского алфавита и Unicode (ISO 10646)



0x0	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
400		Ё	Ъ	Ґ	Є	Ѕ	І	Ї	Ј	Љ	Њ	Ћ	К		У	Ц
410	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
420	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
430	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
440	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
450		ё	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	к		у	ц
490	Ґ	ғ	Ғ	ғ			Җ	җ			Қ	қ				
4A0			Ң	ң										Ү	ү	
4B0	Ҙ	ұ	Ҫ	ч							Һ	һ				
4D0									Ә	ә						
4E0									Ө	ө						

Вспомогательная таблица для ручного преобразования кодировок символов



Китайские иероглифы – 20992 позиций

The screenshot shows the 'Table of Symbols' (Таблица символов) dialog box in Windows. The font is set to 'SimSun'. The main grid displays Chinese characters grouped by Pinyin initials: A, Ai, An, and Ang. A smaller 'Grouping' (Группировка) dialog box is open over the main grid, showing a grid of Pinyin initials (A-Z) with 'U' selected. The main dialog also includes a search field, a 'Find Unicode' (Найти Юникод) button, and a search result field showing 'U+0041 (0x41): Latin Capital Letter A'.

Group	Character 1	Character 2	Character 3	Character 4	Character 5	Character 6	Character 7	Character 8	Character 9	Character 10	Character 11	Character 12	Character 13	Character 14	Character 15	Character 16	Character 17	Character 18	Character 19	Character 20
A	啊	阿	呵	吡	嘎	腌	钢													
Ai	埃	挨	哎	唉	哀	皑	癌	蔼	矮	艾	碍	爱	隘	呆	阨	乃	奇	割	捱	呢
An	鞍	氨	安	俺	按	暗	岸	胺	案	厂	干	广	盒	钳	谗	掩	揞	犴	庵	桉
Ang	肮	昂	盎	仰																

Универсальная система кодирования (**Юникод**) представляет собой набор графических символов и способ их кодирования для компьютерной обработки текстовых данных



Графические символы - это символы, имеющие видимое изображение. Графическим символам противопоставляются управляющие символы и символы форматирования.

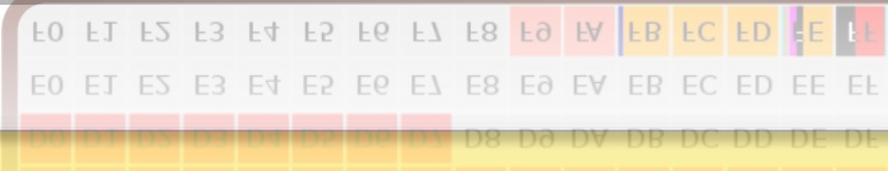
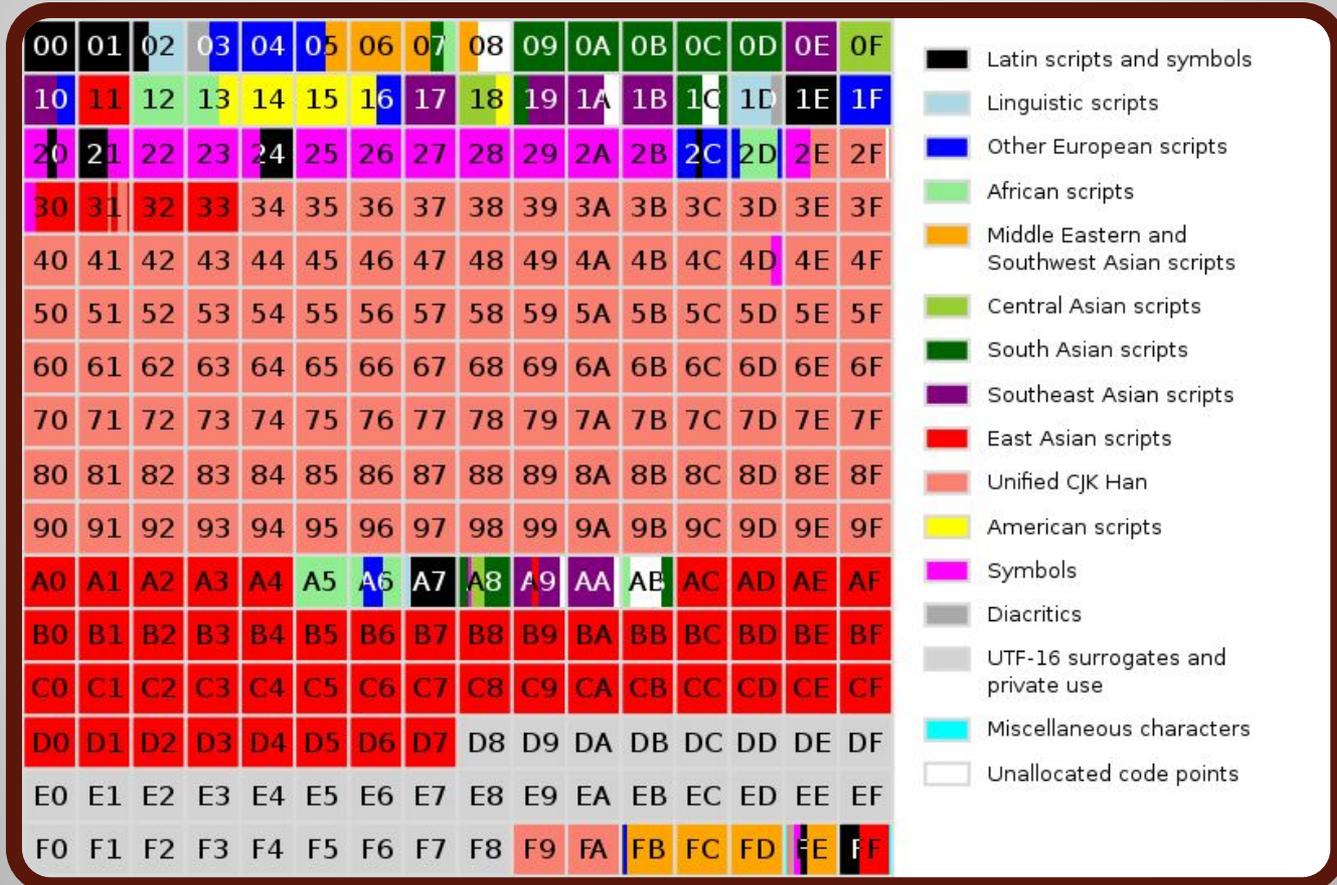
Графические символы

Графические символы включают в себя следующие группы:

- буквы, содержащиеся хотя бы в одном из обслуживаемых алфавитов;
- цифры;
- знаки пунктуации;
- специальные знаки (математические, технические, идеограммы и пр.);
- Разделители.



Распределение представленных в Unicode символов



Проблемы Unicode

Не решенными остаются такие вопросы:

- Не предусмотрено переключение горизонтального и вертикального написания для японского и китайского языков;
- Перевод из строчных букв в заглавные зависит от языка;
- До сих пор не всё прикладное программное обеспечение поддерживает корректную работу с ним.

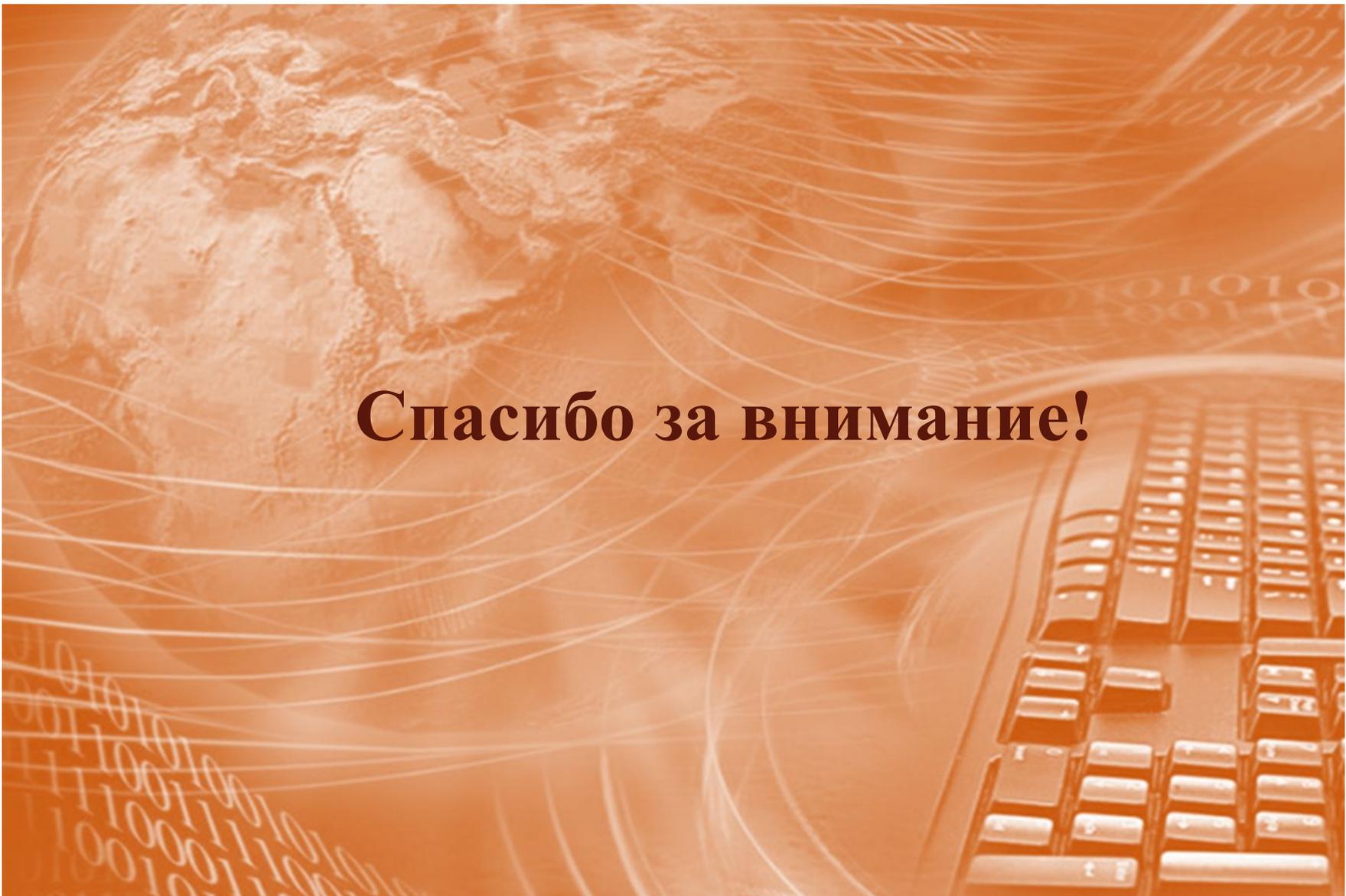
Вывод

Таким образом:

- Unicode стал настоящим спасением, заменив собой все разнообразие локальных кодировок;
- Несмотря на это, Unicode все же не смог разрешить все проблемы, связанные с интернационализацией кодирования.

Список литературы

1. American National Standards Institute <http://www.ansi.org>
2. UTF-8, a transformation format of ISO 10646
<http://tools.ietf.org/html/rfc2279>
3. Internet Assigned Numbers Authority – организация, отвечающая за административное управление в Интернете – <http://www.iana.org>



Спасибо за внимание!