



КОДИРОВАНИЕ ТЕКСТОВОЙ ИНФОРМАЦИИ

ПРЕДСТАВЛЕНИЕ ИНФОРМАЦИИ В КОМПЬЮТЕРЕ

10 класс



ИЗДАТЕЛЬСТВО

БИНОМ

Ключевые слова

- текстовая информация
- кодирование
- кодовые таблицы



Компьютерное представление текстовой информации

Для компьютерного представления текстовой информации достаточно:



...	...
64	01000000
65	01000001
66	01000010
67	01000011
68	01000100

Определить алфавит
(множество всех
символов)

Присвоить каждому
символу алфавита
порядковый номер

Перевести номер
символа в двоичную
систему счисления

Кодировка ASCII

American Standard Code for Information Interchange – американский стандартный код для обмена информацией, разработанный в 1960-х годах в США.

	0	0	0	0	0	0	0	0	0	5						
0	NUL	SOH	STX	ETX	EOT	ENC										
1	0	0	1	0	0	0	0	0	0							
2																
3	0															
4	@															
5	P															
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

Изображаемые символы
(буквы латинского алфавита, цифры, знаки препинания и арифметических операций, скобки и некоторые специальные символы)

Первые 32 символа и 128-й – управляющие
(при выводе текста они не отображаются графически)

A

0 1 0 0 0 0 0 0 1

0 0 0 1 1 1 1 1

0 1 1 1 1 1 1 0

Расширение кодировки ASCII

	0 0 0 0 0 0 0 0							5										
0	NUL	SOH	STX	ETX	EOT	ENC												
1	DLE	DC1	DC2	DC3	DC4	NAK												
2		!	"	#	\$	%												
3	0	1	2	3	4	5												
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	1 0 0 0 0 0 0 0							u	v	w	x	y	z					
								КОИ-8										
	Ъ	Ґ	,	ґ	„	…	†	‡	€	‰	Љ	0	1	1	1	1	1	1
9	Ѓ	Ѕ	Ї	ґ	ґ	•	√	≈	≤	≠	Љ	Ѓ	Ѕ	Ї	ґ	ґ	•	√
A	=	ґ	ґ	€	‰	Љ	Ѓ	Ѕ	Ї	ґ	ґ	•	√	≈	≤	≠	Љ	Ѓ
B	Ѓ	Ѕ	Ї	ґ	ґ	•	√	≈	≤	≠	Љ	Ѓ	Ѕ	Ї	ґ	ґ	•	√
C	Ѓ	Ѕ	Ї	ґ	ґ	•	√	≈	≤	≠	Љ	Ѓ	Ѕ	Ї	ґ	ґ	•	√
D	Ѓ	Ѕ	Ї	ґ	ґ	•	√	≈	≤	≠	Љ	Ѓ	Ѕ	Ї	ґ	ґ	•	√
E	Ѓ	Ѕ	Ї	ґ	ґ	•	√	≈	≤	≠	Љ	Ѓ	Ѕ	Ї	ґ	ґ	•	√
F	Ѓ	Ѕ	Ї	ґ	ґ	•	√	≈	≤	≠	Љ	Ѓ	Ѕ	Ї	ґ	ґ	•	√
								1 1 1 1 1 1 1 1										

Стандартная часть кода (0 ... 127)

Расширение ASCII (128 ... 255)
 (буквы национального алфавита,
 символы национальной валюты и т.п.)

Расширение кодировки ASCII

	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	Windows-1251					KOI-8				

	Ъ	Ѓ	,	ѓ	„	…	†	‡	€	%	Љ	‹	Њ	Ќ	Љ	Ў
9	ђ	‘	’	“	”	•	–	√	≈	≤	™	≥	љ	›	ј	њ
A	џ	џ	џ	џ	џ	џ	џ	џ	џ	џ	џ	џ	џ	џ	џ	џ
B	°	±	І	і	Є	є	µ	¶	·	ё	№	є	»	ј	ѕ	ѕ
C	А	Б	а	б	Г	г	Д	д	Е	е	Ж	ж	З	з	И	и
D	Р	р	С	с	Т	т	У	у	Ф	ф	Х	х	Ц	ц	Ч	ч
E	а	ю	б	а	в	б	г	ц	д	д	е	е	ж	ф	з	г
F	р	п	с	я	т	р	у	с	ф	т	х	у	ц	ж	ч	в

Стандарт Unicode



Unicode — это «уникальный код для любого символа, независимо от платформы, независимо от программы, независимо от языка» (www.unicode.org).

Стандарт Unicode был разработан в 1991 году и описывает алфавиты всех известных, в том числе и «мертвых», языков. Для языков, имеющих несколько алфавитов или вариантов написания (японского и индийского), закодированы все варианты. В кодировку Unicode внесены все математические и иные научные символы и обозначения и даже некоторые придуманные языки (язык эльфов из трилогии Дж. Р. Р. Толкина «Властелин колец»).



Клавиатуры некоторых стран мира



РУССКАЯ



АМЕРИКАНСКАЯ



АРАБСКАЯ



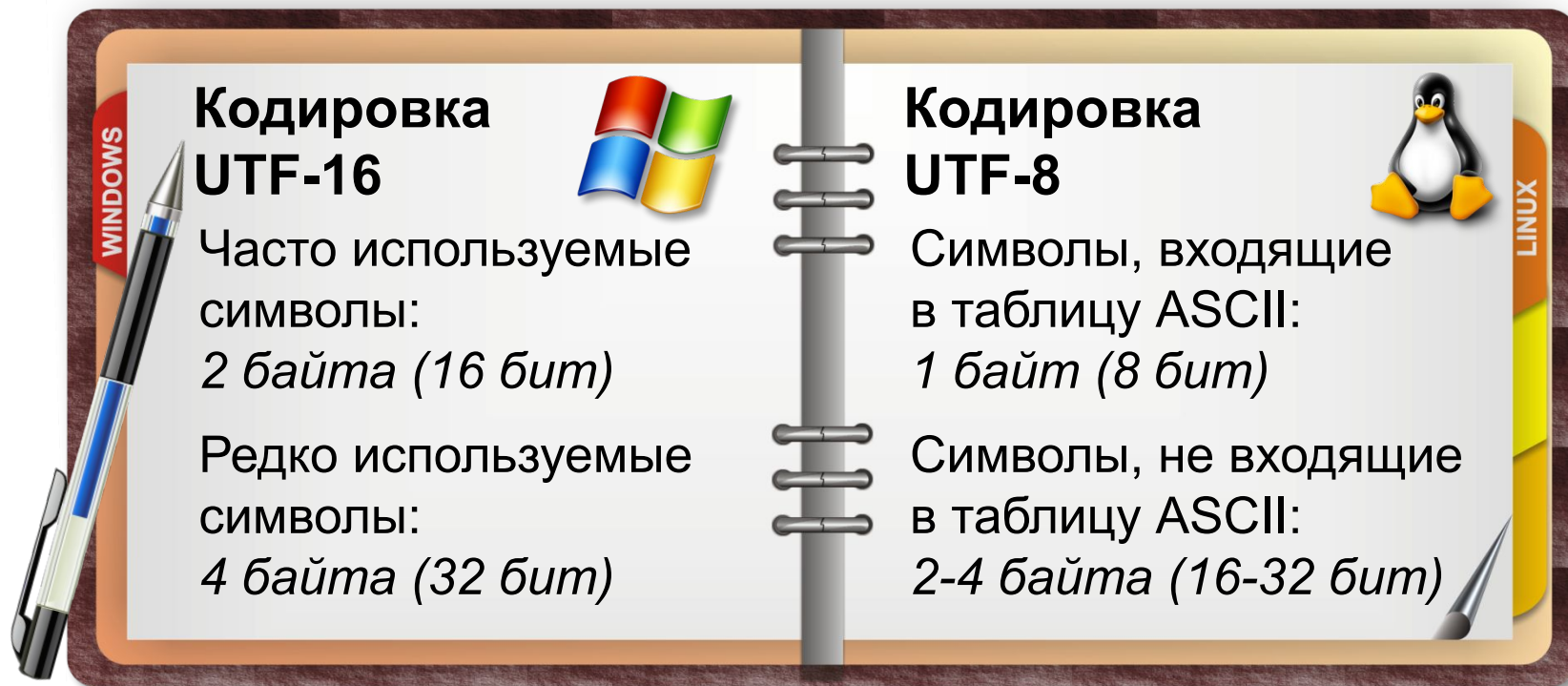
АРМЯНСКАЯ



ЯПОНСКАЯ

Кодировки стандарта Unicode

Для представления символов в памяти компьютера в стандарте Unicode имеется несколько кодировок.



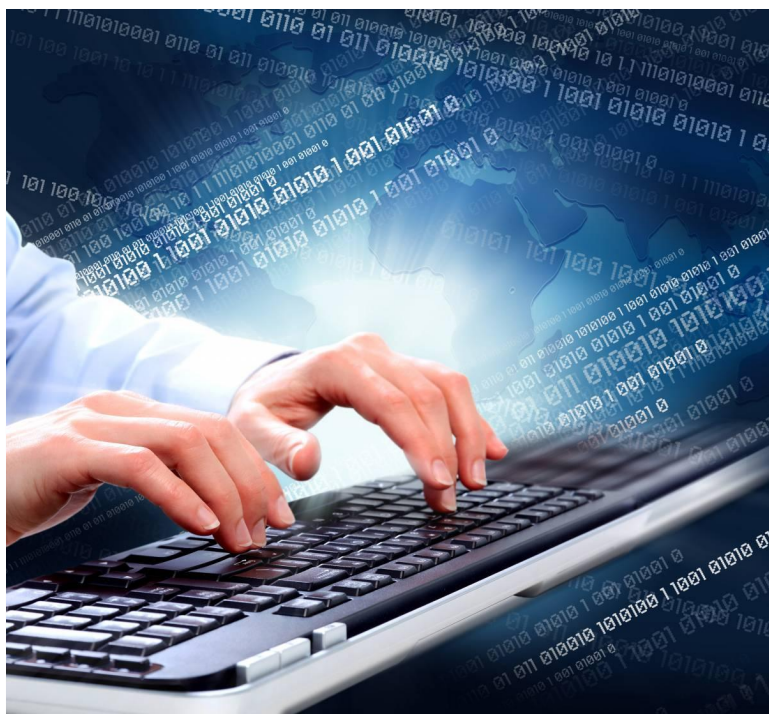
Кодировки Unicode позволяют включать в один документ символы самых разных языков, но их использование ведёт к увеличению размеров текстовых файлов.



Информационный объем сообщения



Информационным объёмом текстового сообщения называется количество бит (байт, килобайт, мегабайт и т. д.), необходимых для записи этого сообщения путём заранее оговоренного способа двоичного кодирования.



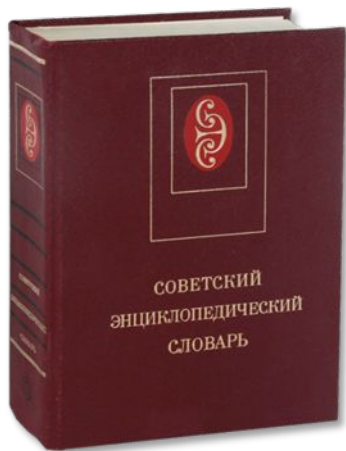
Количество символов в сообщении

$$I = K \cdot i$$

ASCII, KOI-8,
Windows-1251, ...
1 символ = 1 байт

Unicode
1 символ = 2 байта

Вопросы и задания



В Советском энциклопедическом словаре (1983 года издания) 1600 страниц. На одной странице размещается в среднем 100 строк по 140 символов (включая пробелы) в каждой. Найдите объем (в Мбайтах) текстовой информации в словаре, если при записи используется кодировка «*один символ — один байт*».

Дано:

$$i = 1 \text{ байт}$$

$$K = 1600 \cdot 100 \cdot 140$$

$I = ?$

$$I = K \cdot i$$

$$I = \frac{1600 \cdot 100 \cdot 140}{1024 \cdot 1024} \text{ Мб} \approx 21,36 \text{ Мб}$$

Ответ: 21,36 Мбайта

Самое главное

Текстовая информация по своей природе дискретна, так как представляется последовательностью отдельных символов.

В памяти компьютера хранятся специальные кодовые таблицы, в которых для каждого символа указан его двоичный код. Все кодовые таблицы, используемые в любых компьютерах и любых операционных системах, подчиняются международным стандартам кодирования символов.

Основой для компьютерных стандартов кодирования символов послужил код ASCII, рассчитанный на передачу только английского текста. Расширения ASCII-кодировки, в которых первые 128 символов кодовой таблицы совпадают с кодировкой ASCII, а остальные (с 128-го по 255-й) используются для кодирования букв национального алфавита, символов национальной валюты и т. п.



Самое главное

В 1991 году был разработан новый стандарт кодирования символов, получивший название Unicode (Юникод), позволяющий использовать в текстах любые символы любых языков мира. Кодировки Unicode позволяют включать в один документ символы самых разных языков, но их использование ведёт к увеличению размеров текстовых файлов.



Вопросы и задания



Задание 1. Представьте в кодировке ASCII текст
Happy New Year!

а) шестнадцатеричным кодом

48 61 70 70 79 20 4E 65 77 20 59 65 61 72 21

б) десятичным кодом

72 97 112 112 121 32 78 101 119 32 89 101 97 114 33

	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

ОТВЕТ

Подходы к расположению русских букв в различных кодировках



Задание 2. Сравните подходы к расположению русских букв в кодировках Windows-1251 и КОИ-8.

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	
4	@ю	А а	В б	С с	Д д	Е е	Ф ф	Г г	Н н	Х х	И и	Ј ј	К к	Л л	М м	Н н	О о
5	Р р	Qя	Rр	Sс	Tт	Uу	Vж	Wв	Xь	Yы	Zз	[ш	\э]щ	^ч	_ъ	
6	`Ю	а А	б Б	с С	д Д	е Е	ф Ф	г Г	х Х	и И	ј Ј	к К	л Л	м М	н Н	о О	
7	p П	q Я	r Р	s С	t Т	u У	v Ж	w В	x Ъ	y Ы	z З	{ Ш	Э	} Щ	~ Ч		Ъ
С	Аю	Ба	Вб	Гц	Дд	Ее	Жф	Зг	Их	Йи	Кй	Лк	Мл	Нм	Он	По	
D	Рп	Ся	Тр	Ус	Фт	Ху	Цж	Чв	Шь	Щы	Ъз	Ыш	Ьэ	Эщ	Юч	Яъ	
E	аЮ	бА	вБ	гЦ	дД	еЕ	жФ	зГ	иХ	йИ	кЙ	лК	мЛ	нМ	оН	пО	
F	рП	сЯ	тР	уС	фТ	хУ	цЖ	чВ	шь	щы	ъЗ	ыШ	ьЭ	эЩ	юЧ	яЪ	

ПОДСКАЗКА - 1

ПОДСКАЗКА - 2

Вопросы и задания



Задание 3. В 15-м издании энциклопедии Britannica 32 тома, в каждом из которых порядка 1000 страниц. На одной странице размещается в среднем 70 строк по 120 символов (включая пробелы) в каждой. Найдите объем текстовой информации в энциклопедии, если при записи используется кодировка Unicode («один символ — два байта»).

Дано:

$i = 2$ байта

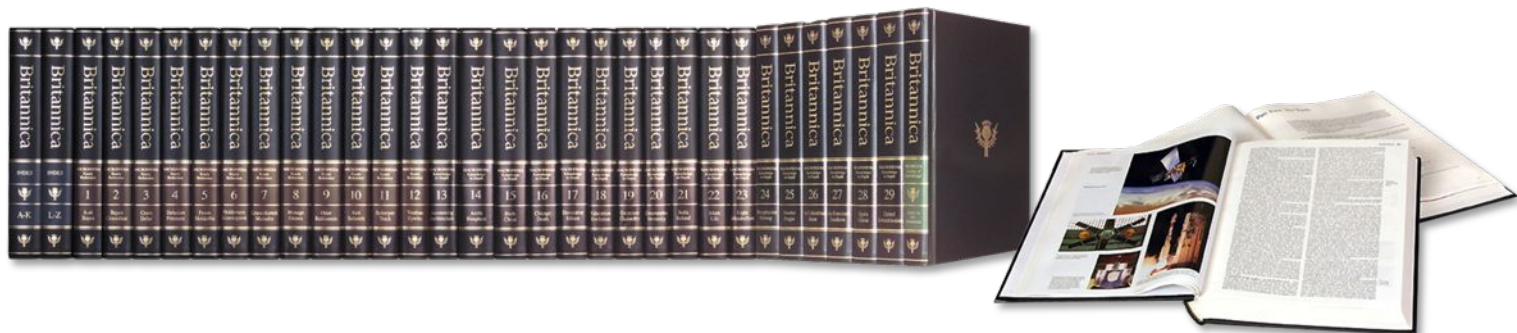
$K =$

$32 \cdot 1000 \cdot 70 \cdot 120$

$I = K \cdot i$

$$I = \frac{32 \cdot 1000 \cdot 70 \cdot 120 \cdot 2}{1024 \cdot 1024} \text{ Мб} \approx 513 \text{ Мб}$$

Ответ: 513 Мбайт



Информационные источники

- <http://dev.bowdenweb.com/a/i/cons/utilities/unicode/unicode-2000px.png>
- https://openclipart.org/image/2400px/svg_to_png/177279/Blank-Generic-Keyboard-Remix-by-Merlin2525.png
- http://arstyle.org/uploads/posts/2010-07/1278744192_1274782943_dreamstime_9113949-converted.jpg
- <http://www.businesstoday.net.my/wp-content/uploads/2015/04/Computer-Programmer-Coding-Camp-shutterstock.jpg>
- <http://static.ozone.ru/multimedia/1005976053.jpg>
- <http://gimnnik.narod.ru/open-office/TextProcessor/p5aa1.html>
- http://media.washtimes.com.s3.amazonaws.com/media/image/2012/03/14/encyclopaedia-britann_lea.jpg
- <http://www.novilist.hr/var/novilist/storage/images/sci-tech/tehnologija/encyclopaedia-britannica-prekida-tiskanje-postaje-digitalna/1306075-1-cro-HR/Encyclopaedia-Britannica-prekida-tiskanje-postaje-digitalna.jpg>