Кодирование и обработка текстовой информации

Автор: учитель информатики Гилева Елена Евгеньевна

Начиная с 60-х годов, компьютеры все больше стали использовать для обработки текстовой информации. В настоящее время большая часть ПК в мире занято обработкой именно текстовой информации.







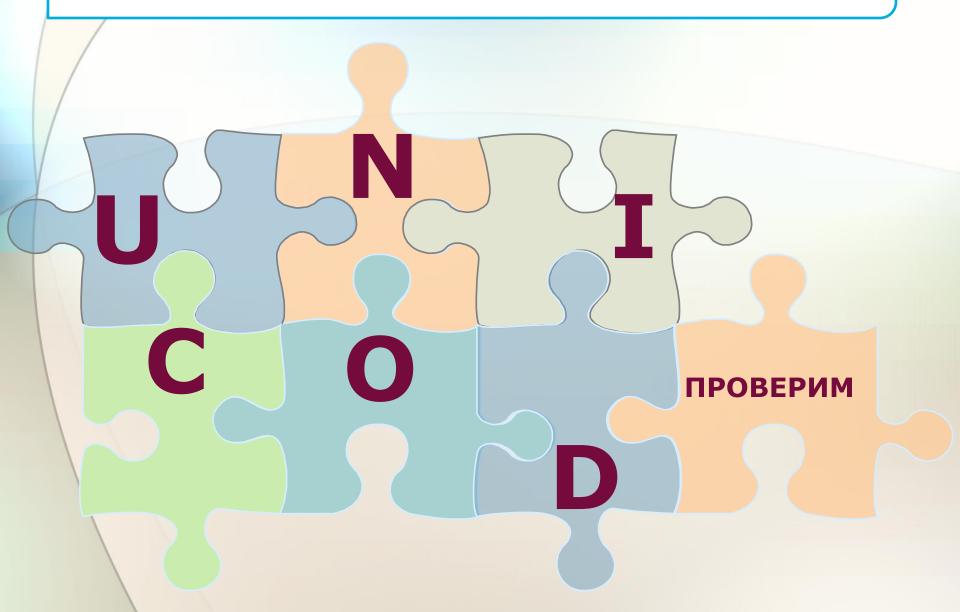


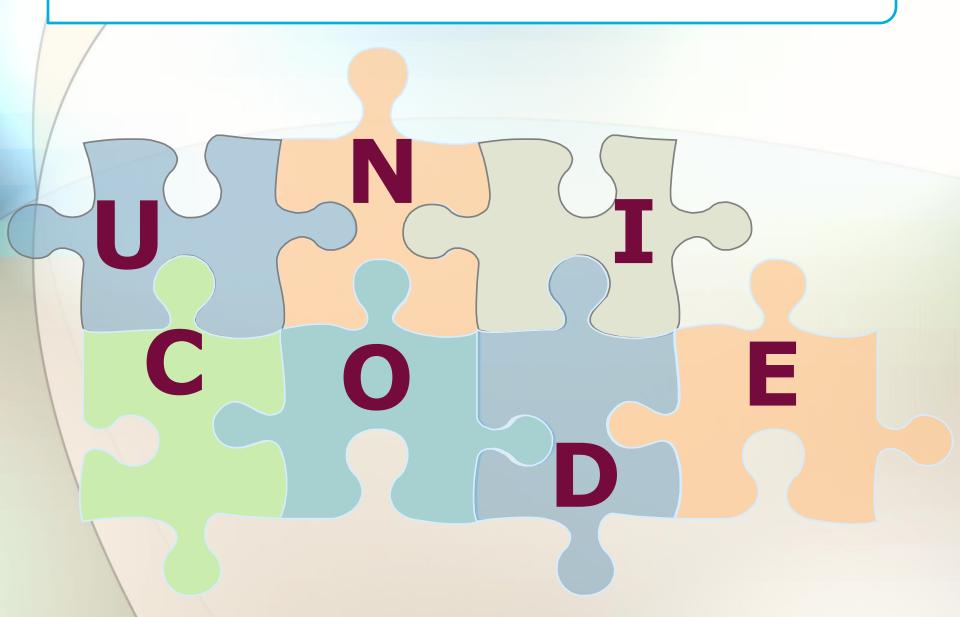














UNICODE

В конце 90-ых годов появился новый международный стандарт Unicode, который отводит под один символ не один байт, а два (с его помощью можно закодировать не 256, а 65536 различных символов).

Полная спецификация стандарта Unicode включает в себя все существующие, вымершие и искусственно созданные алфавиты мира, а также множество математических, музыкальных, химических и прочих символов

САМОЕ ГЛАВНОЕ

- Для кодирования одного символа используется количество информации, равное 1 байту.
- □ Таблица, в которой всем символам компьютерного алфавита поставлены в соответствие порядковые номера (коды), называется таблицей кодировки.
- □ Существуют различные кодовые таблицы: ASC II, КОИ8, СР1251 и др.
- Unicodeмеждународный стандарт, который отводит под один символ два байта.

ДОМАШНЕЕ ЗАДАНИЕ

- 1. п.1.1.1, учебник Угринович Н.Д. Информатика и ИКТ, для 10 класса, М., 2010
- 2. ЭОР, Представление текста в различных кодировках.

http://fcior.edu.ru/card/28666/predstavlenie-teksta-v-razlichnyh-kodirovkah.html-

3. ЭОР, Представление текста в различных кодировках, проверь себя.

http://fcior.edu.ru/card/28605/predstavlenie-teksta-v-razlichnyh-kodirovkah.html

Вспомним известные факты

Процесс преобразования информации в форму, воспринимаемую компьютером называется...

Процесс обратный кодированию называется ...

Каким шифром закодировано словосочетание?

3	е	Л	ë	н	a	Я	ë	Л	К	a
К	3	0	И	p	Г	В	И	0	н	Г

Шифр Цезаря Каждая буква исходного текста заменяется третьей после нее буквой в алфавите. Любая информация кодируется в компьютере с помощью последовательностей двух цифр - 0 и 1. Последовательности из 0 и 1

называются ...,

а цифры 0 и 1 - ... или

Такое кодирование информации на компьютере называется

Повторим основные понятия

- Множество символов, с помощью которых записывается текст, называется...
- Число символов в алфавите это его ...
- Формула определения количества информации: $N = 2^b$, где N это ... b это ...
- Единице измерения 8 бит присвоили название ...
 - С помощью 8 бит можно закодировать ... символов.



Какие символы можно закодировать с помощью 8 бит?

Подсчитаем количество символов

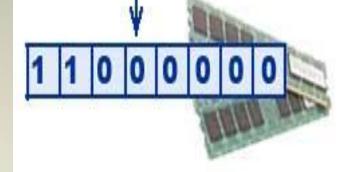
- для русского алфавита 33 строчные буквы
 + 33 прописные буквы = 66;
- для английского алфавита 26 + 26 = 52;
- цифры от 0 до 9.
- Итого 128 символов.
- Остается 128 значений, которые можно использовать для обозначения ...

Человек различает символы по их начертанию, а компьютер - по их коду.



символ

уникальный десятичный код от 0 до 255

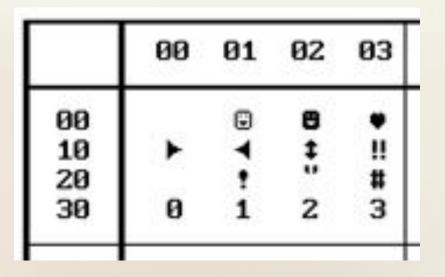


соответствующий двоичный код

от 00000000 до 11111111

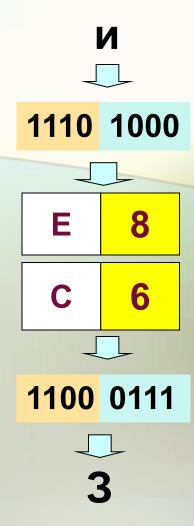
Таблица кодировки -

таблица, в которой устанавливается соответствие между символами и их порядковыми номерами в компьютерном алфавите.



символ - код

		0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
		0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1
		0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
		0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
I	1000	128	129	130	131	132	133	134	5	136	137	138	139	140	141	142	143
		ъ	Ĺ	,	ŕ	"	•••	†	\	€	‰	Љ	(њ	K	ħ	Ţ
	4004	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
	1001	ħ	6	,	"	"	•	-	_		TM	љ	>	њ	K	ħ	Ų
	4040	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
	1010		ÿ	ÿ	J	¤	ľ		§	Ë	©	$\mathbf{\epsilon}$	«	٦		®	Ϊ
	4044	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
	1011	0	±	I	i	ľ	μ	¶	•	ë	Nº	ε	»	j	S	S	ï
Ī	4400	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207
	1100	A	Б	В	Γ	Д	Ε	Ж	3	И	Й	К	Л	M	Н	0	П
Ī	4404	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223
	1101	Р	С	Т	У	Ф	X	Ц	Ч		Щ	Ъ	Ы	Ь	Э	Ю	Я
ĺ	4440	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239
	1110	а	б	В	Γ	Д	е	ж	3	И	Й	K	Л	М	Н	0	П
	4444	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255
	1111	р	С	Т	у	ф	X	ц	ч	Ш	щ	ъ	Ы	Ь	Э	ю	Я



ASCII

(American Standard Code for Information Interchange) – стандартный код информационного обмена США.

sp	1	п	#	\$	%	8	1	()	*	+	,	- E		1
32	33	34	35	36	37	38	33	40	41	42	43	44	45	46	47
0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63
@	Α	В	С	D	E	F	G	Н	-	J	K	L	М	N	0
64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79
Р	Q	R	S	T	U	٧	W	X	Y	Z]	١	1	^	I
80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95
12.00	а	Ь	C	d	е	f	g	h	i	j	k	L	m	n	0
96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111
р	q	r	S	t	u	٧	w	×	у	z	{	1	}	~	
112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	8 8

Основы построения таблицы

- 1. Символы с номерами от нуля (двоичный код 00000000) до 127 (01111111), буквы латинского алфавита, цифры, знаки препинания, скобки и некоторые другие символы.
- 2.От 128 (двоичный код 10000000) до 255 (11111111), используются для кодировки букв национальнх алфавитов, символов псевдографики и научных символов (например ≤, ≥, ≈).

Принцип последовательного кодирования алфавита

В кодовой таблице ASCII располагаются в алфавитном порядке - прописные, строчные с повтором через 32 знака.

Расположение цифр также упорядочено по возрастанию значений.

Благодаря этому и в машинном представлении для символьной информации сохраняется понятие «алфавитный порядок».

КОИ8

Первый стандартов кодирования русских букв на компьютерах -

«Код обмена информацией, 8-битный».

-	100	Γ,	7.	L]	+	+	Τ.	107	+	400			110	110
128	129	130	131 }	132	133	134 \[\(\)	135	136	137	nbsp	139	140	141	142	143 ÷
144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
=		F	ë	П	F	7	П	٦	F	Ш	L	L	П	다	F
160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
l l	l	=	Ë	-11	4	₹	π	īF	土	Ш	兀	+	#	뀨	O
176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
Ю 192	a 193	б 194	Ц 195	Д 196	e 197	ф 198	Г 199	X 200	И 201	Й 202	K 203	л 204	M 205	H 206	0 207
П 208	Я 209	p 210	C 211	T 212	y 213	ж 214	B 215	b 216	Ы 217	3 218	Ш 219	3 220	Щ 221	Ч 222	ъ 223
Ю 224	A 225	Б 226	Ц 227	Д 228	E 229	ф 230	Г 231	X 232	И 233	Й 234	K 235	Л 236	M 237	H 238	0
П	Я	Р	С	Т	у	ж	В	Ь	Ы	3	Ш	Э	Щ	Ч	Ъ
240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255

CP1251

Наиболее распространенная в настоящее время кодировка Microsoft Windows ("CP" означает "Code Page", "кодовая страница").

Á	à	,	è	,,		Ŧ	‡	€	%	É	<	й	Й	ó	ý
128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143
á	•	•	66	"	•	-	_	è	TM	é	>	ò	й	ó	ý
144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
nbsp	ý	Ы	á	Ħ	ы	1	8	Ë	0	Ю́	«	•	shy	8	я́
160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
•	±	ы́	á	,	μ	¶	•	ë	Nº	ю́	>>	à	ю̀	À	я́
176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
A	Б	В	Г	Д	E	Ж	3	И	Й	K	Л	М	Н	0	П
192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207
P 208	C 209	T 210	y 211	ф 212	X 213	Ц 214	215	Ш 216	Щ 217	Ъ 218	Ы 219	ь 220	Э 221	Ю 222	Я 223
а	б	В	Г	Д	е	ж	3	и	й	K	л	М	н	0	п
224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239
Р	C	Т	У	ф	×	ц	ч	ш	щ	ъ	ы	ь	3	ю	Я
240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255



Физкульминутка



Упражнение первое:

резко зажмурить глаза на 2-3 секунды: и широко открыть на 2-3 секунды, повторить упражнение 10 раз.

Упражнение второе:

часто-часто моргать глазами, повторить 10 раз.

Упражнение третье:

поднять глаза вверх, при этом голова остается в одном положении, задержать взгляд на 2-3 секунды, затем опустить глаза вниз и задержать взгляд на 2-3 секунды повторить упражнение 10 раз .

ЗАДАНИЯ



1. В таблице ниже представлена часть кодовой таблицы ASCII.

Каков шестнадцатеричный код символа "q"?

Символ	1	5	Α	В	Q	a
Десятич- ный код	49	53	65	66	81	97
Шест- надцате- ричный код	31	35	41	42	51	61

1) 71 2) 83 3)A1 4) B3

- 2.Буква «і» в таблице кодировки символов имеет десятичный код 105. Что зашифровано последовательностью десятичных кодов: 108 105 110 107?
- 3.С помощью последовательности десятичных кодов: 99 111 109 112 117 116 101 114 зашифровано слово «computer». Какая последовательность десятичных кодов будет соответствовать этому же слову, записанному заглавными буквами?

4. Представьте в форме десятичного кода слово «ЭВМ» в <u>КОИ8</u> и <u>СР1251</u> кодировках.



Ответ

Последовательности десятичных кодов слова «ЭВМ» в различных кодировках составляем на основе кодировочных таблиц:

КОИ8-Р: 252 247 237

CP1251: 221 194 204

сли перевести последовательности кодов из десятичной системы в шестнадцатеричную:

KOИ8-P: FC F7 ED

CP1251: DD C2 CC

ВЫПОЛНИ САМОСТОЯТЕЛЬНО КОНТРОЛЬНЫЕ ЗАДАНИЯ

http://fcior.edu.ru/card/10902/predstavlenie-teksta-v-razlichnyh-kodirovkah.html-Представление текста в различных кодировках, ЭОР