

Кодирование текстовой информации

•

Как вы, наверное, уже знаете, современные компьютерные системы уже давно не работают только с числовой информацией, как это было на начальном этапе их развития. Они наряду с числами могут обрабатывать и буквы, и слова, причем на различных языках, и так называемые **специальные символы**, к которым относят следующие знаки: . , ! ? " # \$ % ^ & * () _ - = + \ > < и др. Очевидно, что любое математическое выражение или строка букв состоит из этих отдельных элементов — букв, цифр, знаков и т. д. Будем в дальнейшем называть их **символами**. Для представления такой информации в компьютерах создают специальные коды, которые так и называют — **символьные коды**.

Это очень существенный момент, не зря еще в 1964 г. в нашей стране был выпущен посвященный этим вопросам государственный общесоюзный стандарт — ГОСТ 10859—64. В последние 15—20 лет эти вопросы определяются международными стандартами. Например, большие вычислительные машины типа IBM 360 и IBM 370 и их аналоги, которые доминировали в мире в 1970—1980-е гг., использовали **специальный код EBCDIC** (Extended Binary Coded Decimal Interchange Code — расширенный двоично-десятичный код обмена информацией), у нас его аналогом был **код ДКОИ-8** (двоичный код для обмена информацией, 8 бит). **Код ASCII** применяется в ПК, совместимых с IBM, работающих под управлением операционной системы MS DOS (Microsoft Disk Operating System — дисковая ОС).

Сейчас, как вы наверняка знаете, **в основном используется ОС Microsoft Windows**. Она применяет так называемую **ANSI-ко-**

дировку. Но эта кодировка ориентирована на английский язык и не содержит символов кириллицы (русских букв), поэтому американская компания **Microsoft** — разработчик Windows — **создала русскую версию ANSI-кодировки**, которая будет приведена далее. **Наша задача** — не заучивать эти многочисленные символы, а **понять механизм формирования кодов**, который, кстати, практически одинаков во всех символьных компьютерных кодах.

Вы уже знаете, что в компьютере каждый символ представлен в виде байта, состоящего из восьми двоичных разрядов, называемых битами. Вы помните также, что содержимым бита может быть либо «0», либо «1», а также то, что 8 бит могут дать 256 комбинаций из «0» и «1». Вы также знаете, что 4 бит, представляющие тетраду байта, «свернуты» в 16-ричную цифру. Если вы внимательно посмотрите на приведенную таблицу русской версии ANSI-кодировки (рис. 2.11), то увидите, что **каждая строка и столбец начинаются с 16-ричной цифры и соответствующей битовой комбинации**, а на пересечении строк и столбцов в соответствующих клетках находятся **кодируемые символы**. Возьмем, например, **знак «+»**. В строке стоит 16-ричная цифра 2 и соответствующая ей кодовая комбинация — тетрада 0010, а в столбце — 16-ричная цифра В и соответствующая ей кодовая комбинация 1011. Следовательно, **знак «+» в компьютере представлен кодом 00101011** (в строках указана первая половина восьмибитового кода, а в столбцах — вторая), в 16-ричном виде этот код можно записать как 2В. Но в каждой клеточке указано еще какое-то десятичное число, в данном случае 43. Это не что иное, как переведенные в десятичную систему или двоичное число 00101011, или 16-ричное число 2В. Действительно,

$$00101011 = 1 + 2 + 8 + 32 = 43;$$

$$2В = 11 + 2 \cdot 16 = 11 + 32 = 43.$$

		0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
		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
0	0000	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	0001	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
2	0010	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47
			!	“	#	\$	%	&	‘	()	*	+	,	-	.	/
3	0011	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63
		0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	0100	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79
		@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	0101	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95
		P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6	0110	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111
		`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	0111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127
		p	q	r	s	t	u	v	w	x	y	z	{		}	~	□
8	1000	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143
		Ђ	Ѓ	„	ѓ	„	…	†	‡	□	‰	Љ	«	Њ	Ћ	Ќ	Џ
9	1001	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
		ђ	‘	’	“	”	•	—	—		™	љ	›	њ	ќ	ћ	џ
A	1010	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
			Ў	ў	Ј	Љ	Ѓ	Ѕ	Ѕ	Ё	©	€	«	¬		®	Ї
B	1011	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
		°	±	І	і	г	μ	¶	·	ë	№	е	»	ј	ѕ	ѕ	ї
C	1100	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207
		А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
D	1101	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223
		Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
E	1110	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239
		а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
F	1101	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255
		р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я

Рис.

2.11. Структура кода ANSI

Если мы спросим вас, какой комбинацией кодируется заглавная буква, то, посмотрев в таблицу, вы легко определите, что русская буква «Ю» имеет код 11011110, или в 16-ричном виде — DE, а в десятичном — 222.

Итак, этот вариант ANSI-кодировки применяется в Windows для всех текстовых шрифтов, содержащих русские буквы. При работе с компьютером многие пользователи используют как MS DOS, так и Windows, отличающиеся друг от друга, помимо прочего, системами кодирования символьной информации. Windows

содержит стандартные средства для перехода от одной кодировки к другой, которые часто выполняются автоматически и не доставляют особых забот пользователю. Тем не менее об этом необходимо помнить, чтобы осознанно анализировать те внештатные ситуации, которые могут возникнуть при использовании компьютера.

Ответить на вопросы и решить задания:

1. Как определить сколько байт содержит текст последнего абзаца параграфа в кодировке ANSI?
2. Для кодирования букв Р, С, Н, О, Г решили использовать двоичное представление чисел 0, 1, 2, 3 и 4 соответственно (с сохранением одного незначащего нуля в случае одноразрядного представления). Закодируйте последовательность букв НОСОРОГ таким способом и результат запишите восьмеричным кодом.
3. Измеряется температура воздуха, которая может быть целым числом от -30 до 34 градусов. Какое наименьшее количество бит необходимо, чтобы закодировать одно измеренное значение?
4. Определить максимальное количество страниц текста, содержащего по 80 символов в каждой строке и 64 строки на странице, которое может содержать файл, сохраненный на гибком магнитном диске объемом 10Кбайт. (кодировка ASCII)
5. Автоматическое устройство осуществило перекодировку информационного сообщения на русском языке, первоначально записанного в коде Windows-1251, в кодировку Unicode. При этом информационное сообщение увеличилось на 400 бит. Какова длина сообщения в символах?