

Введение в информатику

Е. А. Яревский

**физический факультет
СПбГУ
2022**

ЛЕКЦИЯ 4

Представление информации в компьютере

Представление символьной информации

60е годы – появление 7 битовой таблицы **ASCII**
(**A**merican **S**tandard **C**ode for **I**nformation **I**nterchange)

Всего 128 символов:

33 кода управляющих символов

95 кодов знаков (английский алфавит, цифры, и т.д.)

Только английский!

Русский аналог: КОИ-7 (**К**од для **О**бмена **И**нформацией **7**-битный)
Только русский, естественная транслитерация.

Наследник 7-битовых таблиц – 8-битовые таблицы кодировки.
Фактическая замена: **ISO/IEC 8859-1**.

Годится для почти всех западноевропейских языков (кроме французского и финского).

Их включение, а также символ евро, привело к появлению в 1999 году **ISO/IEC 8859-15**.

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

Всего в ISO/IEC 8859 – четырнадцать таблиц.
ISO/IEC 8859-7 – греческий,
ISO/IEC 8859-6 – арабский,
ISO/IEC 8859-5 – русский.

Распространенные/признанные кодировки русских букв

- 1) **ISO/IEC 8859-5** – используется редко(!). Без псевдографики.
- 2) **CP866** (code page) – так называемая альтернативная кодировка, с псевдографикой, MS-DOS. Сейчас используется мало.
- 3) **CP1251** – предложена Microsoft для Windows. Без псевдографики, другие спецсимволы. Очень распространена.
- 4) **koi8-r** (CP878) – много используется в Linux, естественная связь с КОИ-7. На базе ГОСТ КОИ-8.
- 5) **CP10007** – предложена Apple, очень похожа на CP1251.

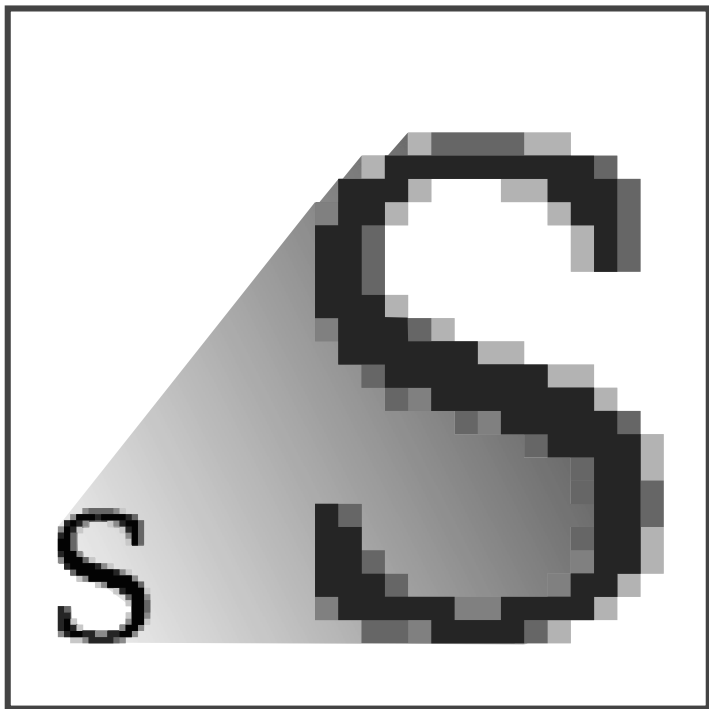
CP1251 (Windows)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
8	Ъ	Ѓ	,	ѓ	„	…	†	‡	□	‰	Љ	<	Њ	Ќ	Ѝ	Ў
9	ђ	‘	’	“	”	•	—	—		™	љ	>	њ	ќ	ћ	џ
A		Ў	ў	Ј	Ќ	Ѓ	!	§	Ё	©	Є	«	¬	·	®	Ї
B	°	±	І	і	г	д	ђ	·	ё	№	є	»	ј	ѕ	ѕ	ї
C	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
D	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
E	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
F	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я

КОИ 8

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
8	—		└	┌	└	┌	└	┌	└	┌	└	■	■	■	■	■
9	▫	▫	▫	┌	■	·	√	≈	≤	≥		┘	°	2	·	÷
A	=		ƒ	ё	п	г	г	п	п	г	ц	ц	г	ц	г	г
B			≠	Ё			≠	п	п	≠	ц	ц	≠			©
C	ю	а	б	ц	д	е	ф	г	х	и	й	к	л	м	н	о
D	п	я	р	с	т	у	ж	в	ь	ы	з	ш	э	щ	ч	ъ
E	Ю	А	Б	Ц	Д	Е	Ф	Г	Х	И	Й	К	Л	М	Н	О
F	П	Я	Р	С	Т	У	Ж	В	Ь	Ы	З	Ш	Э	Щ	Ч	Ъ

Растровые и векторные шрифты



РАСТР

.jpeg .gif .png



ВЕКТОР

.svg

UNICODE

Проблемы традиционных 8-битовых кодировок

- 1) Проблема отображения документов в неправильной кодировке.
- 2) Проблема ограниченности набора символов
- 3) Проблема преобразования одной кодировки в другую
- 4) Проблема дублирования шрифтов

Unicode (Юникод) — стандарт кодирования символов, позволяющий представить знаки практически всех письменных языков.

1.1 - 1991-1996 годы (ISO/IEC 10646—1:1993)

16-битовые символы, $2^{16} = 65536$ символов

2.0 - 1996 год, UTF-16, $2^{20} + 2^{16} - 2048 = 1\,112\,064$ символов

.....

6.2 - 2012 год (кол-во письменностей 100, кол-во символов 110 182)

.....

15.0 - 2022 год (кол-во письменностей 161, кол-во символов 149 186)

Стандарт состоит из двух основных разделов:

1) **универсальный набор символов** (UCS, universal character set).

Задаёт однозначное соответствие символов кодам — элементам кодового пространства, представляющим неотрицательные целые числа (UCS-2, UCS-4).

2) **семейство кодировок** (UTF, Unicode transformation format).

Определяет машинное представление последовательности кодов UCS.

Изначально 16 бит, обозначают 4мя 16-ричными цифрами

коды от U+0000 до U+007F содержит символы набора ASCII
символы кириллицы: коды от U+0400 до U+052F, от U+2DE0 до U+2DFF, от U+A640 до U+A69F

U+0410 **А** CYRILLIC CAPITAL LETTER A

U+0462 **Ъ** CYRILLIC CAPITAL LETTER YAT

U+046D **Ѣ** CYRILLIC SMALL LETTER IOTIFIED YUS

Различные типы представления UTF

UTF-7 - для передачи по 7-битным каналам, не совместима с ASCII, не включена в стандарт

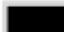

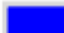
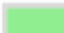












UTF-8 – обеспечивает наилучшую совместимость с 8-битными символами. 1 байт со значением до 128 – символ ASCII (0xxx xxxx). Остальные символы – последовательности от 2 до 6 (фактически, 4) байт. Первый байт имеет вид 11xx xxxx, остальные — 10xx xxxx.

UTF-16 (UTF-16BE, UTF-16LE) – двухбайтовые слова, произвольные значения.
17 плоскостей: от 00 0000 до 10 FFFF (65 536 символов каждая).
U+0000...U+D7FF и U+E000...U+FFFF – обычные символы.
Исключенный диапазон U+D800...U+DFFF – **суррогатные пары**.
Кодируют символы U+010000...U+10FFFF **парой** 16 битных слов, по 10 бит в слове.
Всего $2^{20} + 2^{16} - 2048 = 1\,112\,064$ символов. 2 или 4 байта.

UTF-32 (UTF-32BE, UTF-32LE) – представление 32 битами, 2^{32} символов, избыточно.

Базовая плоскость UNICODE (2¹⁶ СИМВОЛОВ)

00	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	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
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

-  Latin scripts and symbols
-  Linguistic scripts
-  Other European scripts
-  African scripts
-  Middle Eastern and Southwest Asian scripts
-  Central Asian scripts
-  South Asian scripts
-  Southeast Asian scripts
-  East Asian scripts
-  Unified CJK Han
-  American scripts
-  Symbols
-  Diacritics
-  UTF-16 surrogates and private use
-  Miscellaneous characters
-  Unallocated code points

Базовая плоскость UNICODE (2^{16} СИМВОЛОВ)

00	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	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
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

Чёрный — расширенный латинский алфавит;

Голубой — лингвистические символы международного фонетического алфавита IPA;

Синий — другие европейские алфавиты;
Оранжевый — письменности Ближнего Востока;

Светло-оранжевый — письменности Африки;

Зелёный — письменности Южной Азии;

Фиолетовый — письменности ЮВА;

Красный — письменности Восточной Азии;

Розовый — унифицированные

китайско-японско-корейские символы;

Жёлтый — письменности аборигенов

Северной Америки;

Пурпурный — символы;

Тёмно-серый — диакритики;

Светло-серый — суррогатные пары UTF-16 и области для частного использования;

Циан — другие знаки;

Белый — не используется.

ЛИТЕРАТУРА

- 1) Haralambous Y. *Fonts and encodings: from Unicode to advanced typography* (OReilly, 2007), глава 1, 2.
- 2) Д.М. Хэррис и С.Л. Хэррис, *Цифровая схемотехника и архитектура компьютера* (Morgan Kaufman, 2-е издание, 2013).