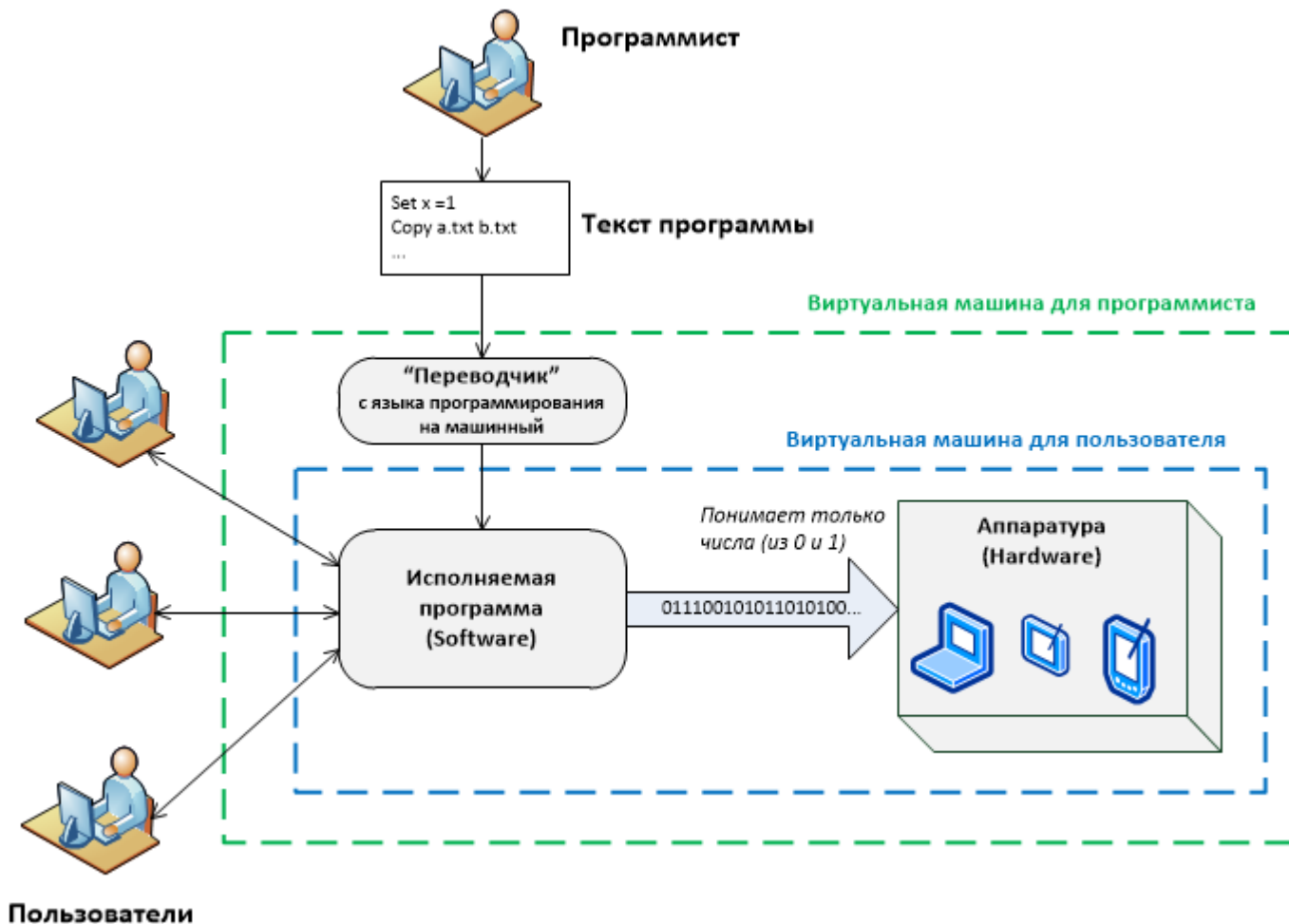


Лекция 2.

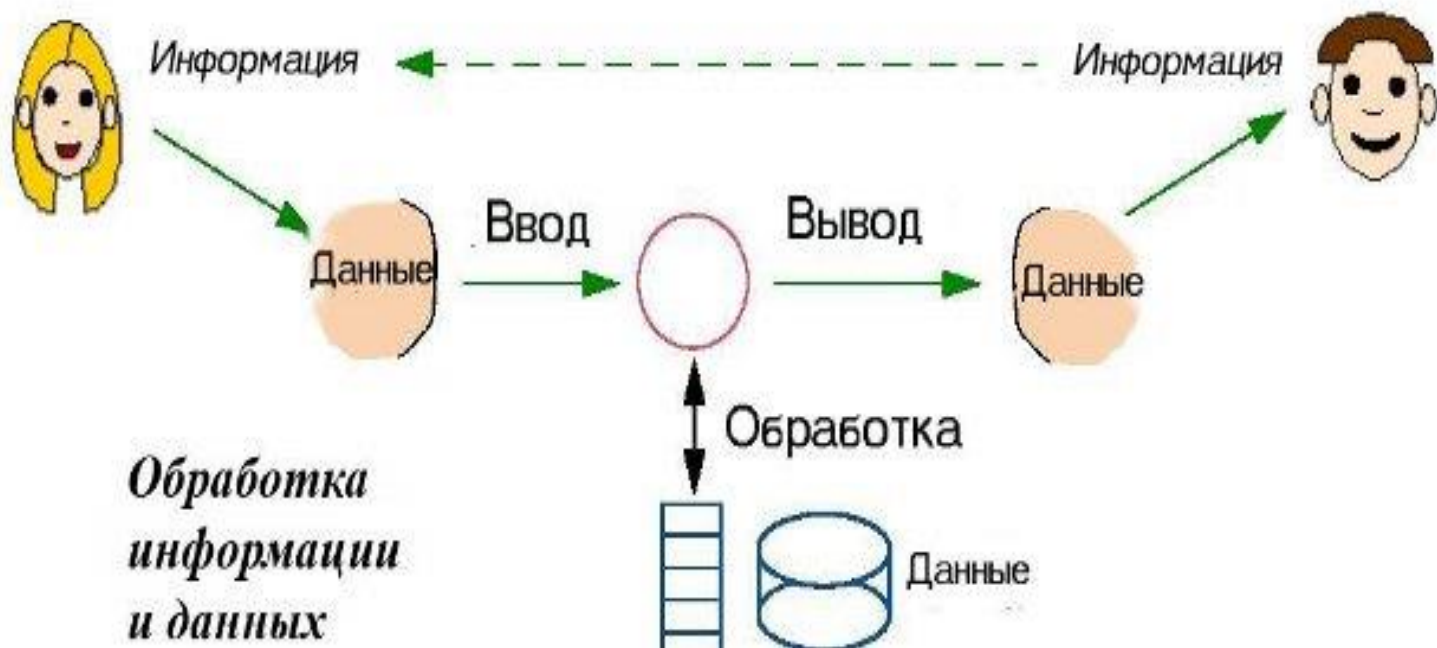
Архитектура фон Неймана. Числовые и символьные данные.



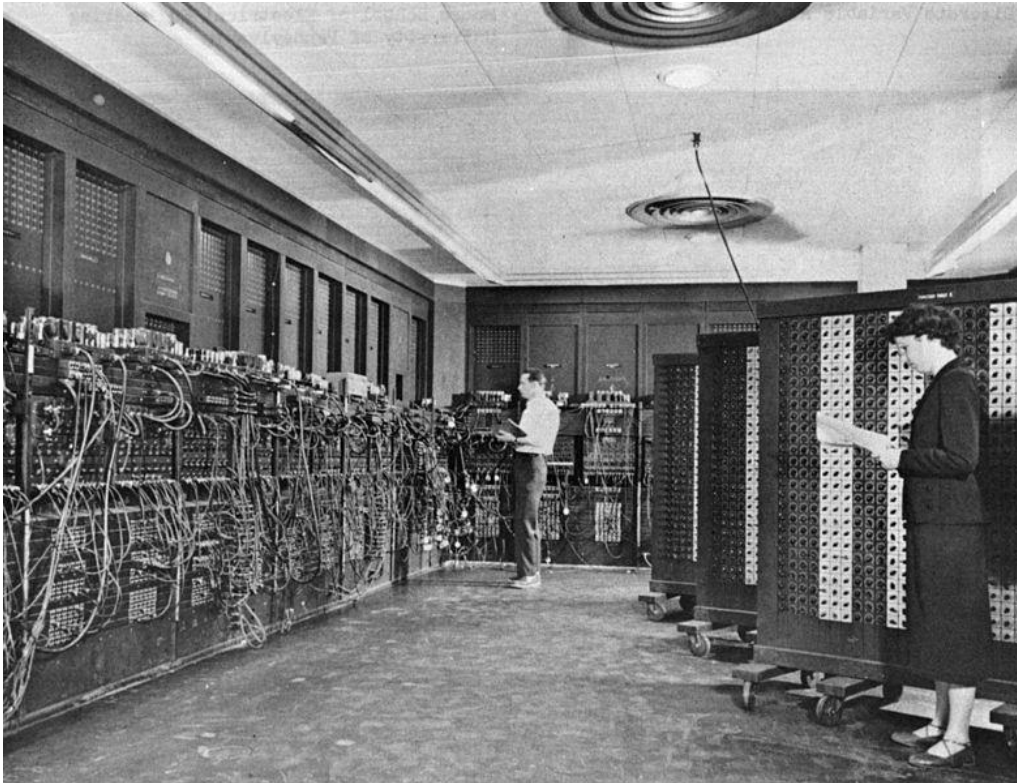
Программа на любом языке программирования выполняется компьютером в машинном коде.

Компьютеры могут: ...

- Хранить и искать информацию в своей памяти.
- Выполнять (очень быстро) базовые операции над информацией: сравнение, замена, арифметические операции.
- Обмениваться данными с другими устройствами (компьютерами, мониторами, принтерами, телефонами, датчиками и т.д.).



Первые компьютерные системы: ограниченный набор исполняемых команд и программ



1945 год. Компьютер
ENIAC (Electronic Numerical
Integrator and Computer)

Площадь - 167 кв.м.

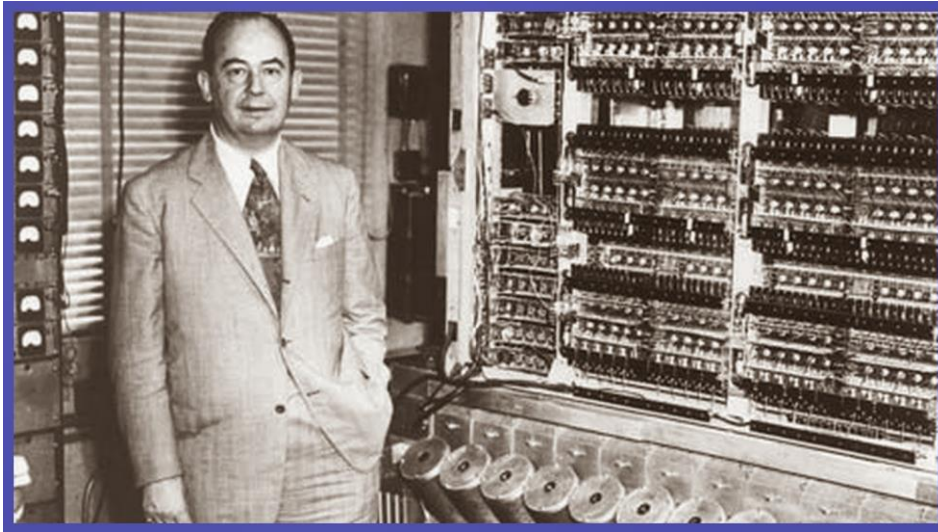
Вес - 27 тонн

Десятичная арифметика

Ввод/вывод данных – перфокарты

Программирование — ручная
установка переключателей

Идея хранения программ в общей памяти



Джон фон Нейман
(1903-1957)

1946 год, статья Дж. фон Нейман, Г. Голдстайна и А. Беркса «Предварительное рассмотрение логической конструкции электронного вычислительного устройства» в журнале Nature

Принципы фон Неймана

- **Принцип двоичного кодирования.** Вся информация кодируется с помощью двоичных сигналов (двоичных цифр, битов) и разделяется на единицы, называемые словами.

Принципы фон Неймана

- **Принцип двоичного кодирования.** Вся информация кодируется с помощью двоичных сигналов (двоичных цифр, битов) и разделяется на единицы, называемые словами.
- **Принцип однородности памяти.** Программы и данные хранятся в одной и той же памяти. Компьютер не различает, что хранится в данной ячейке памяти — число, текст или команда.

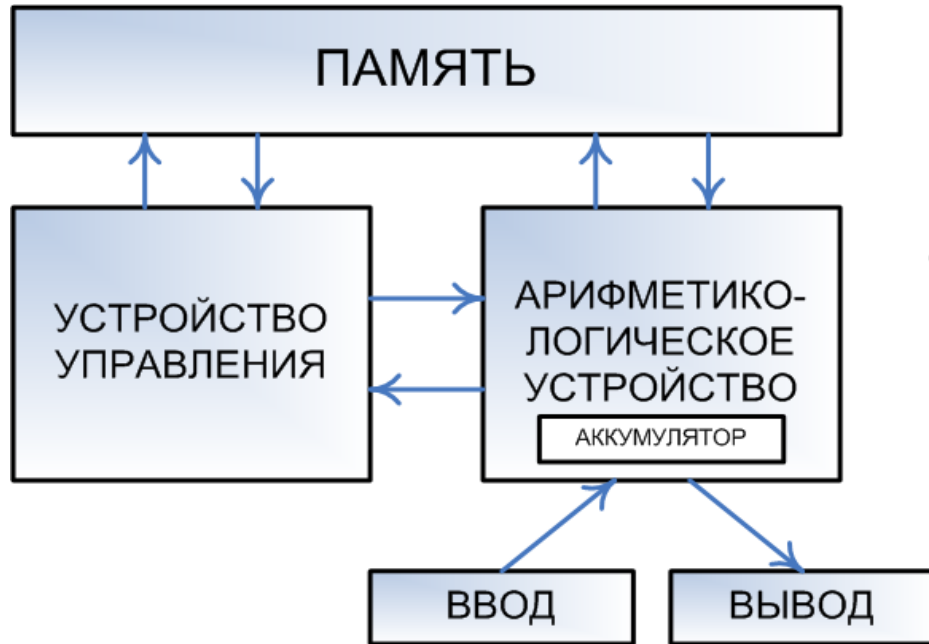
Принципы фон Неймана

- **Принцип двоичного кодирования.** Вся информация кодируется с помощью двоичных сигналов (двоичных цифр, битов) и разделяется на единицы, называемые словами.
- **Принцип однородности памяти.** Программы и данные хранятся в одной и той же памяти. Компьютер не различает, что хранится в данной ячейке памяти — число, текст или команда.
- **Принцип адресуемости памяти.** Память состоит из пронумерованных ячеек; процессору в произвольный момент времени доступна любая ячейка. Можно давать имена областям памяти.

Принципы фон Неймана

- **Принцип двоичного кодирования.** Вся информация кодируется с помощью двоичных сигналов (двоичных цифр, битов) и разделяется на единицы, называемые словами.
- **Принцип однородности памяти.** Программы и данные хранятся в одной и той же памяти. Компьютер не различает, что хранится в данной ячейке памяти — число, текст или команда.
- **Принцип адресуемости памяти.** Память состоит из пронумерованных ячеек; процессору в произвольный момент времени доступна любая ячейка. Можно давать имена областям памяти.
- **Принцип последовательного программного управления.** Предполагает, что программа состоит из набора команд, которые выполняются процессором автоматически друг за другом в определенной последовательности.

Архитектура фон Неймана



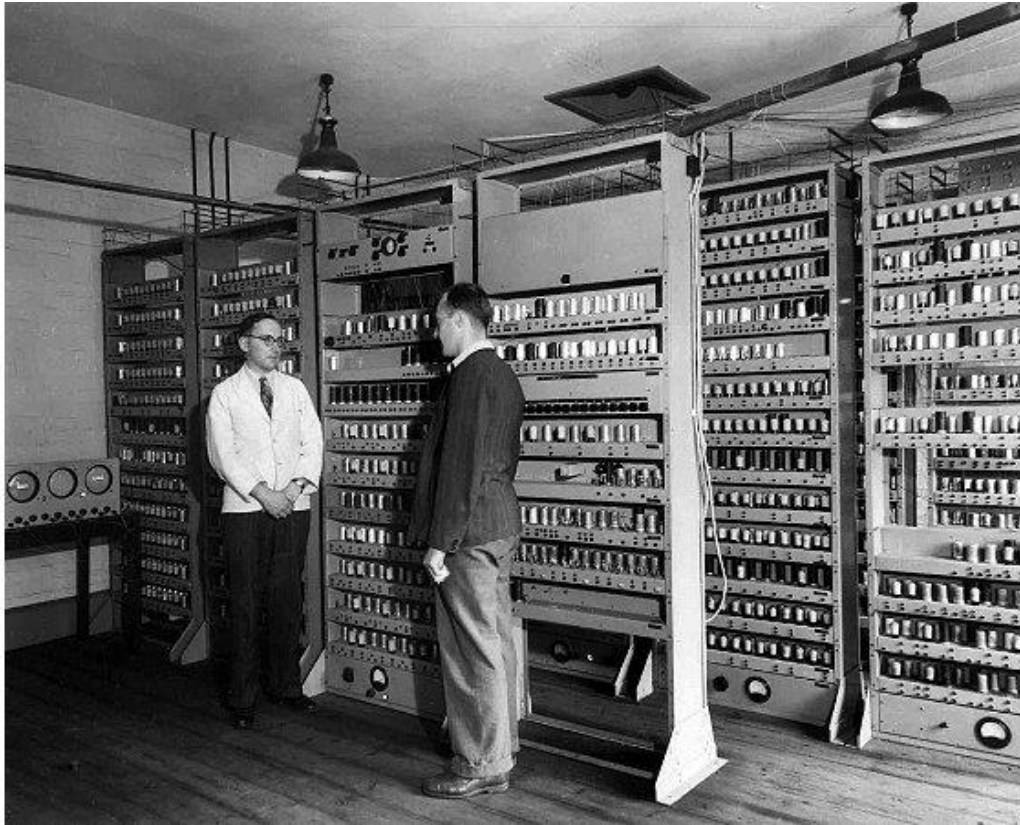
Машина работает в двоичном коде.

Программы и данные хранятся совместно в одном и том же пространстве памяти.

Программа для фон-неймановской машины – последовательность нулей и единиц, одни участки которой кодируют команды, а другие – данные.

Подобная последовательность называется **машинным кодом**.

Первый компьютер с архитектурой фон Неймана



1949 год. Компьютер EDSAC
(Electronic Delay Storage
Automatic Computer)
Кембриджский университет

Поддерживал язык ассемблера и подпрограммы.
В 1952 году на нем была написана первая игра «Крестики-нолики».

Основы двоичной системы

1 бит = 0, 1

1 байт = 8 бит

1 слово = 8 байт = 64 бит (раньше 4 байт = 32 бит). Равно разрядности регистров процессора и/или шины данных.

Основы двоичной системы

1 бит = 0, 1

1 байт = 8 бит

1 слово = 8 байт = 64 бит (раньше 4 байт = 32 бит). Равно разрядности регистров процессора и/или шины данных.

Сколько разных чисел можно закодировать 1 байтом?

Основы двоичной системы

1 бит = 0, 1

1 байт = 8 бит

1 слово = 8 байт = 64 бит (раньше 4 байт = 32 бит). Равно разрядности регистров процессора и/или шины данных.

Сколько разных чисел можно закодировать 1 байтом?

$$00000000_2 = 0_{10}$$

$$00000001_2 = 1_{10}$$

$$00000010_2 = 2_{10}$$

$$00000011_2 = 3_{10}$$

...

$$11111111_2 = 255_{10}$$

Целые положительные числа

Перевод из десятичной в двоичную:

1. Если текущее число меньше 2, то записать его в младший разряд результата, выполнение прекратить. Иначе, разделить текущее число с остатком на 2.
2. Остаток записать в младший разряд результата.
3. Применить пункт 1 к частному.

$$\begin{array}{r} 567 \mid \frac{2}{566} \\ \hline 1 \end{array} \quad \begin{array}{r} 283 \mid \frac{2}{282} \\ \hline 1 \end{array} \quad \begin{array}{r} 141 \mid \frac{2}{140} \\ \hline 1 \end{array} \quad \begin{array}{r} 70 \mid \frac{2}{70} \\ \hline 0 \end{array} \quad \begin{array}{r} 35 \mid \frac{2}{34} \\ \hline 1 \end{array} \quad \begin{array}{r} 17 \mid \frac{2}{16} \\ \hline 1 \end{array} \quad \begin{array}{r} 8 \mid \frac{2}{8} \\ \hline 0 \end{array} \quad \begin{array}{r} 4 \mid \frac{2}{4} \\ \hline 0 \end{array} \quad \begin{array}{r} 2 \mid \frac{2}{2} \\ \hline 0 \end{array} \quad \begin{array}{r} 2 \mid \frac{2}{2} \\ \hline 1 \end{array}$$

$$567_{10} = 1000110111_2$$

Перевод из двоичной в десятичную:

$$1 \cdot 2^9 + 0 \cdot 2^8 + 0 \cdot 2^7 + 0 \cdot 2^6 + 1 \cdot 2^5 + 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = 512 + 0 + 0 + 0 + 32 + 16 + 0 + 4 + 2 + 1 = 567$$

Целые числа со знаком

Как учесть знак числа, как хранить отрицательные числа?

Целые числа со знаком

Как учесть знак числа, как хранить отрицательные числа?

Под знак числа отводят один бит (двоичный разряд). Остальные разряды определяют модуль числа.

Целые числа со знаком

Как учесть знак числа, как хранить отрицательные числа?

Под знак числа отводят один бит (двоичный разряд). Остальные разряды определяют модуль числа.

Прямой код

Старший бит – знак числа (0, если «+» и 1, если «-»). Остальные биты не меняются.

Два недостатка:

$$+0 = 0.0000000_2 \text{ и } -0 = 1.0000000_2$$

$$a = 65_{10} = 0.1000001_2$$

$$-a = -65_{10} = 1.1000001_2$$

$$a+(-a) = 0.0110010 + 1.0110010 = 1.1100100 \neq 0$$

Целые числа со знаком

Как учесть знак числа, как хранить отрицательные числа?

Под знак числа отводят один бит (двоичный разряд). Остальные разряды определяют модуль числа.

Прямой код

Старший бит – знак числа (0, если «+» и 1, если «-»). Остальные биты не меняются.

Два недостатка:

$$\begin{aligned} +0 &= 0.0000000_2 \text{ и } -0 = 1.0000000_2 \\ a = 65_{10} &= 0.1000001_2 \\ -a = -65_{10} &= 1.1000001_2 \\ a+(-a) &= 0.0110010 + 1.0110010 = 1.1100100 \neq 0 \end{aligned}$$

Обратный код

Старший бит – знак числа. Остальные биты инвертируются (0 -> 1, 1 -> 0).

$$\begin{aligned} +0 &= 0.0000000_2 \text{ и } -0 = 1.1111111_2 \\ a = 65_{10} &= 0.1000001_2 \\ -a = -65_{10} &= 1.0111110_2 \end{aligned}$$

В сложении участвуют все биты, единица переноса из знакового разряда прибавляется к младшему разряду суммы.

$$a+(-a) = 0.1000001_2 + 1.0111110_2 = 1.1111111_2 = -0$$

Целые числа со знаком

Дополнительный код

Старший бит – знак числа. Остальные биты инвертируются, к результату прибавляется 1.

Единственная запись нуля: $0 = 0.0000000$

В сложении участвуют все биты, единица переноса из знакового разряда отбрасывается.

$$a = 65_{10} = 0.1000001_2$$

$$-a = -65_{10} = 1.0111111_2$$

$$a + (-a) = 0.1000001_2 + 1.0111111_2 = 0.0000000_2 = 0 \text{ (1 в 9 разряде отбросили)}$$

Целые числа со знаком

Дополнительный код

Старший бит – знак числа. Остальные биты инвертируются, к результату прибавляется 1.

Единственная запись нуля: $0 = 0.0000000$

В сложении участвуют все биты, единица переноса из знакового разряда отбрасывается.

$$a = 65_{10} = 0.1000001_2$$

$$-a = -65_{10} = 1.0111111_2$$

$$a + (-a) = 0.1000001_2 + 1.0111111_2 = 0.0000000_2 = 0 \text{ (1 в 9 разряде отбросили)}$$

Пример: $145 - 102 = 43$

$$\begin{array}{r} [x_1]_o = 0.10010001 \\ + [x_2]_o = 1.10011001 \\ \hline 10.00101010 \\ \begin{array}{l} | \\ \hline \longrightarrow +1 \end{array} \\ \hline 0,00101011 \end{array}$$

Обратный код

$$\begin{array}{r} [x_1]_д = 0.10010001 \\ + [x_2]_д = 1.10011010 \\ \hline 10.00101011 \\ \hline 0,00101011 \end{array}$$

Дополнительный код

Вещественные числа

How to draw an owl

1.



1. Draw some circles

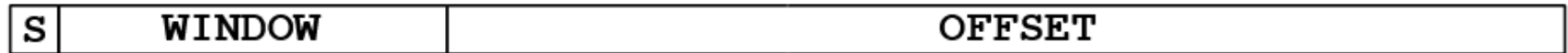
2.



2. Draw the rest of the fucking owl

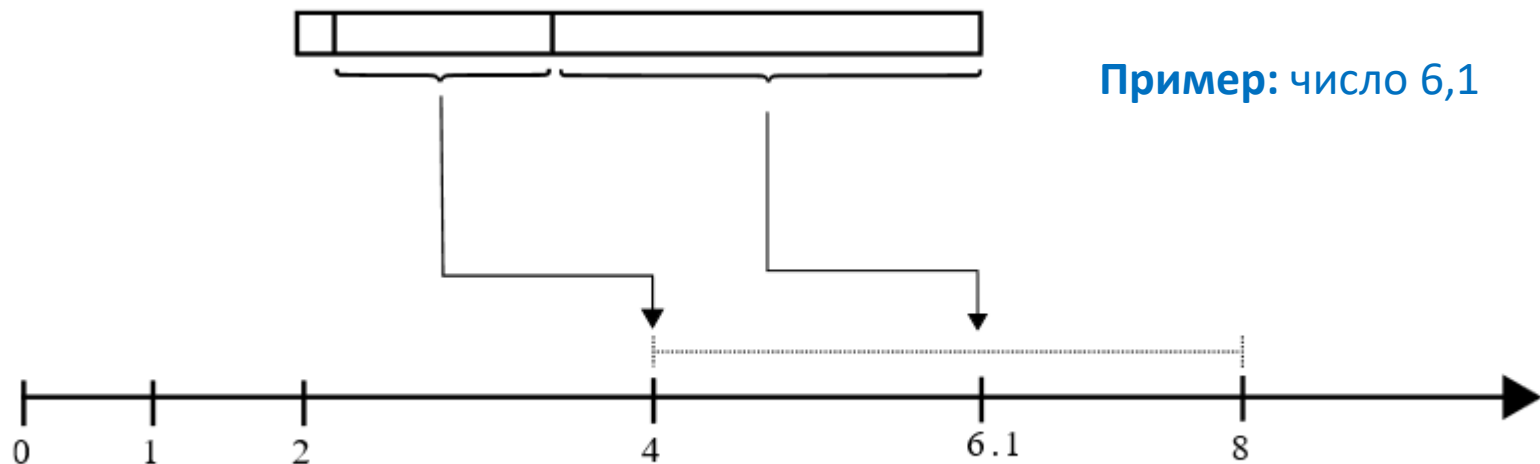
Вещественные числа

«Наглядное объяснение чисел с плавающей запятой» <https://habrahabr.ru/post/337260/>



Три части числа с плавающей запятой.

- **Окно** – число, определяющее интервал $[0, 1]$, $[1, 2]$, $[2, 4]$, $[4, 8]$, ..., $[2^{127}, 2^{128}]$, в который попадает целая часть числа.
- Отрезок Window делится на $2^{23} = 8\,388\,608$ сегментов.
- **Смещение** определяет номер сегмента, в которые попадает число.



Значение 6,1 аппроксимированное с помощью числа с плавающей запятой.

Шестнадцатеричная система

Для представления двоичных данных в удобном виде.

Двоичная тетрада = один 16-ичный знак:

$$0000_2 = 0_{16}$$

$$0001_2 = 1_{16}$$

$$0010_2 = 2_{16}$$

$$0011_2 = 3_{16}$$

$$0100_2 = 4_{16}$$

$$0101_2 = 5_{16}$$

$$0110_2 = 6_{16}$$

$$0111_2 = 7_{16}$$

$$1000_2 = 8_{16}$$

$$1001_2 = 9_{16}$$

$$1010_2 = A_{16}$$

$$1011_2 = B_{16}$$

$$1100_2 = C_{16}$$

$$1101_2 = D_{16}$$

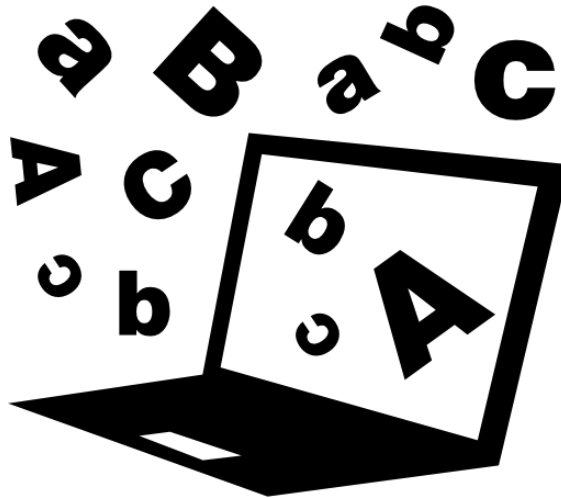
$$1110_2 = E_{16}$$

$$1111_2 = F_{16}$$

Символы, буквы, знаки

Для работы с символами в компьютере нужно:

1. Определить множество символов, которые должны отображаться.
2. Сопоставить каждому символу число (код).
3. Определить способ хранения этих чисел в памяти компьютера (кодирование).



Первые кодировки

Вначале требовалось кодировать минимум символов:

- 10 цифр
- 26 заглавных букв английского алфавита
- 26 строчных букв
- математические знаки $+ - / * = > < %$
- знаки препинания $., ! ? : ; ' "$
- различные скобки $() \{ \} []$
- служебные символы $_ ^ \% \$ @ |$
- 32 непечатных управляющих символов

Первые кодировки

Вначале требовалось кодировать минимум символов:

- 10 цифр
- 26 заглавных букв английского алфавита
- 26 строчных букв
- математические знаки $+ - / * = > < \%$
- знаки препинания $., ! ? : ; ' "$
- различные скобки $() \{ \} []$
- служебные символы $_ \wedge \% \$ @ |$
- 32 непечатных управляющих символов

Для этого хватало **128 символов**.

Для хранения информации использовались **7 бит** в байте (старший бит – для контроля или других целей).

Кодовая таблица ASCII (American Standard Code for Information Interchange)

символ	10- Б код	2-Б код	символ	10- Б код	2-Б код	символ	10-Б код	2-Б код	символ	10-Б код	2-Б код
	32	00100000	8	56	00111000	P	80	01010000	h	104	01101000
!	33	00100001	9	57	00111001	Q	81	01010001	i	105	01101001
"	34	00100010	:	58	00111010	R	82	01010010	j	106	01101010
#	35	00100011	;	59	00111011	S	83	01010011	k	107	01101011
\$	36	00100100	<	60	00111100	T	84	01010100	l	108	01101100
%	37	00100101	=	61	00111101	U	85	01010101	m	109	01101101
&	38	00100110	>	62	00111110	V	86	01010110	n	110	01101110
'	39	00100111	?	63	00111111	W	87	01010111	o	111	01101111
(40	00101000	@	64	01000000	X	88	01011000	p	112	01110000
)	41	00101001	A	65	01000001	Y	89	01011001	q	113	01110001
*	42	00101010	B	66	01000010	Z	90	01011010	r	114	01110010
+	43	00101011	C	67	01000011	[91	01011011	s	115	01110011
,	44	00101100	D	68	01000100	\	92	01011100	t	116	01110100
-	45	00101101	E	69	01000101]	93	01011101	u	117	01110101
.	46	00101110	F	70	01000110	^	94	01011110	v	118	01110110
/	47	00101111	G	71	01000111	_	95	01011111	w	119	01110111
0	48	00110000	H	72	01001000	`	96	01100000	x	120	01111000
1	49	00110001	I	73	01001001	a	97	01100001	y	121	01111001
2	50	00110010	J	74	01001010	b	98	01100010	z	122	01111010
3	51	00110011	K	75	01001011	c	99	01100011	{	123	01111011
4	52	00110100	L	76	01001100	d	100	01100100		124	01111100
5	53	00110101	M	77	01001101	e	101	01100101	}	125	01111101
6	54	00110110	N	78	01001110	f	102	01100110	~	126	01111110
7	55	00110111	O	79	01001111	g	103	01100111	□	127	01111111

Кодовая таблица ASCII (American Standard Code for Information Interchange)

символ	10- Б код	2-Б код	символ	10- Б код	2-Б код	символ	10-Б код	2-Б код	символ	10-Б код	2-Б код
	32	00100000	8	56	00111000	P	80	01010000	h	104	01101000
!	33	00100001	9	57	00111001	Q	81	01010001	i	105	01101001
"	34	00100010	:	58	00111010	R	82	01010010	j	106	01101010
#	35	00100011	;	59	00111011	S	83	01010011	k	107	01101011
\$	36	00100100	<	60	00111100	T	84	01010100	l	108	01101100
%	37	00100101	=	61	00111101	U	85	01010101	m	109	01101101
&	38	00100110	>	62	00111110	V	86	01010110	n	110	01101110
'	39	00100111	?	63	00111111	W	87	01010111	o	111	01101111
(40	00101000	@	64	01000000	X	88	01011000	p	112	01110000
)	41	00101001	A	65	01000001	Y	89	01011001	q	113	01110001
*	42	00101010	B	66	01000010	Z	90	01011010	r	114	01110010
+	43	00101011	C	67	01000011	[91	01011011	s	115	01110011
,	44	00101100	D	68	01000100	\	92	01011100	t	116	01110100
-	45	00101101	E	69	01000101]	93	01011101	u	117	01110101
.	46	00101110	F	70	01000110	^	94	01011110	v	118	01110110
/	47	00101111	G	71	01000111	_	95	01011111	w	119	01110111
0	48	00110000	H	72	01001000	`	96	01100000	x	120	01111000
1	49	00110001	I	73	01001001	a	97	01100001	y	121	01111001
2	50	00110010	J	74	01001010	b	98	01100010	z	122	01111010
3	51	00110011	K	75	01001011	c	99	01100011	{	123	01111011
4	52	00110100	L	76	01001100	d	100	01100100		124	01111100
5	53	00110101	M	77	01001101	e	101	01100101	}	125	01111101
6	54	00110110	N	78	01001110	f	102	01100110	~	126	01111110
7	55	00110111	O	79	01001111	g	103	01100111	□	127	01111111

Почему большие и маленькие буквы идут не сразу друг за другом?

Кодовые страницы

Как кодировать символы других алфавитов?

Кодовые страницы

Как кодировать символы других алфавитов?

Нужно использовать все 8 бит в байте. Числам от 0 до 255 сопоставляются графические образы символов – появляются разные кодовые страницы.

А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143
Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
▒	▓	█		┌	┐	└	┘	┌	┐	└	┘	┌	┐	└	┘
176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
┌	┐	└	┘	─	┌	┐	└	┘	┌	┐	└	┘	┌	┐	└
192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207
Ш	Т	П	Ц	Е	Р	П	Т	Т	Л	Г	■	■	■	■	■
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
224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239
Ё	ё	Є	є	Ї	ї	Ў	ў	°	•	•					
240	241	242	243	244	245	246	247	248	249	250					

Шрифты-кодировки

CP866 (DOS)

CP1251 (Windows)

Á	à	,	è	„	...	†	‡	€	%	É	<	й	Й	Ó	Ú
128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143
á	‘	’	“	”	•	–	—	è	™	é	>	ò	í	ó	ú
144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
nbsp	ÿ	Ы	Э	ж	Ы	ı	§	Ë	©	Ю	«	¬	shy	®	Я
160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
°	±	Ы	Э	’	µ	¶	•	ë	№	ю	»	э	ю	я	я
176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207
Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223
а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239
р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255

Кодовые страницы

КОИ-8 (код обмена информацией) – восьмибитовая кодовая страница (исторически первая для кириллических символов).

Пришла из Unix, широко применялась в электронной почте.

Таблица 1.3. Кодировка КОИ-8

128		144	☐	160	—	176	┌	192	ю	208	п	224	Ю	240	П
129		145	▒	161	Ё	177	└	193	а	209	я	225	А	241	Я
130	Г	146	▓	162	г	178	┘	194	б	210	р	226	Б	242	Р
131	Г	147	┌	163	ё	179	Ё	195	ц	211	с	227	Ц	243	С
132	└	148	•	164	г	180	┘	196	д	212	т	228	Д	244	Т
133	┘	149	•	165	Г	181	┘	197	е	213	у	229	Е	245	У
134	┘	150	√	166	Г	182	┘	198	ф	214	ж	230	Ф	246	Ж
135	┘	151	≈	167	Г	183	┘	199	г	215	в	231	Г	247	В
136	┘	152	◁	168	Г	184	┘	200	х	216	ь	232	Х	248	Ь
137	┘	153	▷	169	Г	185	┘	201	и	217	ы	233	И	249	Ы
138	┘	154	┘	170	Г	186	┘	202	й	218	э	234	Й	250	Э
139	■	155	┘	171	Г	187	┘	203	к	219	ш	235	К	251	Ш
140	■	156	•	172	Г	188	┘	204	л	220	э	236	Л	252	Э
141	■	157	≈	173	Г	189	┘	205	м	221	щ	237	М	253	Щ
142	■	158	•	174	Г	190	┘	206	н	222	ч	238	Н	254	Ч
143	■	159	+	175	Г	191	ё	207	о	223	ь	239	О	255	Ь

Кодовые страницы

КОИ-8 (код обмена информацией) – восьмибитовая кодовая страница (исторически первая для кириллических символов).

Пришла из Unix, широко применялась в электронной почте.

Таблица 1.3. Кодировка КОИ-8

128		144	☐	160	—	176	┌	192	ю	208	п	224	Ю	240	П
129		145	▒	161	Е	177	┘	193	а	209	я	225	А	241	Я
130	Г	146	▓	162	г	178	└	194	б	210	р	226	Б	242	Р
131	Г	147	┌	163	ё	179	Е	195	ц	211	с	227	Ц	243	С
132	└	148	•	164	г	180	┘	196	д	212	т	228	Д	244	Т
133	┘	149	•	165	г	181	└	197	е	213	у	229	Е	245	У
134	┌	150	√	166	г	182	┘	198	ф	214	ж	230	Ф	246	Ж
135	┘	151	≈	167	Г	183	└	199	г	215	в	231	Г	247	В
136	└	152	◁	168	Г	184	┘	200	х	216	ь	232	Х	248	Ь
137	┘	153	▷	169	Г	185	└	201	и	217	ы	233	И	249	Ы
138	+	154	┘	170	Г	186	┘	202	й	218	э	234	Й	250	Э
139	■	155	┘	171	Г	187	┘	203	к	219	ш	235	К	251	Ш
140	■	156	•	172	Г	188	+	204	л	220	э	236	Л	252	Э
141	■	157	≈	173	Г	189	┘	205	м	221	щ	237	М	253	Щ
142	■	158	•	174	Г	190	┘	206	н	222	ч	238	Н	254	Ч
143	■	159	+	175	Г	191	ё	207	о	223	ь	239	О	255	Ь

Задача:

1. Записать в КОИ-8 свое имя и фамилию. Прочитать строку с помощью 7-битного “приложения”.

Символьные строки

С-строки (ASCIIZ-строки)

Строка = массив символов, концом строки считается символ с ASCII-кодом 0.

“abc” = [97 98 99 0]

Pascal-строки

Строка = массив символов, длина определяется счетчиком в первом элементе этого массива (1 байт, максимальная длина строки – 255 символов).

“abc” = [3 97 98 99]

Недостатки однобайтовых кодировок

- Одновременно работать можно лишь с 256 символами (первые 128 зарезервированы).

Недостатки однобайтовых кодировок

- Одновременно работать можно лишь с 256 символами (первые 128 зарезервированы).
- Шрифты привязаны к конкретной кодировке.

Недостатки однобайтовых кодировок

- Одновременно работать можно лишь с 256 символами (первые 128 зарезервированы).
- Шрифты привязаны к конкретной кодировке.
- Каждая кодировка представляет свой набор символов и конвертация из одной в другую возможна только с частичными потерями (отсутствующие символы заменяются на графически похожие).

Недостатки однобайтовых кодировок

- Одновременно работать можно лишь с 256 символами (первые 128 зарезервированы).
- Шрифты привязаны к конкретной кодировке.
- Каждая кодировка представляет свой набор символов и конвертация из одной в другую возможна только с частичными потерями (отсутствующие символы заменяются на графически похожие).
- Перенос файлов между устройствами под управлением разных операционных систем затруднителен. Нужно либо иметь программу-конвертер, либо таскать вместе с файлом дополнительные шрифты.

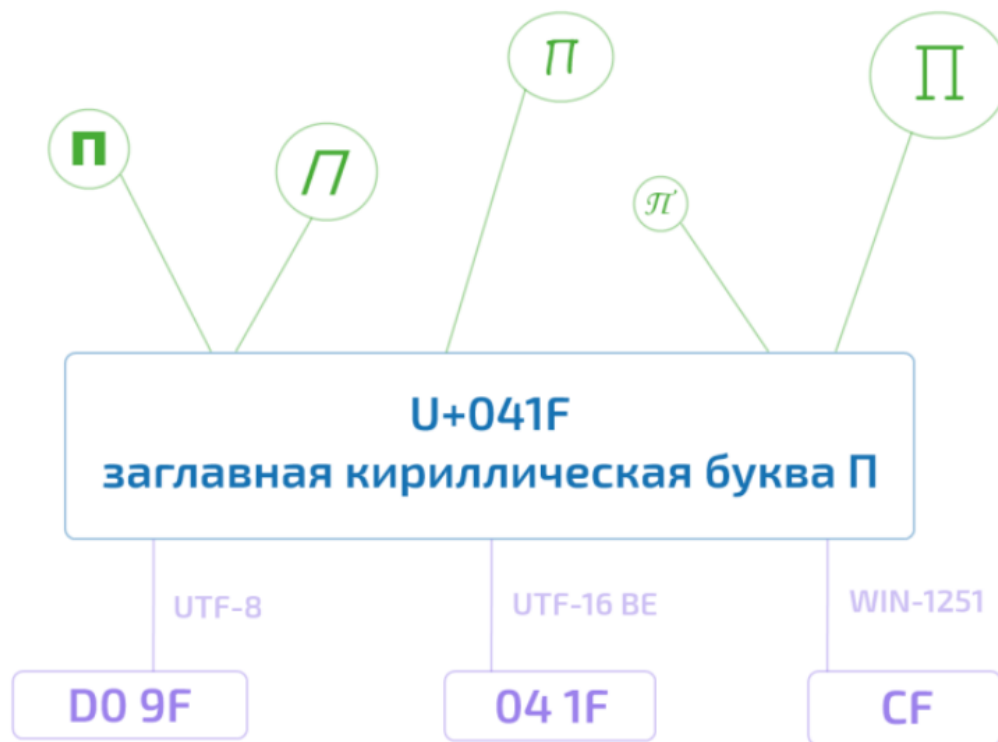
Недостатки однобайтовых кодировок

- Одновременно работать можно лишь с 256 символами (первые 128 зарезервированы).
- Шрифты привязаны к конкретной кодировке.
- Каждая кодировка представляет свой набор символов и конвертация из одной в другую возможна только с частичными потерями (отсутствующие символы заменяются на графически похожие).
- Перенос файлов между устройствами под управлением разных операционных систем затруднителен. Нужно либо иметь программу-конвертер, либо таскать вместе с файлом дополнительные шрифты.
- В мире существуют неалфавитные системы письма (иероглифическая письменность), которые в однобайтной кодировке непредставимы в принципе.

Unicode

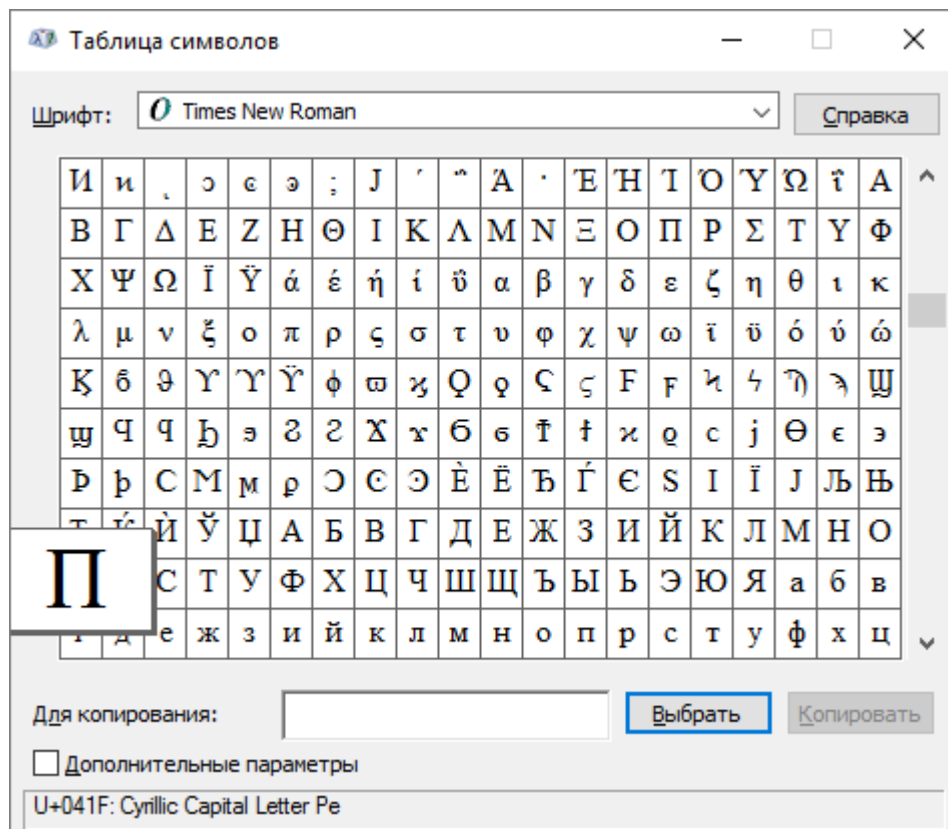
Чёткое разграничение между символами, их представлением в компьютере и их отображением на устройстве вывода.

1. В стандарте определяются абстрактные Unicode-символы.
 - Каждому символу сопоставляется целое число (кодированная позиция, code point). Зарезервировано 1 114 112 кодовых позиций.



Unicode

- Сайт <http://www.unicode.org/charts/>
- Утилита charmap



Unicode

Hello -> U+0048 U+0065 U+006C U+006C U+006F

Как хранить в памяти компьютера?

Unicode

Hello -> U+0048 U+0065 U+006C U+006C U+006F

Как хранить в памяти компьютера?

1. Два байта на каждый символ (65535 символов).

UCS-2 (Universal Character Set, два байта на символ).

Hello -> 00 48 00 65 00 6C 00 6C 00 6F

Hello -> 48 00 65 00 6C 00 6C 00 6F 00

Unicode

Hello -> U+0048 U+0065 U+006C U+006C U+006F

Как хранить в памяти компьютера?

1. Два байта на каждый символ (65535 символов).

UCS-2 (Universal Character Set, два байта на символ).

Hello -> 00 48 00 65 00 6C 00 6C 00 6F

Hello -> 48 00 65 00 6C 00 6C 00 6F 00

2. Мультибайтовая кодировка UTF-8

Unicode Transformation Format, 8 бит

1 байт – 0xxxxxxx

2 байта – 110xxxxx 10xxxxxx

3 байта – 1110xxxx 10xxxxxx 10xxxxxx

4 байта – 11110xxx 10xxxxxx 10xxxxxx 10xxxxxx

Представление первых 127 символов в Unicode совпадает с представлением символов ASCII.

Hello -> 48 65 6C 6C 6F