

Основы программной инженерии (ПОИТ)

Системы программирования

План лекции:

- понятие программного обеспечения;
- системы программирования;
- данные, представление данных, кодировки;
- кодировка ASCII;
- стандарт кодирования Unicode;
- прямой (LE) и обратный (BE) порядок байт;
- маркер последовательности байтов (BOM).

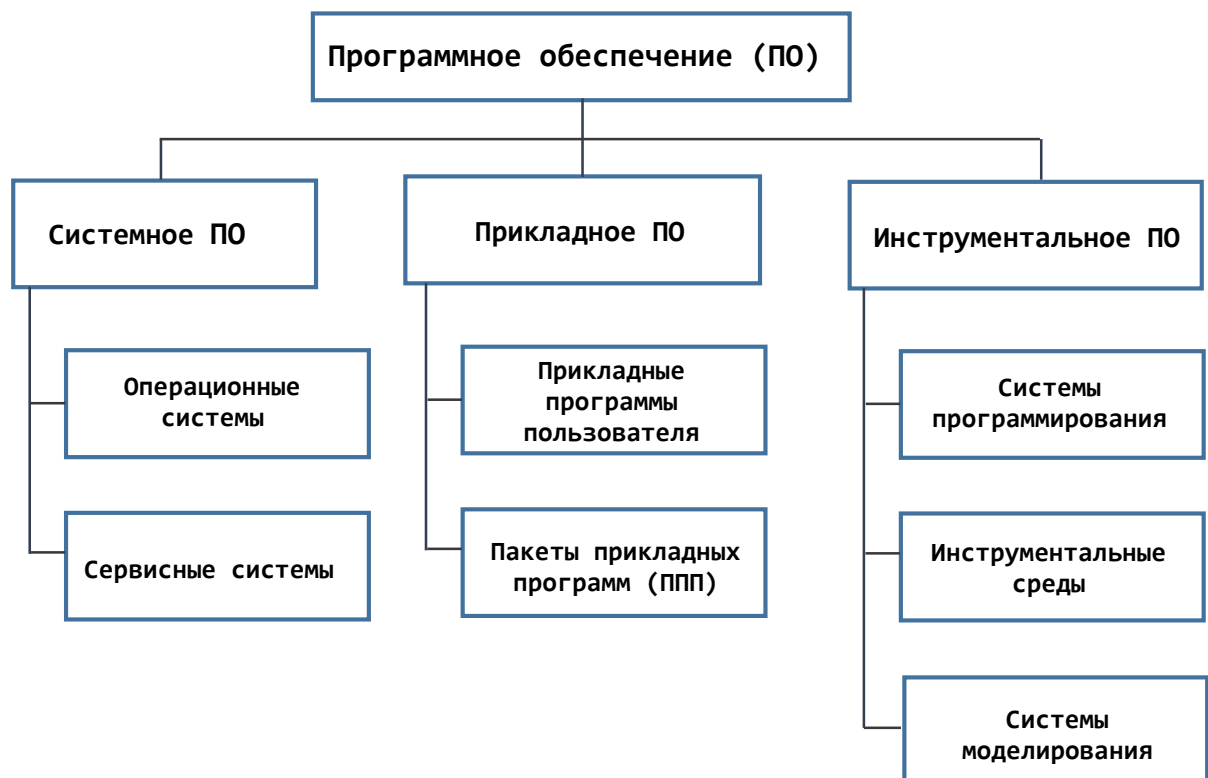
1. Программное обеспечение компьютера

Определение (ГОСТ Р 51904-2002)

Программное обеспечение (ПО) – совокупность компьютерных программ и программных документов, необходимых для эксплуатации этих программ.

Определение (ISO/IEC 26514:2008)

Программное обеспечение (ПО) – программа или множество программ, используемых для управления компьютером.



Классификация программного обеспечения:

Системное ПО – комплекс программ, которые обеспечивают управление компонентами компьютерной системы:

- управление ресурсами компьютера;
- создание копий используемой информации;
- проверка работоспособности устройств компьютера;
- и др.

Прикладное ПО – предназначено для выполнения определённых пользовательских задач и рассчитанная на непосредственное взаимодействие с пользователем.

Инструментальное ПО служит для автоматизации процесса разработки.

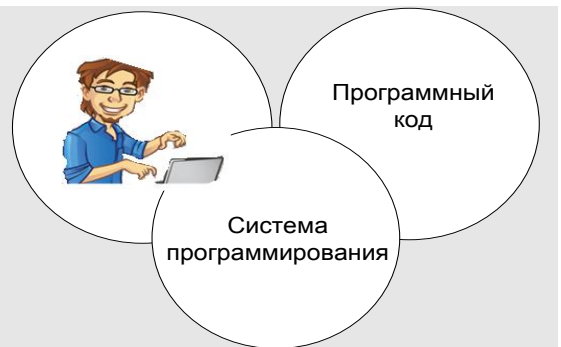
Операционная система – комплекс системных программ, расширяющий возможности вычислительной системы, обеспечивающий управление её ресурсами, загрузку и выполнение прикладных программ, взаимодействие с пользователями.

Системы программирования – системные программы, предназначенные для разработки программного обеспечения.

2. Система программирования

Система программирования:

комплекс программных средств, предназначенных для автоматизации процесса разработки, отладки программного обеспечения и подготовки программного кода к выполнению



Новые требования (тенденции) в современной технологии разработке программного обеспечения:

- распространение промышленных методов организации (планирование трудозатрат, учет, контроль результатов, и т.п.) при проведении работ по работ по разработке программного обеспечения;
- перенос акцента с процесса программирования на более ранние стадии – анализ предметной области, формирование требований.

3. Состав системы программирования:

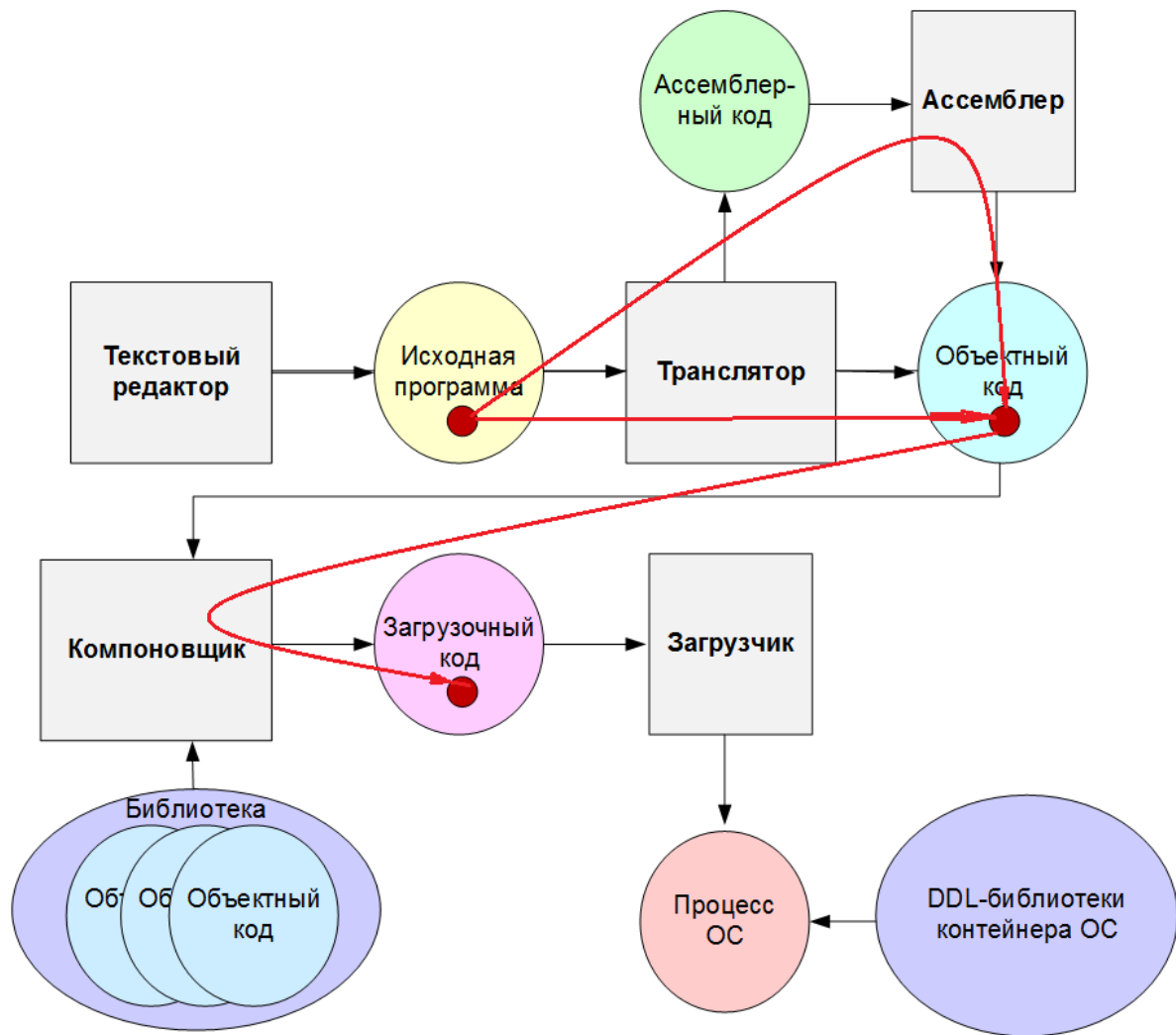
Состав системы программирования

трансляторы
компоновщики
отладчики
профилировщики
программные библиотеки
редакторы кода
системы поддержки версий
и пр.

Система программирования является основным инструментом программиста.



4. Структура классической системы программирования



От исходного кода к исполняемому модулю, основные этапы преобразования:

Классическая схема создания исполняемого файла выполняется для компилируемых языков:

- (1) обработка исходного кода препроцессором,
- (2) компиляция в объектный код и
- (3) компоновка объектных модулей, включая модули из объектных библиотек, в исполняемый файл.

5. Язык программирования:

Язык программирования

формальная знаковая система, предназначенная для записи компьютерных программ.

Знаковая система определяет набор лексических, синтаксических и семантических правил написания программы (программного кода).

Язык программирования представляется в виде набора спецификаций, определяющих его синтаксис и семантику.

Язык программирования представляется в виде набора спецификаций, определяющих его синтаксис и семантику.

Язык программирования определяется не только через спецификации стандарта языка, формально определяющие его синтаксис и семантику, но и через реализации стандарта (программные средства, обеспечивающих трансляцию или интерпретацию программ на этом языке), которые различаются по производителю, версии, времени выпуска, полноте воплощения стандарта, дополнительным возможностям; могут иметь определённые ошибки или особенности реализации.

Спецификация системы программирования: набор требований к системе программирования, достаточный для ее разработки.



6. Кодирование информации

Текстовая информация выражается с помощью естественных или формальных языков в письменной или печатной форме.

Пример:

преподаватель читает лекцию → процесс кодирования студент делает для себя пометки → процесс декодирования студент использует конспект

Системы счисления			Степень двойки		Данные в памяти компьютера							
десятичная	двоичная	шестнадц.			2 ⁷	2 ⁶	2 ⁵	2 ⁴	2 ³	2 ²	2 ¹	2 ⁰
0	0000	0	2 ⁰	1	0	0	0	0	1	0	0	0
1	0001	1	2 ¹	2								
2	0010	2	2 ²	4								
...	2 ³	8								
9	1001	9	2 ⁴	16								
10	1010	A	2 ⁵	32								
11	1011	B	2 ⁶	64								
12	1100	C	2 ⁷	128								
13	1101	D	2 ⁸	256								
14	1110	E										
15	1111	F										

↑							↑
старший							младший
бит							бит

байт – минимальная адресуемая единица
бит – минимальная единица хранения (0 или 1)

Решение проблем с кодировкой текста

Текстовый символ кодируется его **порядковым номером** (0-127), представленным в двоичной системе счисления (1963 ASCII).

Ранние языки, возникшие в эпоху 6-битных символов, использовали более ограниченный набор.

Пример:

алфавит Фортрана включает 49 символов: A B C D E F G H I J K L M N O P Q R S T U V W X Y Z 0 1 2 3 4 5 6 7 8 9 = + - * / () . , \$ ' : пробел

а. Американский стандартный код для обмена информацией. ASCII

ASCII (American Standard Code for Information Interchange) – американский стандартный код для обмена информацией.

ASCII – 8-битная кодировка для представления десятичных цифр, латинского и национального алфавитов, знаков препинания и управляющих символов.

Таблица кодов ASCII делится на две части:

Международным стандартом является первая половина таблицы, т.е. символы с номерами от 0 (00000000), до 127 (01111111).

К концу 1980-х годов стандартом стали 8-битные кодировки.

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

Расширенные таблицы: в расширенных таблицах символы с порядковыми номерами 128-255 представляют символы национальных языков.

Переносимый набор символов

является базовым алфавитом для практически всех современных языков программирования.

Переносимый набор символов (portable character set) – набор из 103 символов, которые должны присутствовать в любой используемой кодировке (стандарт POSIX).

POSIX (англ. Portable Operating System Interface – переносимый интерфейс операционных систем) – набор стандартов, описывающих интерфейсы между операционной системой (ОС) и прикладной программой (системный API), библиотеку языка С и набор приложений и их интерфейсов.

Переносимый набор символов включает в себя все печатные символы US-ASCII и часть управляющих и является базовым алфавитом для практически всех современных языков программирования.

- Альтернативная кодировка **CP866** (операционная система MS-DOS):

Все специфические европейские символы во второй половине таблицы CP866 заменены на кириллицу, а псевдографические символы оставлены без изменения.

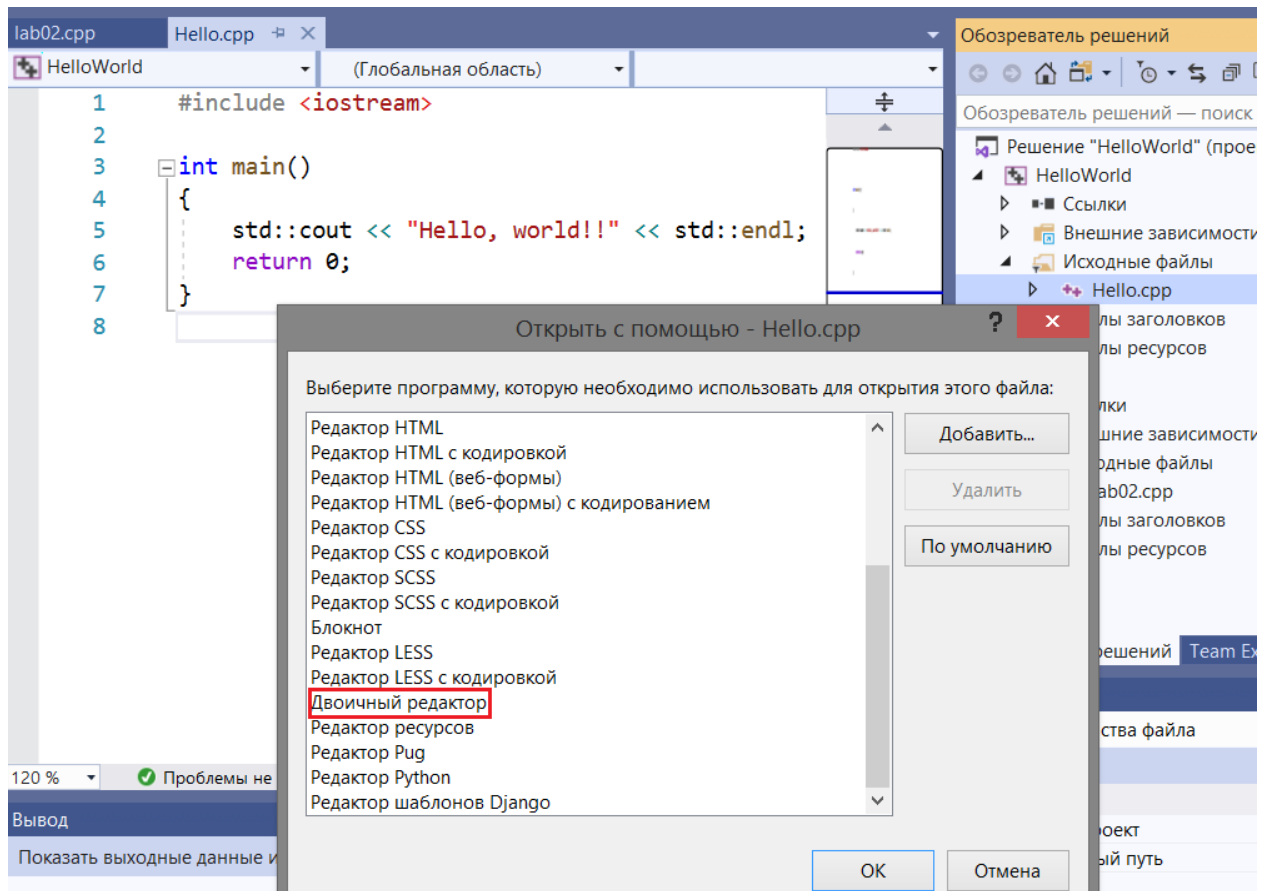
	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
8	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
9	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
A	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
B					┌	┐	└	┘	┌	┐	└	┘	┌	┐	└	┘
C	┐	┌	└	┘	─	┼	┆	┇	┈	┉	┊	┋	┌	┐	└	┘
D	┌	┐	└	┘	┐	┌	└	┘	┐	┌	└	┘	┐	┌	└	┘
E	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
F	Ё	ё	Є	є	Ї	ї	Ў	ў	°	·	·	√	№	□	■	

- русская Windows-кодировка (**Windows-1251**, синоним **CP1251**)

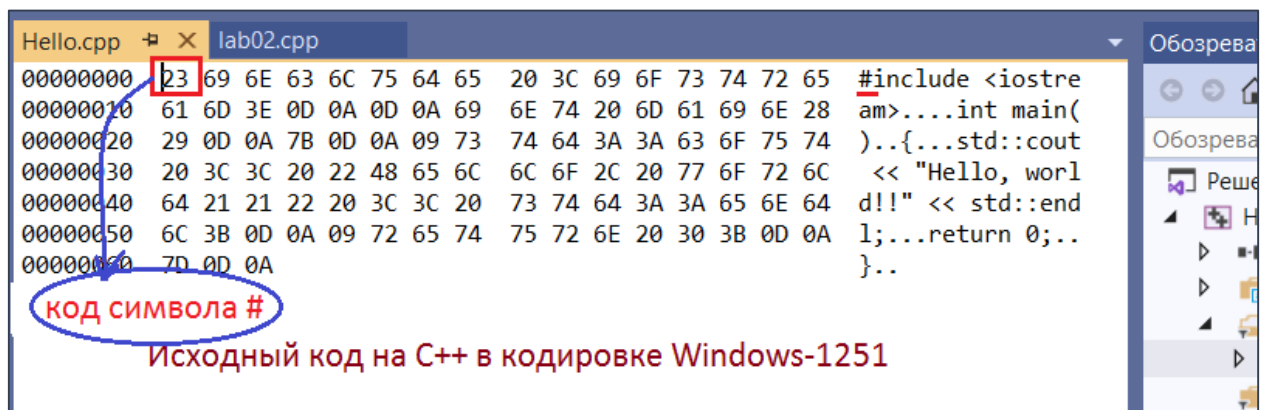
Windows-1251 – набор символов и кодировка, являющаяся стандартной 8-битной кодировкой для русских версий Microsoft Windows до 10-й версии.

	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
00	<u>NUL</u> 0000	<u>STX</u> 0001	<u>SOT</u> 0002	<u>ETX</u> 0003	<u>EOT</u> 0004	<u>ENQ</u> 0005	<u>ACK</u> 0006	<u>BEL</u> 0007	<u>BS</u> 0008	<u>HT</u> 0009	<u>LF</u> 000A	<u>VT</u> 000B	<u>FF</u> 000C	<u>CR</u> 000D	<u>SO</u> 000E	<u>SI</u> 000F
10	<u>DLE</u> 0010	<u>DC1</u> 0011	<u>DC2</u> 0012	<u>DC3</u> 0013	<u>DC4</u> 0014	<u>NAK</u> 0015	<u>SYN</u> 0016	<u>ETB</u> 0017	<u>CAN</u> 0018	<u>EM</u> 0019	<u>SUB</u> 001A	<u>ESC</u> 001B	<u>FS</u> 001C	<u>GS</u> 001D	<u>RS</u> 001E	<u>US</u> 001F
20	<u>SP</u> 0020	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
30	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
40	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
50	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
60	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
70	p	q	r	s	t	u	v	w	x	y	z	{		}	~	<u>DEL</u> 007F
80	Ђ	Ѓ	Ѕ	Ї	Љ	Њ	Ћ	Ќ	Ў	Ъ	Ы	Ь	Э	Ю	Я	а
90	Ђ	Ѓ	Ѕ	Ї	Љ	Њ	Ћ	Ќ	Ў	Ъ	Ы	Ь	Э	Ю	Я	а
A0	<u>NBSP</u> 00A0	Ў	Ў	Ј	Ќ	Љ	Њ	Ћ	Ќ	Ў	Ъ	Ы	Ь	Э	Ю	Я
B0	°	±	І	і	Г	г	Д	д	Е	е	Ж	ж	З	з	И	и
C0	А	В	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
D0	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
E0	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
F0	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я

Интегрированная среда разработки Visual Studio. Открытие файла в двоичном редакторе:



Представление в памяти файла с исходным кодом:



Представление символьной информации в кодировке Windows-1251:

```
lab02.cpp  [icon] [x]
00000000  23 69 6E 63 6C 75 64 65 20 3C 69 6F 73 74 72 65 #include <iostre
00000010  61 6D 3E 0D 0A 0D 0A 69 6E 74 20 6D 61 69 6E 28 am>...int main(
00000020  29 0D 0A 7B 0D 0A 09 63 68 61 72 20 68 65 6C 6C )...{...char hell
00000030  6F 5B 5D 20 3D 20 22 48 65 6C 6C 6F 2C 20 22 3B o[] = "Hello, ";
00000040  0D 0A 09 63 68 61 72 20 66 69 6F 5B 5D 20 3D 20 ...char fio[] =
00000050  22 49 76 61 6E 6F 76 20 49 76 61 6E 20 49 76 61 "Ivanov Ivan Iva
00000060  6E 6F 76 69 63 68 22 3B 0D 0A 09 73 74 64 3A 3A novich";...std::
00000070  63 6F 75 74 20 3C 3C 20 68 65 6C 6C 6F 20 3C 3C cout << hello <<
00000080  20 66 69 6F 20 3C 3C 20 73 74 64 3A 3A 65 6E 64 fio << std::end
00000090  6C 3B 0D 0A 09 72 65 74 75 72 6E 20 30 3B 0D 0A l;...return 0;..
000000a0  7D 0D 0A                                     }..

Исходный код lab02.cpp
```

б. Международный стандарт UNICODE

Решение проблем

неправильного декодирования;
ограниченность набора символов;
преобразования из одной кодировки в другую;
проблема дублирования шрифтов.

Стандарт предложен в 1991 году некоммерческой организацией Unicode Consortium, стандарт ISO/IEC 10646:2020.



Юникод – стандарт кодирования символов, позволяющий представить знаки почти всех письменных языков, состоит из 2х разделов:

- **UCS** – universal character set (универсальный набор символов);
- **UTF** – Unicode transformation format (семейство кодировок).

Принято обозначение символа **U+xxx**, где **xxx**- число в шестнадцатеричном формате.

• UNICODE:

- UCS расположены в 17 плоскостях (0-16);
- в каждой плоскости 2^{16} (65 536) символов;
- плоскость 0 – основная (основные символы);
- 1-14 – дополнительные;
- 15-16 – для частного использования.

- UNICODE: <http://foxtools.ru/Unicode>

Диапазон: 0020-007F: Основная латиница

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
002		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
003	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
004	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
005	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
006	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
007	p	q	r	s	t	u	v	w	x	y	z	{		}	~	

Диапазон: 0400-04FF: Кириллица

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
040	Ё	Ё	Ъ	Ѓ	Є	Ѕ	І	Ї	Ј	Љ	Њ	Ћ	Ќ	Й	Ў	Ц
041	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
042	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
043	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
044	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
045	ё	ё	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	й	ў	ц
046	Ѡ	ѡ	Ѣ	ѣ	Ѥ	ѥ	Ѧ	ѧ	Ѩ	ѩ	Ѫ	ѫ	Ѭ	ѭ	Ѯ	ѯ
047	Ѱ	ѱ	Ѳ	ѳ	Ѵ	ѵ	Ѷ	ѷ	Ѹ	ѹ	Ѻ	ѻ	Ѽ	ѽ	Ѿ	ѿ
048	Ѡ	ѡ	Ѣ	ѣ	Ѥ	ѥ	Ѧ	ѧ	Ѩ	ѩ	Ѫ	ѫ	Ѭ	ѭ	Ѯ	ѯ
049	Ѱ	ѱ	Ѳ	ѳ	Ѵ	ѵ	Ѷ	ѷ	Ѹ	ѹ	Ѻ	ѻ	Ѽ	ѽ	Ѿ	ѿ

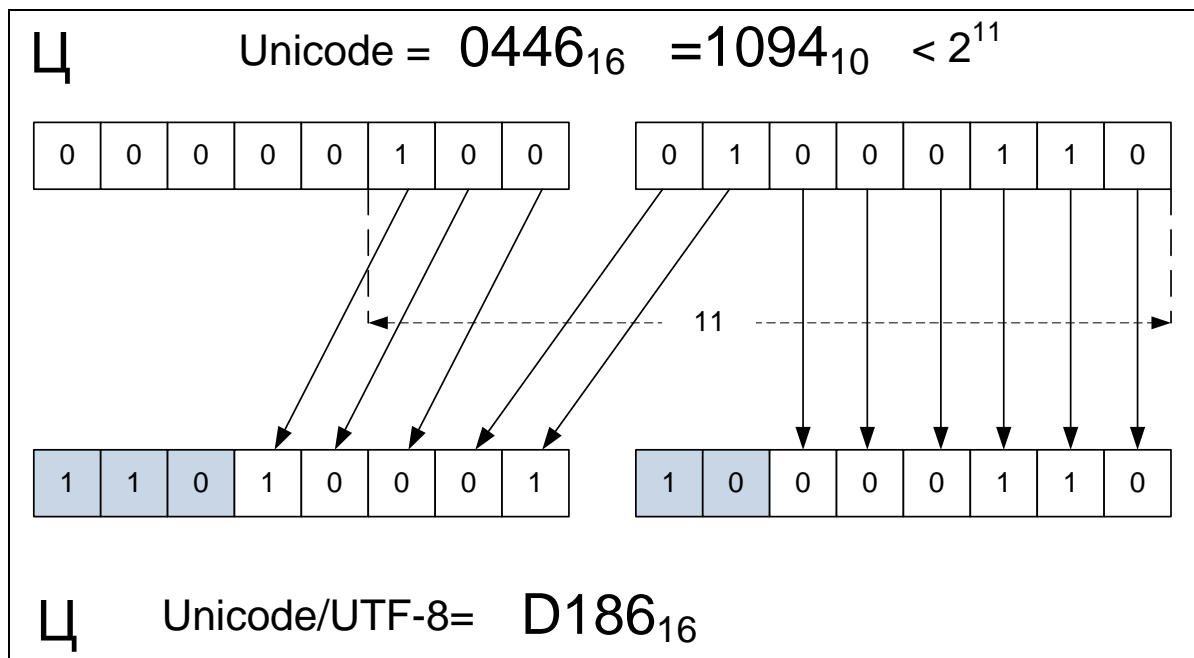
- **UNICODE: кодировка UTF-8**

UTF-8 — представление Юникода, обеспечивающее совместимость со старыми системами, использовавшими 8-битные символы.

Алгоритм кодирования в UTF-8:

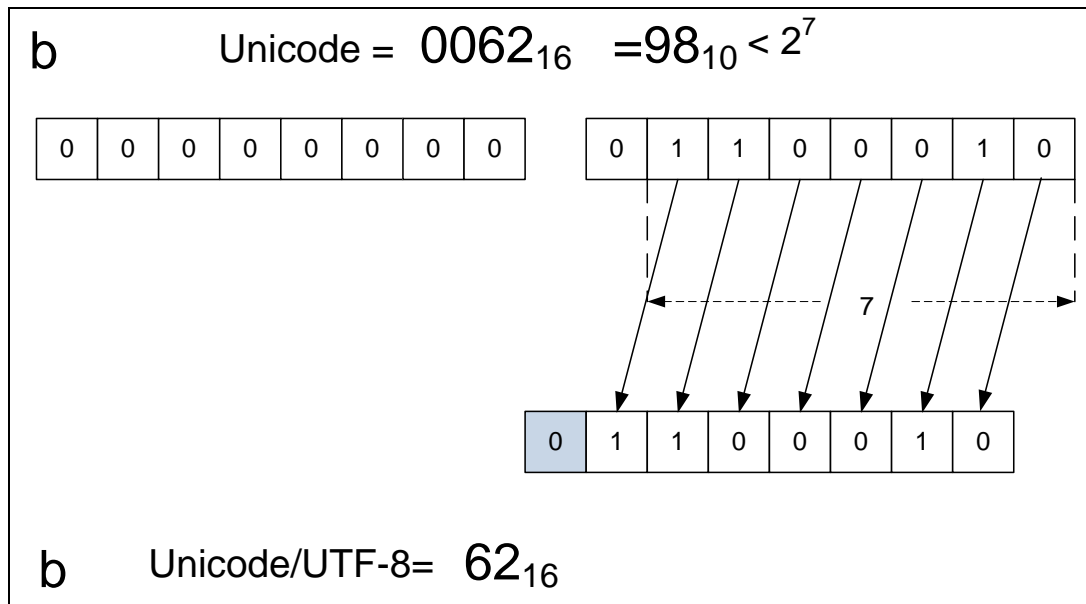
- 1) определить количество октетов (октет: 8 битов или 1 байт) – т.е. в какой диапазон значений попадает количество значащих символов (7, 11, 16, 21, 26, 31);
- 2) подготовить старшие биты первого октета:
 - a. 0xxxxxxx для *одного* октета;
 - b. 110xxxxx – для *двух*;
 - c. 1110xxxx – для *трех* и т.д..
 - d. 10xxxxxx – для *остальных* октетов;
- 3) заполнить оставшиеся биты (выше обозначены как x) в октетах кодом символа Юникода в двоичном виде. Начать с младших битов, поставив их в младшие биты последнего октета кода. И так далее, пока все биты кода символа не будут перенесены в свободные биты октетов.

Пример:



$$0446_{16} = 4 \cdot 16^2 + 4 \cdot 16 + 6 = 1094_{10}$$

Пример:



UNICODE: кодировка UTF-8. Для символов в диапазоне:

$0x00000000 \div 0x0000007F$: **0**xxxxxxx (один октет)

$0x00000080 \div 0x000007FF$: **110**xxxxx **10**xxxxxx (два октета)

$0x00000800 \div 0x0000FFFF$: **1110**xxxx **10**xxxxxx **10**xxxxxx (три октета)

$0x00010000 \div 0x001FFFFF$: **11110**xxx **10**xxxxxx **10**xxxxxx **10**xxxxxx

UNICODE: кодировка UTF-16

В UTF-16 символы кодируются двухбайтовыми словами (16 битов) с использованием всех возможных диапазонов значений (от 0 до FFFF₁₆).

- **Маркер последовательности байтов UNICODE: BOM** (Byte Order Mark)

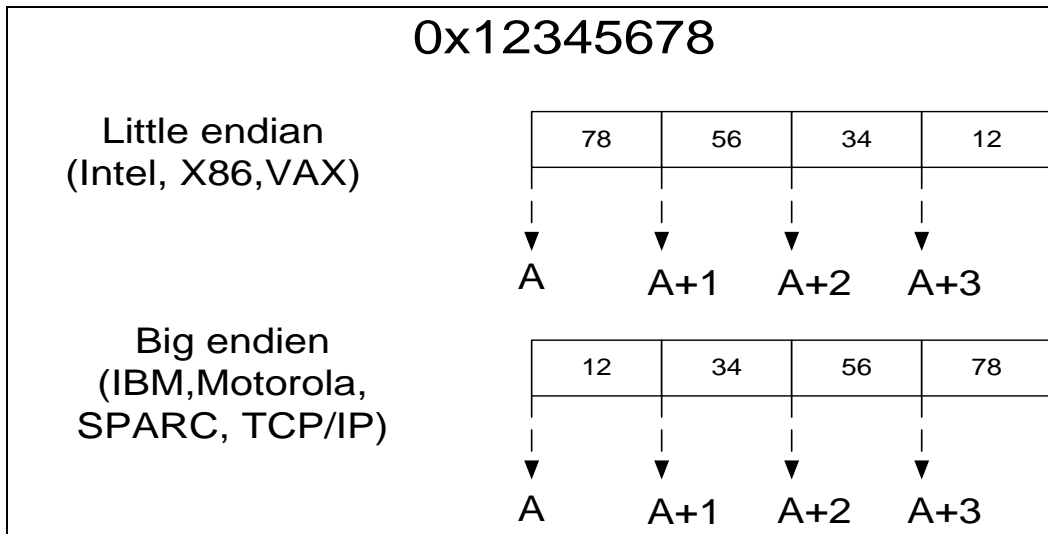
Для определения формата представления Юникода в начало текстового файла записывается **сигнатура** (обозначение) — символ U+FEFF — маркер последовательности байтов.

Шестнадцатеричное представление маркера последовательности байтов для кодировок:

Кодировка	Представление (hex)
UTF-8	EF BB BF
UTF-16 (BE)	FE FF
UTF-16 (LE)	FF FE
UTF-32 (BE)	00 00 FE FF
UTF-32 (LE)	FF FE 00 00

- **Порядок следования байтов:**

- **LE** (Little endian order, прямой порядок, от младшего к старшему);
- **BE** (Big endian order, обратный порядок, от старшего к младшему).



Представление в памяти целочисленного числа на платформе x86:
порядок следования байтов LE
 (Little endian order, прямой порядок, от младшего к старшему)

```

lab02.cpp
Lab_02
1  #include <iostream>
2
3  int main()
4  {
5      int number = 0x12345678;
6      char hello[] = "Hello, ";
7      char fio[] = "Ivanov Ivan Ivanovich";
8      std::cout << hello << fio << std::endl;
9      return 0;
10 }
```

Память 1

Адрес: 0x005DFA14

0x005DFA14	78	56	34	12	cc	cc	cc	cc	1c
0x005DFA18	f0	3d	40	fa	5d	00	33	2e	27
0x005DFA20	01	00	00	00	80	8d	90	00	08
0x005DFA30	90	00	01	00	00	00	80	8d	90
0x005DFA40	08	98	90	00	9c	fa	5d	00	87
0x005DFA50	27	00	a0	91	f0	3d	b6	13	27

Представление шестнадцатиричного числа в памяти компьютера