



HIGHER SCHOOL OF ECONOMICS
NATIONAL RESEARCH UNIVERSITY



Introduction to Programming

Streams and Strings

Sergey Shershakov

#3/17 Jan 2019

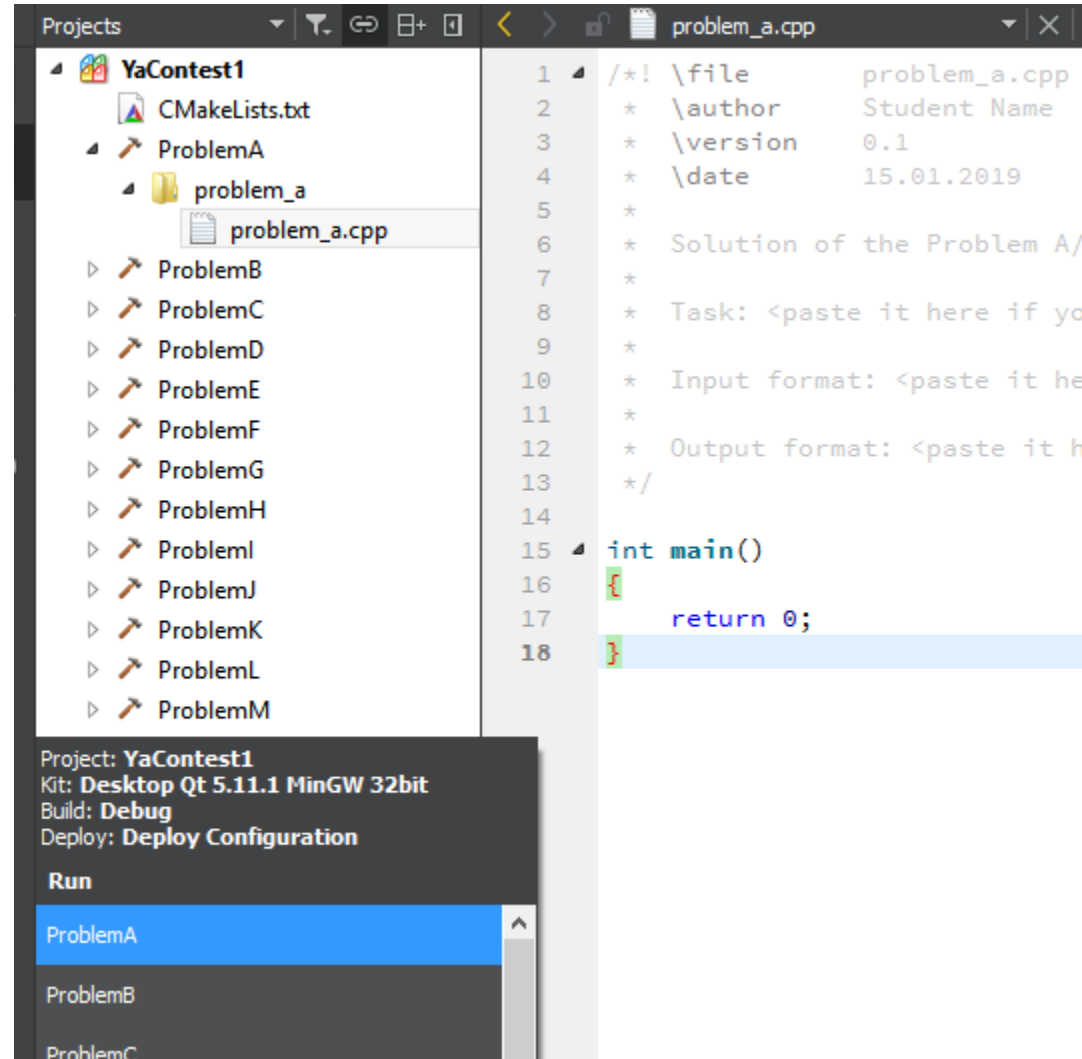
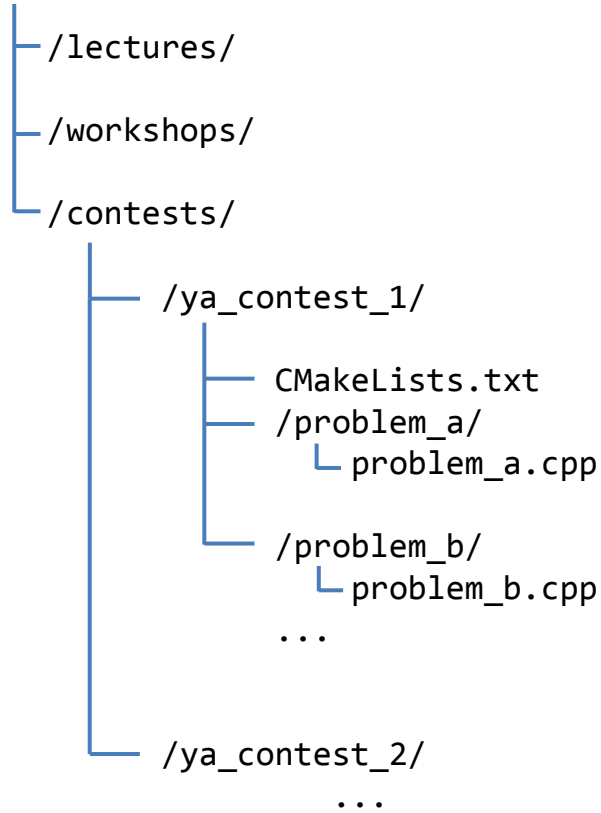
Test 2 (5 pts)



<https://goo.gl/forms/7PqlqR3zSZFgWxMs1>

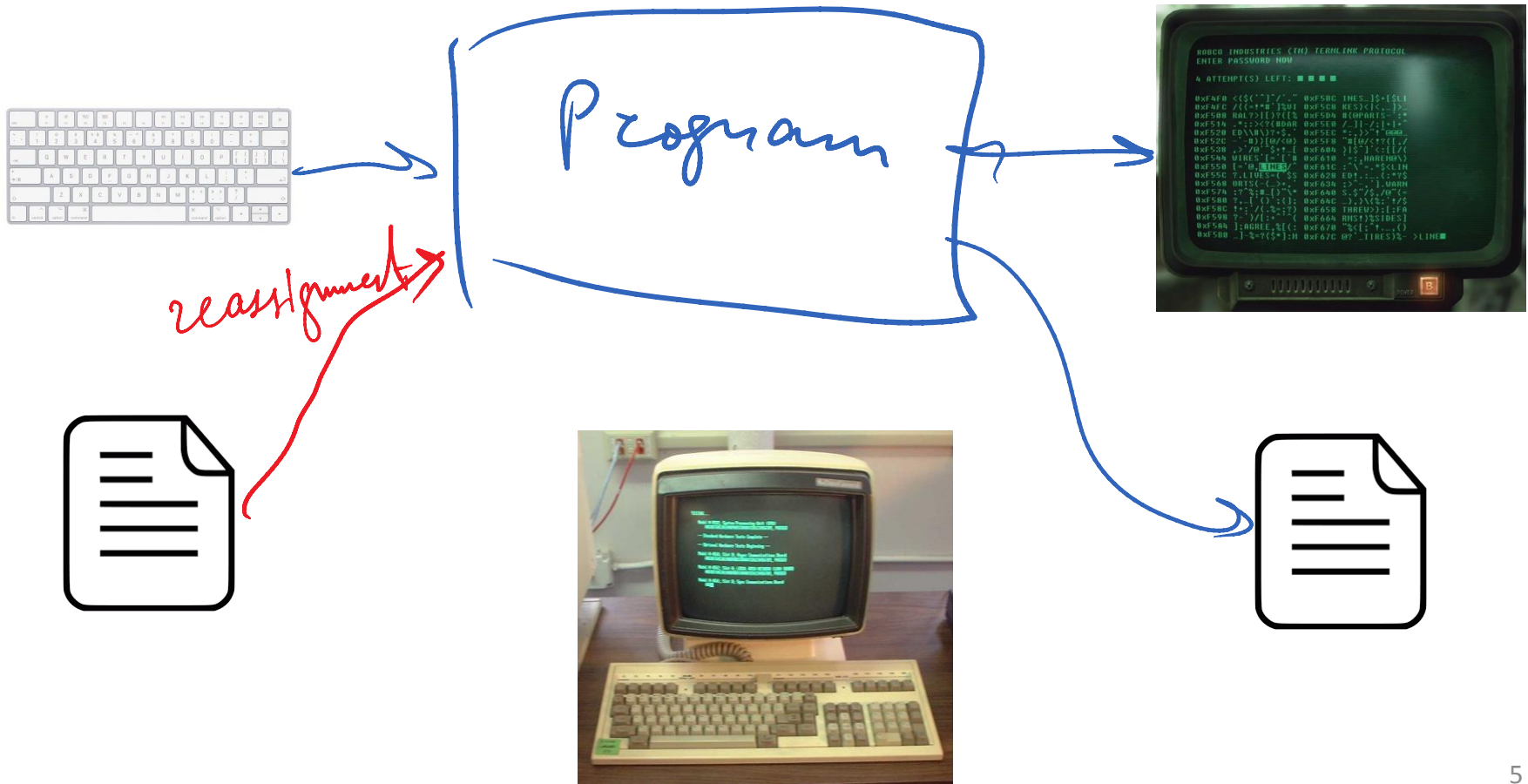
Recommended Contest Project Structure

~/se/progr/



INTRODUCTION TO I/O AND STREAMS

IO Model of a Program



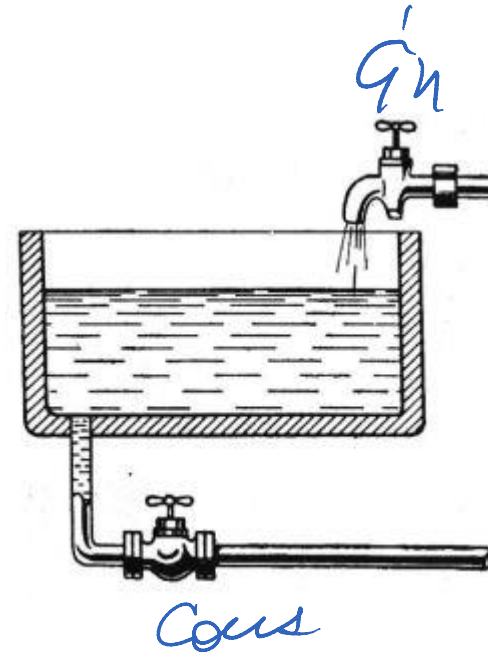
Streams

- A *stream* is an abstraction for input/output.
- Standard input and standard output are represented by global objects `std::cin` (of a type `std::istream`) and `std::cout` (of a type `std::ostream`) respectively.
- One can think of a stream as a source (input) or a destination (output) of characters of an indefinite length.
- Input data to a program from a stream:

```
int x;  
double y;  
std::cin >> x >> y;
```

- Output data from a program to a stream:

```
std::string userName = "Sergey";  
std::cout << "Hello, " << userName;
```



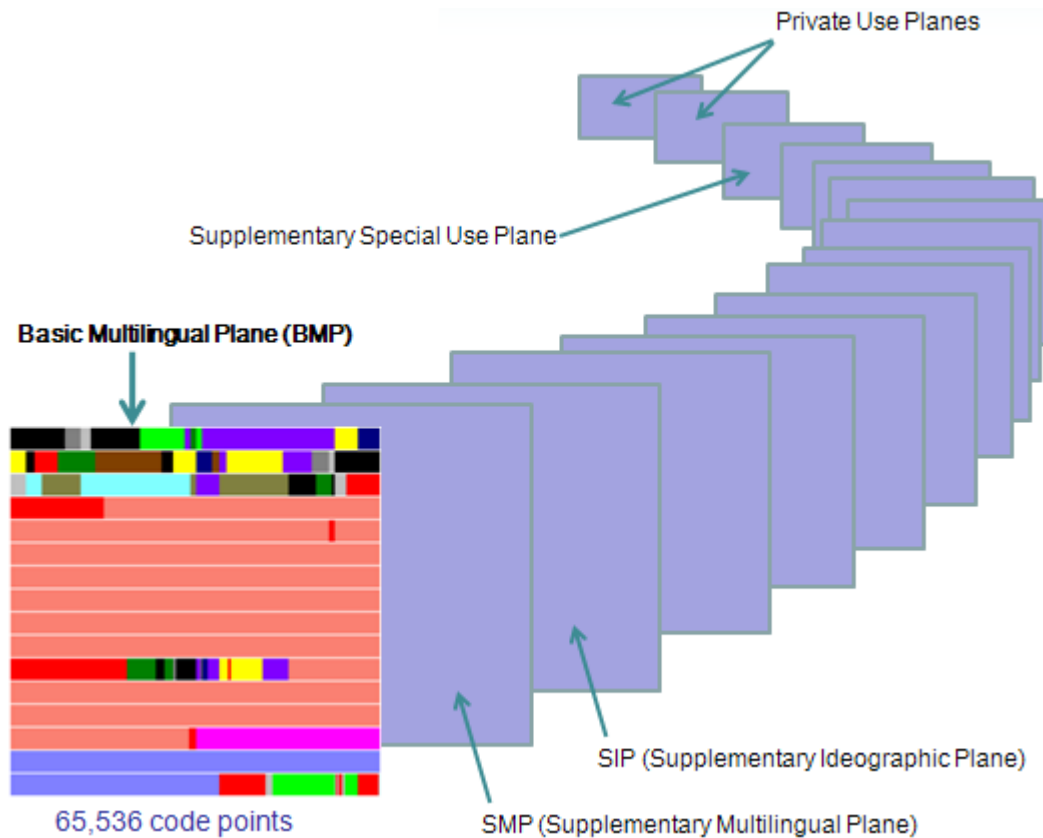
STRINGS AND ENCODING SCHEMES

Representation of Characters

- We have a set of characters we want to use in a computer:
 - A, B, C, D, ..., Z
 - a, b, c, ... z
 - 0, 1, 2, ..., 9
 - !, @, #, \$, %, ^, &, *, (,), -, _, =, +, , ; , ...
 - let's enumerate them!
- What about national alphabets?
 - А, Б, В, ... Я, а, б, в, ... я
 - Α, Β, Γ, Δ, Ε, ..., Ω, α, β, γ, ... ξ, ω
 - ...

Unicode Standard

- Codes almost all known symbols using quite big numbers (millions)
- Structured by using planes and blocks
- Can be implemented by different character encodings



Basic Multilingual Plane

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 script
- Non-Latin European scripts
- African scripts
- Middle Eastern and Southwest Asian scripts
- South and Central Asian scripts
- Southeast Asian scripts
- East Asian scripts
- CJK characters
- Indonesian and Oceanic scripts
- American scripts
- Notational systems
- Symbols
- Private use
- UTF-16 surrogates
- Unallocated code points

As of Unicode 11.0

8-bit code pages

- Each character is coded by a number from 0..255
- Too small to fit even the most basic characters
- Each alphabet has its own coding standard:
 - ASCII — basic set of 127 characters inc. Latin
 - CP866 — obsolete “DOS” Cyrillic codepage
 - CP1251 (ANSI) — Windows-based Cyrillic codepage
 - KOI8-R — obsolete Linux-based Cyrillic codepage
 - ...

ASCII Table

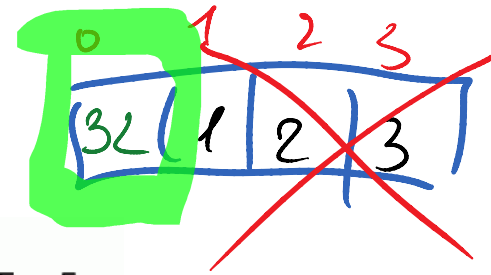
ASCII Code Chart

	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

8 | 80 | 81 | 82 |

- uses only 128 positions (... bits?)
- contains Latin symbols, numbers, basic symbols
- see [Ya.Contest 1 Problems M and N](#)

CP1251 Codepage



Codepage 1251 - Cyrillic Windows

	-0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-A	-B	-C	-D	-E	-F	
0-		0001	0002	0003	0004	0005	0006	0007	0008	0009	000A	000B	000C	000D	000E	000F	
1-		0010	0011	0012	0013	0014	0015	0016	0017	0018	0019	001A	001B	001C	001D	001E	001F
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-		~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~
8-		Ђ	Ѓ	Ѕ	Ї	Љ	Њ	Ћ	Ќ	Ў	Ѕ	Ї	Љ	Њ	Ћ	Ќ	Ў
9-		ђ	ѓ	ѕ	ї	љ	њ	ќ	ќ	ў	ѕ	ї	љ	њ	ќ	ќ	ў
A-		Ў	ў	Ј	Ѡ	Ґ	ґ	Є	©	€	«	¬	-	®	Ї		
B-		°	±	І	і	г	р	¶	·	ё	№	€	»	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

Demo App: Output Characters of the Current Codepage

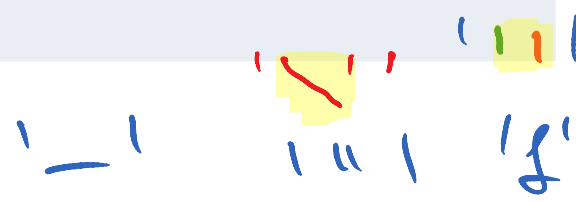


```
char ch = (char)i;  
std::cout << ch << " (" << i << ") \t";
```

Types for Representing Characters

Type	Size	Range	Codepages (examples)
<code>char</code>	8 bit	0..255	Fixed width: ASCII (ext), CP1251, KOI8-r Variable width: UTF-8
<code>char16_t</code>	16 bit	0..65535	Fixed width: UCS-2 “Variable width” (surrogate pairs): UTF-16
<code>char32_t</code>	32 bit	0.. 4294967295	UTF-32
<code>wchar_t</code>	16 or 32 bit		UCS-2

- UCS-2 — 2 bytes per symbol, fixed size
- UTF-16 — 2 or 4 bytes (surrogate pairs) per symbol, variable size
- UTF-32 — 4 bytes per symbol, fixed size
- UTF-8 — 1 to 4 bytes per symbol, fixed size



```
char ch = 'A';
```

UTF-8

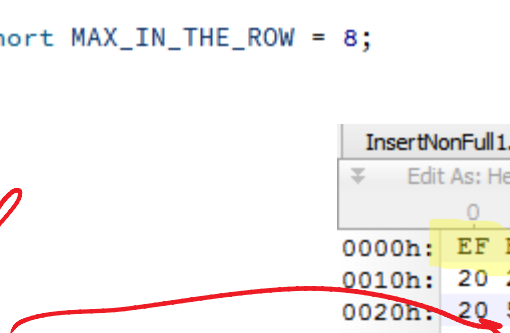
- 1 to 4 bytes per symbol:
 - 1 byte codes ASCII symbols (0..127)
 - 2 bytes code 1 920 symbols including: Latin-script alphabets, and also Greek, Cyrillic, Coptic, Armenian, Hebrew, Arabic, ..., *Combining Diacritical Marks* (é, è, ê, ā, á, ǎ, ï, ü, ë, й, ...)
 - 3 bytes code the rest of the Basic Multilingual Plane including Chinese, Japanese, Korean
 - 4 bytes for other symbols such as math symbols, emoji etc.

Number of bytes	Bits for code point	First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
1	7	U+0000	U+007F	0xxxxxxx			
2	11	U+0080	U+07FF	110xxxxx	10xxxxxx		
3	16	U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	
4	21	U+10000	U+10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx


```
< > ex_1.cpp <Select Symbol>
1  /*! \file      ex_1.cpp
2  *  \author    Sergey Shershakov
3  *  \version   0.1
4  *  \date     15.01.2019
5  *
6  *  Prints printable characters of the current codepage.
7  *  Печатает печатное.
8  */
9
10 #include <iostream>
11
12 int main()
13 {
14     const short MAX_IN_THE_ROW = 8;
```

BOM

std::endl



InsertNonFull1.xibt xiexample2.pyd p1.bin Startup ex_1.cpp																	
Edit As: Hex Run Script Run Template																	
	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	0123456789ABCDEF
0000h:	EF	BB	BF	2F	2A	21	20	5C	66	69	6C	65	20	20	20	20	i»¿/*! \file
0010h:	20	20	20	65	78	5F	31	2E	63	70	70	0D	0A	20	2A	20	ex_1.cpp.. *
0020h:	20	5C	61	75	74	68	6F	72	20	20	20	20	20	53	65	72	\author Ser
0030h:	67	65	79	20	53	68	65	72	73	68	61	6B	6F	76	0D	0A	gey Shershakov..
0040h:	20	2A	20	20	5C	76	65	72	73	69	6F	6E	20	20	20	20	* \version
0050h:	30	2E	31	0D	0A	20	2A	20	20	5C	64	61	74	65	20	20	0.1.. * \date
0060h:	20	20	20	20	20	31	35	2E	30	31	2E	32	30	31	39	0D	15.01.2019.
0070h:	0A	20	2A	0D	0A	20	2A	20	20	50	72	69	6E	74	73	20	. *.. * Prints
0080h:	70	72	69	6E	74	61	62	6C	65	20	63	68	61	72	61	63	printable charac
0090h:	74	65	72	73	20	6F	66	20	74	68	65	20	63	75	72	72	ters of the curr
00A0h:	65	6E	74	20	63	6F	64	65	70	61	67	65	2E	0D	0A	20	ent codepage...
00B0h:	2A	20	20	D0	9F	D0	B5	D1	87	D0	B0	D1	82	D0	B0	D0	* ÐŸÐµÑ+Ð°Ñ,Ð°Ð
00C0h:	B5	D1	82	20	D0	BF	D0	B5	D1	87	D0	B0	D1	82	D0	BD	µÑ, Ð¿ÐµÑ+Ð°Ñ,Ð
00D0h:	D0	BE	D0	B5	2E	0D	0A	20	2A	2F	0D	0A	0D	0A	23	69	Ð%Ðµ... */....#i
00E0h:	6E	63	6C	75	64	65	20	3C	69	6F	73	74	72	65	61	6D	nclude <iostream
00F0h:	3E	0D	0A	0D	0A	69	6E	74	20	6D	61	69	6E	28	29	0D	>....int main().
0100h:	0A	7B	0D	0A	20	20	20	20	63	6F	6E	73	74	20	73	68	{.. const sh
0110h:	6F	72	74	20	4D	41	58	5F	49	4E	5F	54	48	45	5F	52	ort MAX_IN_THE_R
0120h:	4F	57	20	3D	20	38	3B	0D	0A	0D	0A	0D	0A	20	20	20	OW = 8;.....

Introducing Arrays

- *Array* (plain C-style) is a set of objects of the same type placed in adjacent memory cells.
- `int nameOfArray[5];`
- `int nameOfArray[5] = {9, 10, 30, 40, 42};`

Array Initialization

- `int x[4] = {3, 6, 8, 10}; // okay`
- `int y[4]; // okay`
- `y[4] = {5, 6, 7, 9}; // not allowed`
- `y = cards; // not allowed`

- `double points[5] = {3.0, 1.2}; // fewer is ok`
- `short stars[] = {1, 5, 3, 8}; // is ok, compiler counts
// the number of elements for you`

x[0]

- How to find out the number of items?

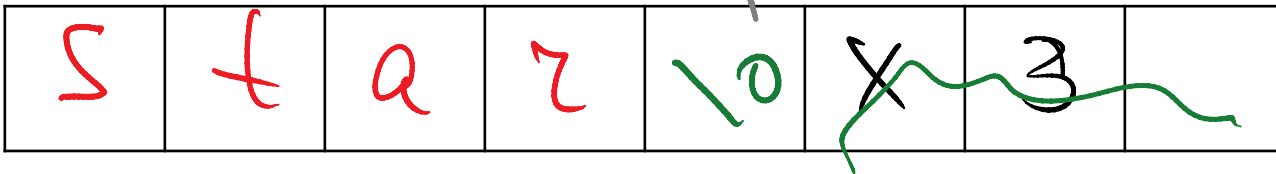
```
short things[] = {1, 7, 2, 10};  
int e1Num = sizeof(things) / sizeof(short);
```

= 4

Is a String an Array of Symbols?

- `char a[8] = {'b', 'e', 'a', 'u', 'x', ' ', 'I', 'I'}; // not a string!`
- `char s[8] = {'f', 'a', 't', 'e', 's', 's', 'a', '\0'}; // a string!`

```
char str1[8] = "Star";
```



Class `std::string`

- `std::string` → `std::basic_string<char>`

8-bit

- represents strings with 8-bit characters:
 - fixed-width: *ASCII, CP1251, CP866,...*
 - variable-width: *UTF-8*


```
std::string s = "I'm a string!";
```

- has an internal buffer to store a continuous sequence of characters
 - `s.c_str()` returns a pointer to the buffer
- has a number of useful methods and operations:
 - `s.size()`, `s.length()` for a string's length;
 - `s.clear()` clears a string;
 - `s.empty()` test if a string is empty;
 - `s[10]` gives an individual (11th) symbol in the sequence;
 - `s + "Abs"` — concatenation, gives a new string.

Reading String Data from a Stream

- Reading strings word by word:

```
std::string s1, s2;  
std::cin >> s1 >> s2; // "Abc Xyz"
```



- Reading a whole line containing spaces (see [Ya.Contest 1 Problems E](#)):

```
std::string bigStr;  
std::getline(cin, bigStr); // "Hello World!"
```



- Reading a text character by character (see [Ya.Contest 1 Problems H](#)):

```
char ch;  
std::cin >> std::noskipws;  
std::cin >> ch;  
std::cin >> ch;  
...  
or  
char ch2 = std::cin.get();  
std::cin.get(ch2)
```

