

Character Encodings and C Programming

The background of the slide is a photograph of a statue in a park, overlaid with a semi-transparent red filter. The statue is a large, classical-style figure, possibly representing a person of historical significance. The scene is set outdoors with trees and foliage visible in the background.

CS 240 - The University of Illinois

Wade Fagen-Ulmschneider

January 20, 2022

Character Encodings

Representing numbers is great -- but what about words? Can we make sentences with binary data?

ASCII



Organization

To begin to create words:

- A letter is _____ binary bits.
_____ hex digits!

(We call this unit a _____.)

Organization

Global standard called the **American Standard Code for Information Interchange (ASCII)** is a _____ for translating numbers to characters.

ASCII

Row	Bit Pattern b7 b6 b5 b4 b3 b2 b1	Column							
		0	1	2	3	4	5	6	7
0	0 0 0 0	NUL	DLE	SP	0	@	P	~	p
1	0 0 0 1	SOH	DC1	!	1	A	Q	a	q
2	0 0 1 0	STX	DC2	"	2	B	R	b	r
3	0 0 1 1	LTX	DC3	#	3	C	S	c	s
4	0 1 0 0	EOT	DC4	\$	4	D	T	d	t
5	0 1 0 1	ENO	NAK	%	5	E	U	e	u
6	0 1 1 0	ACK	SYN	&	6	F	V	f	v
7	0 1 1 1	BEL	ETB	'	7	G	W	g	w
8	1 0 0 0	BS	CAN	(8	H	X	h	x
9	1 0 0 1	HT	EM)	9	I	Y	i	y
10	1 0 1 0	LF	SUB	*	:	J	Z	j	z
11	1 0 1 1	VT	ESC	+	;	K	[k	{
12	1 1 0 0	FF	FS	,	<	L	\	l	!
13	1 1 0 1	CR	GS	-	=	M]	m	}
14	1 1 1 0	SO	RS	.	>	N	^	n	~
15	1 1 1 1	SI	US	/	?	o	_	o	DEL

¹ Change of name

² New character

³ Moved character

Fig. 14.12 ASCII, 1967 and 1968



ASCII

Row	Bit Pattern b7 b6 b5 b4 b3 b2 b1	Column							
		0	1	2	3	4	5	6	7
0	0 0 0 0	NUL	DLE	SP	0	@	P	~	p
1	0 0 0 1	SOH	DC1	!	1	A	Q	a	q
2	0 0 1 0	STX	DC2	"	2	B	R	b	r
3	0 0 1 1	LTX	DC3	#	3	C	S	c	s
4	0 1 0 0	EOT	DC4	\$	4	D	T	d	t
5	0 1 0 1	ENO	NAK	%	5	E	U	e	u
6	0 1 1 0	ACK	SYN	&	6	F	V	f	v
7	0 1 1 1	BEL	ETB	'	7	G	W	g	w
8	1 0 0 0	BS	CAN	(8	H	X	h	x
9	1 0 0 1	HT	EM)	9	I	Y	i	y
10	1 0 1 0	LF	SUB	*	:	J	Z	j	z
11	1 0 1 1	VT	ESC	+	;	K	[k	{
12	1 1 0 0	FF	FS	,	<	L	\	l	!
13	1 1 0 1	CR	GS	-	=	M]	m	}
14	1 1 1 0	SO	RS	.	>	N	^	n	~
15	1 1 1 1	SI	US	/	?	o	_	o	DEL

¹ Change of name

² New character

³ Moved character

Fig. 14.12 ASCII, 1967 and 1968

		Column	0	1	2	3	4	5	6	7
		Bit Pattern	b7	b6	b5					
			0	0	0	0	1	1	1	1
			0	0	1	1	0	0	1	1
			0	1	0	1	0	1	0	1
Row		b4	b3	b2	b1					
0	0 0 0 0	NUL	DLE	SP	0	@	P	~	P	
1	0 0 0 1	SOH	DC1	!	1	A	Q	a	q	
2	0 0 1 0	STX	DC2	"	2	B	R	b	r	
3	0 0 1 1	LTX	DC3	#	3	C	S	c	s	



0b 0100 0001 = 0x41 = A

0b 0100 0001 = 0x41 = A

0b 0100 00**10** = 0x4**2** = **B**

0b 0100 0001 = 0x41 = A

0b 0100 0010 = 0x42 = B

= 0x43 =

= 0x44 =

0b 0100 0001 = 0x41 = A

0b 01**1**0 0001 = 0x**61** = **a**

Shortcomings with ASCII



Other Character Encodings

Character Encodings

There are many other character encodings beyond ASCII.

Character Encodings

One of the most common is the **Unicode Transformation Format (8-bit)**, commonly called:

ISO/IEC 10646



L2/10-038

**ISO/IEC International
Standard
ISO/IEC 10646**

Final Committee Draft

**Information technology –
Universal Coded
Character Set (UCS)**

*Technologie de l'information – Jeu
universel de caractères codés (JUC)*

Second edition, 2010



Length	Byte #1	Byte #2	Byte #3	Byte #4
1-byte:	0 _ _ _ _			
2-bytes:	110 _ _ _ _	10 _ _ _ _		
3-bytes:	1110 _ _ _ _	10 _ _ _ _	10 _ _ _ _	
4-bytes:	1111 0 _ _ _	10 _ _ _ _	10 _ _ _ _	10 _ _ _ _

Characters in UTF-8

a

Characters in UTF-8

€

U+03b5

0100 1000 0110 1001 1111 0000

1001 1111 1000 1110 1000 1001

0100 1000 | 0110 1001 | 1111 0000 |

1001 1111 | 1000 1110 | 1000 1001

Length	Byte #1	Byte #2	Byte #3	Byte #4
1-byte:	0 _ _ _ _			
2-bytes:	110 _ _ _ _	10 _ _ _ _		
3-bytes:	1110 _ _ _ _	10 _ _ _ _	10 _ _ _ _	
4-bytes:	1111 0 _ _ _	10 _ _ _ _	10 _ _ _ _	10 _ _ _ _

0100 1000 | 0110 1001 | 1111 0000 |

1001 1111 | 1000 1110 | 1000 1001

0100 1000 | 0110 1001 | 1111 0000 |
1001 1111 | 1000 1110 | 1000 1001

Hi



Programming in C

The image features a semi-transparent orange overlay on a photograph of a crowd of people gathered around a statue. The statue is the Alma Mater of the University of California, Berkeley, depicting a woman in classical robes. The text 'Programming in C' is centered in white, bold, sans-serif font.

You already know C++!

You already know C++!

Programming in C **is a simplification of C++.**

1. Program Starting Point:

2. Printing to `stdout`

```
1 #include <stdio.h>
2
3 int main() {
4     int i = 42;
5     char *s = "Hello, world!";
6     float f = 3.14;
7
8     printf("%d %s %f\n", i, s, f);
9     printf("%d\n", s[0]);
10    printf("%d\n", s);
11    printf("%d\n", f);
12
13    return 0;
14 }
```

3. Pointers

4. Heap Memory Allocation

```
1 #include <stdlib.h>
2 #include <stdio.h>
3
4 int main() {
5     char *s = malloc(10);
6     int *num = malloc( sizeof(int) );
7
8     printf("%p %p\n", s, num);
9     return 0;
10 }
```


5. Strings

5. Strings

There is no “data type” in C known as a string. Instead, we refer to “C Strings” as a sequence of characters:

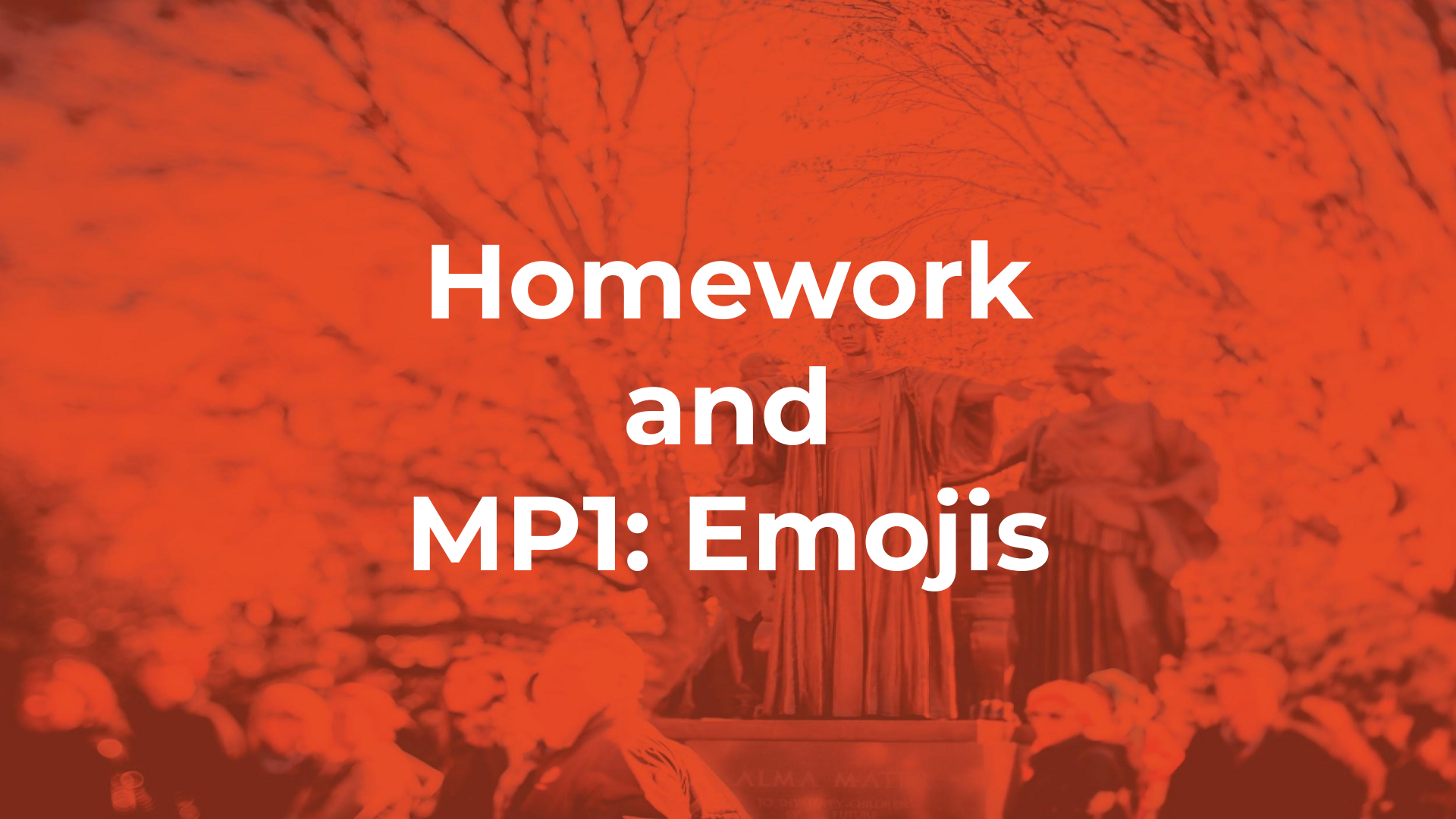
- A “C string” is just a character pointer: _____.
- The string continues until it reaches a _____ byte.

```
6 char *s = malloc(6);
7 strcpy(s, "cs240");
8 printf("s[0]: 0x%x == %d == %c\n", s[0], s[0], s[0]);
9 printf("s[4]: 0x%x == %d == %c\n", s[4], s[4], s[4]);
10 printf("s[5]: 0x%x == %d == %c\n", s[5], s[5], s[5]);
11 printf("s == \"%s\", strlen(s): %ld\n\n", s, strlen(s));
12
13 char *s2 = s + 2;
14 printf("s2[0]: 0x%x == %d == %c\n", s2[0], s2[0], s2[0]);
15 printf("s2 == \"%s\", strlen(s2): %ld\n\n", s2, strlen(s2));
16
17 *s2 = 0;
18 printf("s2[0]: 0x%x == %d == %c\n", s2[0], s2[0], s2[0]);
19 printf("s2 == \"%s\", strlen(s2): %ld\n\n", s2, strlen(s2));
20
21 printf("s == \"%s\", strlen(s): %ld\n", s, strlen(s));
```

```
1 #include <stdio.h>
2 #include <string.h>
3 #include <stdlib.h>
4
5 int main() {
6     char *s = malloc(5);
7     s[0]=0xF0; s[1]=0x9F; s[2]=0x8E; s[3]=0x89; s[4]=0x00;
8
9     char *s1 = "\xF0\x9F\x8E\x89";
10    char *s2 = "👉";
11    char *s3 = "\U0001f389"; // \U - must be 8 bytes
12
13    printf("%s %s %s %s\n", s, s1, s2, s3);
14    printf("strlen(): %ld %ld %ld %ld\n", strlen(s),
15           strlen(s1), strlen(s2), strlen(s3));
16 }
```

Some extremely useful built in string functions:

- `strcmp(char *s1, char *s2)` -- Compares two strings
- `strcat(char *dest, char *src)` -- Concatenate two strings
- `strcpy(char *dest, char *src)` -- Copies a string
- `strlen(char *s)` -- Returns the length of the string

A photograph of a crowd of people gathered around a statue of Alma Mater, overlaid with a semi-transparent orange filter. The statue is the central focus, with people in the foreground and background. The text is centered over the image.

Homework and MP1: Emojis