

C and Linux Programming  
Eastern Washington University  
Computer Science  
March 30th – June 12th, 2020



## Lecture 3

### C String Processing and Files

---

Joe Dumoulin

Mar 30, 2020

Eastern Washington University

# C String Processing and Files

---

# The Story So Far

## Basic C syntax

- Variable Types
- Statements
- Functions

Next we will use these concepts to see how to work with characters and functions in C.

# Basic Character Representations - ASCII

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b> (null)	32	20	040	&#32;	<b>Space</b>	64	40	100	&#64;	<b>@</b>	96	60	140	&#96;	<b>`</b>
1	1	001	<b>SOH</b> (start of heading)	33	21	041	&#33;	<b>!</b>	65	41	101	&#65;	<b>A</b>	97	61	141	&#97;	<b>a</b>
2	2	002	<b>STX</b> (start of text)	34	22	042	&#34;	<b>"</b>	66	42	102	&#66;	<b>B</b>	98	62	142	&#98;	<b>b</b>
3	3	003	<b>ETX</b> (end of text)	35	23	043	&#35;	<b>#</b>	67	43	103	&#67;	<b>C</b>	99	63	143	&#99;	<b>c</b>
4	4	004	<b>EOT</b> (end of transmission)	36	24	044	&#36;	<b>\$</b>	68	44	104	&#68;	<b>D</b>	100	64	144	&#100;	<b>d</b>
5	5	005	<b>ENQ</b> (enquiry)	37	25	045	&#37;	<b>%</b>	69	45	105	&#69;	<b>E</b>	101	65	145	&#101;	<b>e</b>
6	6	006	<b>ACK</b> (acknowledge)	38	26	046	&#38;	<b>&amp;</b>	70	46	106	&#70;	<b>F</b>	102	66	146	&#102;	<b>f</b>
7	7	007	<b>BEL</b> (bell)	39	27	047	&#39;	<b>'</b>	71	47	107	&#71;	<b>G</b>	103	67	147	&#103;	<b>g</b>
8	8	010	<b>BS</b> (backspace)	40	28	050	&#40;	<b>(</b>	72	48	110	&#72;	<b>H</b>	104	68	150	&#104;	<b>h</b>
9	9	011	<b>TAB</b> (horizontal tab)	41	29	051	&#41;	<b>)</b>	73	49	111	&#73;	<b>I</b>	105	69	151	&#105;	<b>i</b>
10	A	012	<b>LF</b> (NL line feed, new line)	42	2A	052	&#42;	<b>*</b>	74	4A	112	&#74;	<b>J</b>	106	6A	152	&#106;	<b>j</b>
11	B	013	<b>VT</b> (vertical tab)	43	2B	053	&#43;	<b>+</b>	75	4B	113	&#75;	<b>K</b>	107	6B	153	&#107;	<b>k</b>
12	C	014	<b>FF</b> (NP form feed, new page)	44	2C	054	&#44;	<b>,</b>	76	4C	114	&#76;	<b>L</b>	108	6C	154	&#108;	<b>l</b>
13	D	015	<b>CR</b> (carriage return)	45	2D	055	&#45;	<b>-</b>	77	4D	115	&#77;	<b>M</b>	109	6D	155	&#109;	<b>m</b>
14	E	016	<b>SO</b> (shift out)	46	2E	056	&#46;	<b>.</b>	78	4E	116	&#78;	<b>N</b>	110	6E	156	&#110;	<b>n</b>
15	F	017	<b>SI</b> (shift in)	47	2F	057	&#47;	<b>/</b>	79	4F	117	&#79;	<b>O</b>	111	6F	157	&#111;	<b>o</b>
16	10	020	<b>DLE</b> (data link escape)	48	30	060	&#48;	<b>0</b>	80	50	120	&#80;	<b>P</b>	112	70	160	&#112;	<b>p</b>
17	11	021	<b>DC1</b> (device control 1)	49	31	061	&#49;	<b>1</b>	81	51	121	&#81;	<b>Q</b>	113	71	161	&#113;	<b>q</b>
18	12	022	<b>DC2</b> (device control 2)	50	32	062	&#50;	<b>2</b>	82	52	122	&#82;	<b>R</b>	114	72	162	&#114;	<b>r</b>
19	13	023	<b>DC3</b> (device control 3)	51	33	063	&#51;	<b>3</b>	83	53	123	&#83;	<b>S</b>	115	73	163	&#115;	<b>s</b>
20	14	024	<b>DC4</b> (device control 4)	52	34	064	&#52;	<b>4</b>	84	54	124	&#84;	<b>T</b>	116	74	164	&#116;	<b>t</b>
21	15	025	<b>NAK</b> (negative acknowledge)	53	35	065	&#53;	<b>5</b>	85	55	125	&#85;	<b>U</b>	117	75	165	&#117;	<b>u</b>
22	16	026	<b>SYN</b> (synchronous idle)	54	36	066	&#54;	<b>6</b>	86	56	126	&#86;	<b>V</b>	118	76	166	&#118;	<b>v</b>
23	17	027	<b>ETB</b> (end of trans. block)	55	37	067	&#55;	<b>7</b>	87	57	127	&#87;	<b>W</b>	119	77	167	&#119;	<b>w</b>
24	18	030	<b>CAN</b> (cancel)	56	38	070	&#56;	<b>8</b>	88	58	130	&#88;	<b>X</b>	120	78	170	&#120;	<b>x</b>
25	19	031	<b>EM</b> (end of medium)	57	39	071	&#57;	<b>9</b>	89	59	131	&#89;	<b>Y</b>	121	79	171	&#121;	<b>y</b>
26	1A	032	<b>SUB</b> (substitute)	58	3A	072	&#58;	<b>:</b>	90	5A	132	&#90;	<b>Z</b>	122	7A	172	&#122;	<b>z</b>
27	1B	033	<b>ESC</b> (escape)	59	3B	073	&#59;	<b>;</b>	91	5B	133	&#91;	<b>[</b>	123	7B	173	&#123;	<b>{</b>
28	1C	034	<b>FS</b> (file separator)	60	3C	074	&#60;	<b>&lt;</b>	92	5C	134	&#92;	<b>\</b>	124	7C	174	&#124;	<b> </b>
29	1D	035	<b>GS</b> (group separator)	61	3D	075	&#61;	<b>=</b>	93	5D	135	&#93;	<b>]</b>	125	7D	175	&#125;	<b>}</b>
30	1E	036	<b>RS</b> (record separator)	62	3E	076	&#62;	<b>&gt;</b>	94	5E	136	&#94;	<b>^</b>	126	7E	176	&#126;	<b>~</b>
31	1F	037	<b>US</b> (unit separator)	63	3F	077	&#63;	<b>?</b>	95	5F	137	&#95;	<b>_</b>	127	7F	177	&#127;	<b>DEL</b>

Source: [www.LookupTables.com](http://www.LookupTables.com)

# What is ASCII?

ASCII stands for "American Standard Code for Information Interchange"

ASCII encodes a symbol into a number and a number to a symbol. The ASCII table is a codebook for changing numbers into characters.

The first 32 ASCII characters are called **Control Characters**. They don't get printed. Instead they cause actions to occur in the terminal.

the characters from 32-127 are characters used to print words, numbers, punctuation and white space for writing english.

**ASCII is insufficient for most languages.** Modern languages and libraries rely on **Unicode**.

# A Brief Interlude About Unicode

Modern applications must support more than english. For this we use **Unicode**.

Unicode uses a 32-bit **code point** to represent each possible character in any language (including made-up languages like elvish and klingon).

東京

Москва - столица России

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

We will not work with unicode, but we should recognize that for modern programming it is very important.

# Categories of ASCII codes

## Categories of ASCII Codes that we Care About

Numbers	ASCII Code
48-57	Digit symbols ('0' - '9')
65-90	Upper-case letters ('A' - 'Z')
97-122	Lower-case letters ('a' - 'z')
33-47 58-64 133-140 123-126	Punctuation
10-13, 32	Whitespace

We can write code to recognize these characters.

# Finding Character Types in a String

The following function takes a character as a parameter and tests if the character is an upper-case letter. It returns 1 if the character is an upper-case letter and zero if it is anything else.

```
int isupperalphabetic(char c)
{
    if (c >= 65 && c <= 90) {
        return 1;
    }
    // else
    return 0;
}
```



# Finding Character Types in a String

We can make the function easier to understand by adding some constants for 1 and 0. This makes it possible for us to talk about the true or false nature of the function. We can also better understand the function by using the ASCII character symbols instead of numbers for comparison.

```
#define TRUE 1
```

```
#define FALSE 0
```

```
int isupperalphabetic(char c)
{
    if (c >= 'A' && c <= 'Z') {
        return TRUE;
    }
    // else
    return FALSE;
}
```

# Predicate Functions

Functions and statements that evaluate to True or False are called **Predicates**. Predicates can be used to make decisions. Here is another predicate function that identifies a lower-case letter if it is found:

```
int isloweralphabetic(char c)
{
    if (c >= 'a' && c <= 'z') {
        return TRUE;
    }
    // else
    return FALSE;
}
```

# Composing Functions

Functions can be combined to make more complex functions. Here's an example of a predicate function created from combining the two previous predicates.

```
int isalphabetic(char c)
{
    if (isupperalphabetic(c) || isloweralphabetic(c)) {
        return TRUE;
    }
    // else
    return FALSE;
}
```

What does this function do?

# Counting Characters in a String

Another type of function we can make is one that counts the number of characters in a string. We can use the fact that c strings are terminated with '\0'. We can find the length of a string by counting all the characters in the string until we find the terminating character.

```
int strlen(char* s)
{
    int x = 0;
    while(s[x] != '\0') {
        x++;
    }
    return x;
}
```

# Writing a Program With Functions

We can use these functions to write a program that prints each letter in the string.

```
// include the functions defined above ...
int main(void)
{
    char testStr[] = "When shall we three meet again \
In thunder, lightning, or in rain?";
    printf("string = %s\n", testStr);
    int i;
    int len = stringlength(testStr);
    printf("%d\n", len);

    for (i = 0; i < len; i++){
        char c = testStr[i];
        if (isalphabetic(c)) {
            printf("LETTER: %c\n", c);
        }
    }
    return EXIT_SUCCESS;
}
```

Question: What is the output of this program?

## Another function: Lower-case a Character

Each upper case and lower case letter in the ASCII table is separated by a constant. We can calculate that constant by subtracting two of the codes for a letter like this:

'a' - 'A'

The code for this is:

```
// if a character is upper case, then change it to lower case
char lower(char c)
{
    if (isupperalphabetic(c)) {
        return c + ('a' - 'A');
    }
    // else
    return c;
}
```

# Counting Letters in a String

We can use the functions we have developed so far to do many things. Let's do something a little more complicated. we will write a program that calculates the number of each letter that appears in a string.

The main function for this task is longer than we have seen so far.

We will need to:

- Define an array of 26 numbers (one for each letter) and initialize it
- Define a string to process
- For each character *c* in the string,
  - If *c* is an upper case letter, change it to lower case
  - If *c* is a letter, add one to that letter's entry in the array
  - Ignore all other characters
- print the contents of the array

# 1. Define and initialize an array

```
int main(void)
{
    int char_counts[26];  // one array slot for each letter

    int i;    // an iterator
    // initialize counts to zero.
    for (i = 0; i < 26; i++) {
        char_counts[i] = 0;
    }
    ....
}
```



## 2. Define the string to process

Define the string, print it, and then print its length.

```
char testStr[] = "Accepting the absurdity of everything \
around us is one step, a necessary experience: it should not \
become a dead end. It arouses a revolt";
printf("string = %s\n", testStr);

int len = stringlength(testStr);
printf("%d\n", len);
```

### 3. For each character in the string

```
...  
    // count characters in the string  
    for (i = 0; i < len; i++){  
        char c = testStr[i];  
        // check for uppercase. if found then lowercase  
        if (isupperalphabetic(c)) {  
            c = lower(c);  
        }  
  
        // if c is a character, count it.  
        if (isloweralphabetic(c)) {  
            int j = c - 'a';  
            // char_counts[j] = char_counts[j] + 1;  
            char_counts[j]++; // The same as above  
        }  
    }  
}
```

...

## 4. Print the array

For each element in the array, we print the character for that array element and the count.

...

```
// print each character followed by its count
for (i = 0; i < 26; i++) {
    printf("character: %c appears %d times.\n", i + 'a', char_co
}
return EXIT_SUCCESS;
}
```

# Assignment 2

We have used the C Programming Language concepts to build some programs that can process ASCII strings. We can use these functions to do other things as well. For example:

- Write a function to recognize if a character is a digit.
- Write a function to recognize if a character is whitespace.
- Write a function that can recognize a number (a consecutive set of digits).
- Write a function to recognize words (consecutive letters between spaces and punctuation).