# Computer Science I

## Strings

Dr. Chris Bourke
cbourke@cse.unl.edu

---

## Outline

1. Introduction
2. Manipulating Strings
3. String Processing
4. Data Processing
5. Exercises

---

# Part I: Introduction

---

## Strings

- A *string* is a sequence of *characters*
- ASCII, but more generally may be Unicode
- ASCII, CJK characters, emojis, etc.
- As of June 2017, Unicode 10.0: 136,755 characters
- Support for 1,112,064 characters
- We'll stick with ASCII

---

## Strings

- Languages represent strings differently
- In C: a string is an array of `char` elements
- Huge caveat: all strings in C must end (terminate) with a null character, `\0`
- Failure to ensure that all strings are *null-terminated* will result in undefined behavior; seg faults, bus errors, etc.
- ASCII value of 0, but it is *NOT*:
  - `'\0'` $\neq$ `'0'`
  - `'\0'` $\neq$ `'\n'`
  - `'\0'` $\neq$ `NULL`
  - it is the *null terminating character*

---

## Strings in C

- So far: single characters, `char`
- Character literals denoted with single quotes: `'A'`
- String literals denoted with double quotes: `"Hello World"`
- String variables:
  - Static strings: `char s[100];`
  - Static string and initialization: `char s[] = "Hello";`
  - Dynamic strings: `char *s;`
  - Constant Strings: `char *s = "Hello";` Don't do this

## Static Strings

- Static string: `char s[100];`
- Creates a string that can hold *up to 99* characters
- Room is needed for the null terminating character
- It may hold *shorter* strings, but it may not exceed 99 characters
- At declaration, *contents are undefined*
- May or may not contain the null terminating character

## Static Strings With Initialization

- Static string + initialization:
  `char s[] = "Hello";`
- Creates a character array of size 6
- Automatically includes the null-terminating character for you
- Contents can be changed, strings in C are *mutable*
- Static strings are allocated on the stack

## Dynamic Strings

- Dynamic string:
  `char *s = (char *) malloc(sizeof(char) * 100);`
- Creates a character array of size 100
- May only hold a string that can hold *up to 99* characters
- Room is needed for the null terminating character
- *Contents are undefined*; may or may not contain the null terminating character
- Allocated on the heap
- Most of your strings will be of this type

## Constant Strings

- Constant string declaration:
  `char *s = "Hello";`
- Dynamic, allocated on the heap, but
- Creates a *read-only, immutable* string
- Actually uses `const char *s = "Hello";`
- But the compiler generally doesn't catch it!
- Avoid unless you *really want* a dynamically allocated, immutable string for some reason

## Manipulating Strings by Character

- Strings are simply character arrays
- Each individual element can be *indexed* and
- a value can be assigned to it.
- Demonstration

## Manipulating Strings by Character

```
 1   char message[] = "hello World!"; //size is 13
 2
 3   message[0] = 'H';
 4   //bad:
 5   message[0] = "H";
 6   printf("message = %s\n", message);
 7
 8   //cut the string short:
 9   message[5] = '\0';
10   printf("message = %s\n", message);
11
12   //restore it: the rest of the contents were unchanged
13   message[5] = ' ';
14   printf("message = %s\n", message);
15
16   message[11] = '?';
17   printf("message = %s\n", message);
18
19   //bad:
20   message[12] = '!';
21   printf("message = %s\n", message);
22
23   //really bad:
24   message = "Goodbye World!";
```

# Part II: String Manipulation

## Copying Strings

- You cannot *assign* a string after it has been declared

```
1   char s1[] = "hello";
2   char *s2 = (char *) malloc(sizeof(char) * 6);
```

- Compiler error:
  ```
  s1 = "Hello";
  ```
- Memory leak:
  ```
  s2 = "World";
  ```
- A string is a character array which is a *memory address*

## Copying Strings

- You must *copy* the contents of one string into another
- String library: `string.h` provides a copy function and many others
- `char * strcpy(char * dest, const char * src);`
- Copies the contents of the *source* string into the *destination* string
- Assumes `src` is properly null-terminated
- It is *your* responsibility to ensure that `dest` is large enough to hold the contents
- Demonstration

## Copying Strings

```
1   char *name = (char *) malloc(sizeof(char) * 10);
2
3   strcpy(name, "Chris");
4   strcpy(name, "Bourke");
5   //invalid:
6   strcpy(name, "Chris Bourke");
```

## String Length

- Essential to know how many characters are stored in a string
- `size_t strlen(const char *s);`
- Result does *not* include the null terminating character!
- Demonstration

## String Length

```
1   char message[] = "Hello World!";
2   int n = strlen(message);
3   printf("n = %d\n", n);
4   message[5] = '\0';
5   n = strlen(message);
6   printf("n = %d\n", n);
7
8   char * stringCopy(const char *s) {
9     char *copy = (char *) malloc(sizeof(char) * (strlen(s) + 1));
10    strcpy(copy, s);
11    return copy;
12  }
```

## String Concatenation

- Copying overwrites a string's contents
- Alternative: append or *concatenation* the contents of one string onto the end of another
- `char *strcat(char *dest, const char *src);`
- Assumes both are null-terminated
- Assumes `dest` is large enough to hold both
- Demonstration

## String Concatenation

```
1   char *firstName = (char *) malloc((5+1) * sizeof(char));
2   char *lastName = (char *) malloc((6+1) * sizeof(char));
3   char *str = (char *) malloc(101) * sizeof(char));
4
5   strcpy(firstName, "Chris");
6   strcpy(lastName, "Bourke");
7
8   strcpy(str, lastName);
9   strcat(str, ", ");
10  strcat(str, firstName);
11  //str contains "Bourke, Chris"
```

## Byte-Limited Versions

- String library provides *byte-limited* versions of both copy and concatenation functions
- `char *strncat(char *dest, const char *src, size_t n);`
- `char *strncpy(char *dest, const char *src, size_t n);`
- Only copies *at most* first `n` bytes/characters
- Stops early if it sees `\0`
- Includes `\0` *only if* it is within the first `n` bytes!
- Demonstration

## Byte-Limited Versions

```
1   char fullName[] = "Christopher";
2   char *nickName = (char *) malloc(6 * sizeof(char));
3
4   strncpy(nickName, fullName, 5);
5   //don't forget:
6   firstName[5] = '\0';
```

# Part III: String Processing

## Iterating Over Strings

- Common to iterate over a string character-by-character
- Straightforward solution: for-loop using `strlen()`
- Generally better ways
- Demonstration

## Iterating Over Strings

```
1   //straightforward code:
2   for(int i=0; i<strlen(s); i++) {
3     printf("%c\n", s[i]);
4   }
5
6   //strlen works like:
7   int i = 0;
8   while(s[i] != '\0') {
9     i++;
10  }
11  //the value of i equals strlen(s)
12
13  //optimized:
14  int n = strlen(s);
15  for(int i=0; i<n; i++) {
16    printf("%c\n", s[i]);
17  }
18
19  //even better:
20  for(int i=0; s[i] != '\0'; i++) {
21    printf("%c\n", s[i]);
22  }
```

## ctype library

- ▶ May want to process or manipulate individual characters
- ▶ The `ctype.h` library provides many useful character functions
- ▶ Functions use ASCII `int` values
- ▶ Automatically type casted

## ctype library

- ▶ `int isdigit(int c)` – returns true if `c` is a digit character, `0` thru `9`
- ▶ `int islower(int c)` – returns true if `c` is a lowercase letter character, `a` thru `z`
- ▶ `int isupper(int c)` – returns true if `c` is an uppercase letter character, `A` thru `Z`
- ▶ `int isspace(int c)` – returns true if `c` is a whitespace character: space, tab, newline, etc.
- ▶ `int tolower(int c)` , `int toupper(int c)` – return the ASCII text value of the lowercase/uppercase version of `c`
- ▶ Demonstration: write a code snippet to count the number of spaces and the total number of whitespace characters.

## ctype library

```
1   char str[] = "Hello how \n\t are you today? \n\n";
2   int numSpaces = 0;
3   int numWhiteSpaces = 0;
4
5   for(int i=0; s[i] != '\0'; i++) {
6     if(isspace(s[i])) {
7       numWhiteSpaces++;
8     }
9     if(s[i] == ' ') {
10      numSpaces++;
11    }
12  }
13  printf("number of spaces: %d\n", numSpaces);
14  printf("total whitespace: %d\n", numWhiteSpaces);
```

## String Comparisons

- ▶ Often need to compare *entire strings* for equality
- ▶ You *cannot* use the equality operator!
- ▶ `s1 == s2` compares memory addresses!
- ▶ Need to use:
  `int strcmp(const char *str1, const char *str2)`
- ▶ Returns:
  - ▶ < 0 if `str1` comes before `str2`
  - ▶ 0 if contents are equal
  - ▶ > 0 if `str1` comes after `str2`
- ▶ Order is determined by ASCII text values
- ▶ Demonstration

## String Comparisons

```
1   int result;
2
3   result = strcmp("apple", "apple"); //0
4   result = strcmp("apple", "apples"); //negative
5   result = strcmp("apples", "apple"); //positive
6   result = strcmp("Apple", "apple"); //negative
7   result = strcmp("apples", "oranges"); //negative
8
9   result = strcmp("100", "99"); //negative!
10
11  result = strncmp("apple", "apples", 5); //zero
12
13  result = strcasecmp("ApPlE", "apple"); //zero
```

## Substrings

- It is possible to reference a *substring* of a string
- Reference a part of the string starting at a particular index
- Demonstration

## Substrings

```
1  char name[] = "Margaret Hamilton";
2  char *lastName = &name[9];
3  printf("Greetings, Ms. %s\n");
```

## Formatting Strings

- Already familiar with `%s` placeholder to print a string to the standard output
- `atoi` and `atof` convert strings to numbers
- Possible to convert numbers to strings
- "Print" to a string instead of the standard output
- `sprintf()` : print to a string
- Demonstration

## Formatting Strings

```
1  char s[100]; //buffer that is "big enough"
2  char state[] = "Nebraska";
3  int numCounties = 93;
4  double population = 1.92;
5  sprintf(s, "%s has %d counties and a population of %.2f million.\n", state,nu
```

## Arrays of Strings

- Arrays of strings are simply 2-D arrays of `char` elements
- Each "row" is a string that must be null-terminated
- Each row/string need not be the same size
- Easy extension of 2-D arrays: `char **arrayOfStrings`
- Demonstration

## Arrays of Strings

```
1  char **names = (char **) malloc(sizeof(char*) * 5);
2  names[0] = stringCopy("Margaret Hamilton");
3  names[1] = stringCopy("Grace Hopper");
4  names[2] = stringCopy("Alan Turing");
5  names[3] = stringCopy("Ada Lovelace");
6  names[4] = stringCopy("Dennis Ritchie");
7  for(int i=0; i<5; i++) {
8      printf("Famous Computer Scientist: %s\n", names[i]);
9  }
```

## Part IV: Data Processing

## String Tokenization

- Strings may contain formatted data: CSV, TSV
- *Tokenization* is the process of splitting a string along some *delimiter* and processing each *token* separately
- Example: `"Hedy,Lamarr,UNL,Avery Hall,Lincoln,NE"`
- Tokens:
  `"Hedy"` `"Lamarr"` `"UNL"` `"Avery"` `"Hall"` `"Lincoln"` `"NE"`
- Generally ignore the delimiter

## String Tokenization

- `char * strtok(char *str, const char *delim);`
- Tokenizes `str` along instances of `delim`
- Usage:
  - First time you call it: pass the string to be tokenized
  - Subsequent calls: pass `NULL` to continue with the same string
  - Returns a pointer to the next token
  - It *modifies your string!*
- Demonstration

## String Tokenization

```
1  char str[] = "Hedy,Lamarr,UNL,Avery Hall,Lincoln,NE";
2  char *token = NULL;
3  token = strtok(str, ",");
4  while(token != NULL) {
5    printf("token = %s\n", token);
6    token = strtok(NULL, ",");
7  }
```

## Part V: Exercises

- Write a string function to change a string's characters to uppercase letters
- Write a string function that returns a new copy of a string with all characters converted to uppercase
- Write a string function to "double space" a paragraph
- Write a "split"-style function: it takes a string and a delimiter and returns an array of strings of the tokens. Ensure no memory leaks!