

Programming Project 3B

EE312 Fall 2014

4 points, plus up to 1 point of additional credit

General: In this project, we'll improve upon the standard C method for representing strings. We're going to do a little work with malloc and free, and along the way, get some practice with pointers, structs and even some macros.

Honors: this project is identical to Project 3R, except for the additional requirement to write `utstrcmp`. Part of your grade will be based on whether your `utstrcmp` function is faster than the `strcmp` function in the standard C library.

The Big Picture: Strings in C are stored using an array of ASCII-encoded characters with a zero on the end. This is well and good, most of the time, but has some drawbacks. For one, in order to determine how long a string is, we have to start counting at the beginning, examining each character until we get to the zero. Even more significant, strings are a common source of buffer overflow errors in programs because programmers often forget how large the array actually is when filling it up with characters.

So, we're going to create an improved version of the string. Our string will be "backwards compatible" with C-style strings. However, we'll store some extra information along with every string. Specifically, in the four bytes immediately preceding the string we'll store the length (number of characters) in the string. In the four bytes immediately preceding the length, we'll store the size of the array (i.e., the capacity for the string, how many characters can it hold before buffer overflow would result). As a small amount of extra safety, in the four bytes immediately preceding the capacity, we'll store a not-so-random looking byte sequence that we can check to confirm that nothing bad has happened (yet). These three extra `uint32_t`s worth of information can be accessed using pointers (`uint32_t*` pointers specifically), or using a *String* struct that we define for you.

Your Mission: You are to implement each of the functions declared inside the `String.h` file. Hopefully I've remembered to describe each of those functions in this handout. But, the official list of what you need to do is contained in the `String.h` file.

You may use (but not modify) the `String` struct that is defined inside `String.h`. This struct is designed specifically for Visual Studio (although it may work on other compilers too), to align the 12 bytes immediately preceding a character string into three integers. You'll note that this struct defines four components to each and every string. These components are:

- **length** – obvious enough, this is the length of the string. Note that the length is the number of characters in the string, and does not include the zero at the end. The length does not include any bogus values that happen to be array after the zero either. Keep in mind that the length of a string may be shorter than the array in which the string is stored.

- **capacity** – the capacity is (one less than) the size of the array in which the characters are stored. So, for example, if we had the string “apple” stored in an array with 20 characters, then capacity would be 19, length would be 5. The first six elements of that array would be { ‘a’, ‘p’, ‘p’, ‘l’, ‘e’, 0 }, and the remaining elements would be random junk. In other words, capacity is the maximum number of characters that can be safely stored inside our string array (and still leave room for the zero on the end).
- **check** – The signature value stored in a string should always be ~0xdeadbeef. This signature serves an unusual purpose. It’s a marker that we place in all our strings. We’ll look for this signature every time we’re asked to work with a string. If signature has exactly the right value (~0xdeadbeef), then the data almost certainly comes from one of our strings. If the signature has the wrong value, then the string probably wasn’t created by us – i.e., someone made a mistake.
- **data** – As a cute trick, our struct won’t contain a pointer, it will contain an actual array. We’re declaring this array to have a size of zero. But, when we allocate space for a String, we’ll always be sure to allocate an extra big chunk.

Here’s some C code that will create and initialize a string for “yo!” using an array with capacity 10.

```
// Example only, this is not great code
char* s; // should be "yo!" when we're done
String* new_string = malloc(sizeof(String) + 10 + 1);
(*new_string).length = 3; // 3 characters in "yo!"
(*new_string).capacity = 10; // malloc'd 10 bytes
(*new_string).check = ~0xdeadbeef;
(*new_string).data[0] = 'y';
(*new_string).data[1] = 'o';
(*new_string).data[2] = '!';
(*new_string).data[3] = 0;
s = (*new_string).data;
printf("the string is %s\n", s);
```

The first executable statement in the example calls malloc and requests a chunk with enough space to store a String (that’s what sizeof(String) does) plus an extra 10 bytes for the character array at the end (our capacity of 10), plus one extra byte for the zero that we will place at the end of the string. We initialize the length, capacity and signature values for this struct, and then start copying our sequence of characters into the array at the end. Since we requested an extra 10 bytes, we know there’s room to spare for these three characters and the zero. Finally, we make the char* pointer s point at the first byte in that array and we’re done.

Here's some functionally equivalent C code without the use of the String struct.

```
// Example only, this is not great code
char* s; // should be "yo!" when we're done
uint32_t* p; // access the ints before the string
p = (uint32_t*) malloc(3 * sizeof(uint32_t) + 10 + 1);
s = (char*) (p + 3);
s[0] = 'y';
s[1] = 'o';
s[2] = '!';
s[3] = 0;
p = (uint32_t*) s;
p = p - 1;
*p = 3; // length immediately before the string
P = p - 1;
*p = 10; // capacity immediately before the length
P = p - 1;
*p = ~0xdeadbeef; // signature before capacity
printf("the string is %s\n", s);
```

Using Strings: Our strings can be used any place that an ordinary string can be used. Anyone choosing to use our string library should declare variables of type `char*` and work with our strings almost as if they were ordinary strings. There are a few exceptions.

- Our strings can only be stored on the heap. Declaring a local or global variable of type "String" will not work. We can, of course, declare local variables of type "String*" (pointers to String). We'll point these pointers at chunks from the heap.
- Our "clients" (programmers who decide to use our strings) agree to create strings calling the *utstrdup* function. Our clients further agree not to poke around inside our struct, or to mess with the characters inside the array. Our clients promise to call *utstrfree* with their strings when they're done with them. Note that clients will never declare any variables of type String or type String*. Client programmers will access strings using variables of type `char*`.

The Functions: You are to write each of the following six functions. Note that the parameters and return values for these functions are **char*** types. You will, of course, actually be creating structs on the heap – which leads to some good practice for us.

`char* utstrdup(char* source);` -- This function is a mirror of the standard C library function called *strdup*. The parameter is a string that we wish to duplicate. The returned pointer will point onto the heap. When you write this function, create a String struct on the heap that holds a copy of source. Set the length and capacity of your String equal to the number of characters in source. Be sure to return the address to the first character in your String (and not the address of the String struct itself). The argument source may or may not be a utstring (i.e., it might be an ordinary C string without our extra info in front).

`void utstrfree(char* p);` -- This function is used by our clients when they are finished working with one of our strings. The function should deallocate the chunk on the heap that holds this String struct (by calling `free`). Just be sure to compute the correct address before calling `free`. Our clients are required to never call this function unless the argument `p` is a utstring. They're also required to always call this function (eventually) for every utstring that they create.

`uint32_t utstrlen(char* str);` -- This function is a mirror of the standard C library function called *strlen*. NOTE: `str` must actually be a String struct. Your function should perform the pointer arithmetic and return the length stored inside this string struct – in other words, this function had better be a whole lot faster than the standard `strlen` function (i.e., no loops (and, please, no recursion)). In this function and all the remaining functions, the first argument must always be a utstring. In the case of `utstrlen`, we knew `str` is a utstring and therefore we can extract the length directly from the bytes in the front of the string.

`char* utstrcpy(char* dest, char* source);` -- This function is a mirror of the standard C library function called *strcpy*. This function should replace the characters in `dest` with the characters from `source`. NOTE: `dest` will actually be a String struct (i.e., the address stored in `dest` will be one of the addresses you returned from `utstrdup`). When you perform the copy, you must not overflow the capacity of this String.

- If the source string is longer than the capacity of `dest`, then copy only as many characters as will fit (i.e., copy capacity characters).
- If the capacity of `dest` is equal to or larger than the length of `source`, then copy all the characters from `source`.
- Be sure to set the length of the `dest` correctly.
- By convention (i.e., for no particularly good reason), this function should return `dest` when it is over. The return value of the function is not actually useful.

`char* utstrcat(char* dest, char* suffix);` -- This function is a mirror of the standard C library function called *strcat*. Similar to `utstrcpy`, `utstrcat` requires that the first parameter be a String struct (well the char array that is part of a String struct). The function should append characters from `suffix`, but must not exceed the total capacity of `dest`.

`char* utstrrealloc(char* str, uint32_t new_capacity);` -- This function is similar to the standard C library function `realloc`. The first parameter is a utstring. We must examine this struct and compare the current capacity to the `new_capacity` parameter.

- If the current capacity is equal to or larger than the `new_capacity`, then do nothing and return `str`.
- If the current capacity is smaller than the `new_capacity`, then create a new String struct with `new_capacity`. Copy all the characters from `str` to this new string. Deallocate the chunk where `str` was stored, and return the newly created string.

`int utstrcmp(const char* s1, const char* s2);` -- This function is similar to the standard C library function `strcmp`, except that both `s1` and `s2` must be utstrings. If the two strings are identical (contain exactly the same characters) you must return 0. If the strings are

different you must return either -1 or +1 depending on whether s1 or s2 is “lexicographically” smaller. To do the comparison, you must compare the first character in s1 to the first character in s2. If the numeric value (e.g., the ASCII code) for s1[0] is smaller than the numeric value for s2[0], then s1 is lexicographically smaller, and you must return -1. If s1[0] > s2[0] then s2 is lexicographically smaller and you must return +1. If the first characters of s1 and s2 are identical, then you must compare s1[1] and s2[1]. So, for example, “apple” is smaller than “ate” (lexicographically, anyway) since the first characters are the same, and ‘p’ is smaller than ‘t’. Note that you’re supposed to do straight numeric comparison using the encoded characters. So, “0abc” is smaller than “abc0” because the first character in “0abc” is ‘0’, which is 48 and that’s smaller than ‘a’.

Please note: your utstrcmp must return the correct value, -1, 0, or +1. However, you’re not required to actually do the comparison in the way I just described. You can compare any character (or characters) to any other character (or characters) you want, in any order you want, provided you return the correct value – the same value (-1, 0 or +1) that you would have returned if you had done the comparison in the order I described.

In order to receive full credit for utstrcmp you must design an implementation that is measurably faster than strcmp from the standard C library. If you want to accomplish that trick, I encourage you to consider introducing additional meta data into your strings. You can do anything you wish provided you do not do anything that more than doubles the amount of heap memory required. So, if the standard utstring would take 24 bytes (for example), then your “improved” utstring cannot require more than 48 bytes to store, including any meta data that you need.

Test program: The stage 1 and stage 2 tests are pretty weak. Hopefully they get the idea across on how our strings will work. The tests might even find some errors in your program. However, please understand that you need to test your own programs. Just because your program happens to pass the (weak) tests provided in main.cpp, please do not assume that your program is correct. Write some additional tests of your own (in particular, utstrcpy is not thoroughly tested).

The stage 3 test is trying to ensure you’ve done nothing silly. If all goes well, this test should run in less than a second and should not have any buffer overflows. If it takes a long time to run (e.g., a minute or more), then you’ve probably done something silly.

To correctly pass stage 4, you must use the assert macro in all of your functions. Stage 4 has two parts, the first part simply checks the heap to make sure you have no memory leaks or buffer overflows. It’s not perfect, but it does a reasonable job of catching a lot of mistakes. The second part calls the utstr functions with incorrect arguments. When this happens, your program is expected to assert(blah) for some expression blah which evaluates to false. A really good assert statement to use could be:

```
assert((*my_string).check == SIGNATURE);
```

where SIGNATURE is #defined to be ~0xdeadbeef. Since this part of stage 4 is supposed to crash your program you'll need to test each function one at a time. Uncomment one of the six lines, compile and run your program. It should crash. Comment out all of the six tests in stage 4, and your program should run without errors.

The stage 5 tests are there to help you evaluate whether your utstrcmp function is in fact faster than the strcmp function. To check, there's a struct call BackwardsCompare with a function inside it (yeah, that's stuff we'll get to later in the semester). Anyway the function is a little 1-liner that calls either utstrcmp or strcmp. The stage 5 test should work with either, and when it runs, the program prints the amount of time taken to sort the entire American-english dictionary, backwards. Compile your program once where BackwardsCompare uses strcmp and run it, then compile your program where BackwardsCompare uses utstrcmp and run it and see which is faster.

NOTE: to test the time your program takes, you MUST enable optimization in the compiler. In g++ compile with the -O flag. In VisualStudio, configure your project for Release mode (instead of the default Debug mode). If you compile your program without optimization, even good utstrcmp functions may appear slow and/or bad utstrcmp functions may appear fast.

Good luck!