

# CS 214: Systems Programming, Spring 2014

## Programming Assignment 1: Tokenizer

### 1 Introduction

In this assignment, you will practice programming with C pointers. Much of the pointer manipulation will come in the form of operating on C strings, although you will be dealing with some pointers to structs as well.

Your task is to write a type and a set of functions, in essence, the equivalent of a Java class that implements a tokenizer. The tokenizer should accept two strings as command-line arguments, the first of which will contain a set of *separator characters* while the second will contain a set of *tokens* separated by one or more separator characters. The tokenizer should return the tokens in the second string one token at a time, hence your program is called a tokenizer.

A string on a command line is a sequence of characters delimited by double quotes ("). Command-line strings can contain newline or double-quote characters, but special syntax is required to contain them and certain other characters. These special characters are represented with *escape characters*. The escape characters may appear in the first or second argument strings to the tokenizer. Each escape character is represented by two regular characters (e.g. "\n" backslash and lower-case 'n' mean a single-byte newline character. The key is the backslash character in the strings. Here are the rules for dealing with the backslash and whatever follows it:

newline (0x0a) \n  
horizontal tab (0x09) \t  
vertical tab (0x0b) \v  
backspace (0x08) \b  
carriage return (0x0d) \r  
form feed (0x0c) \f  
audible alert (0x07) \a  
backslash (0x5c) \\  
double quote (0x22) \"

A backslash at the end of a string should be ignored.

### 2 Implementation

Your implementation needs to export the interface given in the attached `tokenizer.c` file. In particular, you need to define the type needed to represent a tokenizer and three functions for creating and destroying tokenizer objects and getting the next token. Note that we have only defined the minimal interface needed for external code (e.g., our testing code) to use your tokenizer. You will likely need to design and implement additional types and functions.

A token is a sequence of any ASCII characters that does not contain a separator character. Separator characters are provided as a string of one or more ASCII characters. Each pair of tokens are separated by one or more separator characters. Multiple separators may be next to each other (see second example above), and/or at the beginning and/or end of the token string. When this happens, your tokenizer should discard *all* separators.

Your implementation must *not* modify the two original strings in any way. Further, your implementation must return each token as a C string in a character array of the exact right length. For example, the token `usr` should be returned in a character array with 4 elements (the last holds the character `'\0'` to signify the end of a C string).

You may use string functions from the standard C library accessible through `string.h` (e.g, `strlen()`). However, you may not use `strtok()`, `strsep()`, `strpbrk()` or any similar function that already performs the complete tokenization process.

You should also implement a `main()` function that takes 2 string arguments, as defined above. Each character in the first string is a separator. The second string contains zero or more tokens separated by separator characters. Your `main()` function should print out all the tokens in the second string in left-to-right order. Each token should be printed on a separate line. Here is an example invocation of the tokenizer and its output.

```
tokenizer " " "today is sunny"
```

```
today
```

```
is
```

```
sunny
```

Note that many of the escape characters are either unprintable or have undesirable effects on your output. In your output, we want the output of all escape characters (printable and otherwise) to be in bracketed hex of the form `[0xhh]`. So if the `argv` has `""` as the (empty) delimiter string and `"hello\nworld"` as the string to be tokenized, the final output should be `"hello[0x0a]world"` on a single line (16 characters long).

The escape characters are all preceded by a backslash. A backslash may appear before other characters in your program input. For a backslash behind any other character, the effect would be to just remove the backslash, (e.g. `"\k"` would be interpreted as just `"k"`). A backslash at the end of an input string should be discarded and ignored.

Keep in mind that *coding style will affect your grade*. Your code should be well-organized, well-commented, and designed in a modular fashion. In particular, you should design reusable functions and structures, and minimize code duplication. *You should always check for errors*. For example, you should always check that your program was invoked with the minimal number of arguments needed.

Your code should compile correctly (no warnings and errors) with the `-Wall` and either the `-g` or `-O` flags. For example

```
$ gcc -Wall -g -o tokenizer tokenizer.c
```

should compile your code to a debug-able executable named `tokenizer` without producing any warnings or error messages. (Note that `-O` and `-o` are different flags.)

Your code should also be efficient in both space and time. When there are tradeoffs to be made, you need to explain what you chose to do and why.

IMPORTANT NOTE: You may write your code on any machine and operating system you desire, but the code you turn in **MUST** (un)tar (see below), compile and execute on the iLab machines or a zero grade will be given. Be sure to compile and execute your code on an iLab machine before handing it in. This has been clearly stated here and **NO EXCEPTIONS** will be given.

## 3 Examples

### 3.1 Basic Input

- Input:

```
./tokenizer " " "today is a beautiful day"
```

- Output

```
today
is
a
beautiful
day
```

### 3.2 Multiple Delimiters

- Input:

```
./tokenizer "/" "?" "/usr/local?/bin? share"
```

- Output

```
usr
local
bin
share
```

## 4 What to turn in

A tarred gzipped file named `pa1.tgz` that contains a directory called `pa1` with the following files in it:

- A `tokenizer.c` file containing all of your code.
- A file called `testcases.txt` that contains a thorough set of test cases for your code, including inputs and expected outputs.

- A `readme.pdf` file that contains a brief description of the program and any great features you want us to notice.

Suppose that you have a directory called `pa1` in your account (on the iLab machine(s)), containing the above required files. Here is how you create the required tar file. (The `ls` commands are just to help show you where you should be in relation to `pa1`. The only necessary command is the `tar` command.)

```
$ ls pa1 $ ls pa1 readme.pdf testcases.txt tokenizer.c $ tar cfz pa1.tgz pa1
```

You can check your `pa1.tgz` by either untarring it or running `tar tfz pa1.tgz` (see `man tar`).

Your grade will be based on:

- Correctness (how well your code works).
- Quality of your design (did you use reasonable algorithms).
- Quality of your code (how well written your code is, including modularity and comments).
- Efficiency (of your implementation).
- Testing thoroughness (quality of your test cases).