

Project 2: Basic HTML Parsing and Crawling

Instructions

This project is divided into two parts:

1. Build a simple HTML parser that determines if your HTML tags are balanced.
2. Find the number of webpages you can visit from a certain HTML page.

To keep things simple for this project, all pages will be local to your machine.

Part 1: Basic HTML parsing

It is important to know if your HTML tags are balanced. For example:

```
1 <html>
2   <body>
3     <p>
4       <b>Hello World!</b>
5     </p>
6   </body>
7 </html>
```

is balanced since each tag has a begin tag (<tagname>) followed by an end tag (</tagname>) at the same “level depth”. For instance, consider the following:

```
8 <html>
9   <body>
10    <p>
11      <b>Hello World!</b>
12    </body>
13  </p>
14 </html>
```

The above is not balanced because there is a </body> ending a <p> tag. Given the following basic HTML definition, your job is to determine if a given piece of HTML is balanced or not.

HTML	<html>BODY</html>
BODY	<body>TEXT</body>
TEXT	STRING STRING TEXT TAG TAG TEXT
STRING	possibly empty string of printable characters other than < and >
TAG	BOLD ITALICS LINK PARAGRAPH
BOLD	TEXT
ITALICS	<i>TEXT</i>
LINK	TEXT
PARAGRAPH	<p>TEXT</p>
URL	STRING

Part 2: Basic Web Crawler

Now that we have a basic HTML parser, write a web crawler to determine the number of unique pages that can be visited from a certain page. See example below:

Example Output

```
1 ./html-test pages/*
2 Parsing: 'pages/index.html'
3 Parsing: 'pages/unbalanced1.html'
4 Parsing: 'pages/unbalanced2.html'
5 Parsing: 'pages/csu.html'
6 Parsing: 'pages/theend.html'
7
8           Page   Balanced   Visit Count
9   pages/index.html       yes         2
10  pages/unbalanced1.html   no         3
11  pages/unbalanced2.html   no         3
12    pages/csu.html        yes         1
13    pages/theend.html      yes         0
```

`csu.html`, `theend.html`, and `index.html` are balanced, while `unbalanced1.html` and `unbalanced2.html` are not. Also, if the link is **not** visit-able (i.e., the page does not exist), then do not count towards a possible link visit. Here is the explanation of the visit amounts:

File	Visits	Explanation
<code>index.html</code>	2	Can visit <code>csu.html</code> , which is valid, and that page can visit <code>theend.html</code> .
<code>unbalanced1.html</code>	3	Can visit <code>index.html</code> , which links to <code>csu.html</code> , which links to <code>theend.html</code> .
<code>csu.html</code>	1	Can visit <code>theend.html</code> , which is valid and a dead end.
<code>theend.html</code>	0	Links to only one page, but it does not exist.

Remember, the visit count is the number of **unique** pages that can be visited from a page.

Parsing

The first part of this project requires you to read a file and parse out the input. This may seem daunting at first, but here are a couple of things to help out:

- Tags:
 - begin with a “<” and end with a “>”.
 - do **not** have a < or > between the starting < and ending >. Therefore once you see a < parse to a > and that string is your tag and attributes.
 - All tags, except the anchor tag (<a>), will have no attributes or whitespace. That means the body tag will always be “<body>”, there will be no spaces or other characters. The same goes for the end tags.
 - The anchor tag will **strictly** be in the form “” with exactly one space between the “a” and “href”. The URL is in double quotes after “href=”.

- If you encounter a tag that is not valid then the page is not balanced.
- I will ensure all tags are lower case.
- To ignore whitespace, parse everything one character at a time ignoring space, tab, new-line, and carriage return (' ', '\t', '\n', '\r').

Hints

So far we have used arrays, linked lists, pointer, stacks, and queues. Chances are you will be using more than one data structure to solve this problem, but you will only need the ones listed.

Also, if you are familiar with the STL, you may use vector, stack, and queue from the STL library, but nothing else.

The C++ Regular Expression Library may be used for parsing but it isn't required.

Warning

The Data Structure projects are *significantly* more demanding than the labs. Start thinking about the data structure now! If you hit a roadblock (and you should/will), please come talk to me. When you talk to me be prepared to show me what you have thought of so far. If you come in with an idea, I will guide your idea to a good solution. If you come in without an idea, I will tell you to keep thinking about it. Keep in mind, there are many good, fair, and terrible solutions to Data Structure projects.

Write some test cases

Part of this project grade will be how well you can write test cases. You will be in charge of all test cases.

How to turn in

Turn in via GitHub. Ensure the file(s) are in your directory and then:

```
1 git add <files>
2 git commit
3 git push
```

Assigned

October 3, 2022

Due Date

October 31, 2022 11:59 PM

Teamwork

No teamwork, your work must be your own.