







DOM tree traversal

Description:

One of the techniques we use at Medallia for automated regression testing consist on comparing the HTML output of different versions of our code. To do this, we convert the HTML to a "DOM" (Document Object Model) and compare the output. We also filter expected differences from the output via a simple whitelist.

Assume that a "DOM" representation is a tree containing Element and Content nodes:

- A Content node contains:
 - only its content
 - e.g. if the input is a "<body>some content</body>" then "some content" is the content of the content node
- An Element node contains:
 - a tag
 - e.g. if the input is a "<html>...</html>" then "html" is the tag of the element node
 - o an optional id
 - e.g. if the input is "<div id='foo'>" then "foo" is the id of the element node
 - an optional list of children nodes
- The content of an Element node is the concatenation of all its **direct**Content children, separated by spaces
 - e.g. if the input is "<div>some<a>othercontent</div>" then
 the content of "div" is "some content"

Problem:

Given a DOM tree, write out its content subject to the following constraints:

One line per level of the tree

- Each line should include the **tag**, **id** (if present) and **content** (if present) of all Element nodes in that level, in the **same order they appear** in the tree, with each term **separated by spaces**.
- If a whitelisted string is a **substring** of the tag, id or content of an Element node, then **that node and all its children should be ignored** and excluded from the output.

Example

For the input DOM:

```
<html>
     <body id='content'>
          This
          <div id='wrapper1'>
               is a
               <div id='container1'>
                    <div id='container2'>
                         funny
                    </div>
               </div>
               <div id='wrapper2'>
                    enjoyable
               </div>
               little
          </div>
          good
          <div id='wrapper3'>
               harmless
          </div>
          example.
          <div id='wrapper4'>
          </div>
          <a id='link'>
               And a link.
          </a>
     </body>
</html>
```

with a whitelist:

["wrapper4", "a little"]

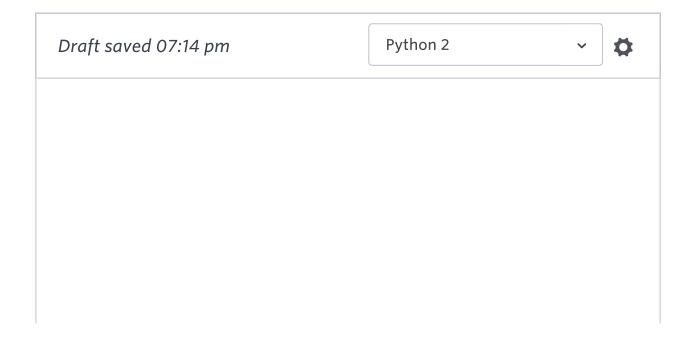
Then the output should be:

html
body content This good example.
div wrapper3 harmless a link And a link.

Explanation:

- "html" is the only element on the top of the tree and it doesn't have any content (everything is below "body").
- "body" has the id "content" and it is the only element on that level. Its content ("This good example") is the concatenation of all its Content nodes.
- "div wrapper3" is the first element in the third level because "div wrapper1" has been white listed (its combined content "is a little" contains the whitelisted string "a little"). Then "a link" is the second element of that level, given that "div wrapper4" has been white listed (its "id" contains "wrapper4"). Finally "And a link." is the content of the "" element.

YOUR ANSWER



```
#!/bin/python
 1
2
 3
   import re
 4
    import collections
 5
    0.00
 6
 7
   An 'Element' node of the DOM tree.
    @param tag: string
 8
9
    @param id: string
    @param children: array of nodes which could be either
10
    'Element' or 'Content' instances
11
12 ▼ class Element:
        def __init__(self, tag, id, children):
13 ▼
14
            self.tag = tag
15
            self.id = id
            self.children = children
16
17
   0.00
18
19
   A 'Content' node of the DOM tree
20
    @param content: string
21
22 ▼class Content:
23
        def init (self, content):
24
            self.content = content
25
    .....
26
27
   A 'DOM' tree
    @param root: an instance of 'Element' which represents to
28
    the root of the tree
    .....
29
30 <del>√</del>class Dom:
31
        def init__(self, root):
32
            self.root = root
33 ▼def format(dom, whiteList):
34
35
        Complete the function below.
36
        @param dom: An instance of the Dom class
        @param whiteList: An array of strings
37
38
        @return the output that will be sent to the STDOUT
39
40
        if dom is not None:
41
            output = ""
42
            queue = []
43
            queue.append(dom.root)
```

```
while queue:
44
45
                cur = queue.pop(0)
46
                if cur.children:
                    output += "\n"
47
                    if cur.tag:
48
                        output += cur.tag + " "
49
                    if cur.id:
50
                        output += cur.id + " "
51
                    for child in cur.children:
52 ▼
                        if isinstance(child, Content):
53
                            output += child.content + " "
54
                        elif isinstance(child, Element):
55
                            queue.append(child)
56
            return output.strip()
57
58
        else:
            return ""
59
```

```
0.00
 60
 61
    The following code is not relevant for this exercise. It
     is only here to help
 62
    running the tests.
 63
 64 ▼def main():
 65
        html = raw input()
 66
         _whiteList = raw_input()
 67
         whiteListArray = [] if not whiteList else
     whiteList.split(",")
 68
        dom = Dom(Parser( html).parse())
 69
 70
        _output = format(_dom, _whiteListArray)
 71
 72
        print output
 73
 74 ▼class Token:
 75 ₹
        def init (self, type, tag, id, content):
 76
             self.type = type
 77
             self.tag = tag
 78
             self.id = id
 79
             self.content = content
 80
 81 ▼class Parser:
 82
        START ELEMENT PATTERN = re.compile('^<(.*?)>')
        END ELEMENT PATTERN = re.compile('^<\/(.*?)>')
 83
        CLASS PATTERN = re.compile("(.*) id='(.*)'")
 84
        CONTENT PATTERN = re.compile('^(.*?)<')
 85
 86
        def init (self, html):
 87 ₹
             self.position = 0
 88
 89
             self.tokens = self.tokenise(html)
 90
 91
        def parse(self):
 92
             return self. parse()
 93
 94
        def parse(self):
 95
             token = self.tokens[self.position]
             _children = []
 96
 97
             self.position += 1
             while self.position < len(self.tokens):
 98 ₹
 99
                 currentToken = self.tokens[self.position]
                 if currentToken.type == 'START':
100
                     children.append(self. parse())
101
                 elif currentToken.type == 'CONTENT':
102
103
```

```
children.append(Content( currentToken.content))
104
                 else:
105
                     break
106
                 self.position += 1
             return Element(_token.tag, _token.id, _children)
107
108
109 ▼
        def tokenise(self, input):
             _tokens = []
110
             while input:
111 \forall
                 _endElement =
112
     self.END ELEMENT PATTERN.match(input)
                 if _endElement:
113 ▼
114
                     tokens.append(Token('END',
    _endElement.group(1), None, None))
115
                     input =
    input[len(_endElement.group(0)):]
116 ▼
                 else:
117
                     startElement =
    self.START ELEMENT PATTERN.match(input)
                     if startElement:
118 ▼
119
                         classElement =
    self.CLASS PATTERN.match( startElement.group(1))
120
                         if classElement:
121
                             tokens.append(Token('START',
     classElement.group(1), classElement.group(2), None))
122
                         else:
123
                             tokens.append(Token('START',
     _startElement.group(1), None, None))
124
                         input =
    input[len( startElement.group(0)):]
125 ₩
                     else:
126
                         content =
     self.CONTENT PATTERN.match(input)
127 ▼
                         if content:
128
                             tokens.append(Token('CONTENT',
    None, None, content.group(1)))
129
                             input =
     input[len( content.group(1)):]
130 ▼
                         else:
131
                             tokens.append(Token('CONTENT',
    None, None, input))
                             input = ''
132
133
             return tokens
134
135
    main()
```

| | - | $\sim -$ | ~ . | |
|---------|-----|----------|-----|-------------|
| Line: | - 1 | ィム | ('\ | |
| TITIE . | | 、フン | CO. | ⊥• / |

☐ Test against custom input

Run Code

Submit code & Continue

L Download sample testcases The input/output files have Unix line endings. Do not use Notepad to edit them on windows.

Status: Compiled successfully. 1/4 sample test cases passed.

Testcase 1: Wrong Answer

Your Output

html
body content This good example.
div wrapper1 is a little
div wrapper3 harmless
div container1
div wrapper2 enjoyable
div container2 funny

Expected Output

html
body content This good example.
div wrapper3 harmless

Testcase 2: Wrong Answer

Your Output

Expected Output

html foo

Testcase 3: Success

Your Output

Testcase 4: Wrong Answer

Your Output

html
body this problem is nice
div foo is an ugly
div bar little
div foo not

Expected Output

html
body this problem is nice
div bar little div bar

About Privacy policy Terms of service