

GCRAWLER: GRAPHICAL WEB CRAWLER

The GCrawler project is a graphical web crawler which crawls a user-supplied website, following links on each page as it goes. It displays a color-coded graph of what pages were crawled and how it reached each page. Try it out here: [GCrawler App](#)

Website

The Website is responsible for:

1. Serving the user with the search form
2. Tracking the user's recent searches
3. Passing input to the Data Transfer program
4. Displaying the results of a search

```
17 router.post('/', function(req, res){
18
19   r = {};
20   var today = new Date();
21   var dd = String(today.getDate()).padStart(2, '0');
22   var mm = String(today.getMonth() + 1).padStart(2, '0');
23   var yyyy = today.getFullYear();
24   today = mm + '/' + dd + '/' + yyyy;
25
26   var link = req.body.link;
27   var search_type = req.body.search_type;
28   var max = req.body.max;
29   var keyword = req.body.keyword;
30
31   if(!link || !search_type || !max)
32   {
33     r.error = "Invalid search...please try again!";
34     r.jsconfigs = ["static_search.js"];
35     r.styles = ["search.css"];
36     res.render('search', r);
37     return;
38 }
```

```
17 /**
18  * Function Name: main
19  * Description: The main controlling program of the crawler.
20  * Inputs: Takes a string representing the starting URL, a char flag for type of
21  * search, an integer representing the link limit and an optional keyword for
22  * halting the search
23  * Outputs: Tree in JSON format
24  */
25 def main(url, search_type, link_limit, keyword = None):
26   #branch based on DFS or BFS (Use different functions for each, general flow below)
27   if search_type == "dfs":
28     tree, robots_flag = cu.depth_search(url, link_limit, keyword)
29   elif search_type == "bfs":
30     tree, robots_flag = cu.breadth_search(url, link_limit, keyword)
31   #add parameter
32   else:
33     return 1
34   #log the resulting output
35   string_list = []
36   for node in tree:
37     string_list.append(node.__str__())
38   #check robots flag
39   if robots_flag:
40     print("robots", end = "")
41   print('[' + ','.join(string_list) + ']')
42   return 0
```

Data Transfer

The Data Transfer program is responsible for:

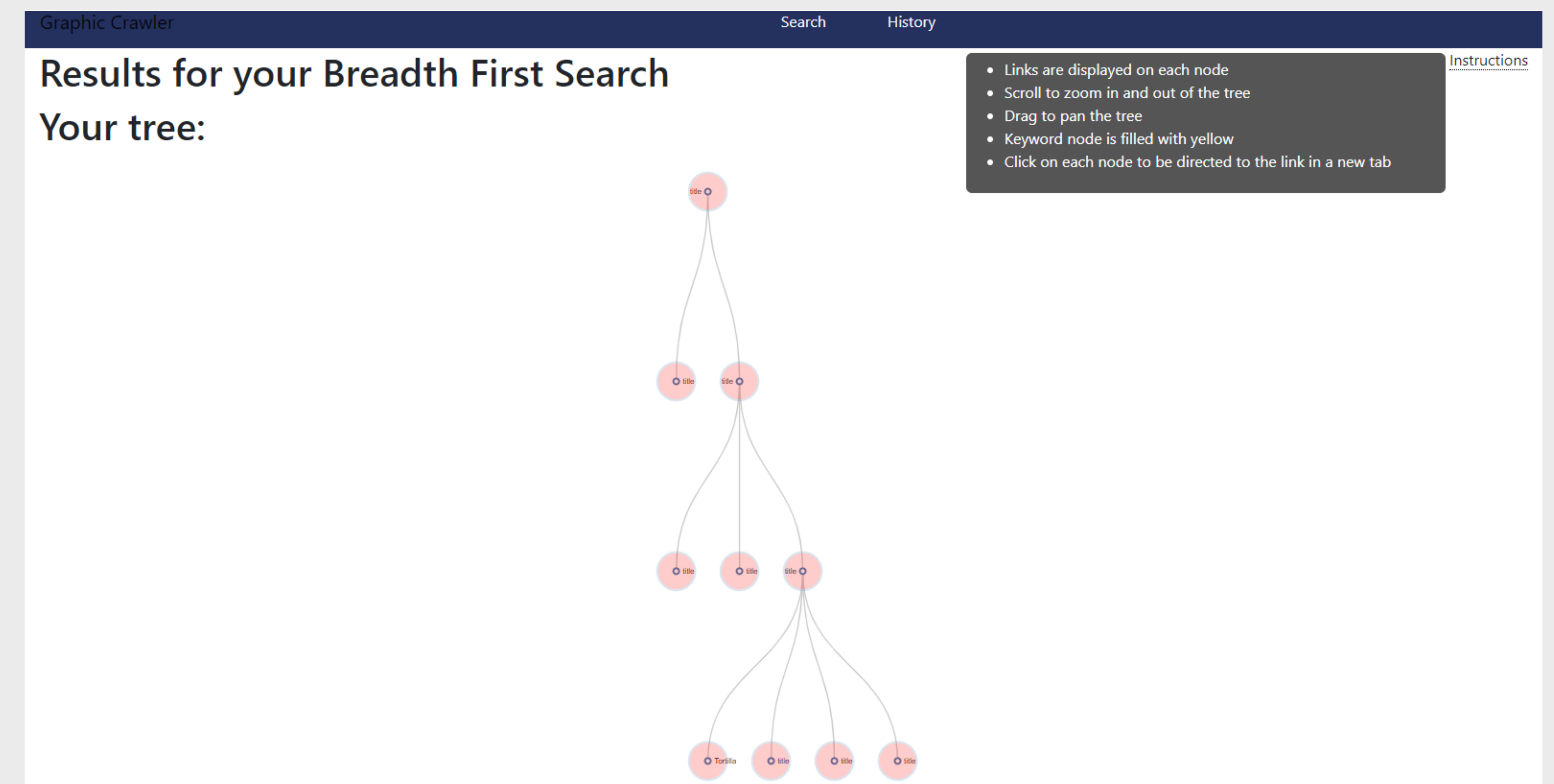
1. Receiving input from the Website
2. Passing the input to the Crawler
3. Receiving the output from the Crawler
4. Passing the output to the Website for rendering

```
66 class Node {
67   constructor(link) {
68     this.name = link;
69     this.children = [];
70     this.parent = "";
71   }
72 }
73
74
75
76
77 function makeLinkTree(jsonArray) {
78   let dict = {};
79   let root;
80
81   jsonArray.forEach(function (urlObject) {
82     let newNode;
83     if (urlObject.URL in dict) {
84       newNode = dict[urlObject.URL];
85     }
86     else {
87       newNode = new Node(urlObject.URL);
88       newNode.title = urlObject.Title;
89       newNode.KeywordFound = urlObject.KeywordFound;
90       dict[urlObject.URL] = newNode;
91     }
92     if (urlObject.Parent === 'null') {
93       root = newNode;
94     }
95     else {
96       let parentNode;
97       if (urlObject.Parent in dict) {
98         parentNode = dict[urlObject.Parent];
99         dict[urlObject.Parent].children.push(newNode);
100       }
101       else {
102         parentNode = dict[urlObject.Parent];
103       }
104       parentNode.children.push(newNode);
105     }
106   });
107 }
```

Crawler

The Crawler is responsible for:

1. Receiving input from the Data Transfer program
2. Crawling each page from the starting link
3. Passing the formatted results to the Data Transfer program for further processing



Enter a website to begin your search!

Full Starting link:

Choose a keyword:

Choose a search type:

- ☐ Depth First Search
☒ Breadth First Search

Page limit (Range 1-3):

GCRAWLER FEATURES:

- **Intuitive:** Complete a simple form to crawl any website. Supply a keyword to search for text on each page. GCrawler will halt when it finds the keyword. (*left*).
- **Retentive:** GCrawler stores your sessions' past searches on the History tab for easy re-use.
- **Flexible:** Choose from Depth First Search or Breadth First Search methods. Depth First Search will follow a random link on each page until the specified page limit (1-10) is reached. Breadth First Search will follow all links on each page until the specified depth limit (1-3) is reached. (*left*)
- **Informative:** GCrawler's color-coded results display the title of each page on the node. The URL can be seen by hovering over the node. A user can click the node to open the page. (*above*)
- **Good Bot:** GCrawler adheres to the Robot Exclusion Protocol by reading a page's robots.txt file and abiding by its rules. It also evaluates pages and links for other indicators that it should not crawl the page or follow the links.

CONNECT WITH THE GCRAWLERS:

[Christopher Beall](#) (Data Transfer): <https://github.com/beallch>

[Helen Jiang](#) (UI/Website): <https://github.com/hyjiang7>

[Brian Metzger](#) (Crawler): <https://github.com/metzgerb>