

WEB CRAWLER PROBLEM

For: Python

Problem

Build a command line script that you can call with a url and it will begin crawling all links on that website and printing out every link it finds to stdout.

For example:

If you put in the url <https://www.amazon.com> it would check the html for that website and then any link it finds, print the link to `stdout` and try to load the html for that page and print any more links it finds on the new page.

This process continues until all links are exhausted. Only links within the same domain should be visited and printed out.

Requirements

- Only crawl links that match the domain of the input url.
 - For example if there was a link from **amazon.com** to **google.com**, it should not be printed out or traversed.
- The cmd line API should also take a flag `-n` that defines the number of “workers” to use for crawling in parallel.
 - These “workers” can be any parallel code execution strategy. *eg: process, thread, or event loop waiting on the network IO, etc.*
 - This should allow for multiple API calls to be performed concurrently using this flag.
- The tool should exit/complete once all links have been visited within that domain.
- You should not crawl the same links twice.
- Optimize for clean code and performance where reasonable.

- Try to just use the standard library if possible or otherwise a minimal set of dependencies. Don't use any third-party libraries or dependencies built for this specific purpose as the goal is for you to build it.

API Example usage

Cmd line:

```
./crawl -n 4 https://www.amazon.com/
```

Output:

```
https://www.amazon.com/gp/help/customer/display.html  
https://www.amazon.com/gp/primecentral
```

etc...

- One link per line of output.
- Should `exit 0` on successful completion.