**CoC News Article Downloader and Viewer**

Miasia Jones

CS 6675-Homework 2-Option 1.2: Write a Web Crawler of your own

## Introduction

For this assignment I decided to crawl the news articles on the College of Computing (CoC) website. The seed URL is the main page of the [CoC news page](). Each page list news articles for a user to choose from and each article listed is linked to a URL to read the full article. I chose this seed URL because I knew I could easily access at least 1,000 URLs from it, being that there are currently 64 pages listing about 20 articles each. The publishing year range that these articles span is 2011-2021.

## Crawler Design

The web crawler is designed to extract the title, date, and URL of every article listed on the CoC news pages. Each article URL is visited, and its HTML content is extracted. All this information is then saved into a .csv file. I created an app using PyQt5 that can crawl the CoC News website and manage the downloaded articles. You do not have to use the app to run the crawler. See README.md for more information. Below is the Pros and Cons for the design of my web crawler alone:
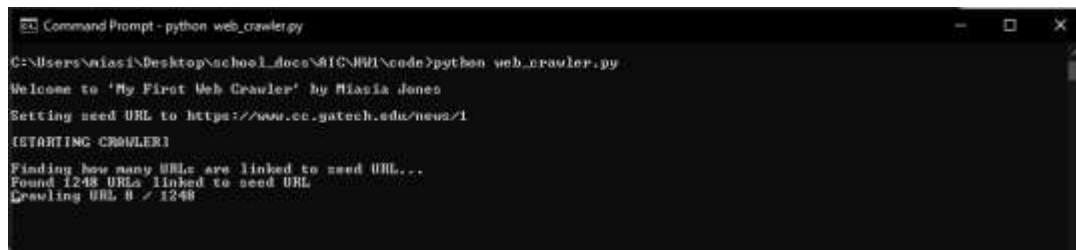
Pros:

- Crawler retrieves all articles from CoC News website, even the undated ones.

- Pages that are visited are not visited again.

- The console output continuously updates the user of its progress.

Cons:

- The crawler calls specific class names in the HTML script of CoC News website to extract certain information, so if these class names were changed, then my crawler may not work.

- The crawler stores the HTML content for each article retrieved in a .csv file, which makes the file hard to read. Saving the HTML content does help within the app I created though.

## Console Screenshots



*Crawling all pages on CoC News page*
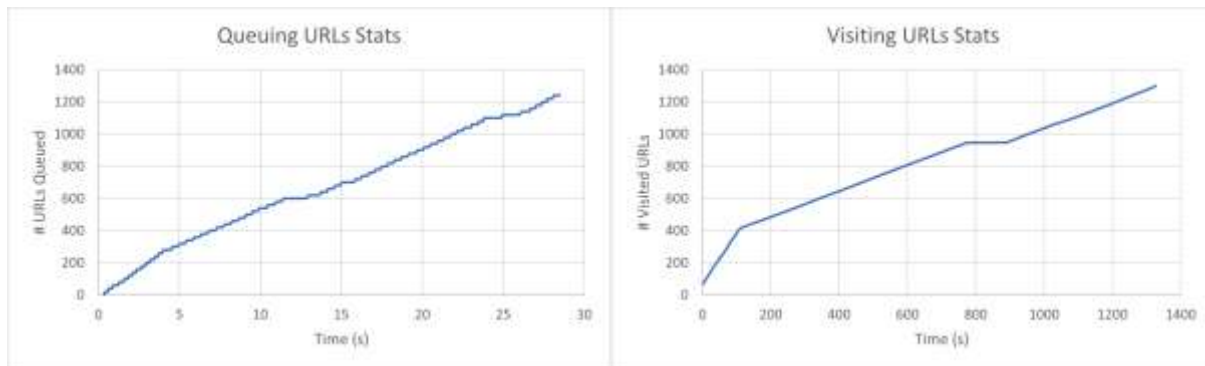


*Crawling only one page on CoC News page*

## App Screenshots



*Using the CoC News Article Downloader and Viewer*

## Crawl Statistics



These are statistics when all the URLs linked to CoC News website is crawled. There were 1,244 URLs found and visited. It took ≈28.47 seconds to find and queue all the URLs linked to the seed URL, in other words, ≈43.74 pages/second. It took ≈1326.76 seconds (≈22.11 minutes) to visit and scrape each URL in the queue, in other words, ≈0.94 pages/second.

## Lessons Learned

- Using the CoC News website was not my first choice, but the structure of the website worked well for what I wanted to do with my crawler. This was my first time implementing a web crawler. My design was simple, but I like the outcome. I was able to brush up on my HTML skill and worked with Python packages I never used before. I was very invested in the assignment and didn't mind the challenge. I'm not sure who would make use of the information I retrieved and I app I created, but I'm very happy with my work. The only thing new that I added that wasn't present on the CoC News page is a news article filtering feature.

- To crawl 10 million pages, I predict that my crawler would take about 4 months. To crawl 1 billion pages, I predict that my crawler would take about 30 years.