



Project Group 10

Cy Heffley, Connor Cole, Eduardo Rodriguez,
Zach McGee



Objective

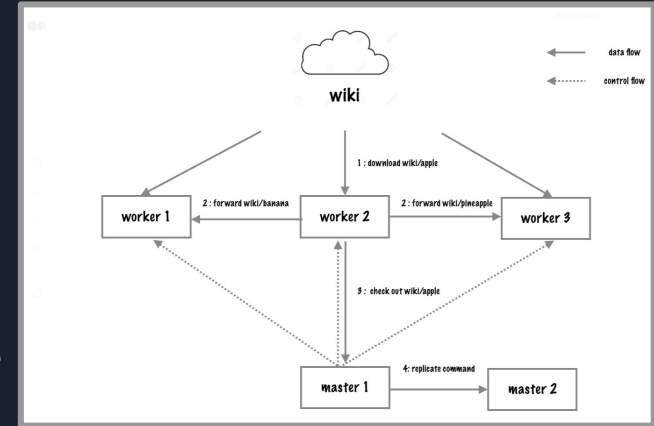
Use modern tools to improve on current internet search technologies, applying a distributed system architecture for scraping websites.

If it is possible to more efficiently crawl the Internet, then a secondary goal is to create a dynamic recommendation system, tailoring results based on the assumption that users with similar interests would wish to receive similar results.

Related Work

In his paper, “Design and Implementation of Distributed Crawler System Based on Scrapy,” Yuhao Fan researched a similar system with quicker persistent memory storage, due to using Redis in conjunction with Scrapy

Issue: False positives in Redis DB increase as time goes on due to Bloom filter used to reduce memory consumption



“Design a Distributed Web Crawler”



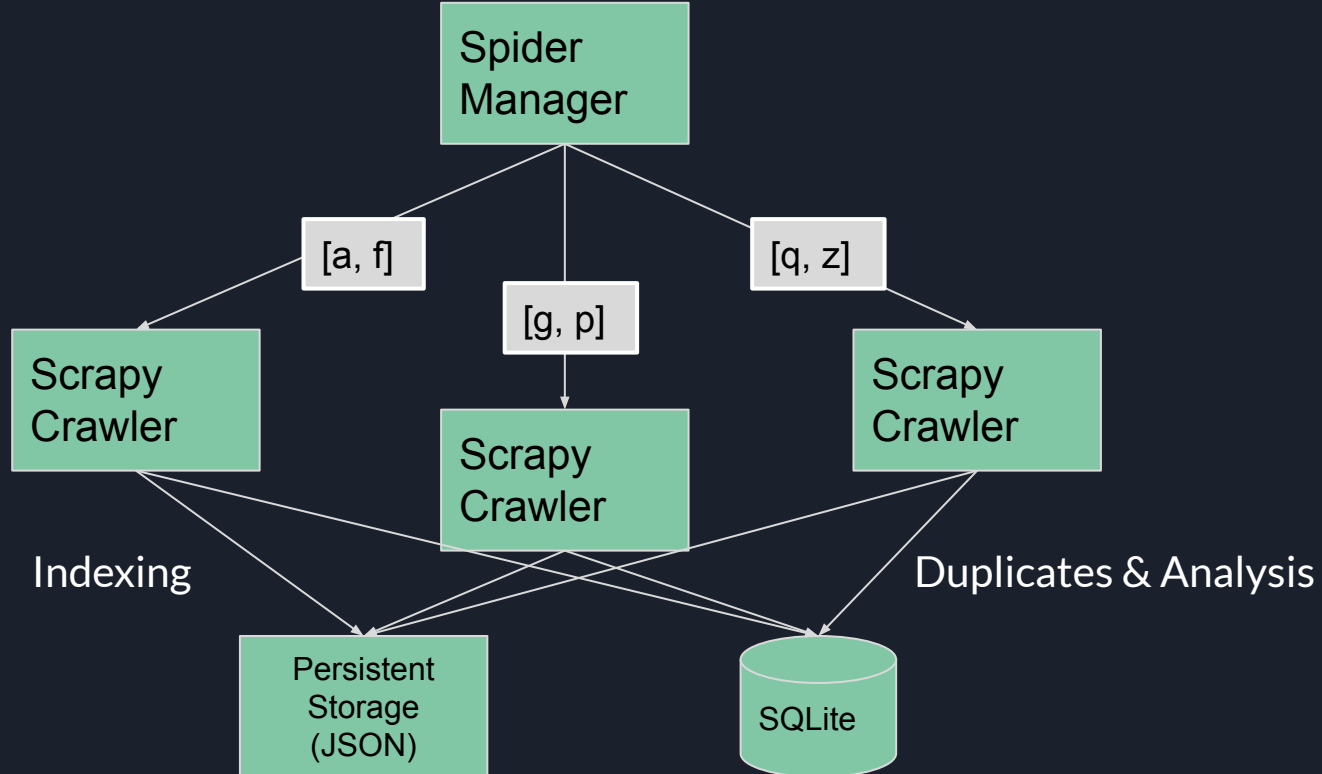
Important definitions

Initial crawling depth - the starting depth of web pages crawled whose links are automatically collected

Filter depth - links are collected past the initial depth, up to the filter depth, only if they have not been previously collected

Maximum depth - links are collected past the filter depth, up to the maximum depth, only if they have not been previously collected and the URL matches the current suffix pattern being crawled

System Architecture





System Architecture and Scrapy

Details:

- Spider manager (`multiple_spiders.py`) is responsible for running multiple individual scrapy processes
 - Scrapy process: the Scrapy framework makes requests/ handles some data scraping
 - We built custom rules/ functionality/ stats monitoring for our designed framework
 - Specifically, we made rules to help deal w/ different levels of depths and to deal w/ distributed systems vs. just the single system



System Architecture and Scrapy

Details:

- Scrapy crawler (post_spyder) is an individual subprocess
 - Run is managed by multiple_spiders, then acts as an independent asynch process
 - Responsible for that processes stack of URLs, filtering/ decisions made internally on if URL should be added to stack
 - Rules for filtering provided by multiple_spiders



High level code execution

Keywords, depths, seed URL defined before start

Stacks = [start_url]

While there are URLs left in the stack,

- See if URL in database, if we are before initial depth pass, else don't add info/ new links from this URL
- If not in DB, add info/ keywords to DB
- If URL is being considered, add all new links (matching suffix criteria, if applicable) to stack

Periodically writes stats ~every 500 new links a subprocess adds



Thought process behind design

Individual processes use DFS when searching URLs

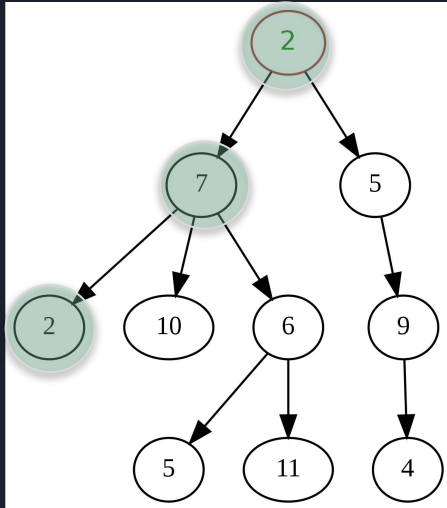
- Similar to traversing down a tree
- Would build up initial stack of URLs to traverse, then new processes would find point in tree to assume responsibility- basically hoping to divide up tree across nodes
- Max depth acts as bottleneck

Tree example- frozen time

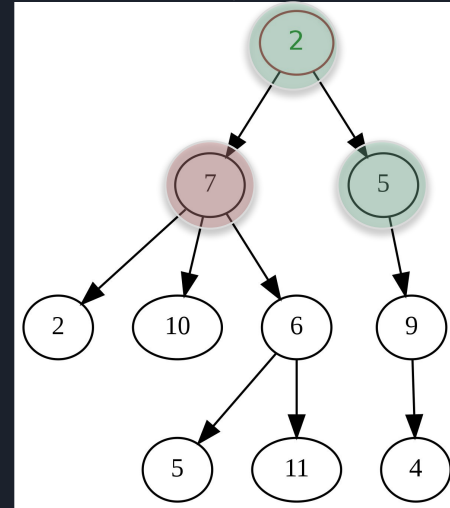
Imagine 2 is a child URL from root URL right at filter depth

- 2 added to booths queues, subprocess 1 takes left side, subprocess 2 sees 7 already in DB so takes right side of tree

Subprocess 1
Working on URL 2



Subprocess 2
Working on URL 5 as 7 is
already in DB



URL Scraping: Initial Depth = 2

Depth

0

Subprocess 1

cc.gatech.edu

1

cc.gatech.edu/xyz

cc.gatech.edu/abc

...

2

cc.gatech.edu/ijk

cc.gatech.edu/ccc

...

Subprocess 2

cc.gatech.edu

cc.gatech.edu/xyz

cc.gatech.edu/abc

...

cc.gatech.edu/ijk

cc.gatech.edu/ccc

...

Adds all links up to filter depth threshold on each scraper to consider- builds up stack of links to work with

Should be small/ quick to get initial stack built up

URL Scraping: Filter Depth = 4

Depth

2

Subprocess 1

cc.gatech.xyz

Checks SQL,
sees
duplicate

Subprocess 2

cc.gatech.abc

3

cc.gatech.edu/abc

cc.gatech.qqq

...

cc.gatech.edu/mmm

cc.gatech.eee

...

4

cc.gatech.edu/rst

cc.gatech.www

...

cc.gatech.edu/qqq

cc.gatech.ccc

...

Queries DB when URL is found to see if exists (if these were perfectly synchronized, cc.gatech.edu/abc has already been added by subprocess 2, not considered in subprocess 1)

URL Scraping: Max Depth = 6

Depth

4

Subprocess 1
Suffix: [a, m]

cc.gatech.xyz

Checks SQL,
sees
duplicate

5

cc.gatech.edu/ooo

cc.gatech.aaa

...

6

cc.gatech.edu/uuu

cc.gatech.ddd

...

Subprocess 2
Suffix: [n, z]

cc.gatech.abc

cc.gatech.edu/eff

cc.gatech.eee

...

cc.gatech.edu/xyz

cc.gatech.ccc

...

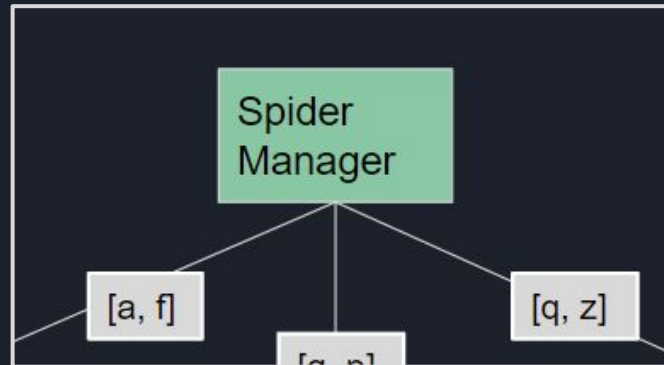
Checks suffix condition, ALSO queries database IF condition is met

Past 6: no links added even if suffix condition met

Spider Manager

Purpose: Divide up different suffixes across multiple scrapers

In true distributed system: manager would be responsible for communication with neighbors, rebalancing suffix rules/ conditions





Technologies

Prototype

Python, Scrapy

SQLite

Production

MySQL, PostgreSQL, Apache
Hadoop, Redis

C++, Rust, or Apache Nutch

Search service: React + Node,
Django, or Spring

Crawling Demonstration

- <https://youtu.be/tVVaFffQCNk>

```
MINGW64/c/Users/ed209/school/CS4675_Group_10/postscrape/postscrape/spi...
2022-04-20 11:06:51 [scraper_1] DEBUG: scraper scraper_1
2022-04-20 11:06:51 [scraper_1] DEBUG: crawling https://www.cc.gatech.edu/people/faculty?page=13
2022-04-20 11:06:51 [scraper_1] DEBUG: current depth: 4
2022-04-20 11:06:51 [scraper_2] DEBUG: scraper scraper_2
2022-04-20 11:06:51 [scraper_2] DEBUG: crawling https://www.cc.gatech.edu/people/faculty?page=13
2022-04-20 11:06:51 [scraper_1] DEBUG: scraper scraper_1
2022-04-20 11:06:51 [scraper_1] DEBUG: crawling https://www.cc.gatech.edu/people/faculty?page=13
2022-04-20 11:06:51 [scraper_1] DEBUG: current depth: 4
2022-04-20 11:06:51 [scraper_2] DEBUG: scraper scraper_2
2022-04-20 11:06:51 [scraper_2] DEBUG: crawling https://www.cc.gatech.edu/people/faculty?page=13
2022-04-20 11:06:51 [scraper_2] DEBUG: current depth: 5
2022-04-20 11:06:51 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.cc.gatech.edu/unit/school-cybersecurity-and-privacy?page=2> (referer: https://www.cc.gatech.edu/unit/school-cybersecurity-and-privacy?page=1)
2022-04-20 11:06:51 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.cc.gatech.edu/unit/school-cybersecurity-and-privacy?page=0> (referer: https://www.cc.gatech.edu/unit/school-cybersecurity-and-privacy?page=1)
2022-04-20 11:06:51 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.cc.gatech.edu/people/haesun-park> (referer: https://www.cc.gatech.edu/person/center/institute-data-engineering-and-science-ideas)
```




Spider Configuration Inputs

Number of spiders in a session: [1, 2, 4]

Initial crawling depth: [2, 3, 5, 10]

Filter depth: [5, 10, 25, 35]

Maximum depth: [10, 45]

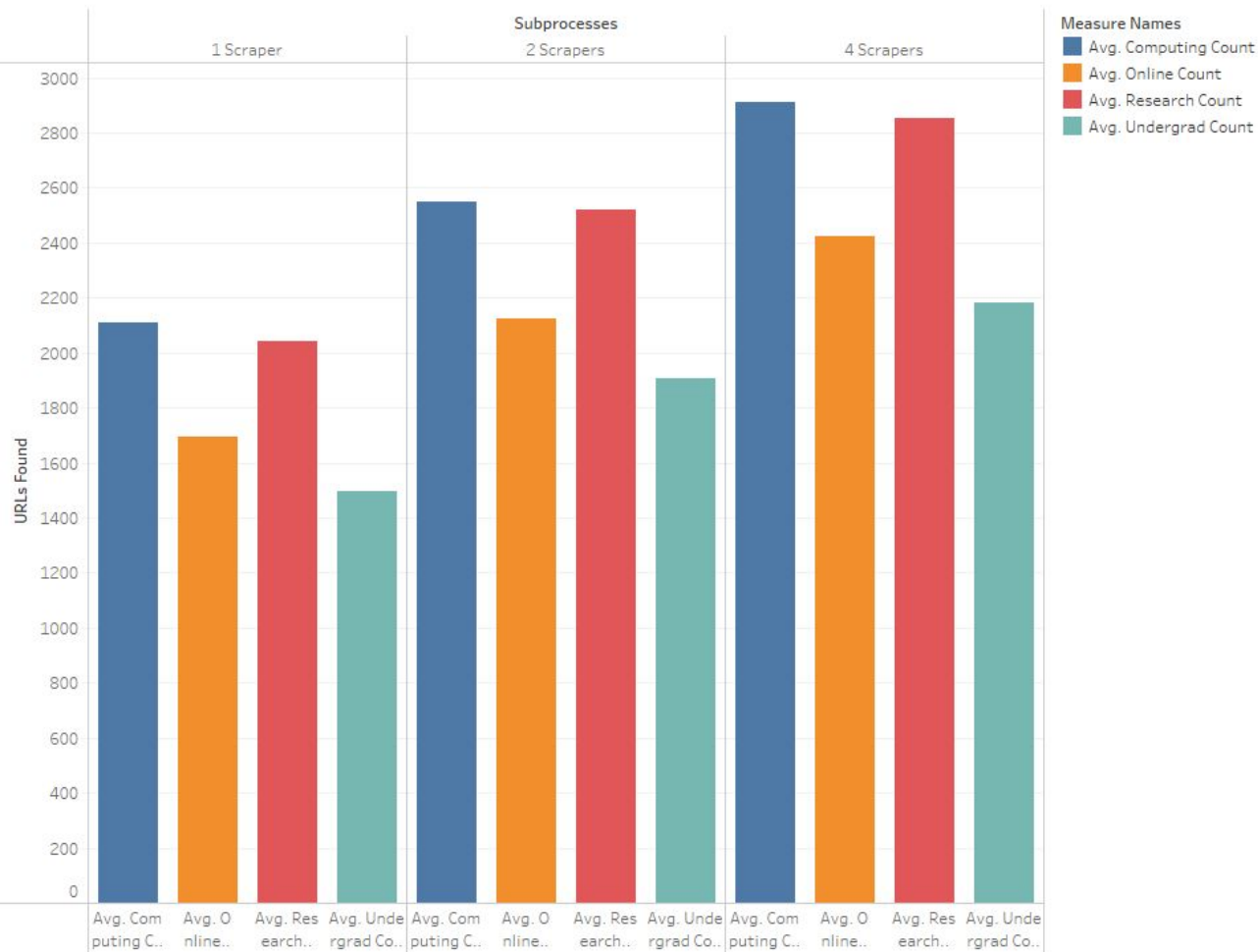


Crawling Results of Our Spider Configurations

Evaluation methods/metrics/datasets and results

- Compare subprocesses, initial, filter, max depths on:
 - Keywords found
 - Duplicate URLs filtered over time
 - Offsite URLs filtered over time
 - GET Requests made over time

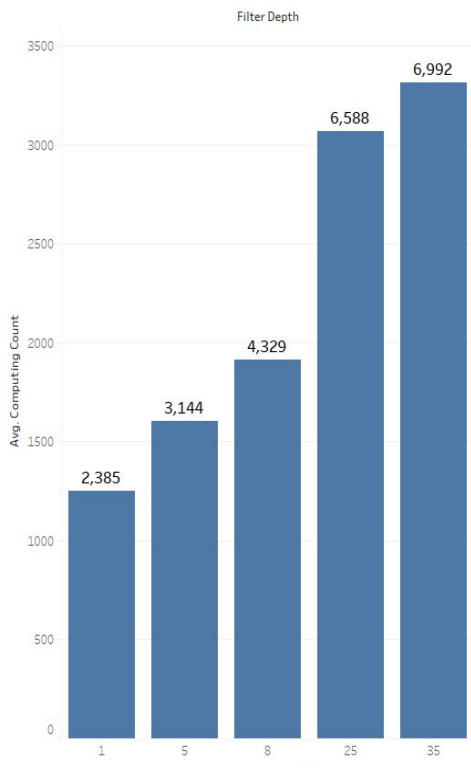
Average Keywords Found By Subprocesses



We see that generally across all run configurations, adding more scrapers gives better results for each keyword.

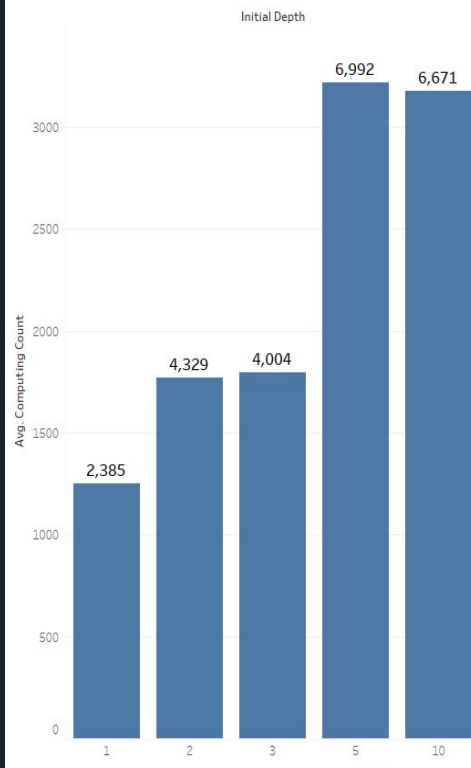
Avg. Computing Count, Avg. Online Count, Avg. Research Count and Avg. Undergrad Count for each Subprocesses. Color shows details about Avg. Computing Count, Avg. Online Count, Avg. Research Count and Avg. Undergrad Count.

Average Keyword Found Across All Runs by Filter Depth



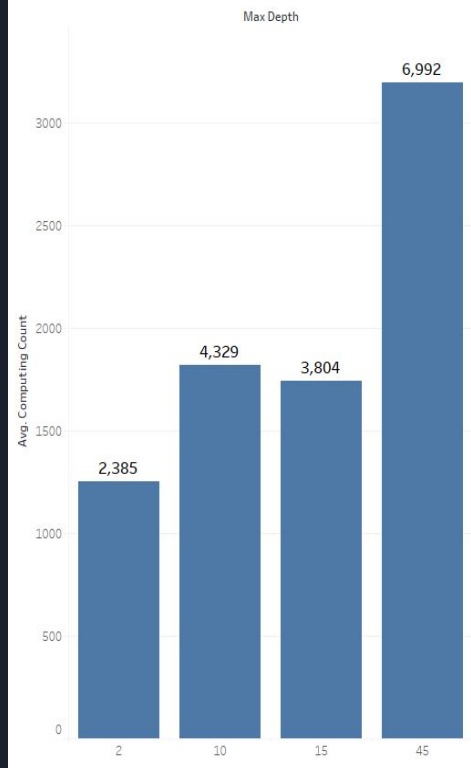
Average of Computing Count for each Filter Depth. The marks are labeled by maximum of Computing Count.

Average Keyword Found Across All Runs by Initial Depth



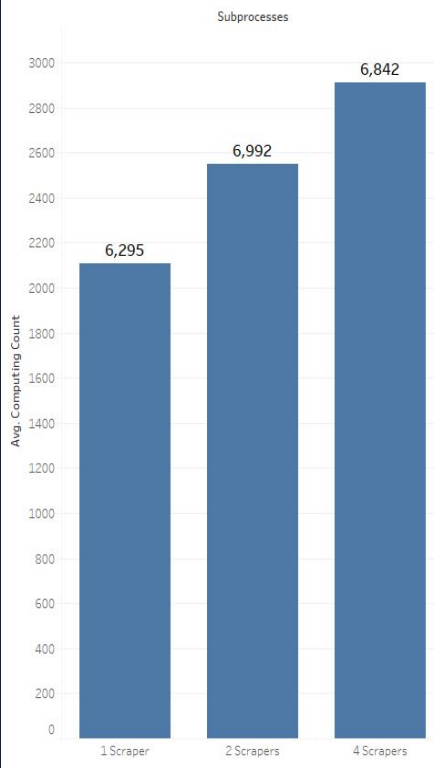
Average of Computing Count for each Initial Depth. The marks are labeled by maximum of Computing Count.

Average Keyword Found Across All Runs by Max Depth



Average of Computing Count for each Max Depth. The marks are labeled by maximum of Computing Count.

Average Keyword Found Across All Runs by Scrapers

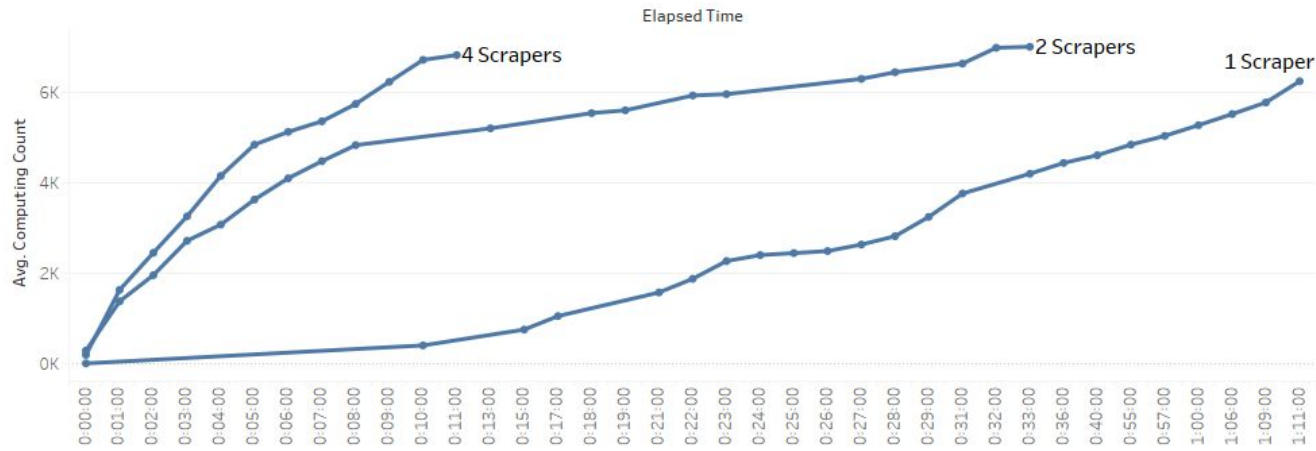


Average of Computing Count for each Subprocesses. The marks are labeled by maximum of Computing Count.

We see a positive correlation between higher depth settings/subprocesses and the number of webpages found containing a given keyword.

Note: This is only the results for the Computing keyword, but all showed the same trends.

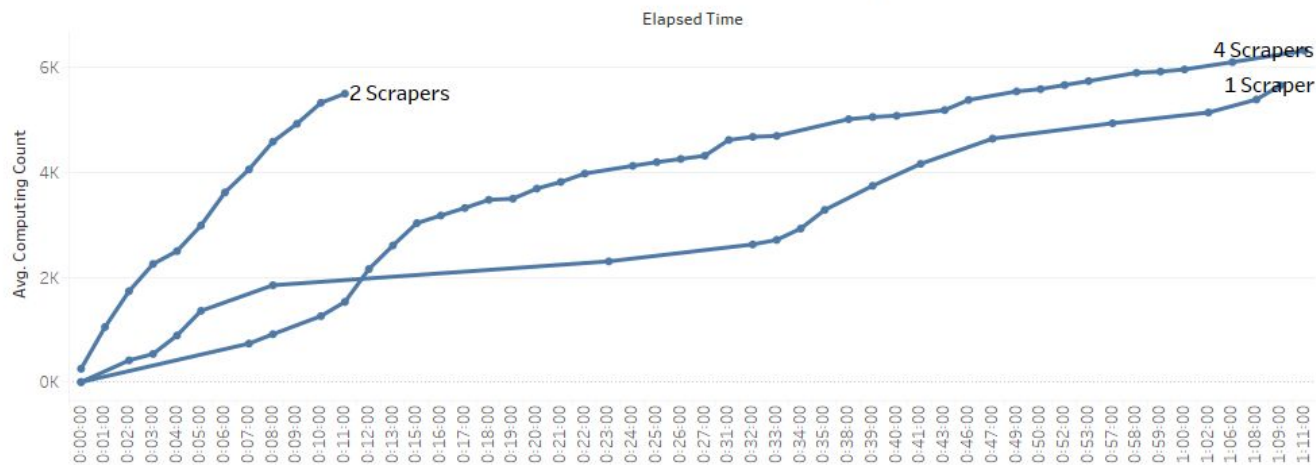
Avg Keyword Found Over Time by Scrapers (5 initial, 35 filter, 45 max)



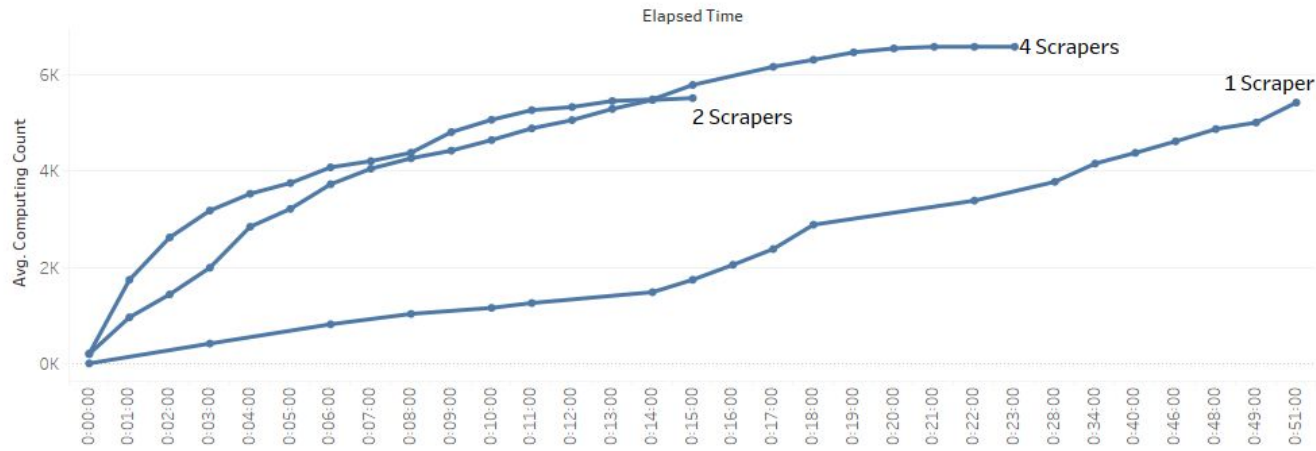
Observations:

- Very different results between filter depths which could indicate an outside factor affecting these quick runs
- The single scraper does consistently worst with a lower initial depth

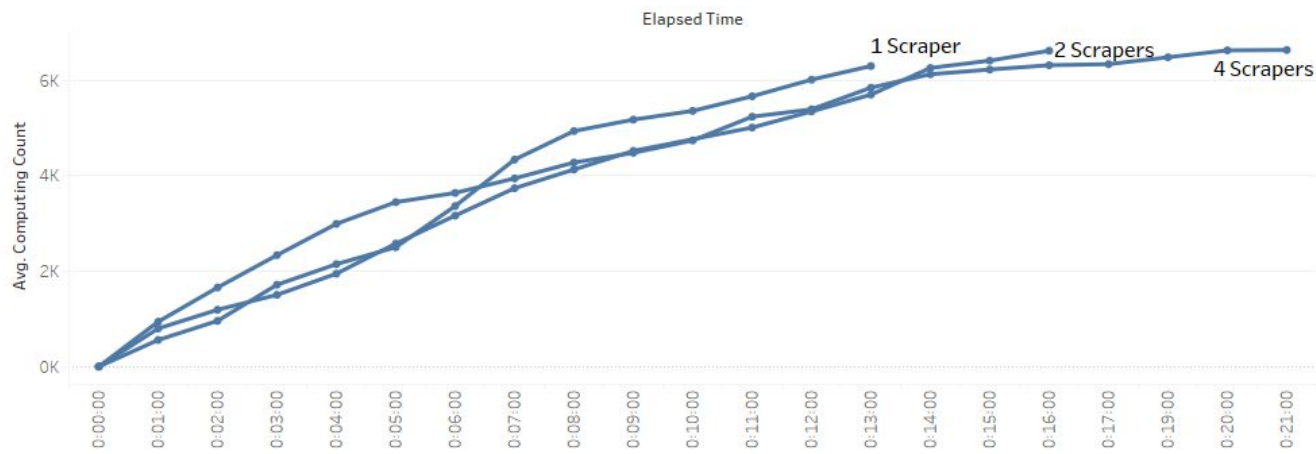
Avg Keyword Found Over Time by Scrapers (5 initial, 25 filter, 45 max)



Avg Keyword Found Over Time by Scrapers (10 initial, 25 filter, 45 max)



Avg Keyword Found Over Time by Scrapers (10 initial, 35 filter, 45 max)



Observations:

- The single scraper does worse with the lower filter depth
- The 10/35/45 configuration was consistently fastest for all configurations of scrapers



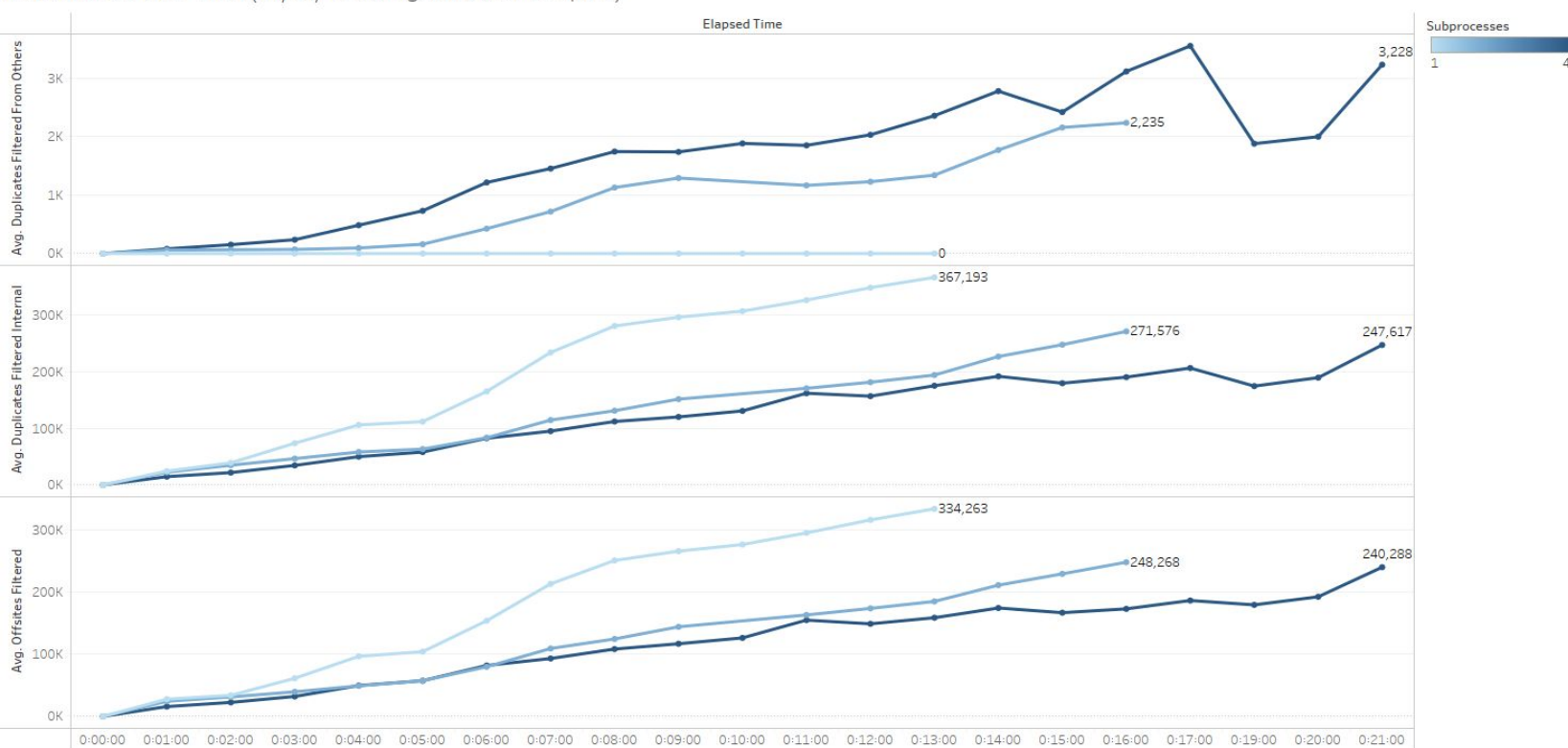
Crawling Results of Our Spider Configurations

To clarify the following graphs

Filtering behavior of spiders:

- **Filtered from others** = URLs filtered because another scraper has already processed it
- **Filtered internal** = URLs that the scraper itself has already encountered
- **Offsites filtered** = URLs that do not match the `cc.gatech.edu/` criteria

URLs Filtered Over Time (10/35/45 Config, AVG over scrapers)



The trends of average of Duplicates Filtered From Others, average of Duplicates Filtered Internal and average of Offsites Filtered for Elapsed Time. Color shows details about Subprocesses. The data is filtered on Configuration, which keeps 10 Initial, 35 Filter, 45 Max. The view is filtered on Subprocesses, which includes everything.

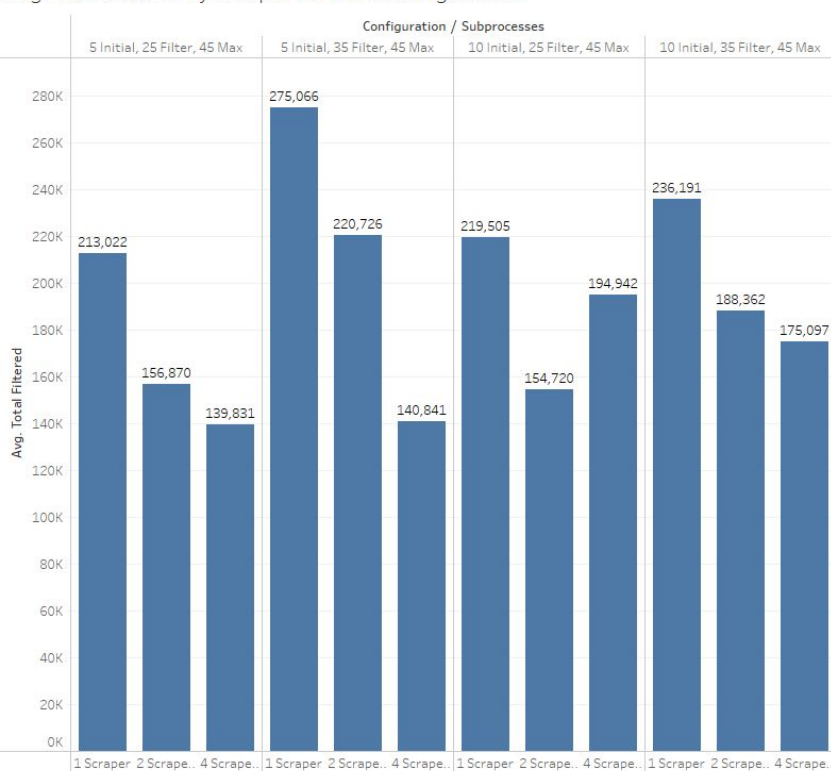
Trends:

- One Scraper's number of filtered results quickly grows, showing a high number of URLs it filters
- Adding more scrapers results in filters less requests but takes longer

Possible Explanations:

- Queries to check conflicts from other scrapers could be slowing down the scraper.

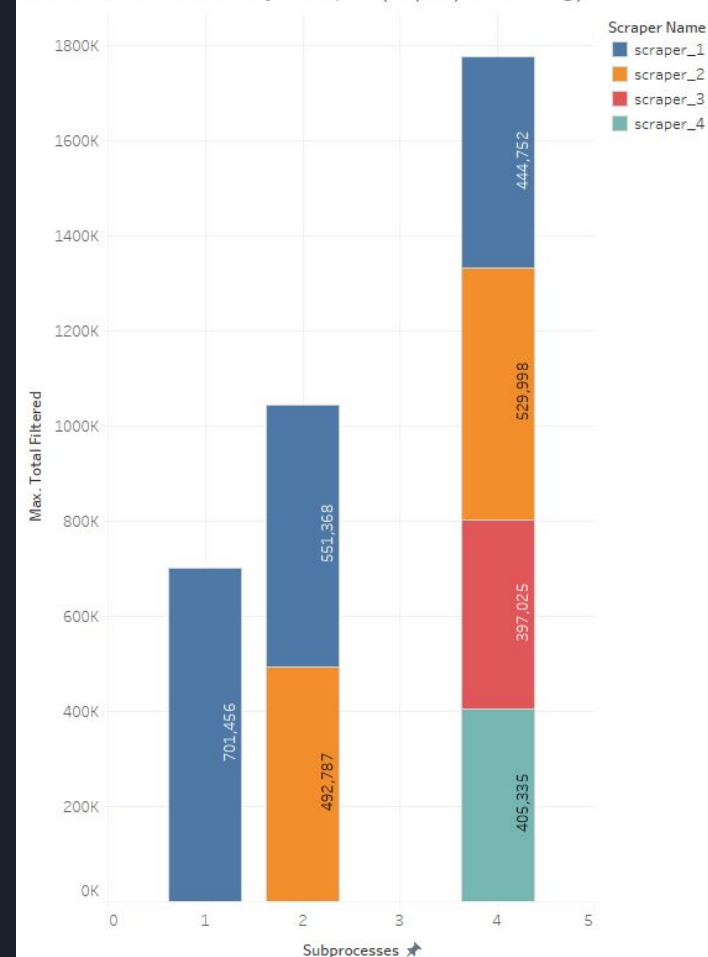
Avg Filtered URLs by Scraper for Each Configuration



Average of Total Filtered for each Subprocesses broken down by Configuration. The view is filtered on Configuration, which keeps 10 Initial, 25 Filter, 45 Max, 10 Initial, 35 Filter, 45 Max, 5 Initial, 25 Filter, 45 Max and 5 Initial, 35 Filter, 45 Max.

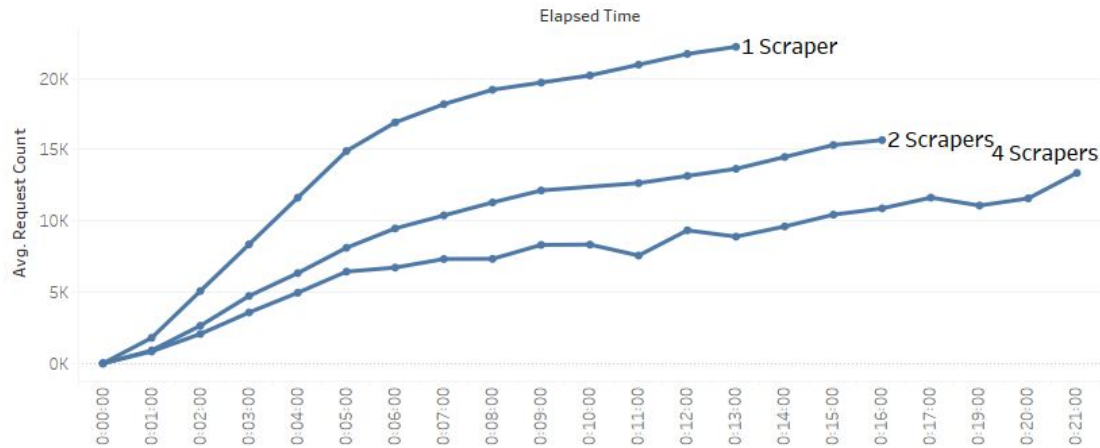
When we add more scrapers, each individual scraper filters less duplicates total, but adding all scrapers together shows that 4 scrapers is likely filtering the same URLs across multiple scrapers.

Total URLs Filtered By Scraper (10/35/45 Config)



The plot of maximum of Total Filtered for Subprocesses. Color shows details about Scraper Name. The data is filtered on Configuration, which keeps 10 Initial, 35 Filter, 45 Max. The view is filtered on Subprocesses, which includes everything.

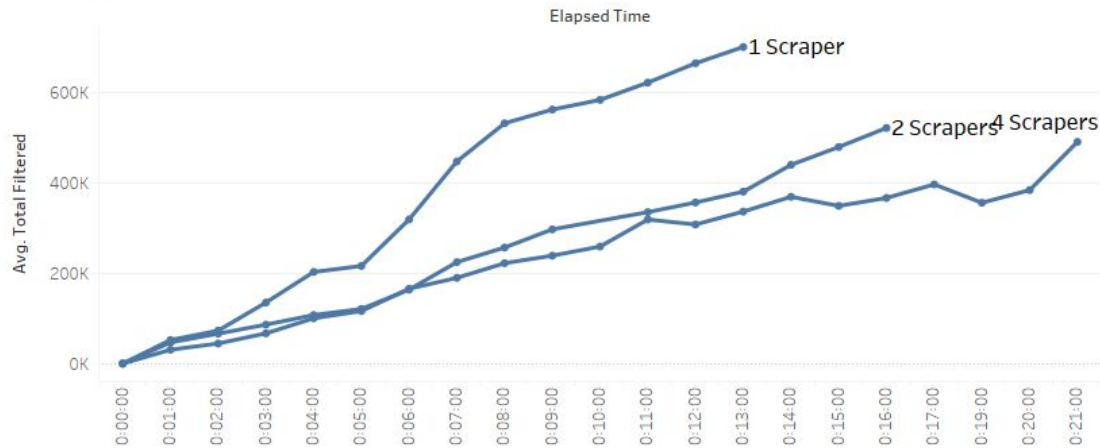
Average Request Count Over Time (10 Initial, 35 Filter, 45 Max)



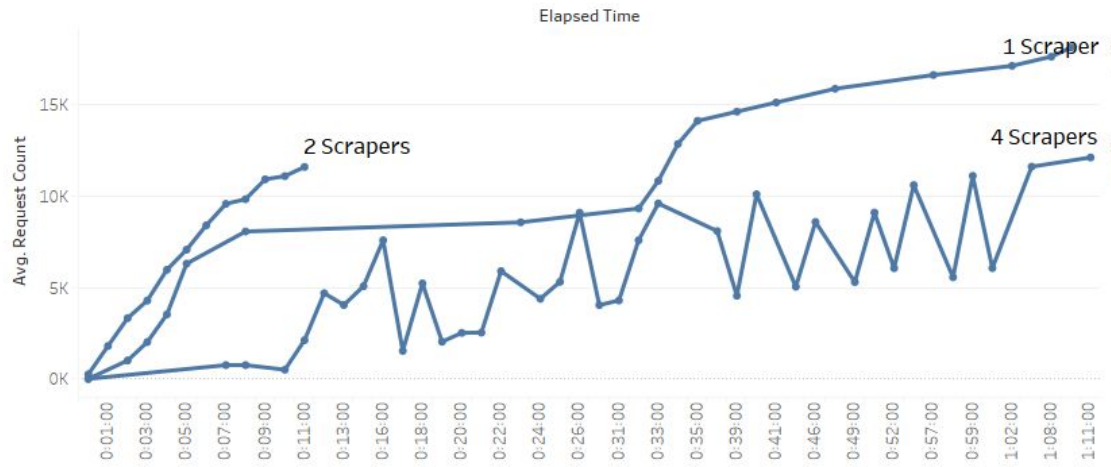
We see a mirroring of average requests made and average total filtered URLs indicating:

- Filtering more URLs doesn't slow down requests.

Average Total Filtered Count Over Time (10 Initial, 25 Filter, 45 Max)



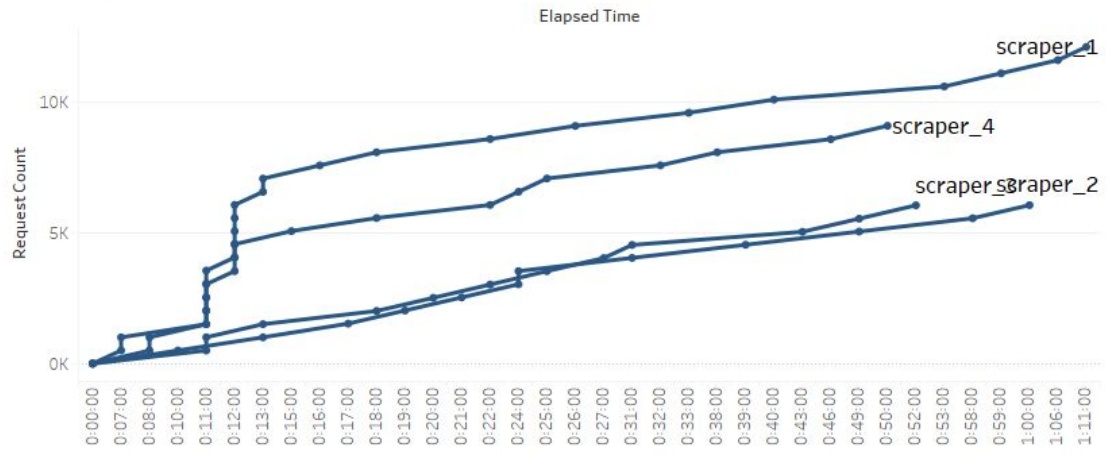
Average Request Count Over Time (5 Initial, 25 Filter, 45 Max)



Some of our data came out a bit sporadic, closer examination shows that some individual data points are missing, resulting in bad data on certain runs.

These runs and missing data points need to be investigated and conducted again.

Average Request Count Over Time for 4 Scrapers (5 Initial, 25 Filter, 45 Max)





Continued work & Aspirations

Setting up the crawling algorithm and getting the results into a useful and presentable form proved to be more difficult than originally anticipated.

Getting the crawlers to work in a truly distributed manner is not complete.

The list of keywords crawled thus far is rather short, so it is not yet useful as part of a search engine.

No real work has been done on the front end, as the data is not stored in a way that can be easily searched.

The dynamic recommendation system has so far eluded us.



Summary and Lessons Learned

Conclusions from the data analysis portion

- Efficiency of multiple crawlers
- Issues dealing with duplicate URLs

Lessons learned about the technologies

- Distributed systems
- Database configurations