

Assignment Project Exam Help

Add WeChat `edu_assist_pro`

Introduction to
Assignment Project Exam Help
Informa |
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`
Lecture 17: Crawling Indexes

Assignment Project Exam Help

Today's lecture

Add WeChat `edu_assist_pro`

- Crawling

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

Basic crawler

Add WeChat edu_assist_pro

- Begin with known “seed” URLs
- Fetch and parse them
 - Extract <https://eduassistpro.github.io/>
 - Place the extracted URLs in a queue
- Fetch each URL on the queue and repeat

Assignment Project Exam Help
<http://www.nature.com/nature/journal/v405/n6783/pdf/405112a0.pdf>

Structure of the Web (2000)

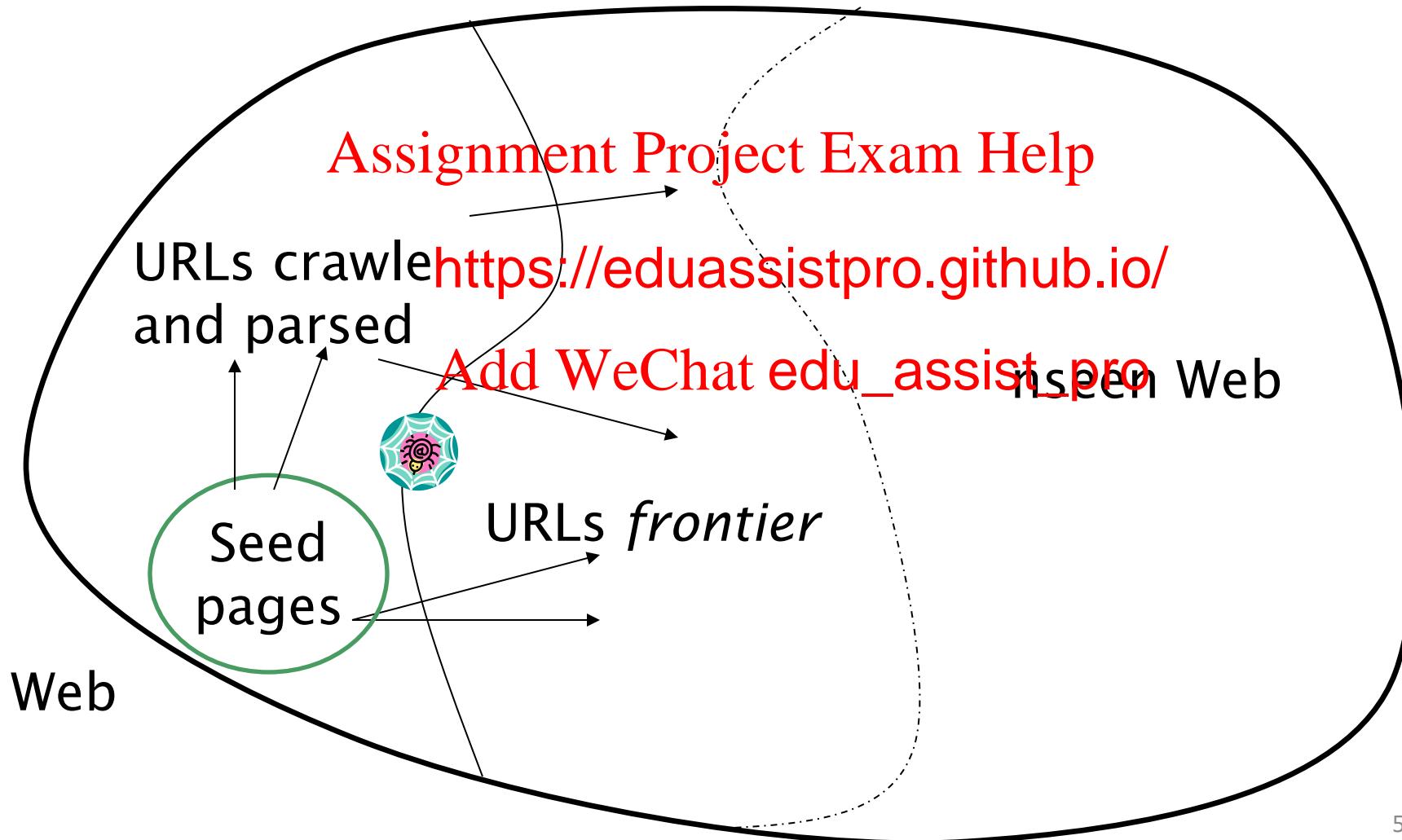
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

Crawling picture



Assignment Project Exam Help

Simple Crawler Thr Add WeChat edu_assist_pro

1

2

3 Assignment Project Exam Help

4

5

6 https://eduassistpro.github.io/

7

8

Add WeChat edu_assist_pro

9

10

11

12

13

14

Assignment Project Exam Help

Simple picture

- Web crawling isn't feasible with one machine
 - All of the above steps distributed
- Malicious
 - Assignment Project Exam Help
 - Spam pages
 - Spider traps<https://eduassistpro.github.io/>
- Even non-malicious
 - Assignment Project Exam Help
 - Latency/bandwidth to remote servers vary
 - Webmasters' stipulations
 - How "deep" should you crawl a site's URL hierarchy?
 - Site mirrors and duplicate pages
- Politeness – don't hit a server too often

Assignment Project Exam Help

What any Add WeChat edu_assist_pro crawler

- Be Polite: Respect implicit and explicit politeness considerations
 - Only cra <https://eduassistpro.github.io/>
 - Respect r this shortly)
- Be Robust: Be immune to r traps and other malicious behavior from web servers

Assignment Project Exam Help

What any crawlers do

- Be capable of distributed operation: designed to run on multiple distributed machines
- Be scalable: Assignment Project Exam Help the crawl rate by adding m <https://eduassistpro.github.io/>
- Performance/efficiency: ppe use of available processing and resources

Assignment Project Exam Help

What any crawlers do

- Fetch pages of “higher quality” first
- Continuous operation: Continue fetching fresh copies of each page
 - https://eduassistpro.github.io/
- Extensible: Add WeChat edu_assist_pro protocols

Assignment Project Exam Help

Freshness

Add WeChat edu_assist_pro

- HTTP protocol has a special request type called HEAD that makes it easy to check for page changes
 - returns information about page, not page itself

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

Freshness

Add WeChat edu_assist_pro

- Web pages are constantly being added, deleted, and modified
- Web crawler must continually revisit pages it has already crawl to maintain th
 - *stale* copies no longer reflect the contents of the web pages

Assignment Project Exam Help

<https://eduassistpro.github.io/>

changed in order

ument colle

Add WeChat edu_assist_pro

ntents of the web

Assignment Project Exam Help

Freshness

Add WeChat edu_assist_pro

- Not possible to constantly check all pages
 - must check important pages and pages that change frequently
- Freshness is the fraction of pages that are fresh
- Optimizing for freshness leads to bad decisions, such as not crawling popular sites
- Age is a better metric

Assignment Project Exam Help

Freshness Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

Age

Add WeChat [edu_assist_pro](#)

- Expected age of a page t days after it was last crawled:

[Assignment Project Exam Help](#)

- Web page update time until the next update is governed by an exponential distribution

<https://eduassistpro.github.io/>

button on

[Add WeChat \[edu_assist_pro\]\(#\)](#)

$$\text{Age}(\lambda, t) = \int_0^t \lambda e^{-\lambda x} (t - x) dx$$

Assignment Project Exam Help

Age

Add WeChat `edu_assist_pro`

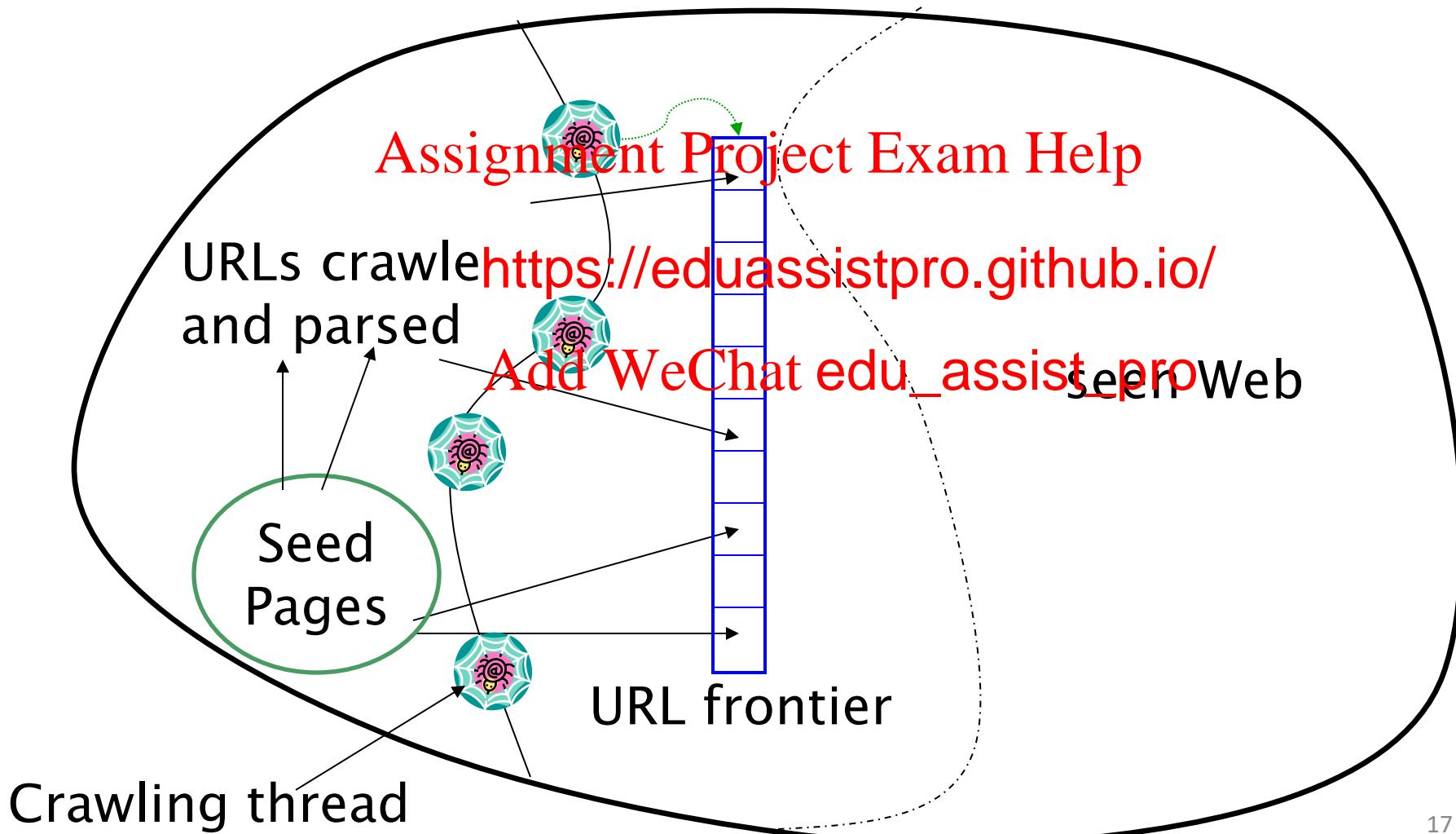
- The older a page gets, the more it costs not to crawl it
 - e.g., expected age with mean change frequency $\lambda = 1/7$ (one change)

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

Updated crawling ~~Add WeChat~~ ~~edu_assist_pro~~



Assignment Project Exam Help

Simple Crawler Thr Add WeChat edu_assist_pro

```
1  
2  
3 Assignment Project Exam Help  
4  
5 https://eduassistpro.github.io/  
6  
7 Add WeChat edu_assist_pro  
8  
9  
10  
11  
12  
13  
14
```

Assignment Project Exam Help

URL frontier

- Can include multiple pages from the same host
 - Add WeChat edu_assist_pro
- Must avoid em all at the same time
 - Assignment Project Exam Help
 - https://eduassistpro.github.io/
- Must try to keep all cr reads busy
 - Add WeChat edu_assist_pro

Assignment Project Exam Help

Explicit and Implicit Politeness

- Explicit politeness: specifications from webmasters on what portions of site can be crawled
 - robots.txt <https://eduassistpro.github.io/>
- Implicit politeness: every specification, avoid hitting any site too often

Assignment Project Exam Help

Robots.txt Add WeChat edu_assist_pro

- Protocol for giving spiders (“robots”) limited access to a website, originally from 1994
 - www.robotstxt.org/wc/norobots.html
- Website ann <https://eduassistpro.github.io/> what can(not) be crawled Add WeChat edu_assist_pro
 - For a URL, create a file URL/robots.txt
 - This file specifies access restrictions

Assignment Project Exam Help

Robots.txt Example

- No robot should visit any URL starting with "/yoursite/temp/", except the robot called "searchengine".

<https://eduassistpro.github.io/>

User-agent: *

Disallow: /yoursite/temp/

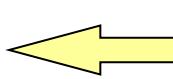
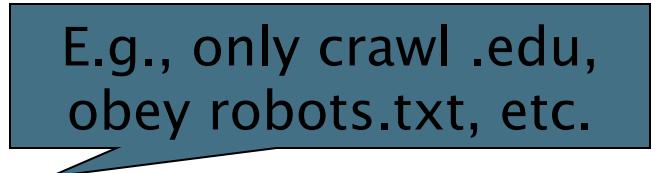
User-agent: searchengine

Disallow:

try <http://www.taobao.com/robots.txt>

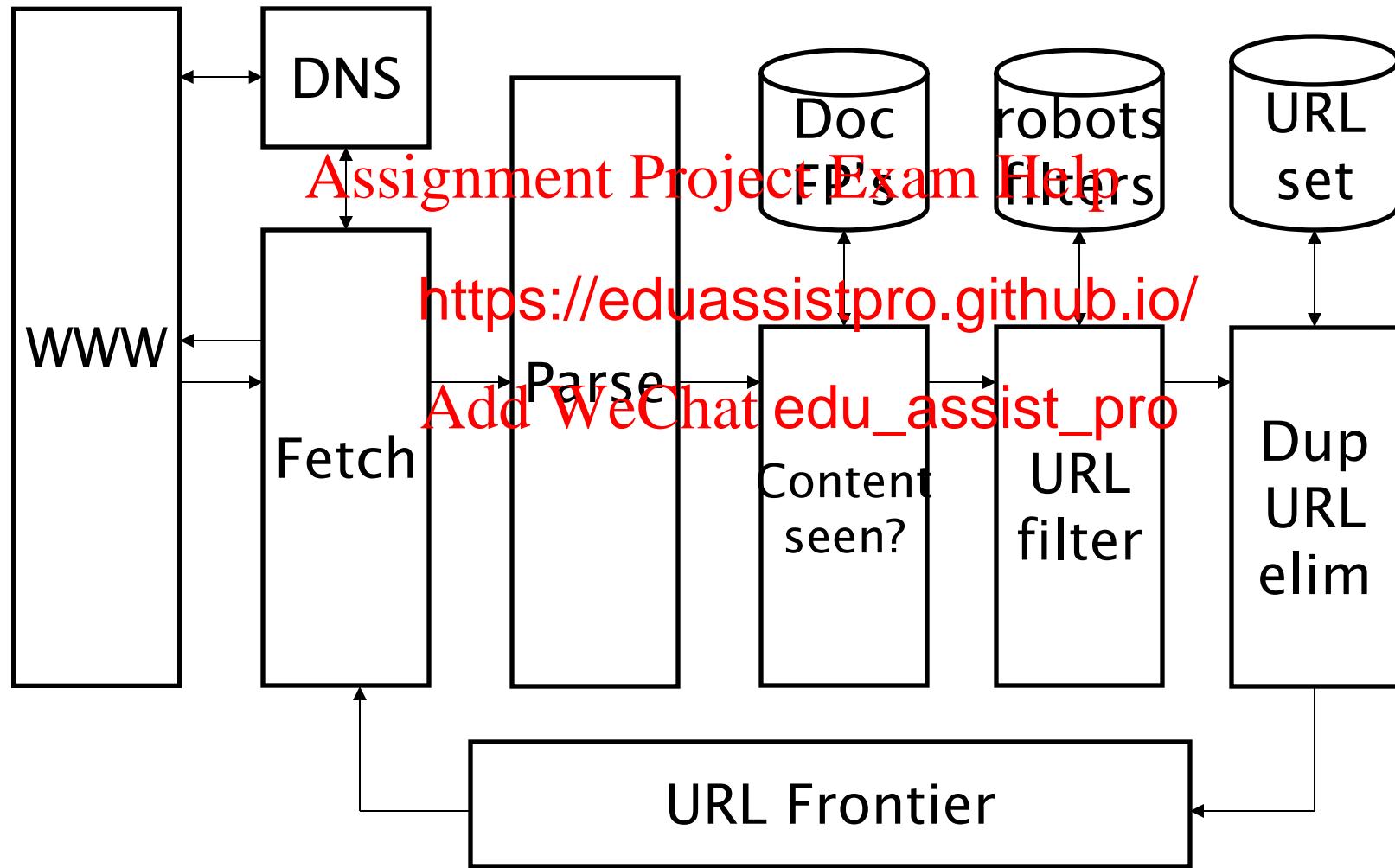
Assignment Project Exam Help

Processing steps in

- Pick a URL from the frontier  Which one?
- Fetch the document at the URL
 Assignment Project Exam Help
- Parse the URL
 - Extract links <https://eduassistpro.github.io/>
- Check if URL has content  Add WeChat to edu_assist_pro
 - If not, add to indexes
- For each extracted URL 
 - Ensure it passes certain URL filter tests
 - Check if it is already in the frontier (duplicate URL elimination)

Assignment Project Exam Help

Basic crawl architecture



Assignment Project Exam Help

DNS (Domain Name System)

- A lookup service on the internet
 - Given a URL, retrieve its IP address
 - Service provided by a distributed set of servers – thus, lookup latency (seconds)
<https://eduassistpro.github.io/>
- Common OS implementations of DNS lookup are *blocking*: only one outstanding request at a time
- Solutions
 - DNS caching
 - Batch DNS resolver – collects requests and sends them out together

Assignment Project Exam Help

Parsing: URL norm

- When a fetched document is parsed, some of the extracted links are *relative URLs*
- E.g., at [https://eduassistpro.github.io/](http://eduassistpro.github.io/)
we have a relative https://eduassistpro.github.io/wikipedia/Wikipedia:General_disclaimer which is the same as the absolute URL
http://en.wikipedia.org/wiki/Wikipedia:General_disclaimer
- During parsing, must normalize (expand) such relative URLs

Assignment Project Exam Help

Content seen? Add WeChat edu_assist_pro

- Duplication is widespread on the web
- If the page just fetched is already in the index, process it
- This is verified using fingerprints or shingles

Assignment Project Exam Help

Removing Noise

- Many web pages contain text, links, and pictures that are not directly related to the main content of the page
Add WeChat edu_assist_pro
- This additional <https://eduassistpro.github.io/> could negatively affect the page
Add WeChat edu_assist_pro
- Techniques have been developed to detect the content blocks in a web page
 - Non-content material is either ignored or reduced in importance in the indexing process

Assignment Project Exam Help

Add WeChat `edu_assist_pro`

Noise Example

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

Finding Content BI

- Cumulative distribution of tags in the example web page

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Main text content of the page corresponds to the “plateau” in the middle of the distribution

Assignment Project Exam Help

Finding Content BI

- Represent a web page as a sequence of bits, where $b_n = 1$ indicates that the n th token is a tag
- **Optimization problem** where we find values of i and j to maximize $\sum_{n=0}^{i-1} b_n + \sum_{n=i}^j (1 - b_n) + \sum_{n=j+1}^{N-1} b_n$
- i.e., maximize

$$\sum_{n=0}^{i-1} b_n + \sum_{n=i}^j (1 - b_n) + \sum_{n=j+1}^{N-1} b_n$$

Assignment Project Exam Help

Finding Content BI

- Other approaches use DOM structure and visual (layout) features

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

Filters and robot

- Filters – regular expressions for URL's to be crawled/not
- Once a robots.txt file from a site, need to do so carefully
 - Doing so burns ban its web server
- Cache robots.txt files

Assignment Project Exam Help

Duplicate URL elim

- For a non-continuous (one-shot) crawl, test to see if an extracted+filtered URL has already been seen
- For a continuous crawl, implement frontier management

Assignment Project Exam Help

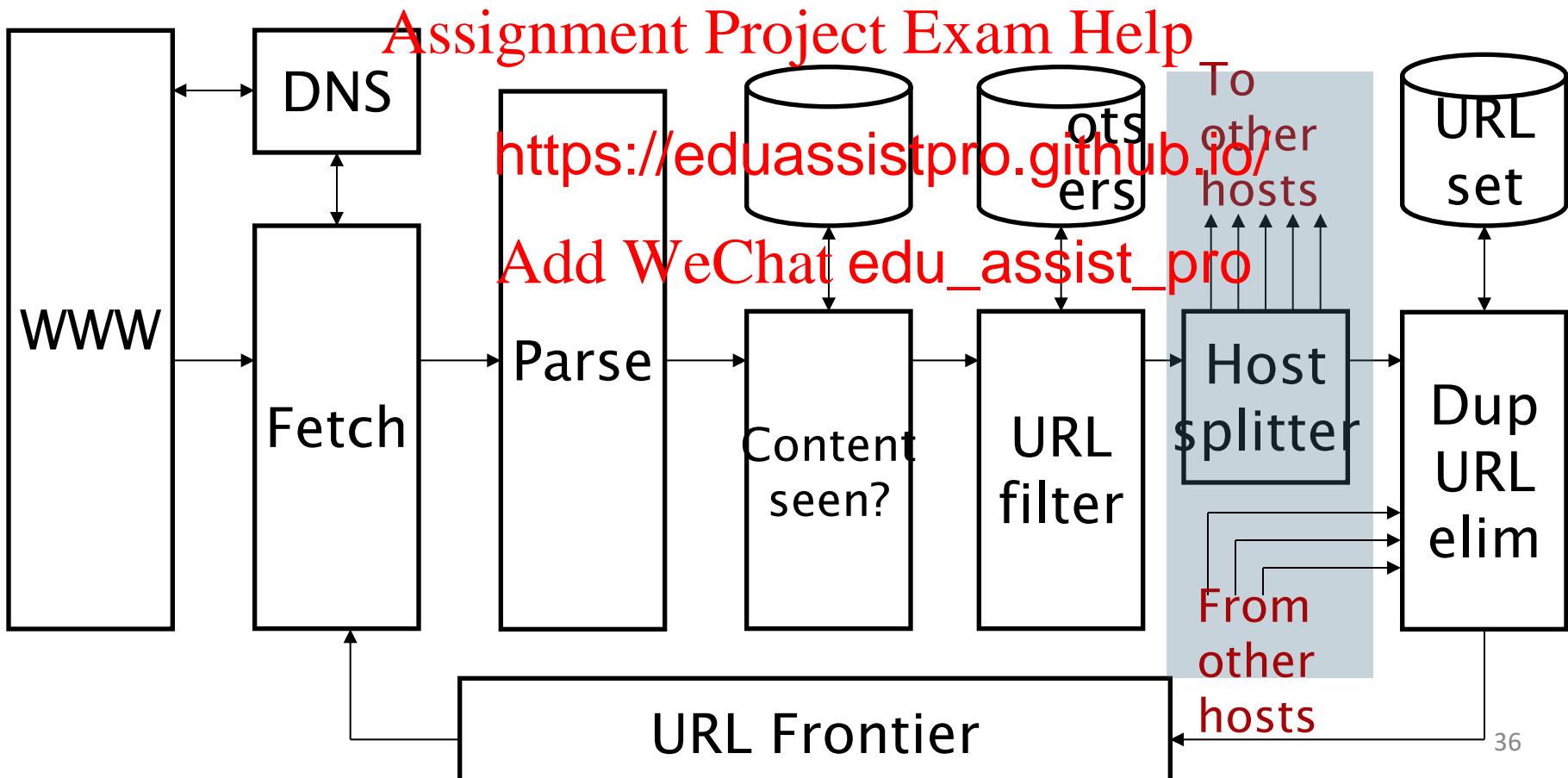
Distributing the crawl Add WeChat edu_assist_pro

- Run multiple crawl threads, under different processes – potentially at different nodes
 - Geographically distributed nodes
- Partition hosts <https://eduassistpro.github.io/nodes/>
 - Hash used for partitioning Add WeChat edu_assist_pro
- How do these nodes communicate?

Assignment Project Exam Help

Communication between nodes

- The output of the URL filter at each node is sent to the Duplicate URL Eliminator at all nodes



Assignment Project Exam Help

URL frontier: two main generations

- Politeness: do not hit a web server too frequently
- Freshness: crawl some pages more often than others
 - E.g., pages <https://eduassistpro.github.io/> those content changes often

These goals may conflict each other

(E.g., simple priority queue fails – many links out of a page go to its own site, creating a burst of accesses to that site.)

Assignment Project Exam Help

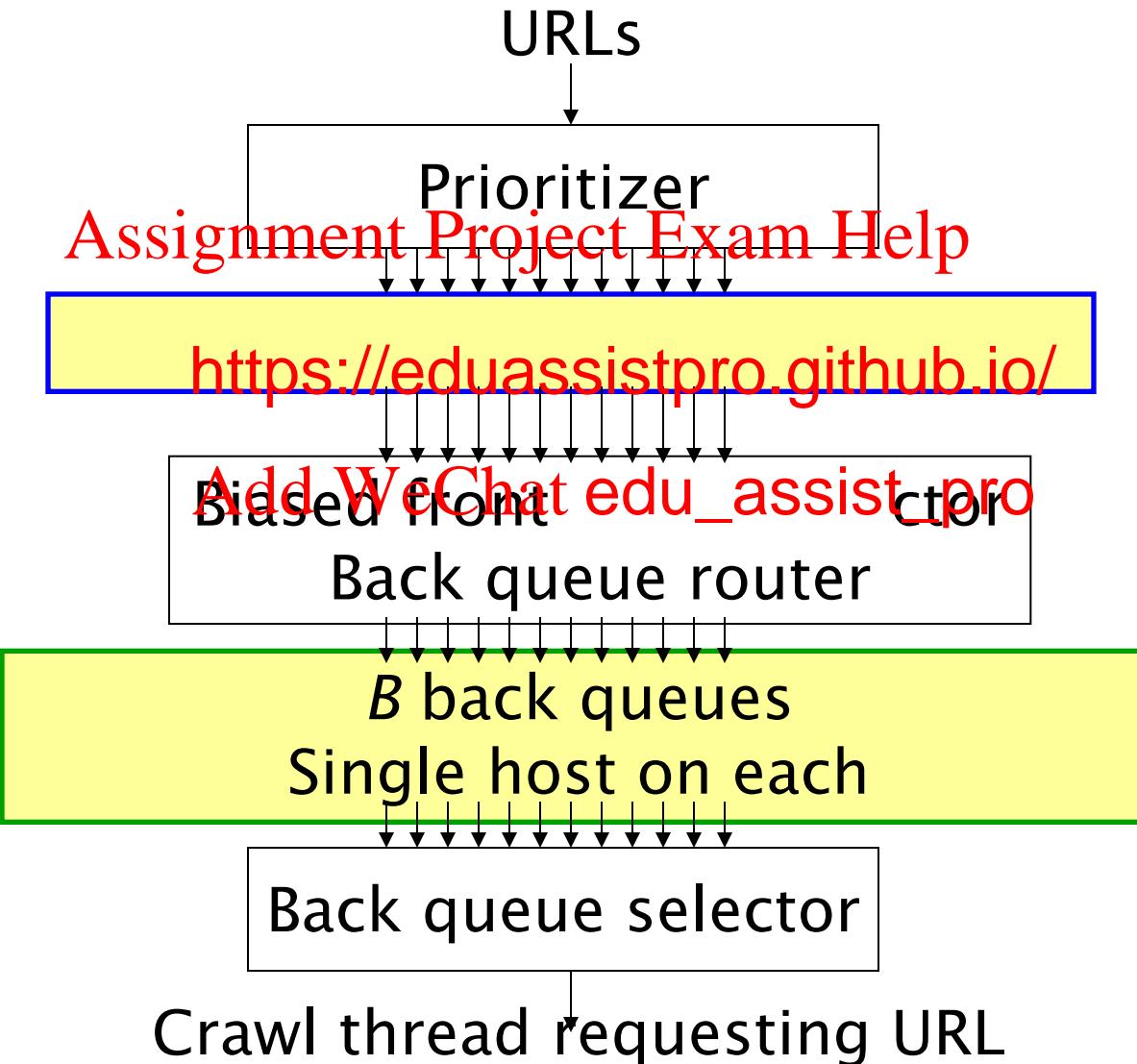
Politeness

Add WeChat edu_assist_pro

- Even if we restrict only one thread to fetch from a host, can hit it repeatedly
- Common https://eduassistpro.github.io/ successive requests have a gap between that is >> time for most recent fetch from host

Assignment Project Exam Help

URL frontier: Merceme



Assignment Project Exam Help

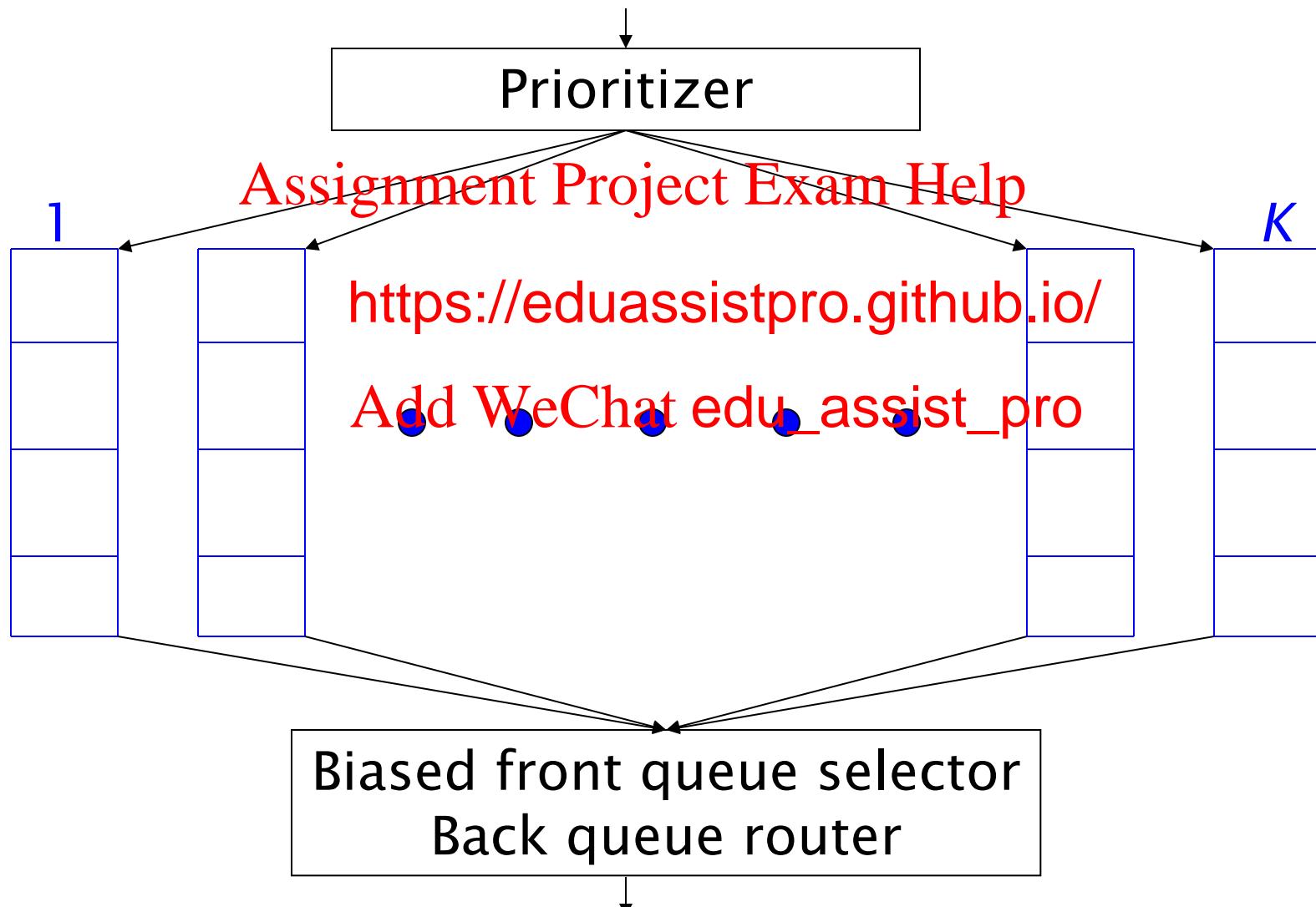
Mercator URL frontier

- URLs flow in from the top into the frontier
- Front queues manage prioritization
 - Assignment Project Exam Help
- Back queues
- Each queue i <https://eduassistpro.github.io/>
 - Add WeChat edu_assist_pro

Assignment Project Exam Help

Front queues

Add WeChat edu_assist_pro



Assignment Project Exam Help

Front queues

Add WeChat edu_assist_pro

- Prioritizer assigns to URL an integer priority between 1 and K
 - Appends URL to corresponding queue
- Heuristics for <https://eduassistpro.github.io/>
 - Refresh rate
 - Application-specific (e.g., “crawl news sites more often”)

Assignment Project Exam Help

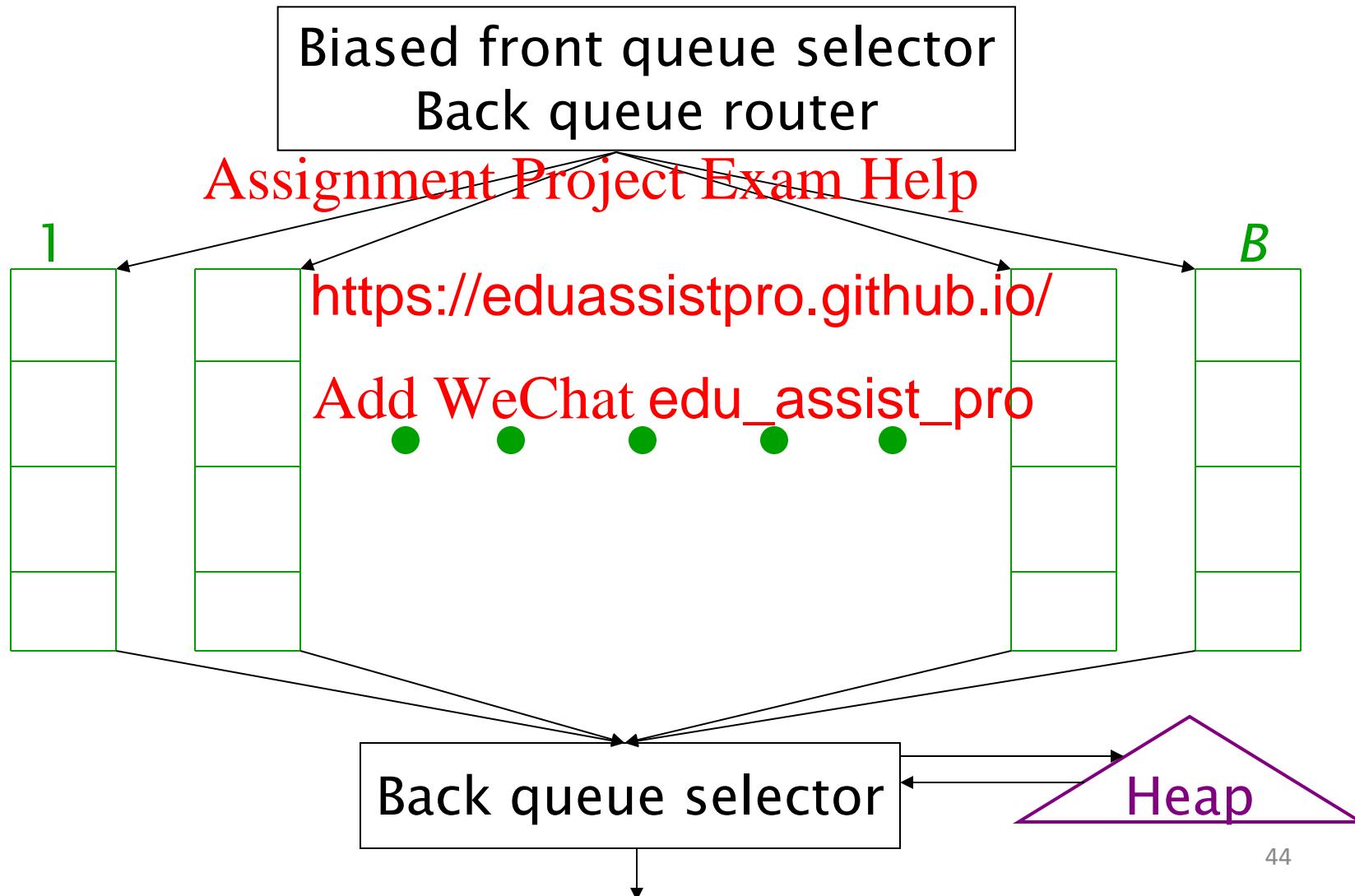
Biased front queue

- When a back queue requests a URL (in a sequence to be described): picks a front queue from which to pull a URL
- This choice can be based on queues of higher priority, or some sophisticated variant
 - Can be randomized

Assignment Project Exam Help

Back queues

Add WeChat edu_assist_pro



Assignment Project Exam Help

Back queue Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

- Each back queue is kept non-empty while the crawl is in progress
- Each back RLS from a single host <https://eduassistpro.github.io/>
 - Maintain a table from back queues

Host name	Back queue
...	3
	1
	B

Assignment Project Exam Help

Back queue ~~Add WeChat edu_assist_pro~~ heap

- One entry for each back queue
- The entry is the earliest time t_e at which the host ~~Assignment Project Exam Help~~ can be hit again corresponding
- This earliest <https://eduassistpro.github.io/>
 - Last access ~~Add WeChat edu_assist_pro~~
 - Any time buffer heuristic we choose

Assignment Project Exam Help

Back queue process

- A crawler thread seeking a URL to crawl:
- Extracts the root of the heap
- Fetches URL (look up from <https://eduassistpro.github.io/>)
- Checks if queue q is non-empty, pulls a URL v from front queues
 - If there's already a back queue for v 's host, append v to q and pull another URL from front queues, repeat
 - Else add v to q
- When q is non-empty, create heap entry for it

Assignment Project Exam Help

Number of back qu

- Keep all threads busy while respecting politeness
- Mercator recommendation: three times as many back queues

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

Resources

- IIR Chapter 20
- [Mercator: A scalable, extensible web crawler \(Heydon et al. 1999\)](#) Assignment Project Exam Help
- [A standard fo](#) <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro