

### Introduction to Information Retrieval <http://informationretrieval.org>

#### IIR 20: Crawling

Hinrich Schütze

Institute for Natural Language Processing, Universität Stuttgart

2008.07.08

- To fetch 1,000,000,000 pages in one month ...
- ... we need to fetch almost 400 pages per second!
- Actually: many more since many of the pages we attempt to crawl will be duplicates, unfetchable, spam etc.

1 / 25

2 / 25

## Basic crawler operation

- Initialize queue with URLs of known seed pages
- Repeat
  - Take URL from queue
  - Fetch and parse page
  - Extract URLs from page
  - Add URLs to queue
- Fundamental assumption: The web is well linked.

## Complications in crawling

- We need many machines – how do we distribute?
- Latency/bandwidth
- How deep should we crawl sites?
- Duplicates
- Spam and spider traps
- Politeness – don't hit a server too often

3 / 25

4 / 25

## Spider trap

- Malicious server that generates an infinite sequence of linked pages
- Sophisticated spider traps generate pages that are not easily identified as dynamic.

## What a crawler must do

- Be polite
  - Don't hit a each site too often
  - Only crawl pages you are allowed to crawl: robots.txt
- Be robust
  - Be immune to spider traps, duplicates, very large pages, very large websites, dynamic pages etc

5 / 25

6 / 25

## Robots.txt

- Protocol for giving crawlers ("robots") limited access to a website, originally from 1994
- Examples:
  - User-agent: \*
  - Disallow: /yoursite/temp/
  - User-agent: searchengine
  - Disallow:
- Important: cache the robots.txt file of each site we are crawling

## What any crawler should do

- Be capable of **distributed** operation
- Be scalable: need to be able to increase crawl rate by adding more machines
- Fetch pages of higher quality first
- Continuous operation: get fresh version of already crawled pages

7 / 25

8 / 25

## URL frontier

- Can include multiple pages from the same host
- Must avoid trying to fetch them all at the same time
- Must keep all crawling threads busy

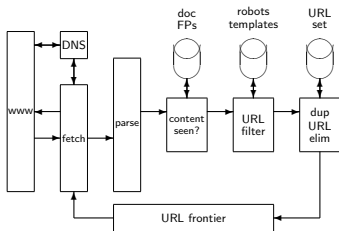
9 / 25

10 / 25

## Processing steps in crawling

- Pick a URL from the frontier
- Fetch the document at the URL
- Check if the document has content already seen (if yes: skip following steps)
- Index document
- Parse the document and extract URLs to other docs
- For each extracted URL:
  - Does it fail certain tests (e.g., spam)? Yes: skip
  - Already in the frontier? Yes: skip

## Basic crawl architecture



11 / 25

12 / 25

- Some URLs extracted from a document are [relative](#) URLs.
- E.g., at <http://mit.edu>, we may have </aboutsitem.html>
  - This is the same as: <http://mit.edu/aboutsitem.html>
- During parsing, we must normalize (expand) all relative URLs.

13 / 25

- For each page fetched: check if the content is already in the index
- Check this using document fingerprints or [shingles](#)
- Skip documents whose content has already been indexed

14 / 25

## Distributing the crawler

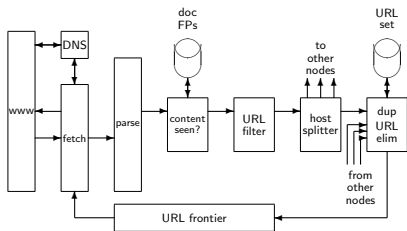
- Run multiple crawl threads, potentially at different nodes
  - Usually geographically distributed nodes
- Partition hosts being crawled into nodes

15 / 25

Google data centers:

<http://www.wayfaring.com/maps/show/48030>

16 / 25

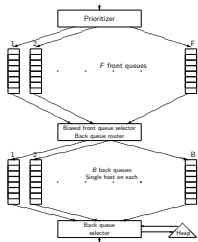


- Politeness: Don't hit a web server too frequently
  - E.g., insert a time gap between successive requests to the same server
- Freshness: Crawl some pages (e.g., news sites) more often than others
- Not an easy problem: simple priority queue fails.
- Why?

17 / 25

18 / 25

## Mercator URL frontier



- URLs flow in from the top into the frontier.
- Front queues manage prioritization.
- Back queues enforce politeness.
- Each queue is FIFO.

19 / 25

## Front queues

- Prioritizer assigns to URL an integer priority between 1 and  $K$ .
- Then appends URL to corresponding queue
- Heuristics for assigning priority: refresh rate, PageRank etc

20 / 25

- Selection from front queues is initiated by back queues (see below)
- Pick a front queue from which to select next URL
  - Round robin
  - Randomly
  - Or more sophisticated variant
  - But with a bias in favor of high-priority front queues

21 / 25

- Each back queue is kept non-empty while the crawl is in progress.
- Each back queue only contains URLs from a single host.
- Maintain a table from hosts to back queues.

22 / 25

- One entry for each back queue
- The entry is the earliest time  $t_e$  at which the host corresponding to the back queue can be hit again.
- This earliest time is determined by (i) last access to that host (ii) time gap heuristic

23 / 25

- Extract the root of the back queue heap
- Fetch the URL at head of corresponding back queue  $q$
- Check if  $q$  is now empty
- If yes: keep (i) pulling URLs from the front queues and (ii) adding them to their corresponding back queues until ...
- ... the URL's host does not have a back queue – then put the URL in  $q$  and create heap entry for it.

24 / 25

## Resources

- Chapter 20 of IIR
- Resources at <http://ifnlp.org/ir>
- Paper on Mercator by Heydon et al.
- Robot exclusion standard