

CIT-496

Chapter-3

Search Engines

Information Retrieval in
Practice

Web Crawler

- Finds and downloads web pages automatically
- provides the collection for searching
- Web is huge and constantly growing
- Web is not under the control of search engine providers
- Web pages are constantly changing
- Crawlers also used for other types of data

Retrieving Web Pages

- Every page has a unique *uniform resource locator* (URL)
- Web pages are stored on web servers that use HTTP to exchange information with client software
- e.g.,

http://www.cs.umass.edu/csinfo/people.html

The diagram illustrates the components of the URL `http://www.cs.umass.edu/csinfo/people.html`. Three double-headed arrows point from the labels below to their corresponding parts in the URL above:

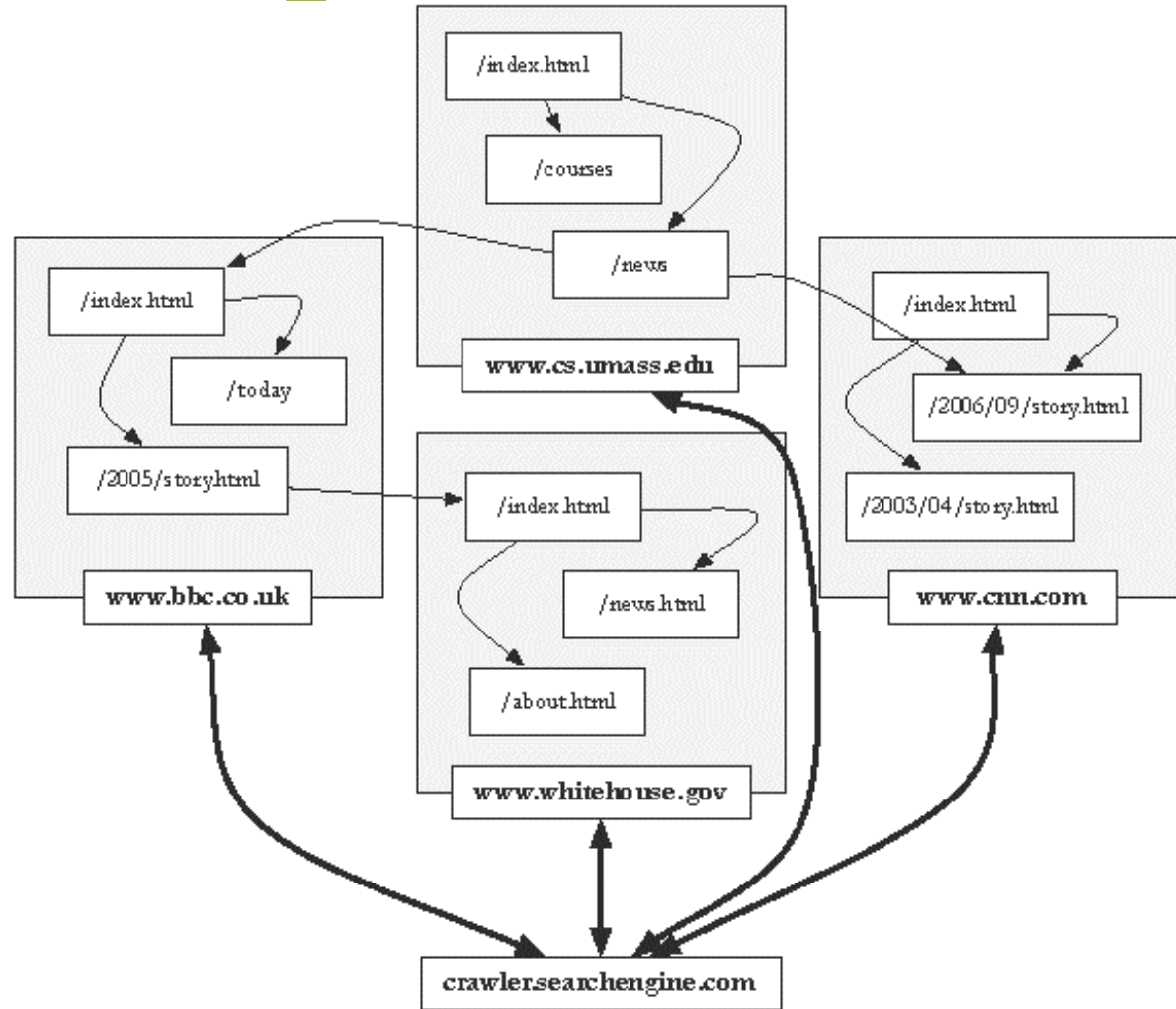
- http** is connected to the `http` part of the URL.
- hostname** is connected to `www.cs.umass.edu`.
- resource** is connected to `/csinfo/people.html`.

scheme **hostname** **resource**

Retrieving Web Pages

- ◉ Web crawler client program connects to a *domain name system (DNS)* server
- ◉ DNS server translates the hostname into an *internet protocol (IP)* address
- ◉ Crawler then attempts to connect to server host using specific *port*
- ◉ After connection, crawler sends an HTTP request to the web server to request a page
 - ◉ usually a GET request

Crawling the Web



Web Crawler

- Starts with a set of *seeds*, which are a set of URLs given to it as parameters
- Seeds are added to a URL request queue
- Crawler starts fetching pages from the request queue
- Downloaded pages are parsed to find link tags that might contain other useful URLs to fetch
- New URLs added to the crawler's request queue, or *frontier*
- Continue until no more new URLs or disk full

Web Crawling

- Web crawlers spend a lot of time waiting for responses to requests
- To reduce this inefficiency, web crawlers use threads and fetch hundreds of pages at once
- Crawlers could potentially flood sites with requests for pages
- To avoid this problem, web crawlers use *politeness policies*
 - e.g., delay between requests to same web server

Controlling Crawling

- Even crawling a site slowly will anger some web server administrators, who object to any copying of their data
- Robots.txt file can be used to control crawlers

```
User-agent: *  
Disallow: /private/  
Disallow: /confidential/  
Disallow: /other/  
Allow: /other/public/
```

```
User-agent: FavoredCrawler  
Disallow:
```

```
Sitemap: http://mysite.com/sitemap.xml.gz
```


Freshness

- Web pages are constantly being added, deleted, and modified
- Web crawler must continually revisit pages it has already crawled to see if they have changed in order to maintain the *freshness* of the document collection

Freshness

- 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

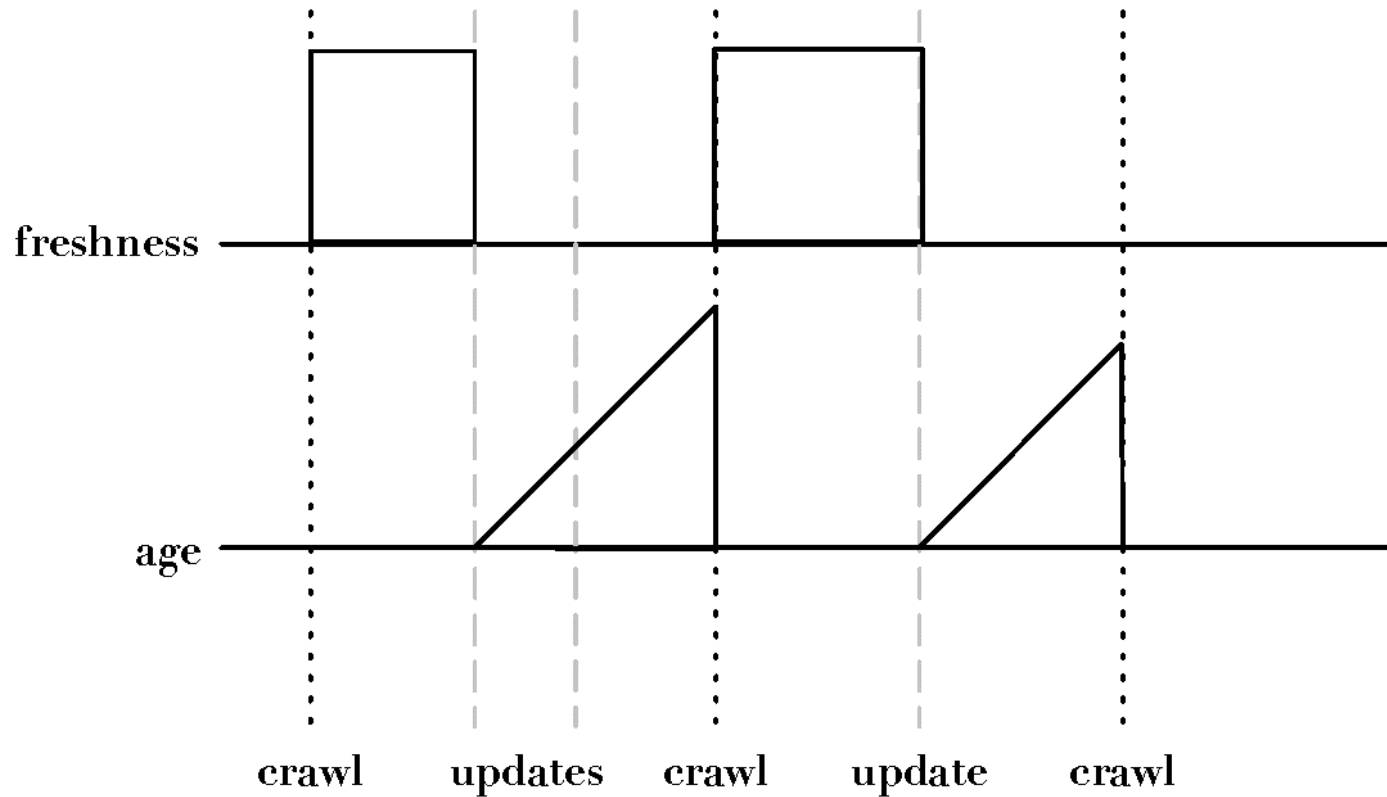
```
Client request: HEAD /csinfo/people.html HTTP/1.1
Host: www.cs.umass.edu

HTTP/1.1 200 OK
Date: Thu, 03 Apr 2008 05:17:54 GMT
Server: Apache/2.0.52 (CentOS)
Last-Modified: Fri, 04 Jan 2008 15:28:39 GMT
Server response: ETag: "239c33-2576-2a2837c0"
Accept-Ranges: bytes
Content-Length: 9590
Connection: close
Content-Type: text/html; charset=ISO-8859-1
```

Freshness

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

Freshness vs. Age



Focused Crawling

- Attempts to download only those pages that are about a particular topic
 - used by *vertical search* applications
- Rely on the fact that pages about a topic tend to have links to other pages on the same topic
 - popular pages for a topic are typically used as seeds
- Crawler uses *text classifier* to decide whether a page is on topic

Deep Web

- Sites that are difficult for a crawler to find are collectively referred to as the *deep* (or *hidden*) Web
 - much larger than conventional Web
- Three broad categories:
 - private sites
 - no incoming links, or may require log in with a valid account
 - form results
 - sites that can be reached only after entering some data into a form
 - scripted pages
 - pages that use JavaScript, Flash, or another client-side language to generate links

Sitemaps

- Sitemaps contain lists of URLs and data about those URLs, such as modification time and modification frequency
- Generated by web server administrators
- Tells crawler about pages it might not otherwise find
- Gives crawler a hint about when to check a page for changes

Sitemap Example

```
<?xml version="1.0" encoding="UTF-8"?>
<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9">
  <url>
    <loc>http://www.company.com/</loc>
    <lastmod>2008-01-15</lastmod>
    <changefreq>monthly</changefreq>
    <priority>0.7</priority>
  </url>
  <url>
    <loc>http://www.company.com/items?item=truck</loc>
    <changefreq>weekly</changefreq>
  </url>
  <url>
    <loc>http://www.company.com/items?item=bicycle</loc>
    <changefreq>daily</changefreq>
  </url>
</urlset>
```


Distributed Crawling

- Three reasons to use multiple computers for crawling
 - Helps to put the crawler closer to the sites it crawls
 - Reduces the number of sites the crawler has to remember
 - Reduces computing resources required
- Distributed crawler uses a hash function to assign URLs to crawling computers
 - hash function should be computed on the host part of each URL

Desktop Crawls

- Used for desktop search and enterprise search
- Differences to web crawling:
 - Much easier to find the data
 - Responding quickly to updates is more important
 - Must be conservative in terms of disk and CPU usage
 - Many different document formats
 - Data privacy very important

Conversion

- Text is stored in hundreds of incompatible file formats
 - e.g., raw text, RTF, HTML, XML, Microsoft Word, ODF, PDF
- Other types of files also important
 - e.g., PowerPoint, Excel
- Typically use a conversion tool
 - converts the document content into a tagged text format such as HTML or XML
 - retains some of the important formatting information

Character Encoding

- A character encoding is a mapping between bits and glyphs
 - i.e., getting from bits in a file to characters on a screen
 - Can be a major source of incompatibility
- ASCII is basic character encoding scheme for English
 - encodes 128 letters, numbers, special characters, and control characters in 7 bits, extended with an extra bit for storage in bytes

Character Encoding

- Other languages can have many more glyphs
 - e.g., Chinese has more than 40,000 characters, with over 3,000 in common use
- Many languages have multiple encoding schemes
 - e.g., CJK (Chinese-Japanese-Korean) family of East Asian languages, Hindi, Arabic
 - must specify encoding
 - can't have multiple languages in one file
- Unicode developed to address encoding problems

Unicode

- Single mapping from numbers to glyphs that attempts to include all glyphs in common use in all known languages
- Unicode is a mapping between numbers and glyphs
 - does not uniquely specify bits to glyph mapping!
 - e.g., UTF-8, UTF-16, UTF-32

Storing the Documents

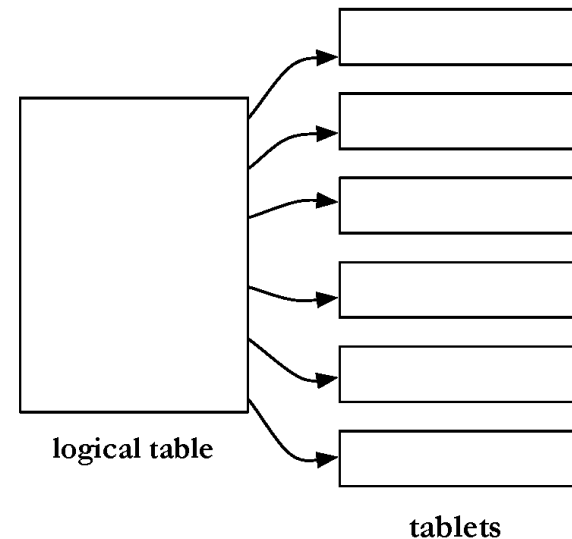
- Many reasons to store converted document text
 - saves crawling time when page is not updated
 - provides efficient access to text for snippet generation, information extraction, etc.
- Database systems can provide document storage for some applications
 - web search engines use customized document storage systems

Compression

- Text is highly redundant (or predictable)
- Compression techniques exploit this redundancy to make files smaller without losing any of the content
- Compression of indexes covered later
- Popular algorithms can compress HTML and XML text by 80%
 - e.g., DEFLATE (zip, gzip) and LZW (UNIX compress, PDF)
 - may compress large files in blocks to make access faster

BigTable

- Google's document storage system
 - Customized for storing, finding, and updating web pages
 - Handles large collection sizes using inexpensive computers

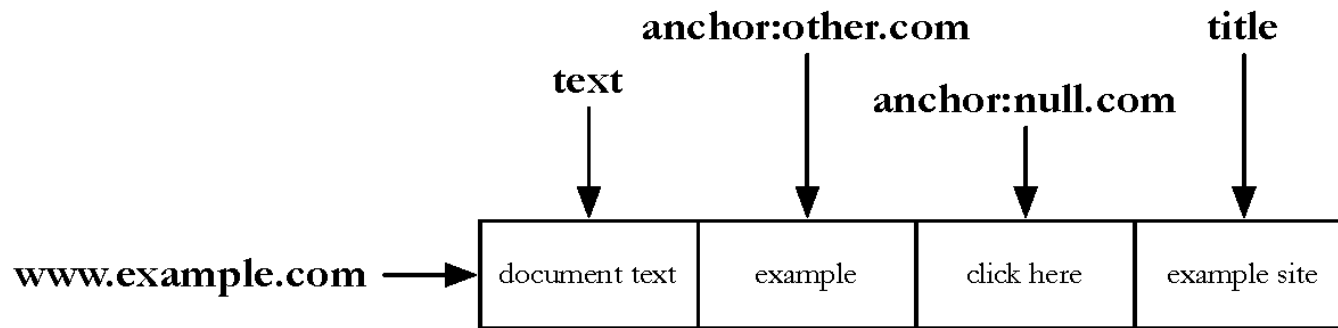


BigTable

- No query language, no complex queries to optimize
- Only row-level transactions
- Tablets are stored in a replicated file system that is accessible by all BigTable servers
- Any changes to a BigTable tablet are recorded to a transaction log, which is also stored in a shared file system
- If any tablet server crashes, another server can immediately read the tablet data and transaction log from the file system and take over

BigTable

- Logically organized into rows
- A row stores data for a single web page



- Combination of a row key, a column key, and a timestamp point to a single cell in the row

Detecting Duplicates

- Duplicate and near-duplicate documents occur in many situations
 - Copies, versions, plagiarism, spam, mirror sites
 - 30% of the web pages in a large crawl are exact or near duplicates of pages in the other 70%
- Duplicates consume significant resources during crawling, indexing, and search
 - Little value to most users

Duplicate Detection

- Exact duplicate detection is relatively easy
- Checksum techniques
 - A checksum is a value that is computed based on the content of the document
 - e.g., sum of the bytes in the document file

T	r	o	p	i	c	a	l		f	i	s	h	<i>Sum</i>
54	72	6F	70	69	63	61	6C	20	66	69	73	68	508

- Possible for files with different text to have same checksum
- Functions such as a *cyclic redundancy check* (CRC), have been developed that consider the positions of the bytes

Near-Duplicate Detection

- More challenging task
 - Are web pages with same text context but different advertising or format near-duplicates?
- A near-duplicate document is defined using a threshold value for some similarity measure between pairs of documents
 - e.g., document $D1$ is a near-duplicate of document $D2$ if more than 90% of the words in the documents are the same

Fingerprints

1. The document is parsed into words. Non-word content, such as punctuation, HTML tags, and additional whitespace, is removed.
2. The words are grouped into contiguous *n-grams* for some *n*. These are usually overlapping sequences of words, although some techniques use non-overlapping sequences.
3. Some of the *n-grams* are selected to represent the document.
4. The selected *n-grams* are hashed to improve retrieval efficiency and further reduce the size of the representation.
5. The hash values are stored, typically in an inverted index.
6. Documents are compared using overlap of fingerprints

Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species.

(a) Original text

tropical fish include, fish include fish, include fish found, fish found in, found in tropical, in tropical environments, tropical environments around, environments around the, around the world, the world including, world including both, including both freshwater, both freshwater and, freshwater and salt, and salt water, salt water species

(b) 3-grams

938 664 463 822 492 798 78 969 143 236 913 908 694 553 870 779

(c) Hash values

664 492 236 908

(d) Selected hash values using $0 \bmod 4$

Removing Noise

- Many web pages contain text, links, and pictures that are not directly related to the main content of the page
- This additional material is mostly *noise* that could negatively affect the ranking of the page
- 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

Noise Example

CNN.com Member Center Sign In | Register International Edition

SEARCH

Home Page World U.S. Weather Business Sports Analysis Politics Law Technology Science & Space Health Entertainment Critique Travel Education Special Reports Video Autos E-Reports

SCIENCE & SPACE

Aquarium plays whale shark matchmaker

Two females flown 8,000 miles for double date in Atlanta

Monday, June 5, 2006, 10:28 p.m. EDT (11:28 GMT)

ATLANTA, Georgia (CNN) — Ralph and Norton, meet Alice and Trixie.

The Georgia Aquarium's two male whale sharks got some female companionship on Saturday, when they were joined by two females transported to Atlanta from Taipei, Taiwan.

Researchers are hoping the sharks will mate.

The females — 11 feet and 14 feet long — were flown more than 8,000 miles by UPS, which incorporated a company B-747 freighter with advanced marine life support systems to carry them. [\(Watch what's hot & get the sharks together — 1:55\)](#)

The pilot said they treated the massive fish like first-class passengers.

"As we were doing the descent, we asked to start down a little sooner to make a nice shallow descent, to not make things too uncomfortable back there for the whale sharks," UPS pilot Capt. Bob Crum said.

The plane's center of balance was carefully planned, according to a statement from the aquarium, and veterinarians accompanied the sharks.

The delivery company also brought the two males to Atlanta, where researchers can study the whale sharks' behavior, breeding and development.

The whale sharks — named after the main characters in the 1950s sitcom "The Honeymooners" — were delivered to the aquarium in special transportation containers.

The Georgia Aquarium, which opened in November, is the world's largest aquarium. It was a \$250 million gift to Georgia from Bernie Marcus, co-founder of The Home Depot and his wife, Bill, through the Marcus Foundation.

It is the only aquarium outside of Asia to showcase whale sharks, which are the largest fish on Earth.

The aquarium's 6.2-million gallon "Ocean Voyager" tank can hold up to six whale sharks at a time.

Story Tools

SPACE **TOP STORIES**

- Astronauts prepare for third spacewalk
- Russians choose Putin's successor
- Astronomers vie to make biggest telescope
- Iran's president makes landmark visit to Iraq
- NASA to beam Beatles song to North Star
- Israel PM: Attacks on militants go on
- U.S. plans for falling satellite
- Cable arrested in abandoned baby case

International Edition Languages **CNN TV** **CNN International** **Headline News** **Transcripts** **Advertise with Us** **About Us**

SEARCH

© 2007 Cable News Network. A Time Warner Company. All Rights Reserved. External sites open in new window; not endorsed by CNN.com on (TV-14) Pay service with fee and archive open. Lastname:

Content block