

# Crawling

T. Yang, UCSB 290N  
 Some of slides from Crofter/Metzler/Strohman's textbook

## Table of Content

- **Basic crawling architecture and flow**
  - Distributed crawling
- **Scheduling: Where to crawl**
  - Crawling control with robots.txt
  - Freshness
  - Focused crawling
- **URL discovery**
  - Deep web, Sitemaps, & Data feeds
- **Data representation and store**

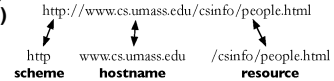
## Web Crawler

- Finds and downloads web pages automatically for search and web mining
- Web is huge and constantly growing

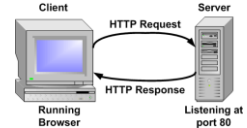


## Downloading Web Pages

- 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
  - HTTP /1.1



## HTTP

```

Terminal
File Edit View Terminal Go Help
~$ telnet java.sun.com 80
Trying 209.249.116.141...
Connected to java.sun.com.
Escape character is '^]'.
GET / HTTP/1.0

HTTP/1.1 200 OK
Server: Netscape-Enterprise/6.0
Date: Thu, 22 Jul 2004 18:27:16 GMT
Content-type: text/html; charset=ISO-8859-1
Set-cookie: JSESSIONID=java.sun.com-1745625344100070182534e493f41ad6f1e;path=/;expires=Thu, 22-Jul-2004 18:57:14 GMT
Connection: close

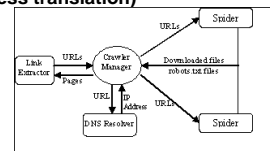
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
<html>
<head>
<title>Java Technology</title>
<meta name="keywords" content="Java, platform" />
<meta name="description" content="Java technology is a portfolio of products that are based on the power of networks and the idea that the same software should run on many different kinds of systems and devices." />
<meta http-equiv="Content-Type" content="text/html; charset=ISO-8859-1"/>
<meta name="date" content="2003-11-23" />
    
```

Figure 3 Using Telnet to Connect to a Web Server

## Downloading Web Pages

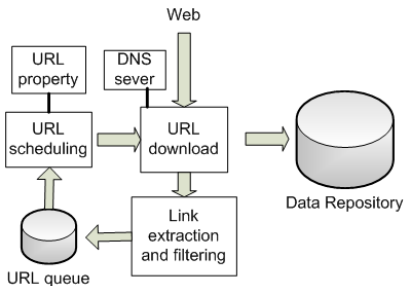
- Need a scalable *domain name system* (DNS) server (hostname to IP address translation)

- Crawler 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

## A Crawler Architecture



## Web Crawler

- Starts with a set of **seeds**
  - 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**
- Scheduler prioritizes to discover new or refresh the existing URLs
- Repeat the above process

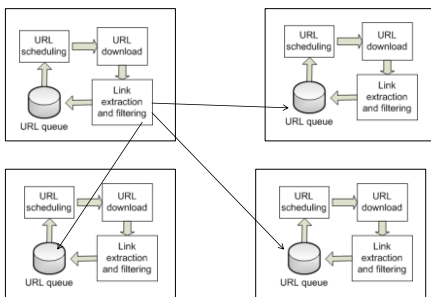
## Distributed Crawling: Parallel Execution

- Crawlers may be running in diverse geographies – USA, Europe, Asia, etc.
  - Periodically update a master index
  - Incremental update so this is “cheap”
- Three reasons to use multiple computers
  - Helps to put the crawler closer to the sites it crawls
  - Reduces the number of sites the crawler has to remember
  - More computing resources

## Variations of Distributed Crawlers

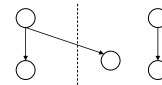
- Crawlers are **independent**
  - Fetch pages oblivious to each other.
- **Static assignment**
  - Distributed crawler uses a hash function to assign URLs to crawling computers
  - hash function can be computed on the host part of each URL
- **Dynamic assignment**
  - Master-slaves
  - Central coordinator splits URLs among crawlers

## A Distributed Crawler Architecture



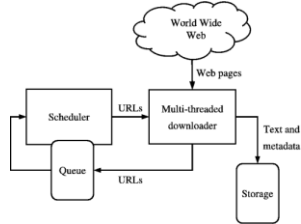
## Options of URL outgoing link assignment

- **Firewall mode**: each crawler only fetches URL within its partition – typically a domain
  - inter-partition links not followed
- **Crossover mode**: Each crawler may following inter-partition links into another partition
  - possibility of duplicate fetching
- **Exchange mode**: Each crawler periodically exchange URLs they discover in another partition



## Multithreaded page downloader

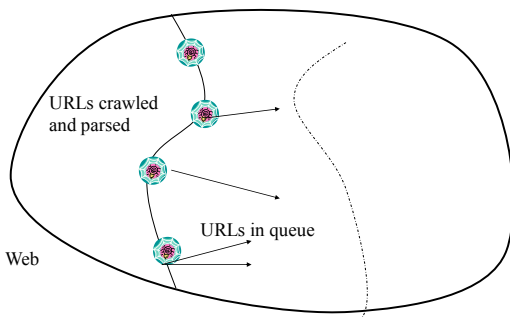
- **Web crawlers spend a lot of time waiting for responses to requests**
  - Multi-threaded for concurrency
  - Tolerate slowness of some sites
- **Few hundreds of threads/machine**



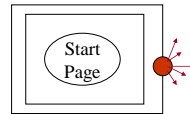
## Table of Content

- **Crawling architecture and flow**
- **Schedule: Where to crawl** ←
  - Crawling control with robots.txt
  - Freshness
  - Focused crawling
- **URL discovery:**
  - Deep web, Sitemaps, & Data feeds
- **Data representation and store**

## Where do we spider next?



## How fast can spam URLs contaminate a queue?

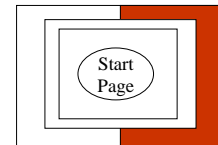


BFS depth = 2

Normal avg outdegree = 10

100 URLs on the queue including a spam page.

Assume the spammer is able to generate dynamic pages with 1000 outlinks



BFS depth = 3

2000 URLs on the queue  
50% belong to the spammer

BFS depth = 4

1.01 million URLs on the queue  
99% belong to the spammer

## Scheduling Issues: Where do we spider next?

- **Keep all spiders busy (load balanced)**
  - Avoid fetching duplicates repeatedly
- **Respect politeness and robots.txt**
  - Crawlers could potentially flood sites with requests for pages
  - use *politeness policies*: e.g., delay between requests to same web server
- **Handle crawling abnormality:**
  - Avoid getting stuck in traps
  - Tolerate faults with retry

## More URL Scheduling Issues

- **Conflicting goals**
  - Big sites are crawled completely;
  - Discover and recrawl URLs frequently
    - Important URLs need to have high priority
      - What's best? Quality, fresh, topic coverage
    - Avoid/Minimize duplicate and spam
  - Revisiting for recently crawled URLs should be excluded to avoid the endless of revisiting of the same URLs.
- **Access properties of URLs to make a scheduling decision.**

## /robots.txt

- Protocol for giving spiders (“robots”) limited access to a website, originally from 1994
  - [www.robotstxt.org/](http://www.robotstxt.org/)
- Website announces its request on what can(not) be crawled
  - For a URL, create a file `robots.txt`
  - This file specifies access restrictions
  - Place in the top directory of web server.
    - E.g. [www.cs.ucsb.edu/robots.txt](http://www.cs.ucsb.edu/robots.txt)
    - [www.ucsb.edu/robots.txt](http://www.ucsb.edu/robots.txt)

## Robots.txt example

- No robot should visit any URL starting with `/yoursite/temp/`, except the robot called “searchengine”:

```
User-agent: *  
Disallow: /yoursite/temp/
```

```
User-agent: searchengine  
Disallow:
```

## More Robots.txt example

```
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
  - *stale* copies no longer reflect the real contents of the web pages

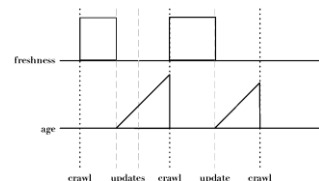
## Freshness

- 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
  - Information is not reliable.

```
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
  - Need to check important pages and pages that change frequently
- Freshness is the proportion of pages that are fresh
- Age as an approximation



## 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

## Table of Content

- Basic crawling architecture and flow
- Schedule: Where to crawl
  - Crawling control with robots.txt
  - Freshness
  - Focused crawling
- Discover new URLs 
  - Deep web, Sitemaps, & Data feeds
- Data representation and store

## Discover new URLs & Deepweb

- Challenges to discover new URLs
  - Bandwidth/politeness prevent the crawler from covering large sites fully.
  - Deepweb
- Strategies
  - Mining new topics/related URLs from news, blogs, facebook/twitters.
  - Identify sites that tend to deliver more new URLs.
  - Deepweb handling/sitemaps
  - RSS feeds

## 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

- Placed at the root directory of an HTML server.
  - For example, <http://example.com/sitemap.xml>.
- 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>
```

## Document Feeds

- **Many documents are published**
  - created at a fixed time and rarely updated again
  - e.g., news articles, blog posts, press releases, email
- **Published documents from a single source can be ordered in a sequence called a document feed**
  - new documents found by examining the end of the feed

## Document Feeds

- **Two types:**
  - A *push feed* alerts the subscriber to new documents
  - A *pull feed* requires the subscriber to check periodically for new documents
- **Most common format for pull feeds is called RSS**
  - Really Simple Syndication, RDF Site Summary, Rich Site Summary, or ...
- **Examples**
  - CNN RSS newfeed under different categories
  - Amazon RSS popular product feeds under different tags

## RSS Example

```
<?xml version="1.0"?>
<rss version="2.0">
  <channel>
    <title>Search Engine News</title>
    <link>http://www.search-engine-news.org/</link>
    <description>News about search engines.</description>
    <language>en-us</language>
    <pubDate>Tue, 19 Jun 2008 05:17:00 GMT</pubDate>
    <ttl>60</ttl>

    <item>
      <title>Upcoming SIGIR Conference</title>
      <link>http://www.sigir.org/conference/</link>
      <description>The annual SIGIR conference is coming!
        Mark your calendars and check for cheap
        flights.</description>
      <pubDate>Tue, 05 Jun 2008 09:50:11 GMT</pubDate>
      <guid>http://search-engine-news.org#500</guid>
    </item>
```

## RSS Example

```
...
  <item>
    <title>New Search Engine Textbook</title>
    <link>http://www.cs.umass.edu/search-book/</link>
    <description>A new textbook about search engines
      will be published soon.</description>
    <pubDate>Tue, 05 Jun 2008 09:33:01 GMT</pubDate>
    <guid>http://search-engine-news.org#499</guid>
  </item>
</channel>
</rss>
```

## RSS

- **A number of channel elements:**
  - Title
  - Link
  - description
  - ttl tag (time to live)
    - amount of time (in minutes) contents should be cached
- **RSS feeds are accessed like web pages**
  - using HTTP GET requests to web servers that host them
- **Easy for crawlers to parse**
- **Easy to find new information**

## Table of Content

- **Crawling architecture and flow**
- **Scheduling: Where to crawl**
  - Crawling control with robots.txt
  - Freshness
  - Focused crawling
- **URL discovery**
  - Deep web, Sitemaps, & Data feeds
- **Data representation and store** ←

## 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**
  - 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

## Software Internationalization with Unicode

- **Search software needs to be able to run for serving different international content**
- **Proliferation of encodings comes from a need for compatibility and to save space**
  - UTF-8 uses one byte for English (ASCII), as many as 4 bytes for some traditional Chinese characters
  - variable length encoding, more difficult to do string operations
  - UTF-32 uses 4 bytes for every character
- **Many applications use UTF-32 for internal text encoding (fast random lookup) and UTF-8 for disk storage (less space)**

## Example of Unicode

Decimal	Hexadecimal	Encoding			
0-127	0-7F	0xxxxxxx			
128-2047	80-7FF	110xxxxx	10xxxxxx		
2048-55295	800-D7FF	1110xxxx	10xxxxxx	10xxxxxx	
55296-57343	D800-DFFF	Undefined			
57344-65535	E000-FFFF	1110xxxx	10xxxxxx	10xxxxxx	
65536-1114111	10000-10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

- e.g., Greek letter pi ( $\pi$ ) is Unicode symbol number 960
- In binary, 00000011 11000000 (3C0 in hexadecimal)
- Final encoding is **11001111 10000000** (CF80 in hexadecimal)

## 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.
- **Data stores used for page repository**
- **Store many documents in large files, rather than each document in a file**
  - avoids overhead in opening and closing files
  - reduces seek time relative to read time
- **Compound documents formats**
  - used to store multiple documents in a file
  - e.g., TREC Web

## TREC Web Format

```
<DOC>
<DOCID>WT001-001-104/DOCID</DOCID>
<DOCLEN>
http://www.example.com/test.html 204.244.59.33 19970101013145 text/html 440
HTTP/1.0 200 OK
Date: Wed, 01 Jan 1997 01:21:13 GMT
Server: Apache/1.0.3
Content-type: text/html
Content-length: 270
Last-modified: Mon, 28 Nov 1996 05:31:24 GMT
</DOCLEN>
<HTML>
<TITLE>Tropical Fish Store</TITLE>
<META>
</META>
</HTML>
<DOC>
<DOCID>WT001-001-1094/DOCID</DOCID>
<DOCLEN>
http://www.example.com/fish.html 204.244.59.33 19970101013149 text/html 440
HTTP/1.0 200 OK
Date: Wed, 01 Jan 1997 01:21:19 GMT
Server: Apache/1.0.3
Content-type: text/html
Content-length: 270
Last-modified: Mon, 28 Nov 1996 05:31:24 GMT
</DOCLEN>
<HTML>
<TITLE>Fish Information</TITLE>
This page will soon contain interesting
information about tropical fish.
</HTML>
</DOC>
```

## Text 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: a separate topic**
- **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