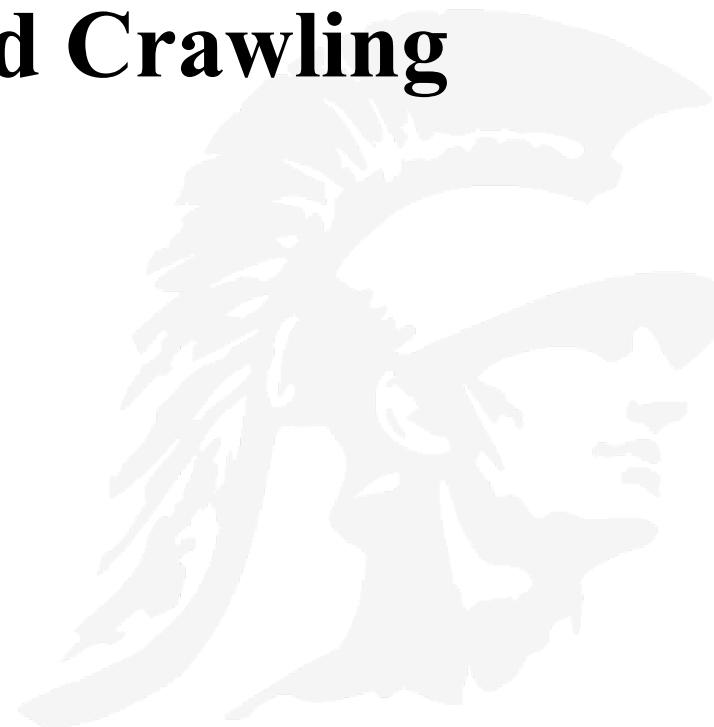


# Crawlers and Crawling



# There are Many Crawlers

- A web crawler is a computer program that visits web pages in an organized way

- Sometimes called a spider or robot

- A list of web crawlers can be found at

- [http://en.wikipedia.org/wiki/Web\\_crawler](http://en.wikipedia.org/wiki/Web_crawler)

- Google's crawler is called googlebot, see

- <http://support.google.com/webmasters/bin/answer.py?hl=en&answer=182072>

- Yahoo's web crawler is/was called Yahoo! Slurp, see

- [http://en.wikipedia.org/wiki/Yahoo!\\_Search](http://en.wikipedia.org/wiki/Yahoo!_Search)

- Bing uses five crawlers

- Bingbot, standard crawler
  - Adidxbot, used by Bing Ads
  - MSNbot, remnant from MSN, but still in use
  - MSNBotMedia, crawls images and video
  - BingPreview, generates page snapshots

- For details see: <http://www.bing.com/webmaster/help/which-crawlers-does-bing-use-8c184ec0>

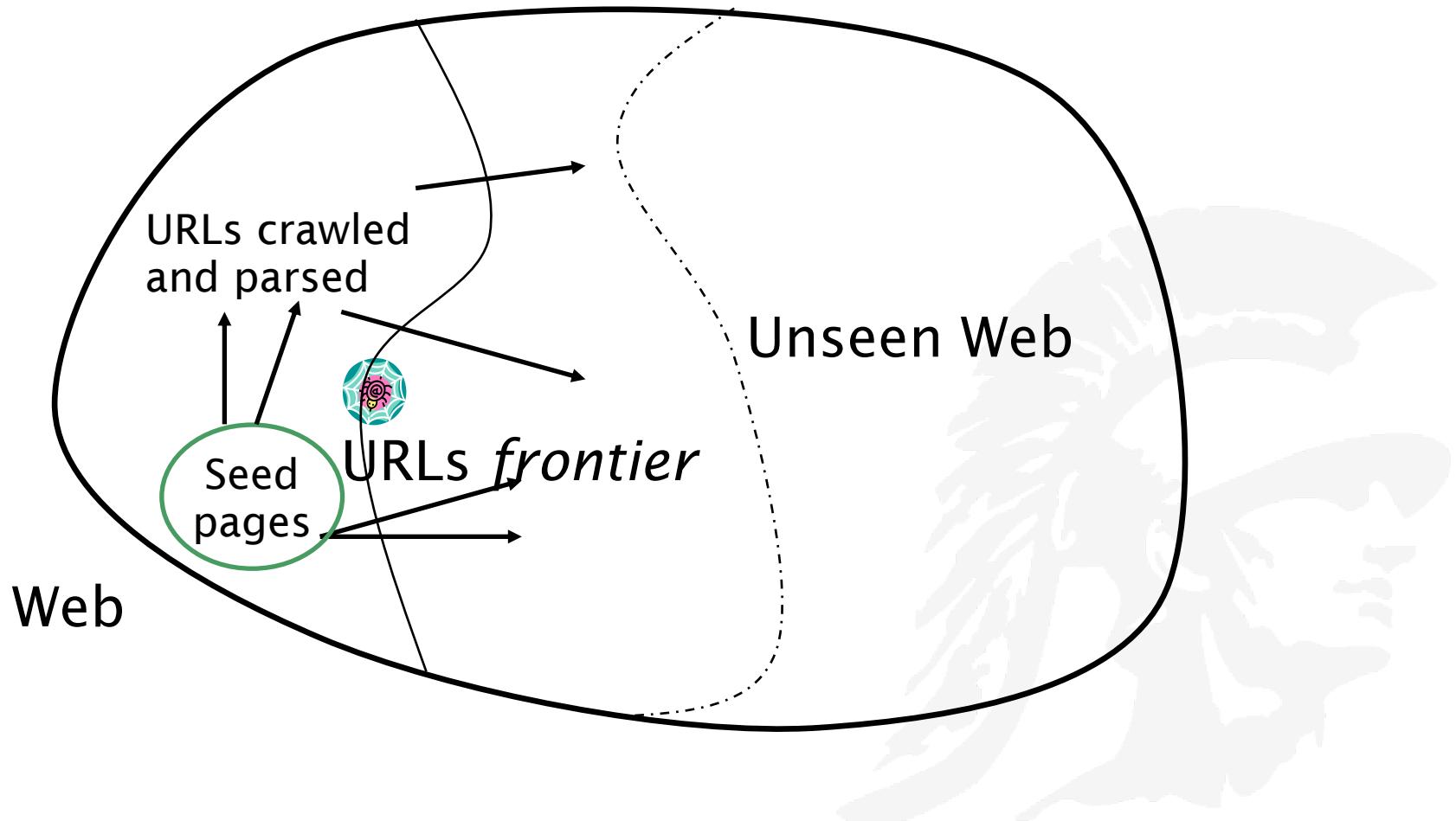
# Web Crawling Issues

- **How to crawl?**
  - *Quality*: how to find the “Best” pages first
  - *Efficiency*: how to avoid duplication (or near duplication)
  - *Etiquette*: behave politely by not disturbing a website’s performance
- **How much to crawl? How much to index?**
  - *Coverage*: What percentage of the web should be covered?
  - *Relative Coverage*: How much do competitors have?
- **How often to crawl?**
  - *Freshness*: How much has changed?
  - How much has really changed?

# Simplest Crawler Operation

- Begin with known “seed” pages
- Fetch and parse a page
  - Place the page in a database
  - Extract the URLs within the page
  - Place the extracted URLs on a queue
- Fetch each URL on the queue and repeat

# Crawling Picture



# Simple Picture – Complications

- **Crawling the entire web isn't feasible with one machine**
  - But all of the above steps can be distributed
- **Handling/Avoiding malicious pages**
  - Some pages contain spam
  - Some pages contain spider traps – especially dynamically generated pages
- **Even non-malicious pages pose challenges**
  - Latency/bandwidth to remote servers can vary widely
  - Robots.txt stipulations can prevent web pages from being visited
  - How can one avoid site mirrors and duplicate pages
- **Maintain politeness – don't hit a server too often**

# Robots.txt

- There is a protocol that defines the limitations for a web crawler as it visits a website; its definition is here
  - <http://www.robotstxt.org/orig.html>
- The website announces its request on what can(not) be crawled by placing a robots.txt file in the root directory
  - e.g. see  
<http://www.ticketmaster.com/robots.txt>

# Robots.txt Example

- **No robot visiting this domain should visit any URL starting with "/yoursite/temp/":**

User-agent: \*

Disallow: /yoursite/temp/

- **Directives are case sensitive**
- **Additional symbols allowed in the robots.txt directives include:**
  - '\*' - matches a sequence of characters
  - '\$' - anchors at the end of the URL string
- **Example of '\*':**

User-agent: Slurp

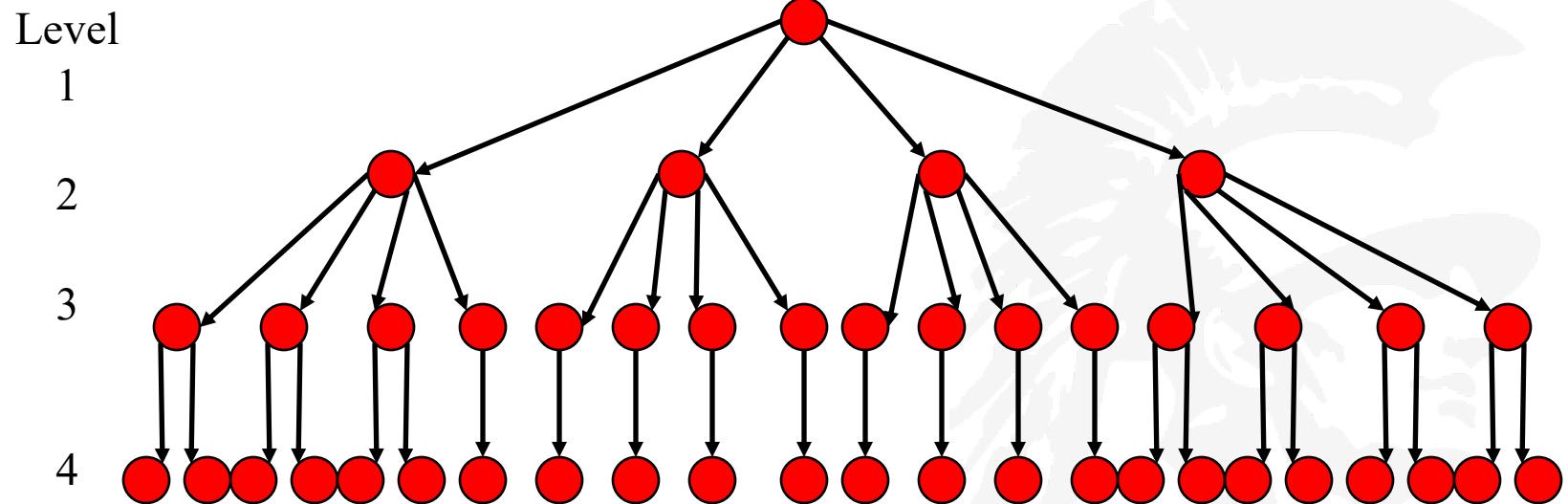
Allow: /public\*/

Disallow: /\*\_print\*.html

Disallow: /\*?sessionid

# Basic Search Strategies

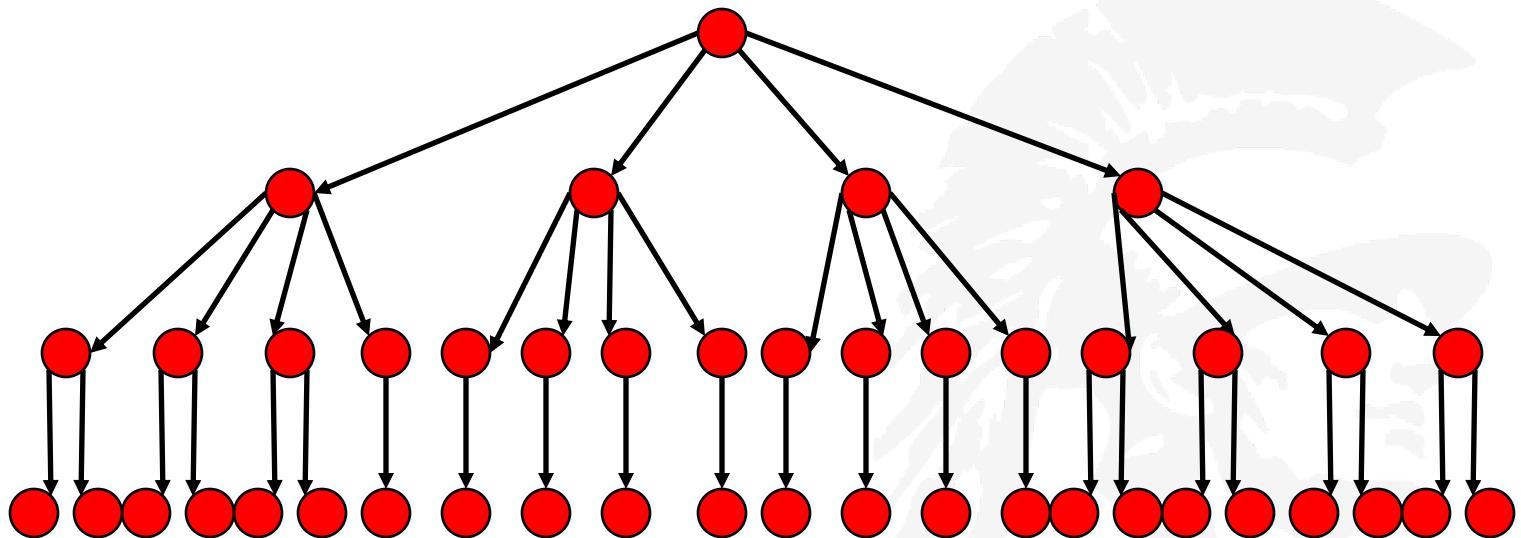
## Breadth-first Search



Examine all pages at level  $i$  before examining pages at level  $i+1$

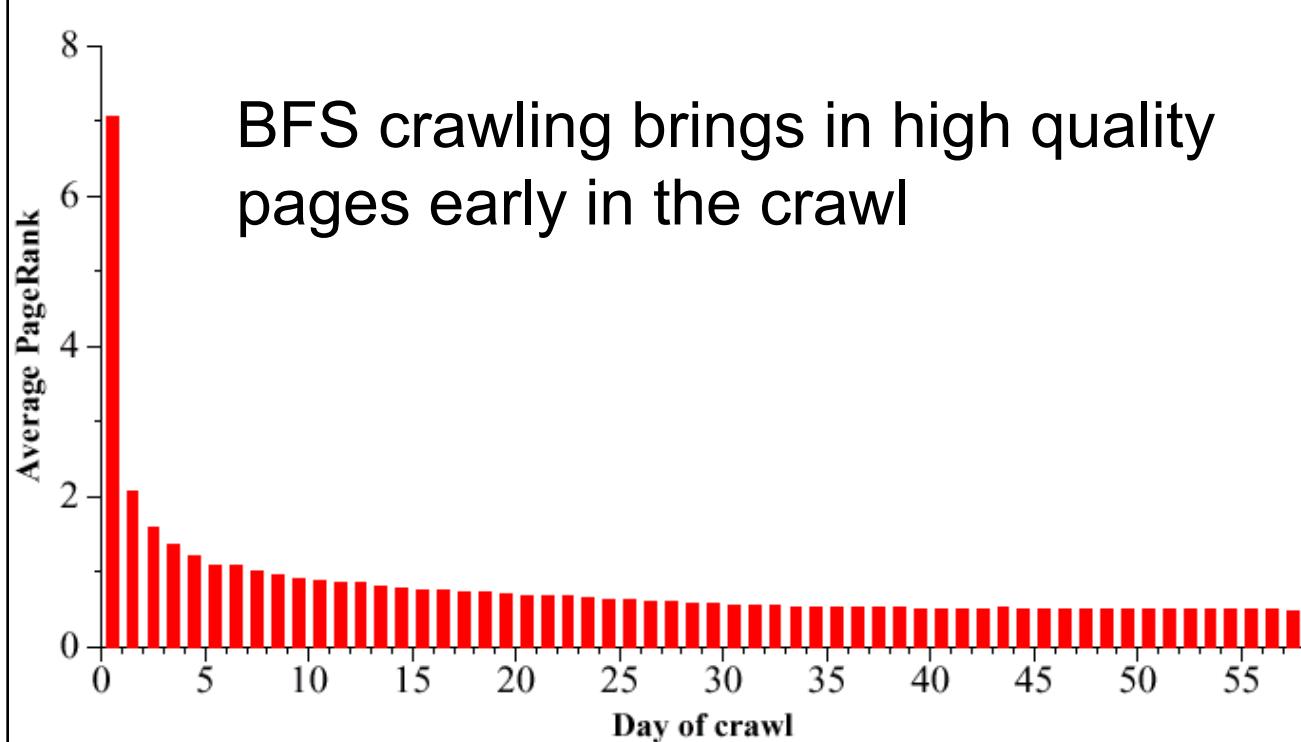
# Basic Search Strategies (cont)

## Depth-first Search



At each step move to a page down the tree

## Web Wide Crawl (328M pages) [Naj01]



Average PageRank score by day of crawl

Page Rank is an algorithm developed by Google for determining the value of a page

## Crawling Algorithm – Version 2

Initialize queue (Q) with initial set of known URL's.

**Loop** until Q empty or page or time limit exhausted:

    Pop a URL, call it L, from the front of Q.

    If L is not an HTML page (e.g. .gif, .jpeg, ....)  
        continue the loop

    If L has already been visited, continue the loop.

    Download page, P, for L

    If cannot download P (e.g. 404 error, robot excluded)  
        continue loop

    Index P (e.g. add to inverted index and store cached copy)

    Parse P to obtain list of new links N.

    Append N to the end of Q

**End** loop

# Queueing Strategy

- How new links are added to the queue determines the search strategy.
- FIFO (append to end of Q) gives breadth-first search.
- LIFO (add to front of Q) gives depth-first search.
- Heuristically ordering the Q gives a “focused crawler” that directs its search towards “interesting” pages; e.g.
  - A document that changes frequently could be moved forward
  - A document whose content appears relevant to some topic can be moved forward
  - e.g. see Focused Crawling: A New Approach by S. Chakrabarti et al  
<https://www.sciencedirect.com/science/article/pii/S1389128699000523>
- One way to re-order the URLs on the queue is to:
  - Move forward URLs whose In-degree is large
  - Move forward URLs whose PageRank is large
    - We will discuss the PageRank algorithm later

# Avoiding Page Duplication

- A crawler must detect when revisiting a page that has already been crawled  
**(Remember: the web is a graph not a tree).**
- Therefore, a crawler must efficiently index URLs as well as already visited pages
- To determine if a URL has already been seen,
  - Must store URLs in a standard format (discussed ahead)
  - Must develop a fast way to check if a URL has already been seen
- To determine if a new page has already been seen,
  - Must develop a fast way to determine if an *identical* page was already indexed
  - Must develop a fast way to determine if a *near-identical* page was already indexed

# Link Extraction

- **Must find all links in a page and extract URLs.**
  - URLs occur in tags other than <a>, e.g.
    - <frame src=“site-index.html”>, <area, href=“...”>, <meta>, <link>, <script>
- **Relative URL’s must be completed, e.g. using current page URL or <base> tag**
  - <a href=“proj.html”> to http://www.myco.com/special/tools/proj.html
  - <a href=“..../outline/syllabus.html”> to http://www.myco.com/special/outline/syllabus.html
- **Two Anomalies**
  1. Some anchors don’t have links, e.g. <a name=“banner”>
  2. Some anchors produce dynamic pages which can lead to looping  
<a href=http://www.mysite.com/search?x=arg1&y=arg2>

# Representing URLs

- URLs are rather long, 80 bytes on the average, implying 1 billion URLs will require 80 Gigabytes

- Recently Google reported finding 30 trillion unique URLs, which by the above would require 2400 terabytes (or 2.4 petabytes) to store

## 1. One Proposed Method: To determine if a new URL has already been seen

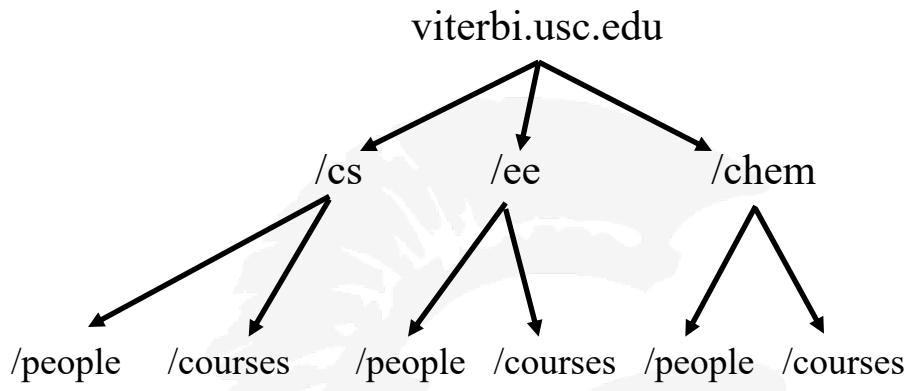
- First hash on host/domain name, then
- Use a trie data structure to determine if the path/resource is the same as one in the URL database

## 2. Another Proposed Method: URLs are sorted lexicographically and then stored as a delta-encoded text file

- Each entry is stored as the difference (delta) between the current and previous URL; this substantially reduces storage
- However, restoring the actual URL is slower, requiring all deltas to be applied to the initial URL
- To improve speed, checkpointing (storing the full URL) is done periodically

# Trie for URL Exact Matching

- Simplest (and worst) algorithm to determine if a new URL is in your set
  - `grep -i <search_url> <url_file>`
  - For  $n$  URLs and maximum length  $k$ , time is  $O(nk)$
- Characteristics of tries
  - They share the same prefix among multiple “words”
  - Each path from the root to a leaf corresponds to one “word”
  - *Endmarker symbol, \$, at the ends of all words*
    - To avoid confusion between words with almost identical elements
      - Assume all words are \$ terminated



If we store  $N$  URLs, each of maximum length  $M$ , in a binary search tree, then the search time is  $O(M * \log N)$ ; however, using a trie, the search time is  $O(M)$ , at the expense of more storage

# Why Normalizing URLs is Important

- For example, all the following URLs have the same meaning, but different hashes:
  - `http://www.google.com`
  - `http://www.google.com/`
  - `https://www.google.com`
  - `www.google.com`
  - `google.com`
  - `google.com.`



# Normalizing URLs (4 rules)

1. Convert the scheme and host to lower case. The scheme and host components of the URL are case-insensitive.
  - HTTP://www.Example.com/ → http://www.example.com/
2. Capitalize letters in escape sequences. All letters within a percent-encoding triplet (e.g., "%3A") are case-insensitive, and should be capitalized.  
**Example:**
  - http://www.example.com/a%c2%b1b →  
http://www.example.com/a%C2%B1b
3. Decode percent-encoded octets of unreserved characters.  
http://www.example.com/%7Eusername/ →  
http://www.example.com/~username/
4. Remove the default port. The default port (port 80 for the “http” scheme) may be removed from (or added to) a URL. Example:
  - http://www.example.com:80/bar.html →  
http://www.example.com/bar.html
  - See [https://en.wikipedia.org/wiki/URL\\_normalization](https://en.wikipedia.org/wiki/URL_normalization)

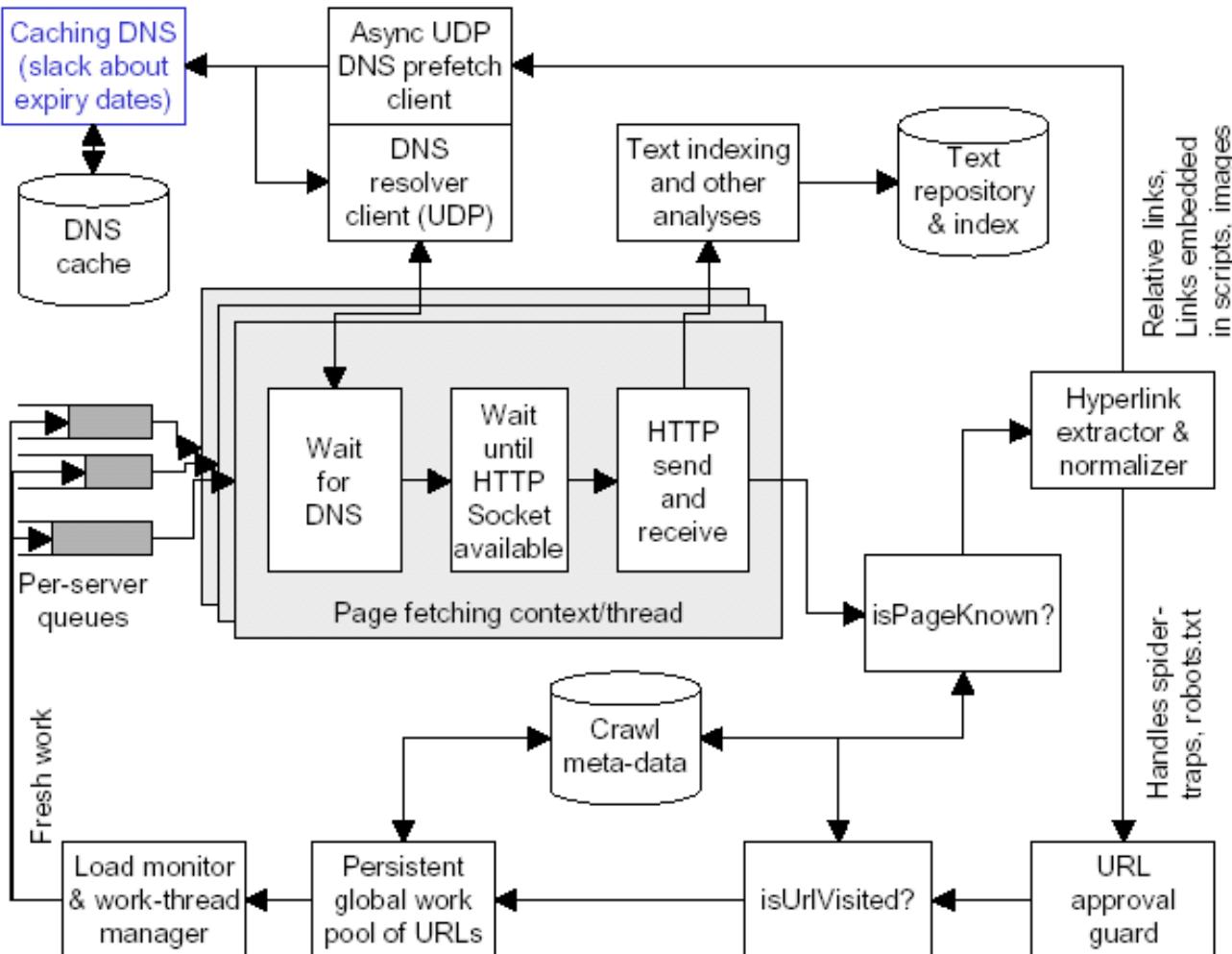
# Avoiding Spider Traps

- A spider trap is when a crawler re-visits the same page over and over again
- The most well-known spider trap is the one created by the use of Session ID's
- A Session ID is often used to keep track of visitors, and some sites puts a unique ID in the URL:
  - An example is [www.webmasterworld.com/page.php?id=264684413484654](http://www.webmasterworld.com/page.php?id=264684413484654) (**Note** this URL doesn't exist).  
Each user gets a unique ID and it's often requested from each page.  
The problem here is when Googlebot comes to the page, it spiders the page and then leaves, it comes back to another page and it finds a link to the previous page, but since it has been given a different session id now, the link shows up as another URL.
- One way to avoid such traps is for the crawler to be careful when the querystring “ID=” is present in the URL
- Another technique is to monitor the length of the URL, and stop if the length gets “too long”

# Handling Spam

- The **first generation** of spam consisted of pages with a high number of repeated terms, so as to score high on search engines that ranked by word frequency
  - Words were typically rendered in the same color as the background, so as to not be visible, but still count
- The **second generation** of spam used a technique called *cloaking*;
  - When the web server detects a request from a crawler, it returns a different page than the page it returns from a user request
  - The page is mistakenly indexed
- A **third generation**, called a doorway page, contains text and metadata chosen to rank highly on certain search keywords, but when a browser requests the doorway page it instead gets a more “commercial” page
- **Cloaking** and **doorway pages** are not permitted according to Google’s webmaster suggestions
  - See <http://support.google.com/webmasters/bin/answer.py?hl=en&answer=66355>
  - (show video, start at 3:40 into the video to see why showing different pages based on geolocation is sometimes permitted; show 3 minutes worth)

# The Mercator Web Crawler



**The diagram points out all of the key elements of a crawler;**  
**Notice**

1. The DNS caching server
2. Use of UDP for DNS
3. Load and thread monitor
4. Parallel threads waiting for a page to download

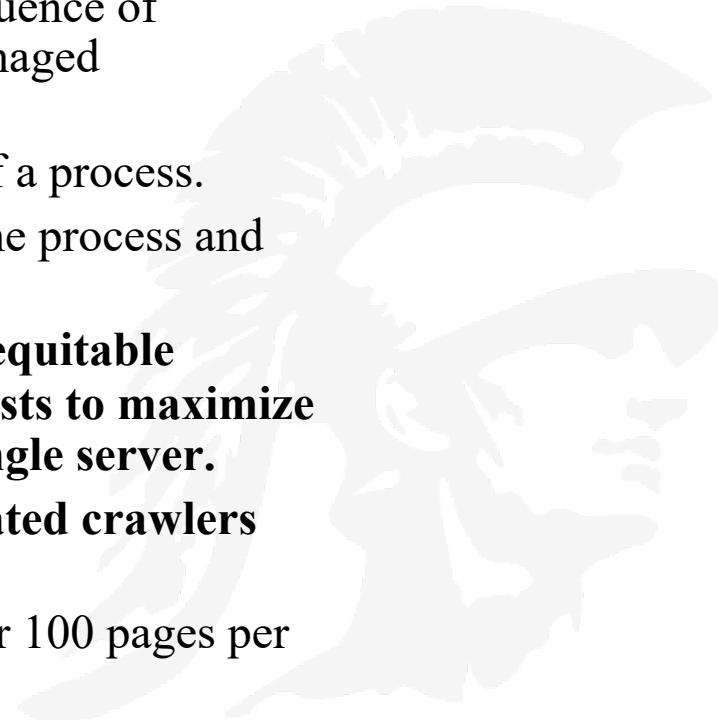
# Measuring and Tuning a Crawler

- **Measuring and tuning a crawler for peak performance eventually reduces to**
  - Improving URL parsing speed
  - Improving network bandwidth speed
  - Improving fault tolerance
- **More Issues (some of which are discussed ahead)**
  - Refresh Strategies: how often is the process re-started
  - Detecting duplicate pages
  - Detecting mirror sites
  - Speeding up DNS lookup (see previous slide)
  - URL normalization (discussed earlier)
  - Handling malformed HTML

- *A common operating system's implementation of DNS lookup is blocking: only one outstanding request at a time; so*
1. **DNS caching:** build a caching server that retains IP-domain name mappings previously discovered
  2. **Pre-fetching client**
    - once a page is parsed,
      - immediately make DNS resolution requests to the caching server; and
      - if unresolved, use UDP (User Datagram Protocol) to resolve from the DNS server
  3. **Customize the crawler so it allows issuing of many resolution requests simultaneously; there should be many DNS resolvers**

# Multi-Threaded Crawling

- One bottleneck is network delay in downloading individual pages.
- It is best to have multiple threads running in parallel each requesting a page from a different host.
  - a **thread** of execution is the smallest sequence of programmed instructions that can be managed independently by a scheduler.
  - In most cases, a thread is a component of a process.
  - Multiple threads can exist within the same process and share resources
- Distribute URL's to threads to guarantee equitable distribution of requests across different hosts to maximize through-put and avoid overloading any single server.
- Early Google spider had multiple coordinated crawlers with about 300 threads each,
  - together they were able to download over 100 pages per second.



# Distributed Crawling Approaches

- Once the crawler program itself has been optimized, the next issue to decide is how many crawlers will be running at any time
- Scenario 1: A *centralized crawler* controlling a set of parallel crawlers all running on a LAN
  - A *parallel crawler* consists of multiple crawling processes communicating via local network (sometimes called an intra-site parallel crawler)
- Scenario 2: A *distributed set of crawlers* running on widely distributed machines, with or without cross communication

# Distributed Model

- If crawlers are running in diverse geographic locations, how do we organize them
  - By country, by region, by available bandwidth
  - Distributed crawlers must periodically update a master index
  - But incremental update is generally “cheap”
    - Why? Because
    - a. you can compress the update, and
    - b. you need only send a differential update
    - both of which will limit the required communication

# Issues and Benefits of Distributed Crawling

- **Benefits:**
  - scalability: for large-scale web-crawls
  - costs: use of cheaper machines
  - network-load dispersion and reduction: by dividing the web into regions and crawling only the nearest pages
- **Issues:**
  - overlap: minimization of multiple downloaded pages
  - quality: depends on the crawl strategy
  - communication bandwidth: minimization

# Coordination of Distributed Crawling

## – Three strategies

### 1. Independent:

- ▶ no coordination, every process follows its extracted links

### 2. Dynamic assignment:

- ▶ a central coordinator dynamically divides the web into small partitions and assigns each partition to a process

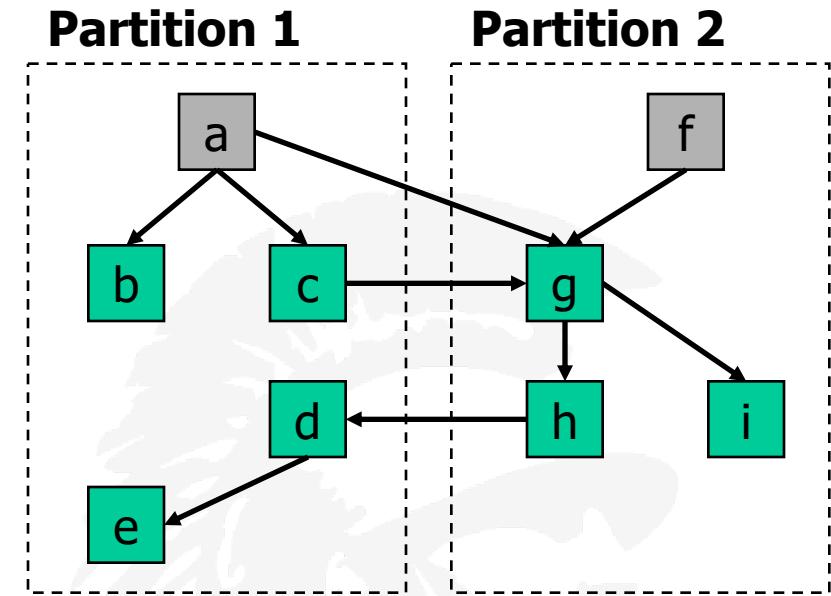
### 3. Static assignment:

- ▶ Web is partitioned and assigned without a central coordinator before the crawl starts

# Static Assignment

Links from one partition to another (inter-partition links) can be handled in one of three ways:

1. *Firewall mode:*  
**a process does not follow any inter-partition link**
2. *Cross-over mode:*  
**a process also follows inter-partition links and possibly discovers also more pages in its partition**
3. *Exchange mode:*  
**processes exchange inter-partition URLs; this mode requires communication**



# Classification of Parallel Crawlers

- If exchange mode is used, communication can be limited by:
  - Batch communication: every process collects some URLs and sends them in a batch
  - Replication: the  $k$  most popular URLs are replicated at each process and are not exchanged (previous crawl or on the fly)
- Some ways to partition the Web:
  - URL-hash based: this yields many inter-partition links
  - Site-hash based: reduces the inter partition links
  - Hierarchical: by TLD, e.g. .com domain, .net domain ...
- General Conclusions of Cho and Garcia-Molina
  - Firewall crawlers attain good, general coverage with low cost
  - Cross-over ensures 100% quality, but suffer from overlap
  - Replicating URLs and batch communication can reduce overhead

- The behavior of a Web crawler is the outcome of a combination of policies:
  - A *selection policy* that states which pages to download.
  - A *re-visit policy* that states when to check for changes to the pages.
  - A *politeness policy* that states how to avoid overloading websites.
  - A *parallelization policy* that states how to coordinate distributed web crawlers.

# Keeping Spidered Pages Up to Date

- Web is very dynamic: many new pages, updated pages, deleted pages, etc.
- Periodically check crawled pages for updates and deletions:
  - Just look at LastModified indicator to determine if page has changed, only reload entire page if needed
- Track how often each page is updated and preferentially return to pages which are historically more dynamic.
- Preferentially update pages that are accessed more often to optimize freshness of more popular pages.

# Implications for a Web Crawler

- **A *steady crawler* runs continuously without pause**
  - Typically search engines use multiple crawlers
- **When a crawler replaces an old version by a new page, does it do it “in-place” or “shadowing”**
  - Shadowing implies a new set of pages are collected and stored separately and all are updated at the same time
  - The above implies that queries need to check two databases, the current database and the database of new pages
  - Shadowing either slows down query processing or decreases freshness
- **Conclusions:**
  - running multiple types of crawlers is best
  - Updating in-place keeps the index current

# Cho and Garcia-Molina, 2000

- **Two simple re-visiting policies**
  - **Uniform policy:** This involves re-visiting all pages in the collection with the same frequency, regardless of their rates of change.
  - **Proportional policy:** This involves re-visiting more often the pages that change more frequently. The visiting frequency is directly proportional to the (estimated) change frequency.
- **Cho and Garcia-Molina proved the surprising result that, in terms of average freshness, the uniform policy outperforms the proportional policy in both a simulated Web and a real Web crawl.**
- **The explanation for this result comes from the fact that, when a page changes too often, the crawler will waste time by trying to re-crawl it too fast and still will not be able to keep its copy of the page fresh.**
- **To improve freshness, we should penalize the elements that change too often**

# Help the Search Engine Crawler Creating a SiteMap

- A sitemap is a list of pages of a web site accessible to crawlers
- This helps search engine crawlers find pages on the site
- XML is used as the standard for representing sitemaps
- Here is an example of an XML sitemap for a three page website

```
<?xml version="1.0" encoding="UTF-8"?>
<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9">
<url>
  <loc>http://www.example.com/?id=who</loc>
  <lastmod>2009-09-22</lastmod>
  <changefreq>monthly</changefreq>
  <priority>0.8</priority> </url>
<url>
  <loc>http://www.example.com/?id=what</loc>
  <lastmod>2009-09-22</lastmod>
  <changefreq>monthly</changefreq>
  <priority>0.5</priority> </url>
<url>
  <loc>http://www.example.com/?id=how</loc>
  <lastmod>2009-09-22</lastmod>
  <changefreq>monthly</changefreq>
  <priority>0.5</priority> </url>
</urlset>
```

Google originally introduced the sitemap format; now Bing, Yahoo, and Ask also support sitemaps

See the Google, Bing, Yahoo, Ask announcement:  
<http://www.google.com/press/pressrel/sitemapsorg.html>

# Google Crawlers

- Google now uses multiple crawlers
  - Googlebot
  - Googlebot News
  - Googlebot Images
  - Googlebot Video
  - Google Mobile Smartphone
  - Google Mobile AdSense
  - Google AdsBot
  - Google app crawler
  - For details see  
<https://support.google.com/webmasters/answer/1061943?hl=en>
  - see also Google's tool for checking how Googlebot sees your website
  - <https://support.google.com/webmasters/answer/6066468?rd=2>

Google crawlers - Search

Secure | https://support.google.com/webmasters/answer/1061943?hl=en

CSCI 571 Home Page CSCI 572 Home Page Other bookmarks

Crawler	User agent token	Full user agent string (as seen in website log files)
Googlebot (Google Web search)	Googlebot	Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html) or (rarely used): Googlebot/2.1 (+http://www.google.com/bot.html)
Googlebot News	Googlebot-News (Googlebot)	Googlebot-News
Googlebot Images	Googlebot-Image (Googlebot)	Googlebot-Image/1.0
Googlebot Video	Googlebot-Video (Googlebot)	Googlebot-Video/1.0
Google Smartphone	Googlebot	Mozilla/5.0 (Linux; Android 6.0.1; Nexus 5X Build/MMB29P) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/41.0.2272.96 Mobile Safari/537.36 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)
Google Mobile AdSense	Mediapartners-Google or Mediapartners (Googlebot)	[various mobile device types] (compatible; Mediapartners-Google/2.1; +http://www.google.com/bot.html)
Google AdSense	Mediapartners-Google Mediapartners (Googlebot)	Mediapartners-Google
Google AdsBot landing page quality check	AdsBot-Google	AdsBot-Google (+http://www.google.com/adsbot.html)
Google app crawler (Used to fetch	AdsBot-Google-Mobile-Apps	AdsBot-Google-Mobile-Apps

# Google's Googlebot

- Begins with a list of webpage URLs generated from previous crawls
- Uses Sitemap data provided by webmasters
- Many versions of Googlebot are run on multiple machines located near the site they are indexing
- Googlebot cannot see within Flash files, audio/video tracks, and content within programs
- Advice
  - To prevent “File not found” in a website’s error log, create an empty robots.txt file
  - To prevent Googlebot from following any links on a page, use “nofollow” meta tag
  - To prevent Googlebot from following an individual link, add “rel=“nofollow”” attribute to the link
- Assertion: Chrome was in fact a repackaging of the Googlebot search crawler;
- Why: browsers don’t just render the DOM Hierarchy of HTML, they include transformations via CSS and JavaScript, and for Googlebot to extract the most meaningful features from a web page it would be necessary to have access to these transformations
- Conclusion: Googlebot and Chrome share a great deal of code

- In a public statement to Forbes magazine Google admitted that Googlebot does evaluate JavaScript as it crawls web pages
- Fundamentally, a browser is just software that provides an implementation of the W3C DOM Specification via a Rendering Engine, and a scripting engine to enable any additional scripting resources.
- Evidence: Googlebot is now requesting links that don't appear directly in JavaScript — links that get put together on the fly
- What would the advantages be?
  - it emulates the user experience including page load time, final markup, positions of elements;
  - it allows detection of hidden links or hidden text

<https://searchengineland.com/tested-googlebot-crawls-javascript-heres-learned-220157>

Read more: <http://ipullrank.com/googlebot-is-chrome/#ixzz2qbmjVGDj>

<http://www.forbes.com/sites/velocity/2010/06/25/google-isnt-just-reading-your-links-its-now-running-your-code/print/>

## More on Googlebot

### Google Isn't Just Reading Your Links, It's Now Running Your Code

 Taylor Buley, Contributor

It's long been observed that Google's search indexer can read JavaScript code, the lingua franca of dynamic Web applications. But for years it's been unclear whether or not the Googlebot actually understood what it was looking at or whether it was merely doing "dumb" searches for well-understood data structures like hyperlinks.

On Friday, a Google spokesperson confirmed to *Forbes* that Google does indeed go beyond mere "parsing" of JavaScript. "Google can parse and understand some JavaScript," said the spokesperson.

Rather than just read a page for links, Google's acknowledgment suggests that it might be able to interact with applications like a human would — and crack open parts of the Web that search engines like Bing might not be able to see. That would mean that Google has redefined what it means to be a search engine.

Very few [JavaScript files](#) exist in Google's results and for its part, the company has kept any JavaScript interpreting abilities pretty close to its chest. Documentation for Google's Site Search product, for example, says that it [cannot index content](#) contained in JavaScript. A [primer on indexing](#) says it "cannot process the content of some rich media files or dynamic pages."

Yet those looking closely at their server logs may notice that Google is now requesting links that don't appear directly in JavaScript — links that get put together on the fly and Google could not possibly know about unless it could execute at least part of that JavaScript code.

Mark Drummond, chief executive of [Wowd](#), a unique search engine company

Above is a portion of the Forbes article

# Testing Googlebot on Crawling JavaScript Code

- For details see <https://searchengineland.com/tested-googlebot-crawls-javascript-heres-learned-220157>
- **JavaScript redirects**
  - The redirect was followed and the original page dropped from the index
- **JavaScript Links**
  - Links in dropdown menus, onchange event handlers were fully crawled and followed
- **Dynamically inserted content including text, images, links**
  - Text was crawled and indexed and the page ranked for the content
- **Dynamically inserted metadata and page elements such as <title>, <meta>**
  - Crawled
- **Example with rel=“nofollow”**
  - Worked as expected, the link was not followed