

Crawling - part II

Coverage

Good coverage is obtained by carefully selecting seed URLs and using a good page selection policy to decide what to crawl next.

Breadth-first search is adequate when you have simple needs, but many techniques outperform it. It particularly helps to have an existing index from a previous crawl.

Coverage Goals

The Internet is too large and changes too rapidly for any crawler to be able to crawl and index it all. Instead, a crawler should focus on strategic crawling to balance coverage and freshness.

A crawler should prioritize crawling high-quality content to better answer user queries. The Internet contains a lot of spam, redundant information, and pages which aren't likely to be relevant to users' information needs.

```
def crawl(seeds):  
    # The frontier is initially the seed set  
    frontier.add_pages(seeds)  
  
    # Iteratively crawl the next item in the frontier  
    while not frontier.is_empty():  
  
        # Crawl the next URL and extract anchor tags from it  
        url = frontier.choose_next()  
        page = crawl_url(url)  
        urls = parse_page(page)  
  
        # Update the frontier and send the page to the indexer  
        frontier.add_pages(urls)  
        send_to_indexer(page)
```

Basic Crawler Algorithm

Selection Policies

A **selection policy** is an algorithm used to select the next page to crawl. Standard approaches include:

- **Breadth-first search:** This distributes requests across domains relatively well and tends to download high-PageRank pages early.
- **Backlink count:** Prioritize pages with more in-links from already-crawled pages.
- **Larger sites first:** Prioritize pages on domains with many pages in the frontier.
- **Partial PageRank:** Approximate PageRank scores are calculated based on already-crawled pages.

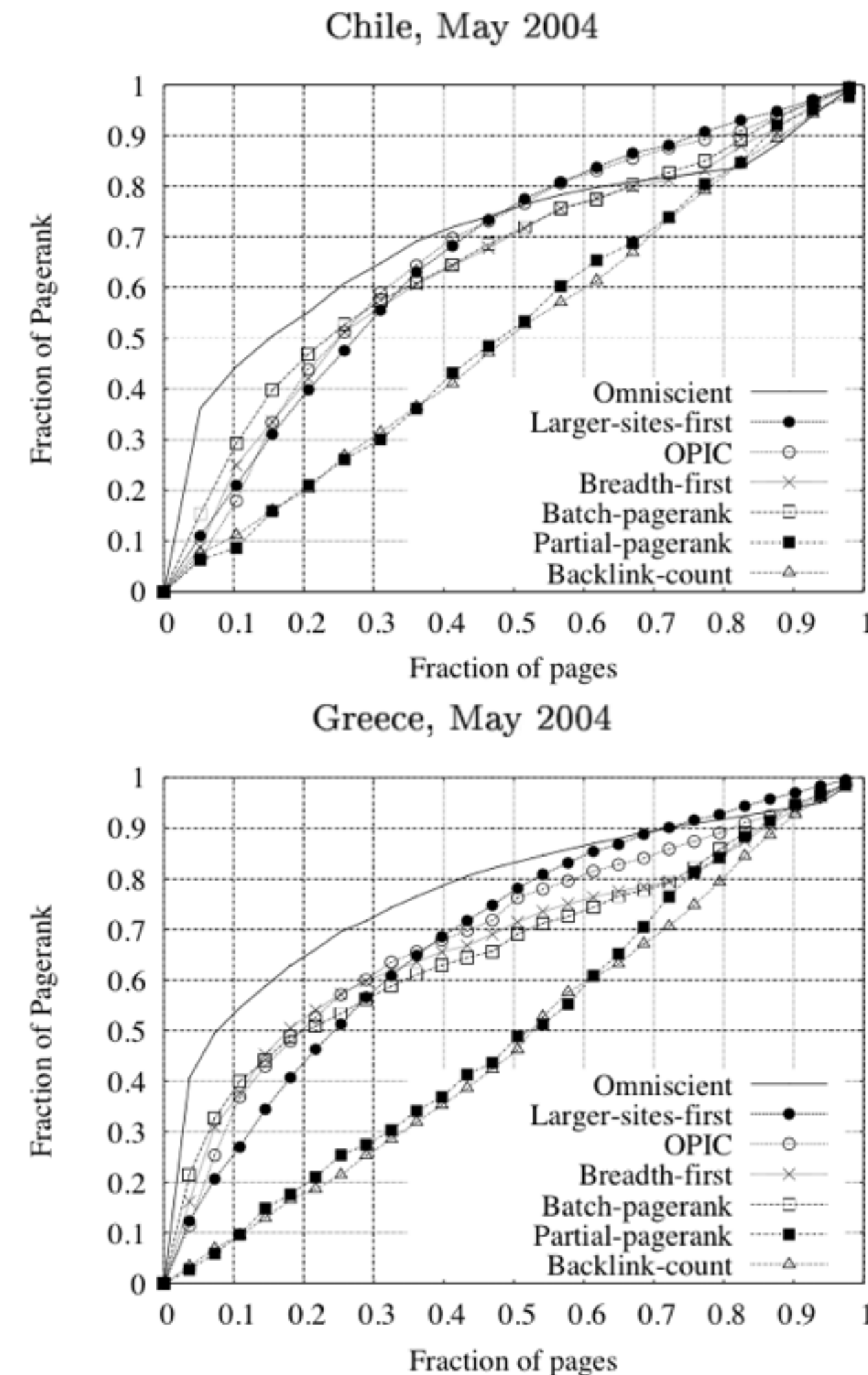
There are also approaches which estimate page quality based on a prior crawl.

Comparing Approaches

Baeza-Yates et al compare these approaches to find out which fraction of high quality pages in a collection is crawled by each strategy at various points in a crawl.

Breadth-first search does relatively poorly. Larger sites first is among the best approaches, along with “historical” approaches which take PageRank scores from a prior crawl into account.

OPIC, a fast approximation to PageRank which can be calculated on the fly, is another good choice. The “omniscient” baseline always fetches the highest PR page in the frontier.



Obtaining Seed URLs

It's important to choose the right sites to initialize your frontier. A simple baseline approach is to start with the sites in an Internet directory, such as <http://www.dmoz.org>.

In general, good hubs tend to lead to many high-quality web pages. These hubs can be identified with a careful analysis of a prior crawl.

dmoz In partnership with **AOL.**

Follow @dmoz | [about dmoz](#) | [dmoz blog](#) | [suggest URL](#) | [help](#) | [link](#) | [editor login](#)

Search [advanced](#)

Arts
[Movies](#), [Television](#), [Music...](#)

Business
[Jobs](#), [Real Estate](#), [Investing...](#)

Computers
[Internet](#), [Software](#), [Hardware...](#)

Games
[Video Games](#), [RPGs](#), [Gambling...](#)

Health
[Fitness](#), [Medicine](#), [Alternative...](#)

Home
[Family](#), [Consumers](#), [Cooking...](#)

Kids and Teens
[Arts](#), [School Time](#), [Teen Life...](#)

News
[Media](#), [Newspapers](#), [Weather...](#)

Recreation
[Travel](#), [Food](#), [Outdoors](#), [Humor...](#)

Reference
[Maps](#), [Education](#), [Libraries...](#)

Regional
[US](#), [Canada](#), [UK](#), [Europe...](#)

Science
[Biology](#), [Psychology](#), [Physics...](#)

Shopping
[Clothing](#), [Food](#), [Gifts...](#)

Society
[People](#), [Religion](#), [Issues...](#)

Sports
[Baseball](#), [Soccer](#), [Basketball...](#)

World
[Català](#), [Česky](#), [Dansk](#), [Deutsch](#), [Español](#), [Esperanto](#), [Français](#), [Galego](#), [Hrvatski](#), [Italiano](#), [Lietuvių](#), [Magyar](#), [Nederlands](#), [Norsk](#), [Polski](#), [Português](#), [Română](#), [Slovensky](#), [Suomi](#), [Svenska](#), [Türkçe](#), [Български](#), [Ελληνικά](#), [Русский](#), [Українська](#), [العربية](#), [עברית](#), [עִבְרִית](#), [ไทย](#), [日本語](#), [简体中文](#), [繁體中文](#), ...

[Become an Editor](#) Help build the largest human-edited directory of the web

Copyright © 1998-2015 AOL Inc.



<http://www.dmoz.org>

The Deep Web

Despite these techniques, a substantial fraction of web pages remains uncrawled and unindexed by search engines. These pages are known as “the deep web.”

These pages are missed for many reasons.

- Dynamically-generated pages, such as pages that make heavy use of AJAX, rely on web browser behavior and are missed by a straightforward crawl.
- Many pages reside on private web sites and are protected by passwords.
- Some pages are intentionally hidden, using robots.txt or more sophisticated approaches such as “darknet” software.

Special crawling and indexing techniques are used to attempt to index this content, such as rendering pages in a browser during the crawl.

Freshness

The web is constantly changing, and re-crawling the latest changes quickly can be challenging.

It turns out that aggressively re-crawling as soon as a page changes is sometimes the wrong approach: it's better to use a cost function associated with the expected age of the content, and tolerate a small delay between re-crawls.

Page Freshness

The web is constantly changing as content is added, deleted, and modified. In order for a crawler to reflect the web as users will encounter it, it needs to recrawl content soon after it changes.

This need for freshness is key to providing a good search engine experience. For instance, when breaking news develops, users will rely on your search engine to stay updated.

It's also important to refresh less time-sensitive documents so the results list doesn't contain spurious links to deleted or modified data.

HTTP HEAD Requests

A crawler can determine whether a page has changed by making an HTTP HEAD request.

The response provides the HTTP status code and headers, but not the document body. The headers include information about when the content was last updated.

However, it's not feasible to constantly send HEAD requests, so this isn't an adequate strategy for freshness.

Request

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

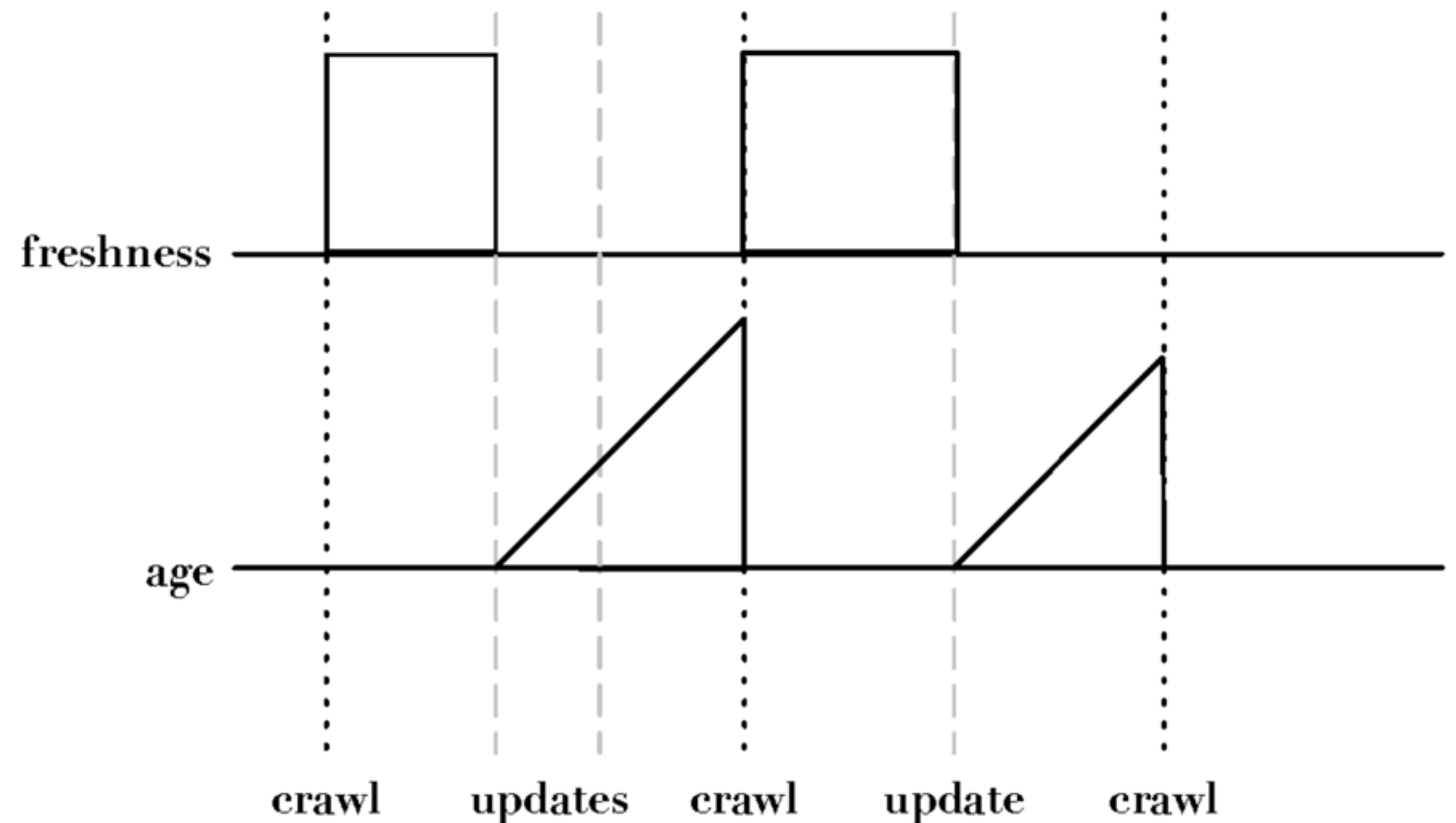
Response

```
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
ETag: "239c33-2576-2a2837c0"
Accept-Ranges: bytes
Content-Length: 9590
Connection: close
Content-Type: text/html; charset=ISO-8859-1
```

Freshness vs. Age

It turns out that optimizing to minimize freshness is a poor strategy: it can lead the crawler to ignore important sites.

Instead, it's better to re-crawl pages when the age of the last crawled version exceeds some limit. The *age* of a page is the elapsed time since the first update after the most recent crawl.



Freshness is binary, age is continuous.

Expected Page Age

The expected age of a page t days after it was crawled depends on its update probability:

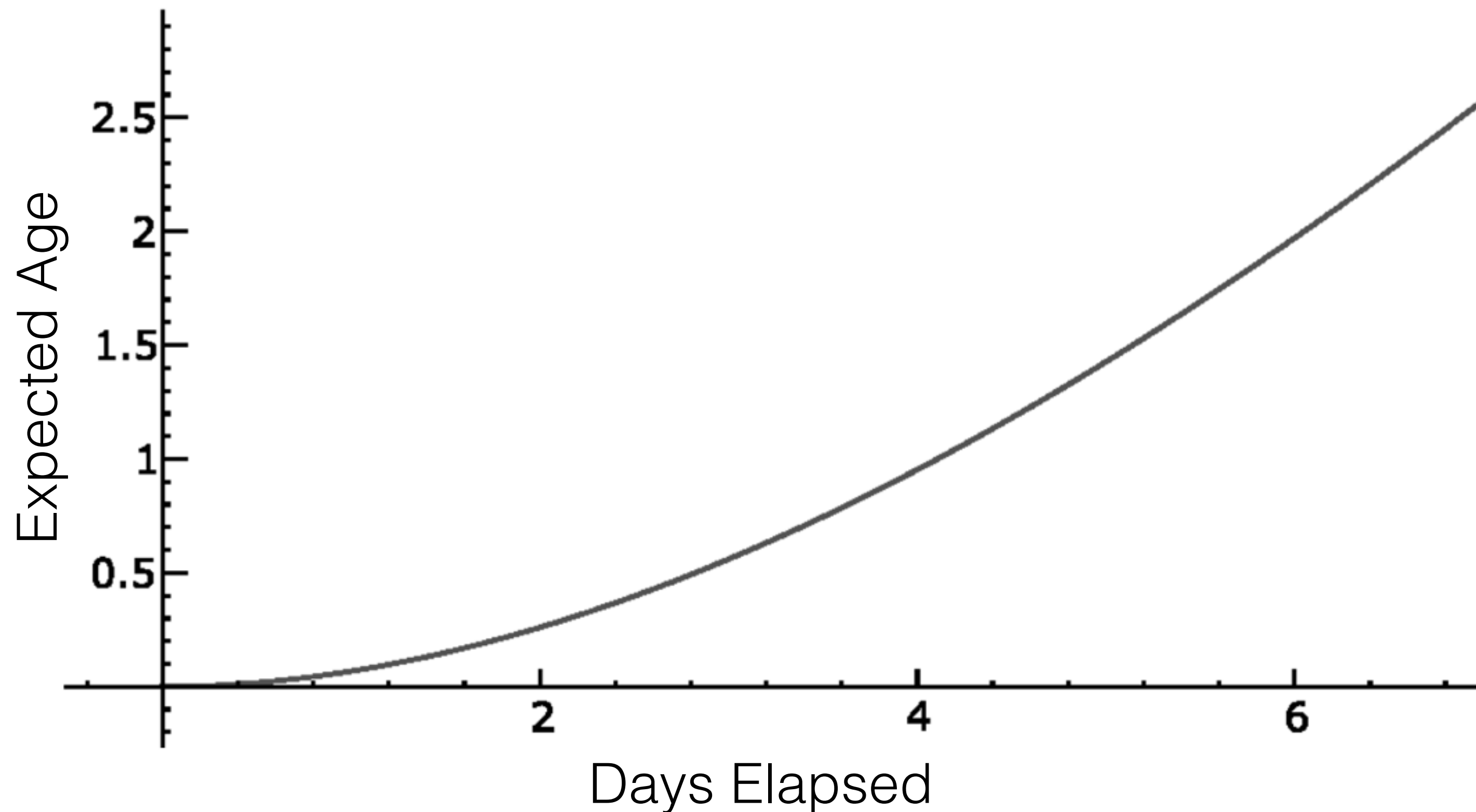
$$age(\lambda, t) = \int_0^t P(\text{page changed at time } x)(t - x)dx$$

On average, page updates follow a Poisson distribution – the time until the next update is governed by an exponential distribution. This makes the expected age:

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

Cost of Not Re-crawling

The cost of not re-crawling a page grows exponentially in the time since the last crawl. For instance, with page update frequency $\lambda = 1/7$ days:



Freshness vs. Coverage

The opposing needs of Freshness and Coverage need to be balanced in the scoring function used to select the next page to crawl.

Finding an optimal balance is still an open question. Fairly recent studies have shown that even large name-brand search engines only do a modest job at finding the most recent content.

However, a reasonable approach is to include a term in the page priority function for the expected age of the page content. For important domains, you can track the site-wide update frequency λ .

Technique	Objectives		Factors considered		
	<i>Coverage</i>	<i>Freshness</i>	<i>Importance</i>	<i>Relevance</i>	<i>Dynamicity</i>
Breadth-first search [43, 95, 108]	✓				
Prioritize by indegree [43]	✓		✓		
Prioritize by PageRank [43, 45]	✓		✓		
Prioritize by site size [9]	✓		✓		
Prioritize by spawning rate [48]	✓				✓
Prioritize by search impact [104]	✓		✓	✓	
Scoped crawling (Section 4.2)	✓			✓	
Minimize obsolescence [41, 46]		✓	✓		✓
Minimize age [41]		✓	✓		✓
Minimize incorrect content [99]		✓	✓		✓
Minimize embarrassment [115]		✓	✓	✓	✓
Maximize search impact [103]		✓	✓	✓	✓
Update capture (Section 5.2)		✓		✓	✓
WebFountain [56]	✓	✓			✓
OPIC [1]	✓	✓	✓		

3.2 Taxonomy of crawl ordering techniques.

Pitfalls of Crawling

A breadth-first search implementation of crawling is not sufficient for coverage, freshness, spam avoidance, or other needs of a real crawler.

Scaling the crawler up takes careful engineering, and often detailed systems knowledge of the hardware architecture you're developing for.

Crawling at Scale

A commercial crawler should support thousands of HTTP requests per second. If the crawler is distributed, that applies for each node. Achieving this requires careful engineering of each component.

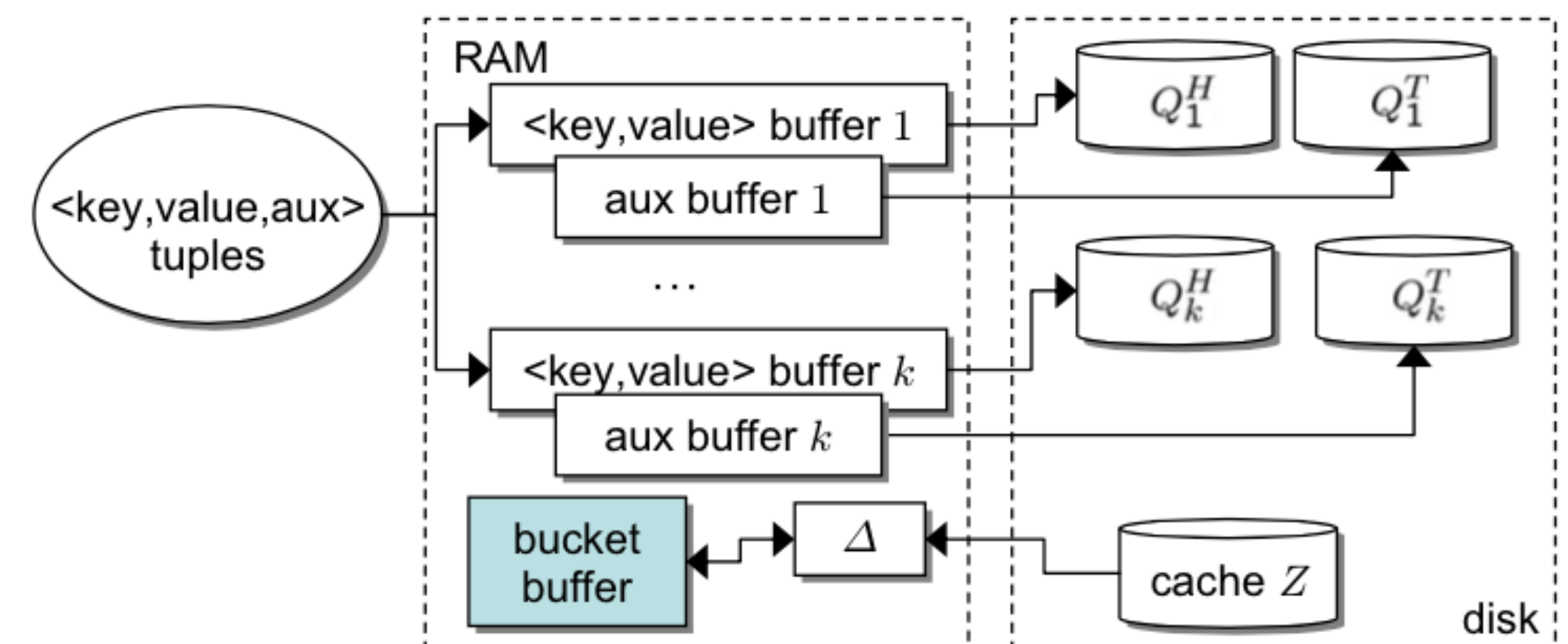
- DNS resolution can quickly become a bottleneck, particularly because sites often have URLs with many subdomains at a single IP address.
- The frontier can grow extremely rapidly – hundreds of thousands of URLs per second are not uncommon. Managing the filtering and prioritization of URLs is a challenge.
- Spam and malicious web sites must be addressed, lest they overwhelm the frontier and waste your crawling resources. For instance, some sites respond to crawlers by intentionally adding seconds of latency to each HTTP response. Other sites respond with data crafted to confuse, crash, or mislead a crawler.

Duplicate URL Detection at Scale

Lee et al's DRUM algorithm gives a sense of the requirements of large scale de-duplication.

It manages a collection of tuples of keys (hashed URLs), values (arbitrary data, such as quality scores), and aux data (URLs). It supports the following operations:

- **check** – Does a key exist? If so, fetch its value.
- **update** – Merge new tuples into the repository.
- **check+update** – Check and update in a single pass.



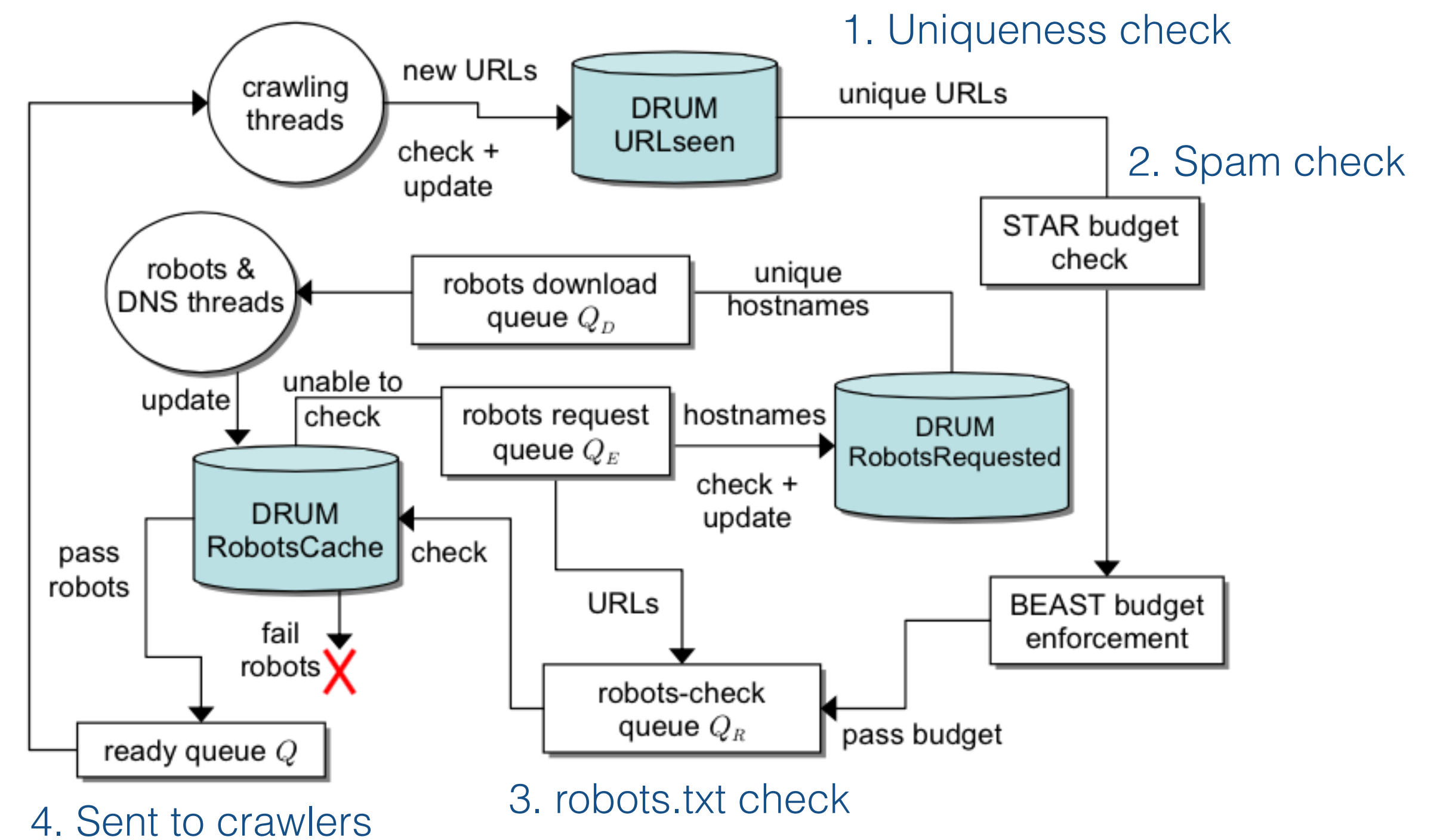
Data flow for DRUM: A tiered system of buffers in RAM and on disk is used to support large-scale operations.

IRLBot Operation

DRUM is used as a storage for the IRLBot crawler. A new URL passes through the following steps.

1. The URLSeen DRUM checks whether the URL has already been fetched.
2. If not, two budget checks filter out spam links (discussed next).
3. Next, we check whether the URL passes its robots.txt. If necessary, we fetch robots.txt from the server.
4. Finally, the URL is passed to the queue to be crawled by the next available thread.

IRLBot Architecture



Link Spam

The web is full of link farms and other forms of link spam, generally posted by people trying to manipulate page quality measures such as PageRank.

These links waste a crawler's resources, and detecting and avoiding them is important for correct page quality calculations.

One way to mitigate this, implemented in IRLBot, is based on the observation that spam servers tend to have very large numbers of pages linking to each other.

They assign a budget to each domain based on the number of in-links from other domains. The crawler de-prioritizes links from domains which have exceeded their budget, so link-filled spam domains are largely ignored.

Spider Traps

A *spider trap* is a collection of web pages which, intentionally or not, provide an infinite space of URLs to crawl.

Some site administrators place spider traps on their sites in order to trap or crash spambots, or defend against malicious bandwidth-consuming scripts.

A common example of a benign spider trap is a calendar which links continually to the next year.



A benign spider trap on
<http://www.timeanddate.com>

Avoiding Spider Traps

The first defense against spider traps is to have a good politeness policy, and always follow it.

- By avoiding frequent requests to the same domain, you reduce the possible damage a trap can do.
- Most sites with spider traps provide instructions for avoiding them in robots.txt.

```
[...]  
User-agent: *  
Disallow: /createshort.html  
Disallow: /scripts/savecustom.php  
Disallow: /scripts/wquery.php  
Disallow: /scripts/tzq.php  
Disallow: /scripts/savepersonal.php  
Disallow: /information/mk/  
Disallow: /information/feedback-save.php  
Disallow: /information/feedback.html?  
Disallow: /gfx/stock/  
Disallow: /bm/  
Disallow: /eclipse/in/*?iso  
Disallow: /custom/save.php  
Disallow: /calendar//index.html  
Disallow: /calendar//monthly.html  
Disallow: /calendar//custom.html  
Disallow: /counters//newyeara.html  
Disallow: /counters//worldfirst.html  
[...]
```

From <http://www.timeanddate.com/robots.txt>

Storing Crawled Content

We need to normalize and store the contents of web documents so they can be indexed, so snippets can be generated, and so on.

Online documents have many formats and encoding schemes. There are hundreds of character encoding systems we haven't mentioned here.

A good document storage system should support efficient random access for lookups, updates, and content retrieval. Often, a distributed storage system like Big Table is used.

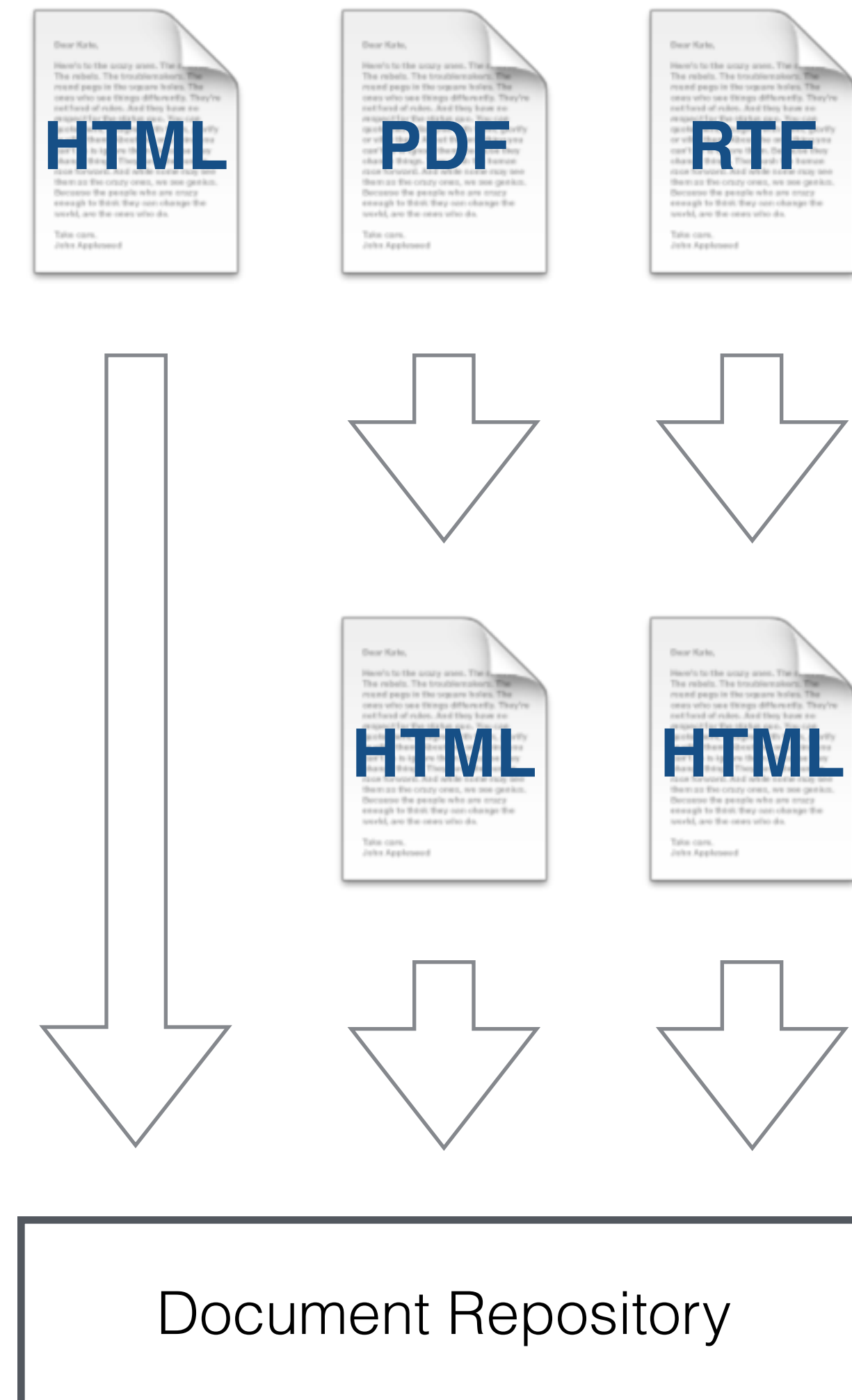
Content Conversion

Downloaded page content generally needs to be converted into a stream of tokens before it can be indexed.

Content arrives in hundreds of incompatible formats: Word documents, PowerPoint, RTF, OTF, PDF, etc.

Conversion tools are generally used to transform them into HTML or XML.

Depending on your needs, the crawler may store the raw document content and/or normalized content output from a converter.



Character Encodings

Crawled content will be represented with many different character encodings, which can easily confuse text processors.

A *character encoding* is a map from bits in a file to glyphs on a screen. In English, the basic encoding is ASCII.

ASCII uses 8 bits: 7 bits to represent 128 letters, numbers, punctuation, and control characters and an extra bit for padding.

USASCII code chart

The chart is a grid with 15 rows and 8 columns. The top-left corner has a 3D box labeled 'Bits' with arrows pointing to bit positions b7, b6, b5, b4, b3, b2, b1. The top row of the grid shows bit patterns for columns 0-7. The left side of the grid shows bit patterns for rows 0-15. The grid contains the following data:

Row	b4	b3	b2	b1	Column	0	1	2	3	4	5	6	7
0	0	0	0	0	0	NUL	DLE	SP	0	@	P	`	p
1	0	0	0	1	1	SOH	DC1	!	1	A	Q	a	q
2	0	0	1	0	2	STX	DC2	"	2	B	R	b	r
3	0	0	1	1	3	ETX	DC3	#	3	C	S	c	s
4	0	1	0	0	4	EOT	DC4	\$	4	D	T	d	t
5	0	1	0	1	5	ENQ	NAK	%	5	E	U	e	u
6	0	1	1	0	6	ACK	SYN	&	6	F	V	f	v
7	0	1	1	1	7	BEL	ETB	'	7	G	W	g	w
8	1	0	0	0	8	BS	CAN	(8	H	X	h	x
9	1	0	0	1	9	HT	EM)	9	I	Y	i	y
10	1	0	1	0	10	LF	SUB	*	:	J	Z	j	z
11	1	0	1	1	11	VT	ESC	+	;	K	[k	{
12	1	1	0	0	12	FF	FS	,	<	L	\	l	
13	1	1	0	1	13	CR	GS	-	=	M]	m	}
14	1	1	1	0	14	SO	RS	.	>	N	^	n	~
15	1	1	1	1	15	SI	US	/	?	O	_	o	DEL

Image courtesy Wikipedia

Unicode

The various Unicode encodings were invented to support a broader range of characters. Unicode is a single mapping from numbers to glyphs, with various encoding schemes of different sizes.

- UTF-8 uses one byte for ASCII characters, and more bytes for extended characters. It's often preferred for file storage.
- UTF-32 uses four bytes for every character, and is more convenient for use in memory.

	ASCII	UTF-8	UTF-32
A	0x41	0x41	0x00000041
&	0x26	0x26	0x00000026
π	N/A	0xCF 0x80	0x000003C0
👍	N/A	0xF0 0x9F 0x91 0x8D	0x0001F44D

UTF-8

UTF-8 uses a variable-length encoding scheme.

If the most significant (leftmost) bit of a given byte is set, the character takes another byte.

The first 128 numbers are the same as ASCII, so any ASCII document could be said to (retroactively) use UTF-8.

UTF-8 is designed to minimize disk space for documents in many languages, but UTF-32 is faster to decode and easier to use in memory.

UTF-8 Encoding Scheme

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	11110xxxx	10xxxxxx	10xxxxxx	
65536–1114111	10000–10FFFF	111110xxx	10xxxxxx	10xxxxxx	10xxxxxx

Document Repositories

What do we need from our document repository?

- **Fast random access** – need to store and obtain documents by their URLs (or a hash of the URL)
- **Fast document updates** – need to associate and update metadata with documents, and replace (or append to) records when documents are re-crawled
- **Compressed storage** – greatly reduces storage needs, and minimizes disk reads for access
- **Large file storage** – multiple documents are stored in a single large file to reduce filesystem overhead

Most companies use custom storage systems, or distributed systems like Big Table.

Large File Storage

TREC Web Format

Placing millions or billions of web pages in individual files results in substantial filesystem overhead for opening, writing, and finding files.

It's important to store many files into larger files, generally with an indexing scheme to give fast random access.

A simple index might store a B-tree mapping document URL hash values to the byte offset to the document contents in the file.

```
<DOC>
<DOCNO>WTX001-B01-10</DOCNO>
<DOCHDR>
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, 25 Nov 1996 05:31:24 GMT
</DOCHDR>
<HTML>
<TITLE>Tropical Fish Store</TITLE>
Coming soon!
</HTML>
</DOC>
<DOC>
<DOCNO>WTX001-B01-109</DOCNO>
<DOCHDR>
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, 25 Nov 1996 05:31:24 GMT
</DOCHDR>
<HTML>
<TITLE>Fish Information</TITLE>
This page will soon contain interesting
information about tropical fish.
</HTML>
</DOC>
```

Vertical Search

Vertical Search depends on crawling a collection on the topic of interest.

General search engines also use topical crawlers to improve their coverage for key topics.

The main trick to topical crawling is finding topical pages which are only reachable by exploring off-topic pages through careful risk-taking.

Vertical Search

Vertical Search engines focus on a particular domain of information.

The primary difference between vertical and general search engines is the set of documents they crawl. Vertical Search engines typically use what are known as *topical crawlers*.

The screenshot shows the CiteSeerX search interface. At the top, there are navigation tabs for 'Documents', 'Authors', and 'Tables', along with links for 'MetaCart', 'Sign up', and 'Log in'. The search bar contains the text 'vertical search' and a 'Search' button. Below the search bar, there is a checkbox for 'Include Citations' and a link for 'Advanced Search'. The search results are displayed in a list format, showing the first 10 of 119,317 results. The first result is 'Comparison of Three Vertical Search Spiders' by Michael Chau, published in IEEE Computer in 2003. The second result is 'Depth first search and linear graph algorithms' by Robert Tarjan, published in SIAM Journal on Computing in 1972. The third result is 'Suffix arrays: A new method for on-line string searches' by Udi Manber and Argo Gene Myerst, published in SIAM J. Comput in 1993. The fourth result is 'A greedy randomized adaptive search procedure for the 2-partition problem' by Thomas A. Feo and M. Pardalos, published in Operations Research in 1994. On the right side of the page, there is a 'Tools' section with a 'Sorted by:' dropdown menu set to 'Relevance' and a 'Try your query at:' section with links to 'Scholar', 'Yahoo!', 'DBLP', 'Bing', 'CSB', and 'Academic'.

CiteSeer, Vertical Search for Research

Topical Crawlers

Topical Crawlers focus on documents related to a particular topic of interest.

These crawlers are useful for improving the collection quality of general search engines, too. Many search engines use a variety of topical crawlers to supplement their primary crawler.

A basic approach uses a topical set of seed URLs and text classifiers to decide whether links appear to be on topic.

```
def crawl(seeds):  
  
    # High quality topical hubs are used as seeds  
    frontier.add_pages(seeds)  
  
    # Iteratively crawl the next item in the frontier  
    while not frontier.is_empty():  
  
        # Crawl the next URL and extract anchor tags from it  
        url = frontier.choose_next()  
        page = crawl_url(url)  
        urls = parse_page(page)  
  
        # The URLs are filtered to stay on topic  
        urls = filter_by_topic(urls)  
  
        # Update the frontier and send the page to the indexer  
        frontier.add_pages(urls)  
        send_to_indexer(page)
```

Basic Topical Crawler

Text Classifiers

Text classification is a Machine Learning task that we'll see later in the course.

The idea is to use properties of the URL, anchor text, and document to predict whether the URL links to a page on the topic of interest.

For example, we could use a unigram language model trained on anchor text for topical links.

Classification with Language Models

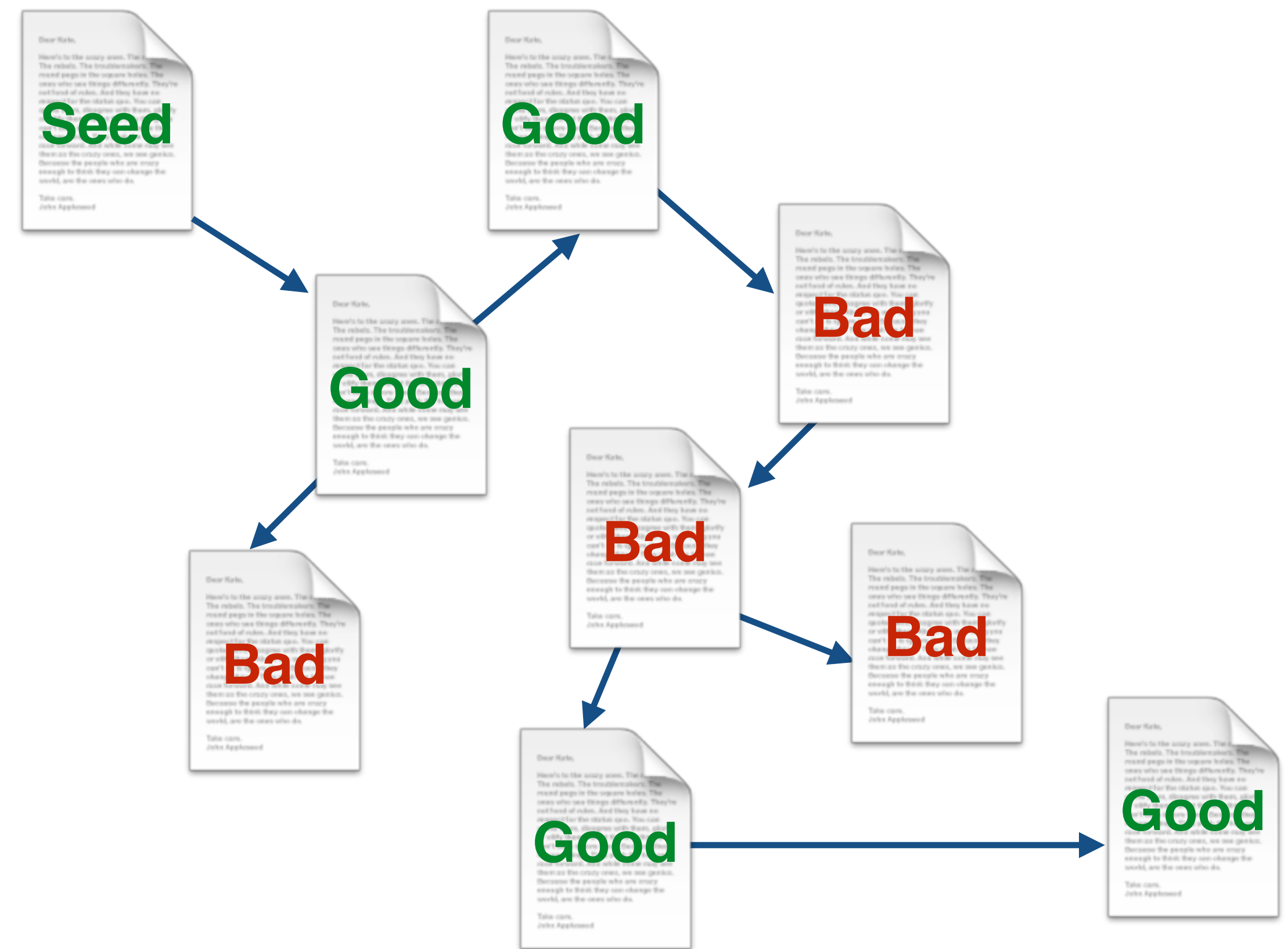
1. Collect anchor text for links to topical and non-topical pages.
2. Train a unigram language model by producing smoothed probability estimates of topicality for each term.
3. Classify new links using the odds ratio from training data for some threshold λ :

$$\prod_{w \in \text{text}} \frac{\Pr(w | \text{topic} = 1)}{\Pr(w | \text{topic} = 0)} \stackrel{?}{>} \lambda$$

Explore vs. Exploit Tradeoff

More sophisticated topical crawlers use machine learning techniques to balance the tradeoff between *exploring* new territory and *exploiting* links which are probably high-quality.

- Exploit-only strategies may miss high quality pages which aren't tightly linked to the seed set.
- Explore-only strategies will ignore high-quality pages we can easily find.



Sometimes bad links must be explored to find good links

Careful Exploration

There are many ways to balance exploration and exploitation, and this topic is actively researched for many applications. Here are some simple ways for this task.

- Adjust the classification threshold to manage your risk threshold.
- Flip a biased coin to decide whether to visit a page which doesn't seem promising.
- If using a document quality score such as PageRank, explore for a while without updating quality scores. Links on crawled pages won't be taken into account, so scores will be somewhat inaccurate and you will explore more.

There are more sophisticated approaches if maximizing performance is important.

Crawling Structured Data

In addition to the obvious content for human readers, the web contains a great deal of structured content for use in automated systems.

- Document feeds are an important way to manage freshness at some of the most frequently-updated web sites.
- Much of the structured data owned by various web entities is published in a structured format. This can provide signals for relevance, and can also aid in reconstructing structured databases.

Structured Web Data

In addition to unstructured document contents, a great deal of structured data exists on the web. We'll focus here on two types:


- Document feeds, which sites use to announce their new content
- Content metadata, used by web authors to publish structured properties of objects on their site

Document Feeds

Sites which post articles, such as blogs or news sites, typically offer a listing of their new content in the form of a document feed.













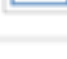

Several common feed formats exist. One of the most popular is RSS, which stands for (take your pick):

- Rich Site Summary
- Really Simple Syndication
- RDF Site Summary
- ...?

CNN.com now offers [podcasting feeds](#). 

Please note that by accessing CNN RSS feeds, you agree to our [terms of use](#).

[What is RSS?](#) | [How do I access RSS?](#)

Title	Copy URLs to RSS Reader		
Top Stories	http://rss.cnn.com/rss/cnn_topstories.rss		
World	http://rss.cnn.com/rss/cnn_world.rss		
U.S.	http://rss.cnn.com/rss/cnn_us.rss		
Business (CNNMoney.com)	http://rss.cnn.com/rss/money_latest.rss		
Politics	http://rss.cnn.com/rss/cnn_allpolitics.rss		
Crime	http://rss.cnn.com/rss/cnn_crime.rss		
Technology	http://rss.cnn.com/rss/cnn_tech.rss		

<http://www.cnn.com/services/rss/>

RSS Format

RSS is an XML format for document listings.

RSS files are obtained just like web pages, with HTTP GET requests.

The `ttl` field provides an amount of time (in minutes) that the contents should be cached.

RSS feeds are very useful for efficiently managing freshness of news and blog content.

```
<?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>
    <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 Example

Structured Data

Many web pages are generated from structured data in databases, which can be useful for search engines and other crawled document collections.

Several schemas exist for web authors to publish their structured data for these tools.

The WHATWG web specification working group has produced several standard formats for this data, such as *microdata* embedded in HTML.

```
<section itemscope itemtype="http://schema.org/Person">
  Hello, my name is
  <span itemprop="name">John Doe</span>,
  I am a
  <span itemprop="jobTitle">graduate research assistant</span>
  at the
  <span itemprop="affiliation">University of Dreams</span>.
  My friends call me
  <span itemprop="additionalName">Johnny</span>.
  You can visit my homepage at
  <a href="http://www.JohnnyD.com" itemprop="url">www.JohnnyD.com</a>.
  <section itemprop="address" itemscope itemtype="http://schema.org/PostalAddress">
    I live at
    <span itemprop="streetAddress">1234 Peach Drive</span>,
    <span itemprop="addressLocality">Warner Robins</span>,
    <span itemprop="addressRegion">Georgia</span>.
  </section>
</section>
```

Source: [http://en.wikipedia.org/wiki/Microdata_\(HTML\)](http://en.wikipedia.org/wiki/Microdata_(HTML))

Web Ontologies

The main web ontology is published at schema.org. These schemas are used to annotate web pages for automated information extraction tools.

As the published information is not necessarily authoritative, the data needs to be carefully validated for quality and spam removal.

Popular schema.org entities

- Creative works: [CreativeWork](#), [Book](#), [Movie](#), [MusicRecording](#), [Recipe](#), [TVSeries](#) ...
- Embedded non-text objects: [AudioObject](#), [ImageObject](#), [VideoObject](#)
- [Event](#)
- [Health and medical types](#): notes on the health and medical types under [MedicalEntity](#).
- [Organization](#)
- [Person](#)
- [Place](#), [LocalBusiness](#), [Restaurant](#) ...
- [Product](#), [Offer](#), [AggregateOffer](#)
- [Review](#), [AggregateRating](#)
- [Action](#)

Crawling - Wrap Up

Goals of Crawling

A good crawler will balance several factors:

- **Coverage:** Pages should be selected to maximize the number of distinct high-quality pages.
- **Freshness:** Pages which have been updated should be re-crawled soon.
- **Performance:** Each machine should crawl thousands of pages per second.
- **Politeness:** Requests to the same domain are infrequent, and site owners' requested crawler policies are respected.

Major Challenges

High-performance data structures, such as IRLbot's DRUM, must be used to efficiently de-duplicate URLs, manage robots.txt caches, etc.

Malicious web content should be carefully avoided, and low-quality content (malformed HTML, unreliable web sites, etc.) should be identified and dealt with as appropriate.

Web site owners generally want their information to be crawled, so they provide assistance in terms of sitemaps, RSS, embedded metadata, etc.

Spam Technologies

- **Cloaking**

- Serve fake content to search engine robot
- *DNS cloaking*: Switch IP address. Impersonate

- **Doorway pages**

- Pages optimized for a single keyword that re-direct to the real target

- **Keyword Spam**

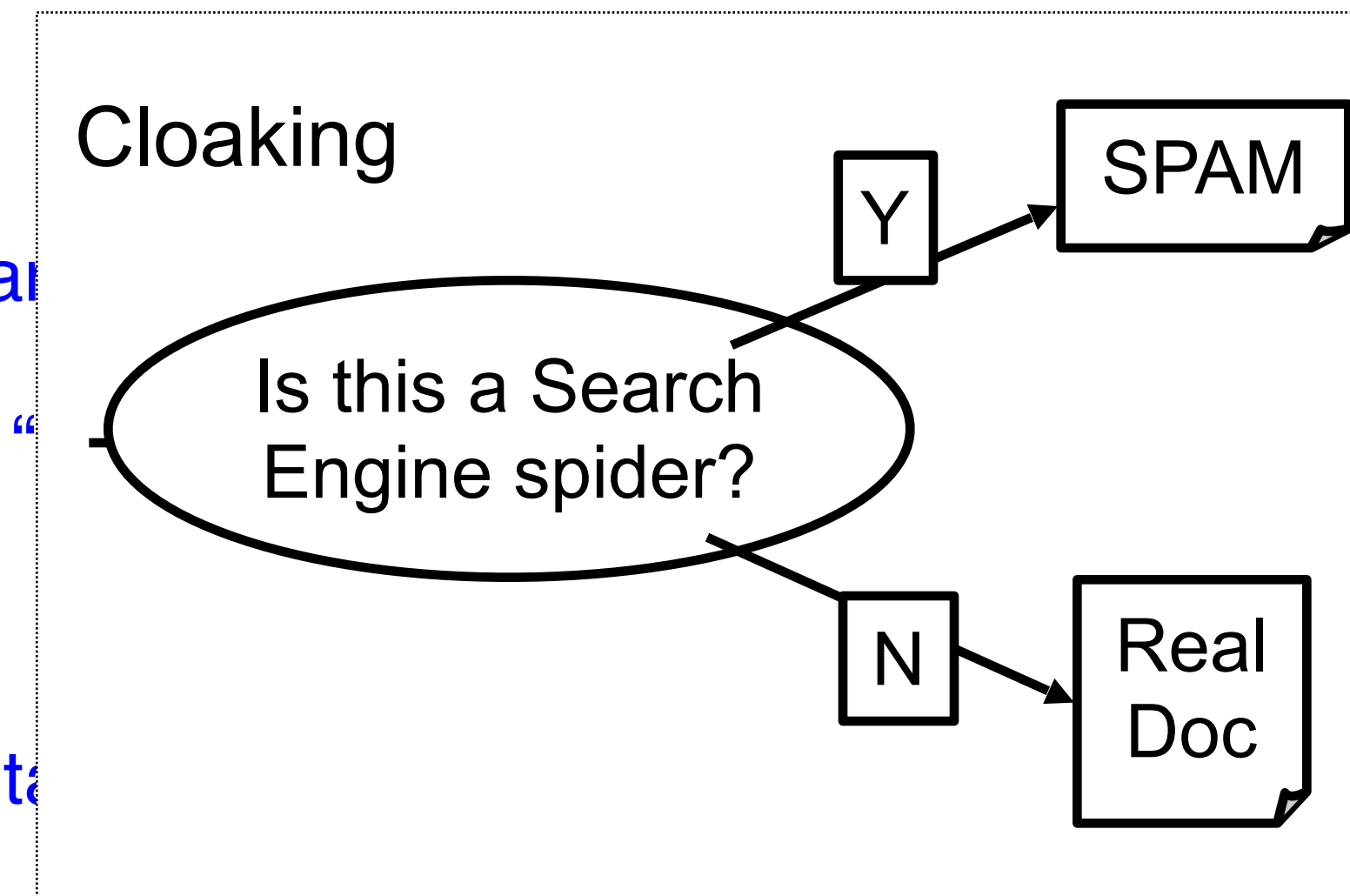
- Misleading meta-keywords, excessive repetition of a term, fake “*nofollow*”
- Hidden text with colors, CSS tricks, etc.

- **Link spamming**

- Mutual admiration societies, hidden links, awards
- *Domain flooding*: numerous domains that point or re-direct to a target

- **Robots**

- Fake click stream
- Fake query stream
- Millions of submissions via Add-Url



Meta-Keywords =

"... London hotels, hotel, holiday inn, hilton, discount, booking, reservation, sex, mp3, britney spears, viagra, ..."