

Web Search: Agenda

- Quick web overview
- How is the web different?
- Types of web information needs
- Crawling
- Link analysis

The WWW was invented in 1989, at CERN.



Sir Tim Berners-Lee
1955 –

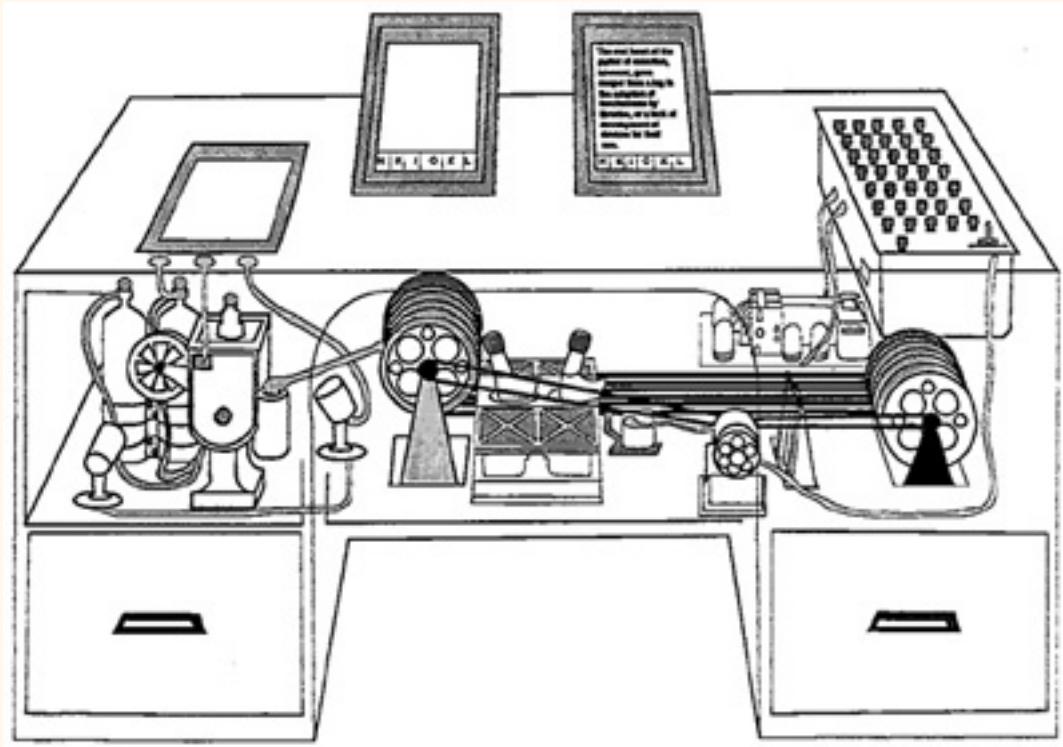
Keys to the WWW: HTTP, HTML, and hypertext.

"I just had to take the hypertext idea and connect it to the [Transmission Control Protocol](#) and [domain name system](#) ideas and—ta-da!—the World Wide Web ... **Creating the web was really an act of desperation**, because the situation without it was very difficult when I was working at CERN later. Most of the technology involved in the web, like the hypertext, like the Internet, multifont text objects, had all been designed already. I just had to put them together. It was a step of generalising, going to a higher level of abstraction, thinking about all the documentation systems out there as being possibly part of a larger imaginary documentation system."

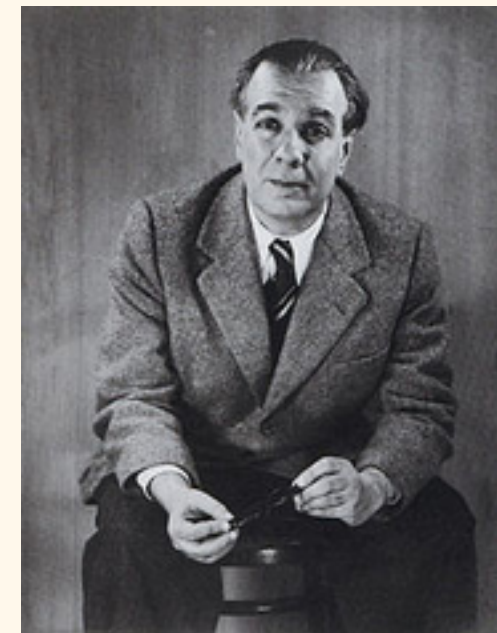
https://commons.wikimedia.org/wiki/File:NeXTcube_first_webserver.JPG

https://en.wikipedia.org/wiki/Tim_Berners-Lee

Depending on who you ask, the concept of hypertext was invented by either:



Vannevar Bush
1890 – 1974



Jorge Luis Borges
1899 – 1986

Regardless, it was not until the 1960s that actual hypertext systems were built...



Ted Nelson
1937 –



Douglas Engelbart
1925 – 2013

... and not until the 1970s that any were really usable.



What distinguished the web from previous information systems?

1. Simplicity

2. Decentralization

3. Hypertext

In other words: anybody could easily add any content and cross-link it any-which-way.

This led to incredibly rapid and uncontrolled growth...
... which in turn posed major challenges for search:

	Standard	Web
<i>Index size</i>	(relatively) small	petabytes+
<i>Updates & Additions</i>	rare	constant
<i>Language</i>	monolingual	aggressively multilingual
<i>Content</i>	consistent (news articles, etc.)	diverse formats
<i>Authors</i>	curated	uncontrolled
<i>Behavior</i>	cooperative	adversarial
<i>Quality</i>	homogeneous	heterogeneous, subjective

Web information needs fall into three broad categories:

Informational:

*Seeking information on a broad topic;
answer usually not found on any single page*
“Leukemia”, “art galleries in Portland”

Navigational:

*Seeking a specific website from a particular
entity*
“alaska airlines”

Transactional:

*A prelude to the user performing a transaction (buying
something, making a reservation, etc.)*
“alaska airlines reservation desk”

Search engines try to infer the category of need from a query (and user behavior: prior queries, clicks, etc.)...

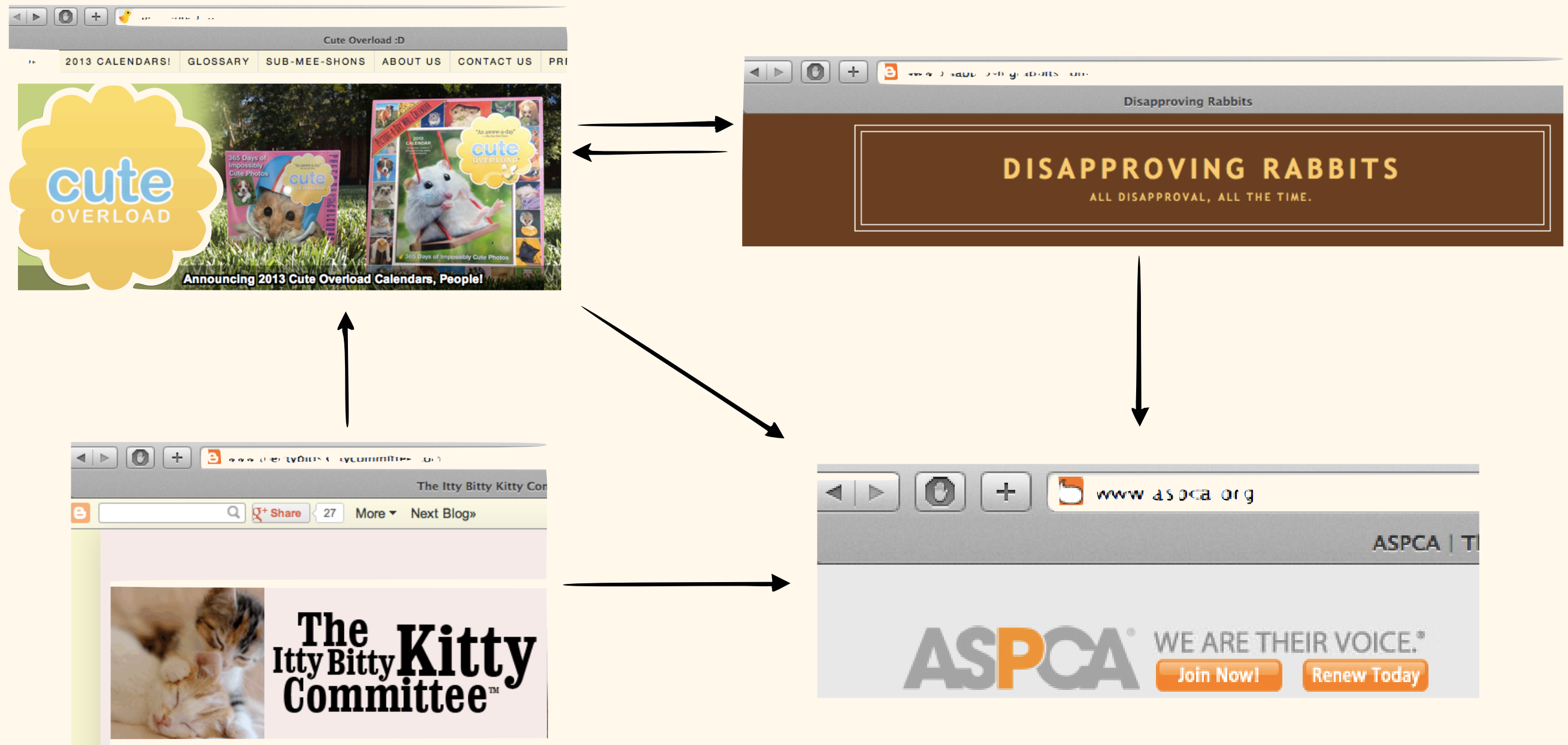
... this is a huge area of research!

Early web search engines took one of two basic approaches:



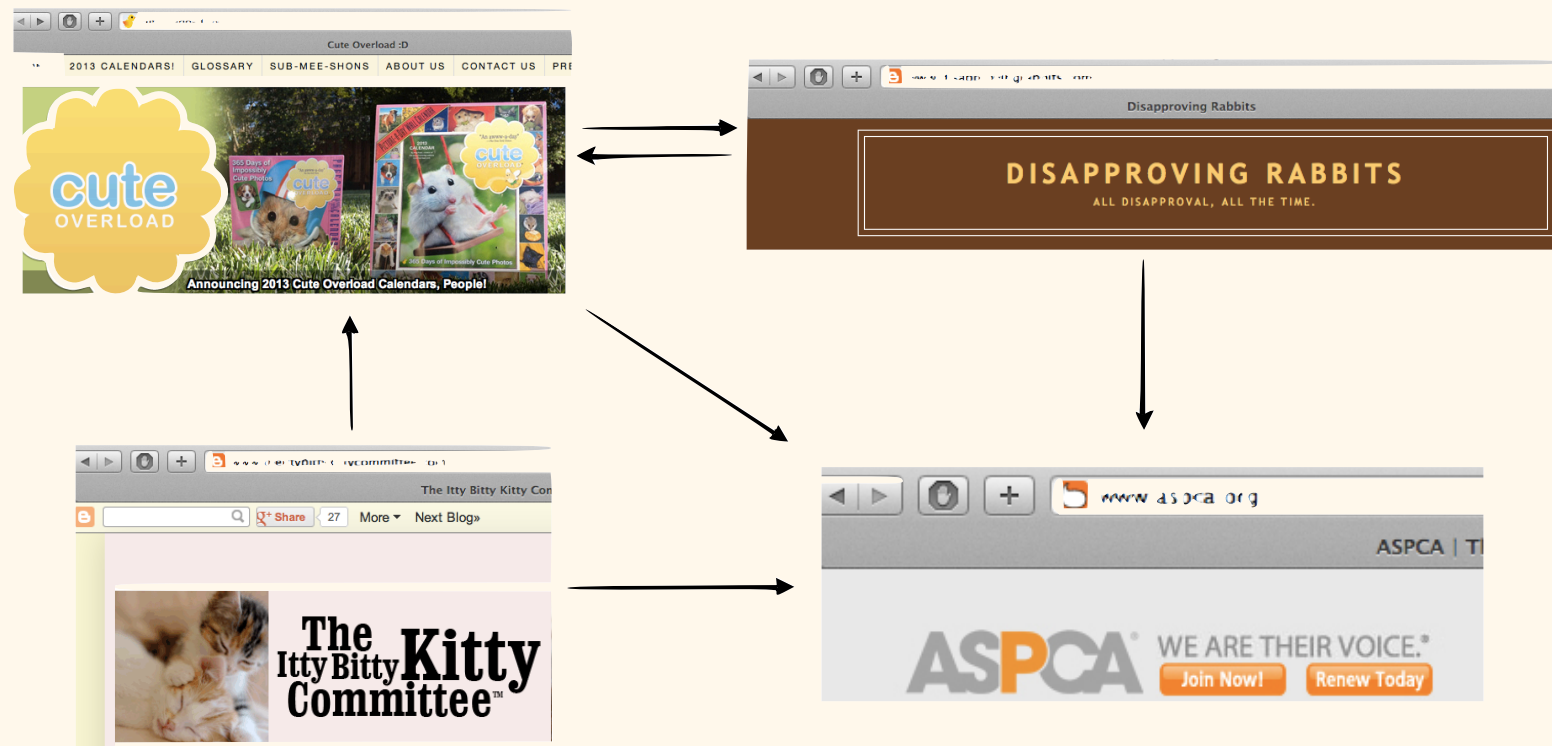
Both were based primarily on the *content* of the indexed web pages.

The web can be thought of as a directed graph...



This is a very useful formulation!

The web can be thought of as a directed graph...



A page's *in-degree*: # incoming links; *out-degree*: # outgoing links

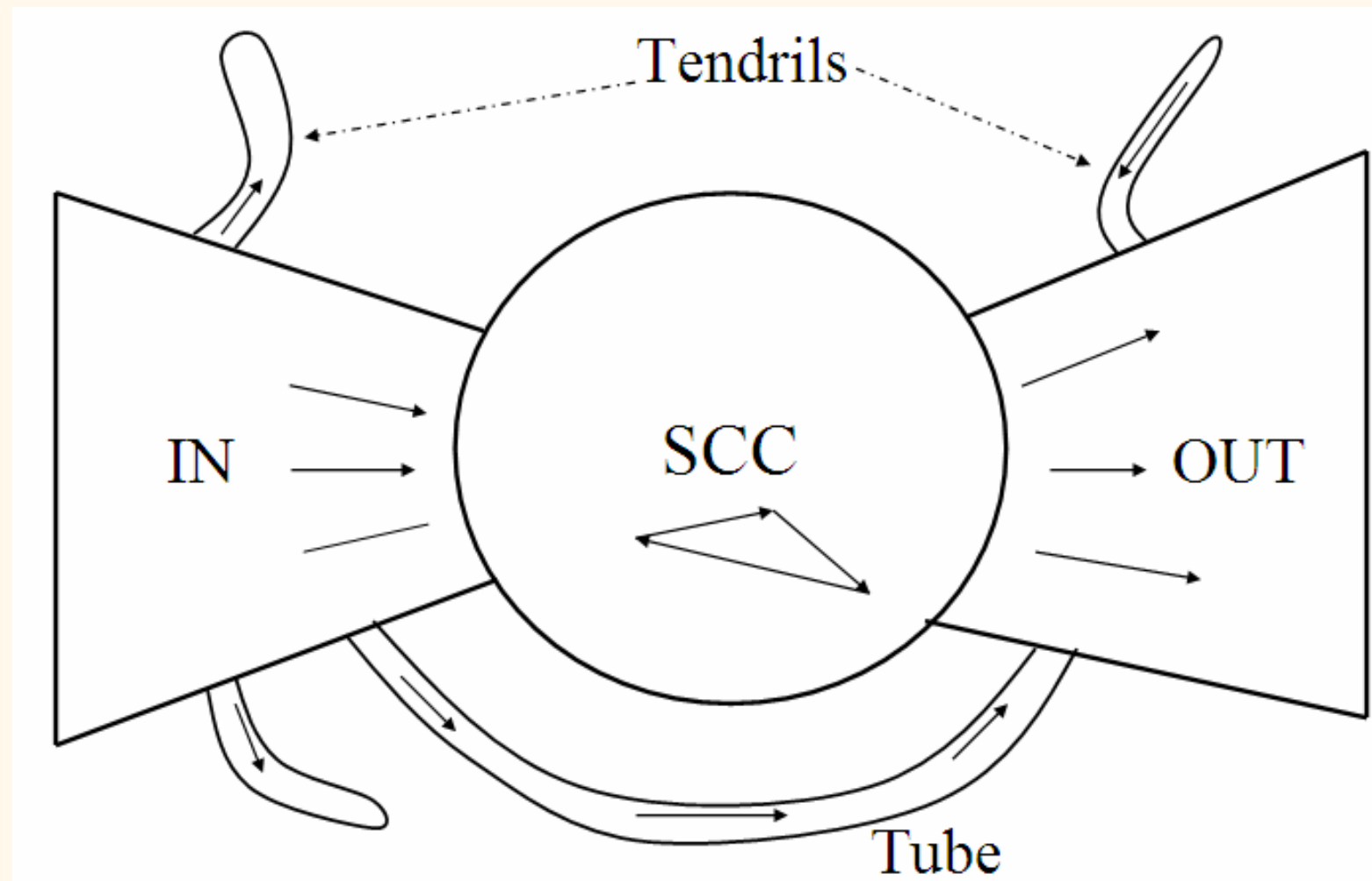
The web graph is *not* strongly connected;

Links are *non-random*;

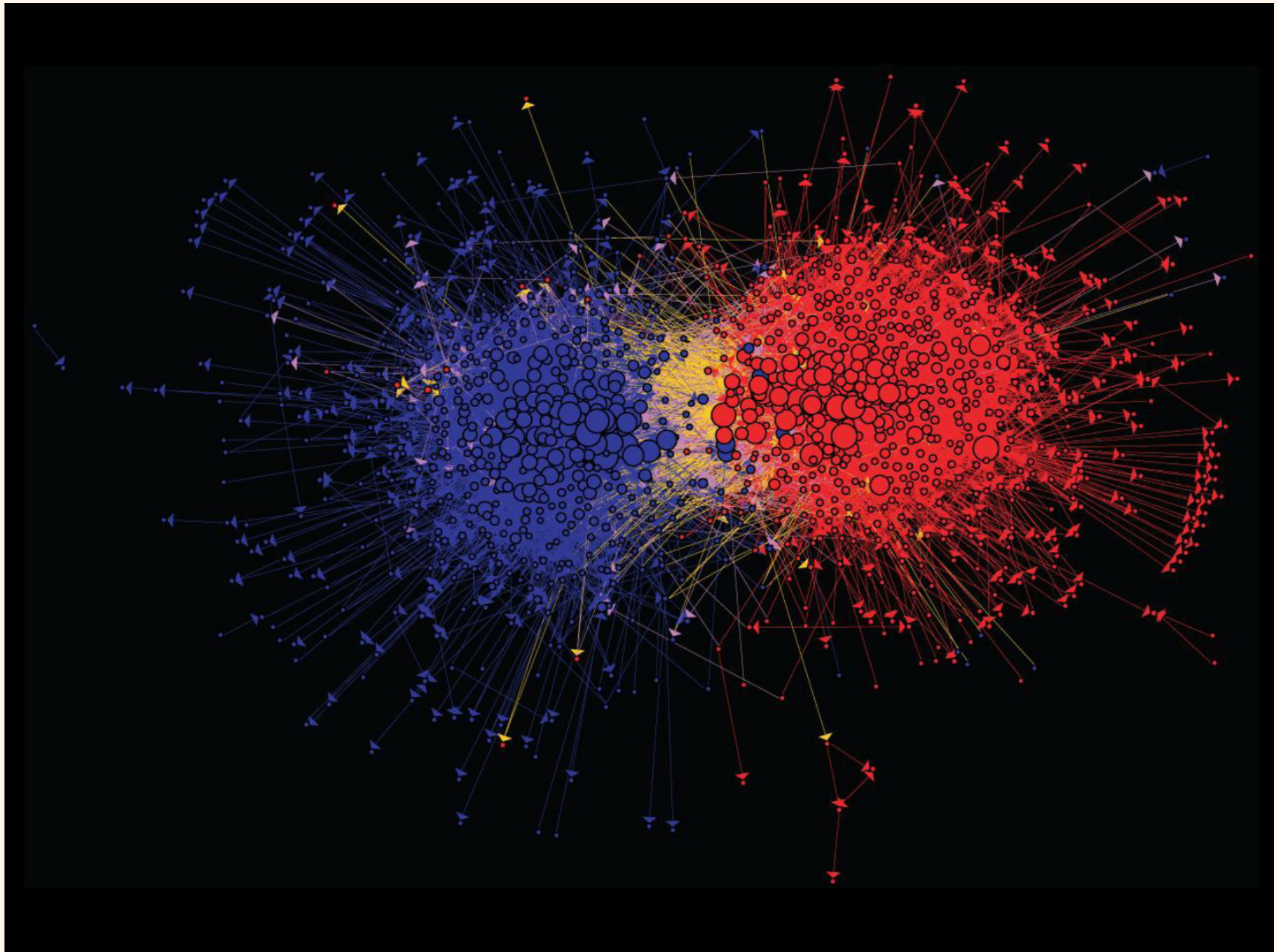
The link distribution is often said to follow a power-law distribution:

num. pages with in-degree of i is proportional to $1/i^\alpha$

The web can be thought of as a directed graph...



► **Figure 19.4** The bowtie structure of the Web. Here we show one tube and three tendrils.

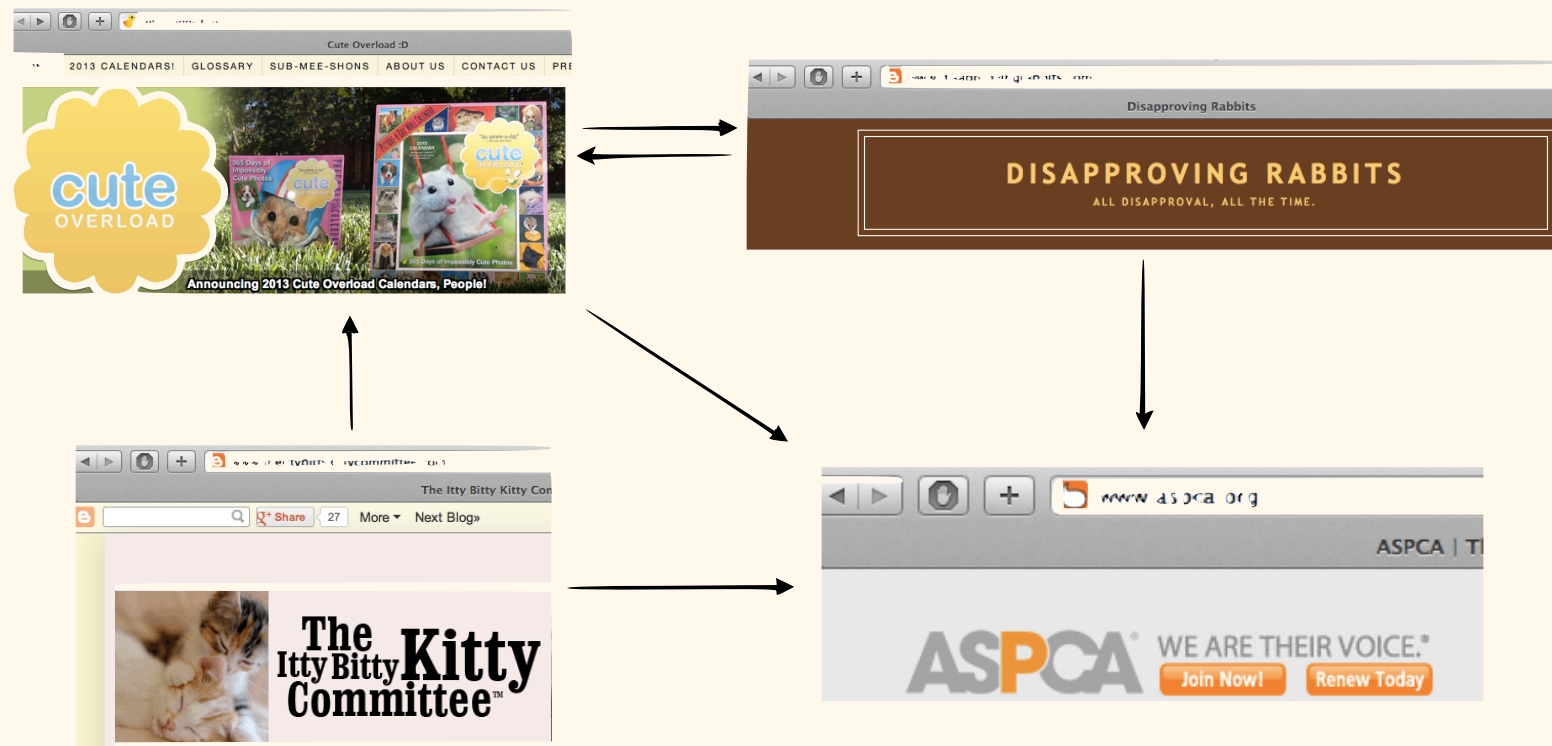


Connections between political blogs

Polarization of the network [Adamic-Glance, 2005]

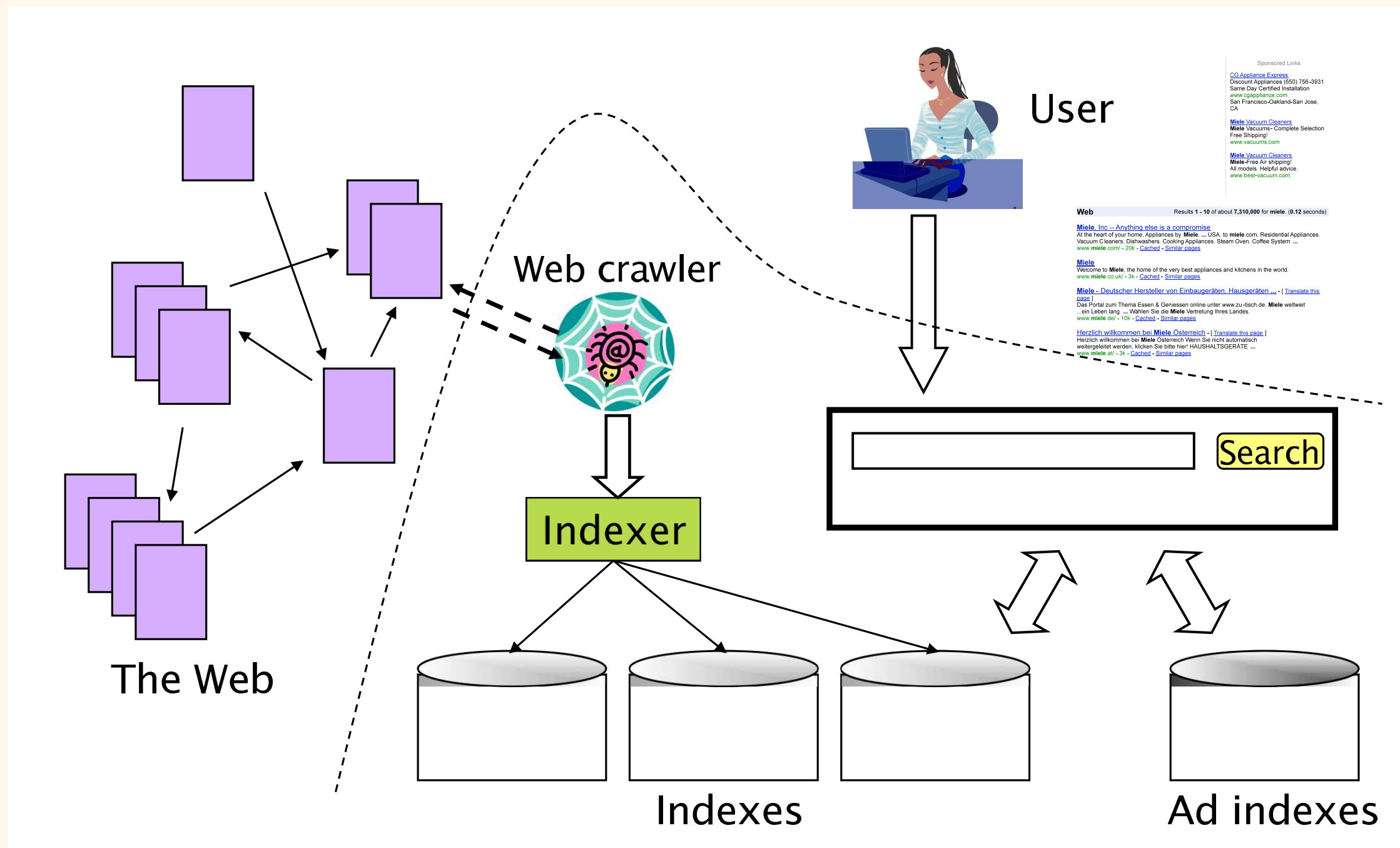
J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmnds.org>

The web can be thought of as a directed graph...



We can use link data in two main ways:

1. Inferring *authority* or *trustworthiness* of a page;
2. Identifying new pages to be added to our index.



► **Figure 19.7** The various components of a web search engine.

“Crawling” is a basic and important piece of any search engine:

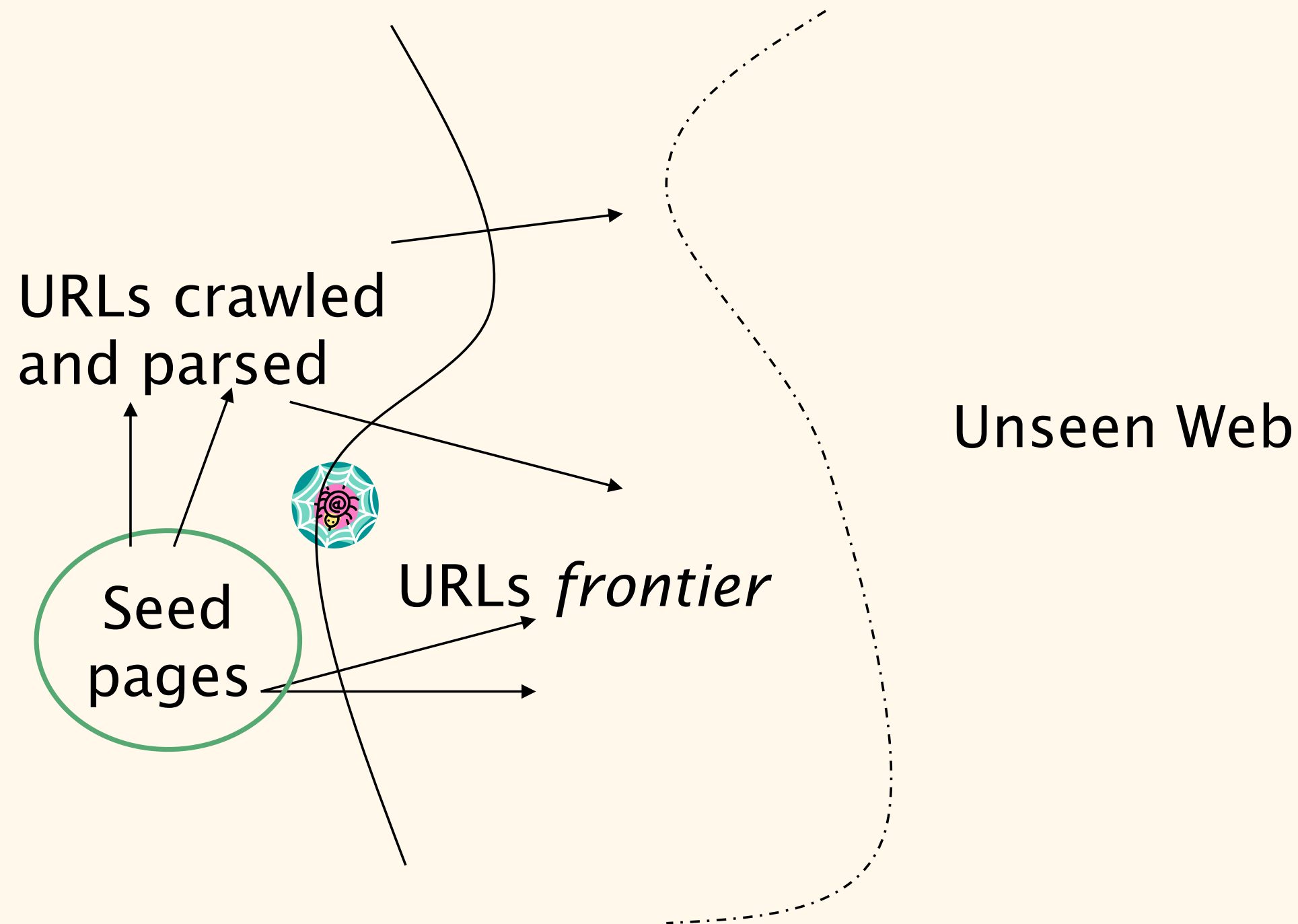
Starting with a set of “seed” URLs...

Fetch & parse, extract URLs...

Add to queue...

Fetch & parse each URL on queue, etc.

“Crawling” is a basic and important piece of any search engine:



This may sound easy, but it actually is not:

Crawling requires parallelism, with all the joy that brings...

Malicious pages seek to confound crawlers:

- Spam, Spider traps (both static & dynamic)

- Latency & Bandwidth vary widely across sites

- Sites have different crawling policies (depth, frequency, etc.)

- Duplicate pages, site mirrors, “dynamic” UIs, etc.

- Politeness

There are two primary things that crawlers *must* do:

Be Robust:

Be immune to spider traps; gracefully handle bandwidth, timeout, and web server issues

Be Polite:

Respect implicit and explicit conventions

There are two primary things that crawlers *must* do:

Be Robust:

Be immune to spider traps; gracefully handle bandwidth, timeout, and web server issues

Be Polite:

Respect implicit and explicit conventions

Implicit politeness: avoid hitting a site too often, etc.

Explicit politeness: honor specifications from webmasters about what portions of the site may be crawled (`robots.txt`), etc.

Robots.txt:

A way to “control” what parts of your site get crawled (and by whom).

Dating from 1994, the protocol has remained (relatively) static.

Well-behaved crawlers look for a text file named “robots.txt” at the root of your website, and follow its directives:

```
User-Agent: *
```

```
Allow: /about_us
```

```
Disallow: /docs/private_files/
```

Robots.txt:

Extensions: directives for crawl frequency & rate, etc.

Sitemap, Host (to specify canonical mirrors, etc.)

Caution: a robots.txt file is *not* a security mechanism!

```
User-agent: *  
Disallow: /_archive/  
Disallow: /_resources/  
Disallow: /academic/som/  
Disallow: /bigbrain_courses_staging/  
...  
Disallow: /itgdba/
```

Quite the opposite, in fact...

Robots.txt:

“During the reconnaissance stage of a web application testing, the tester (or attacker) usually uses a list of known subdirectories to brute force the server and find hidden resources.”

<http://thiébaud.fr/robots.txt.html>

User-agent: *

Disallow: /admin/

Disallow: /stats/

Disallow: /internaljobs/

Disallow: /internaljobsbyorganization/

Disallow: /internaljobsearch/

<http://www.behindthefirewalls.com/2013/07/using-robotstxt-to-locate-your-targets.html>

There are two primary things that crawlers *must* do:

Be Robust:

Be immune to spider traps; gracefully handle bandwidth, timeout, and web server issues

Be Polite:

Respect implicit and explicit conventions

Implicit politeness: avoid hitting a site too often, etc.

Explicit politeness: honor specifications from webmasters about what portions of the site may be crawled (`robots.txt`), etc.

Beyond those requirements, there are several “*shoulds*”:

Be capable of distributed operation;

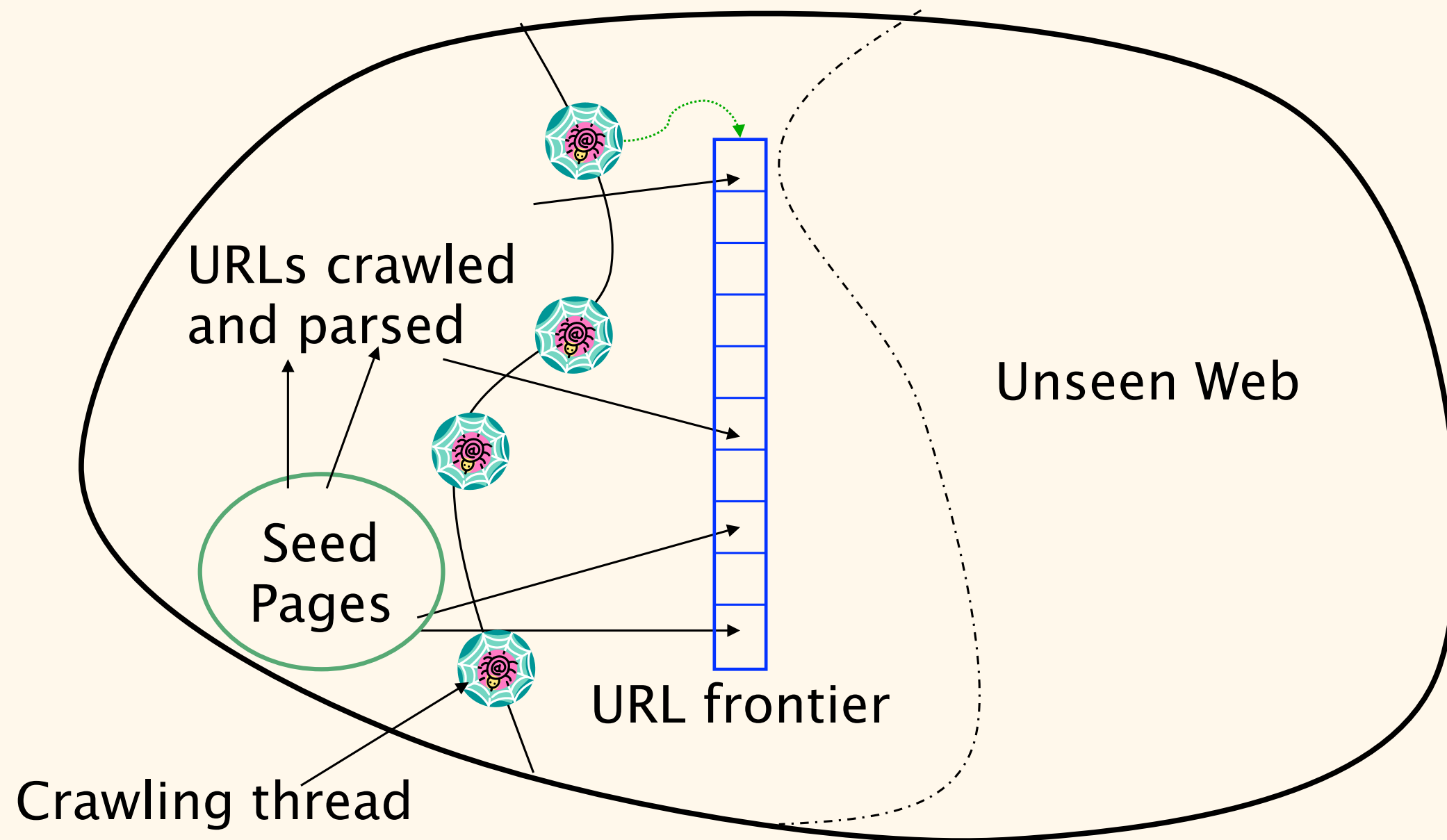
Be scalable (increase crawl rate by adding machines)

Be as efficient as possible (both in terms of CPU and network)

Be clever: fetch “higher quality” pages first, etc.

Operate continuously

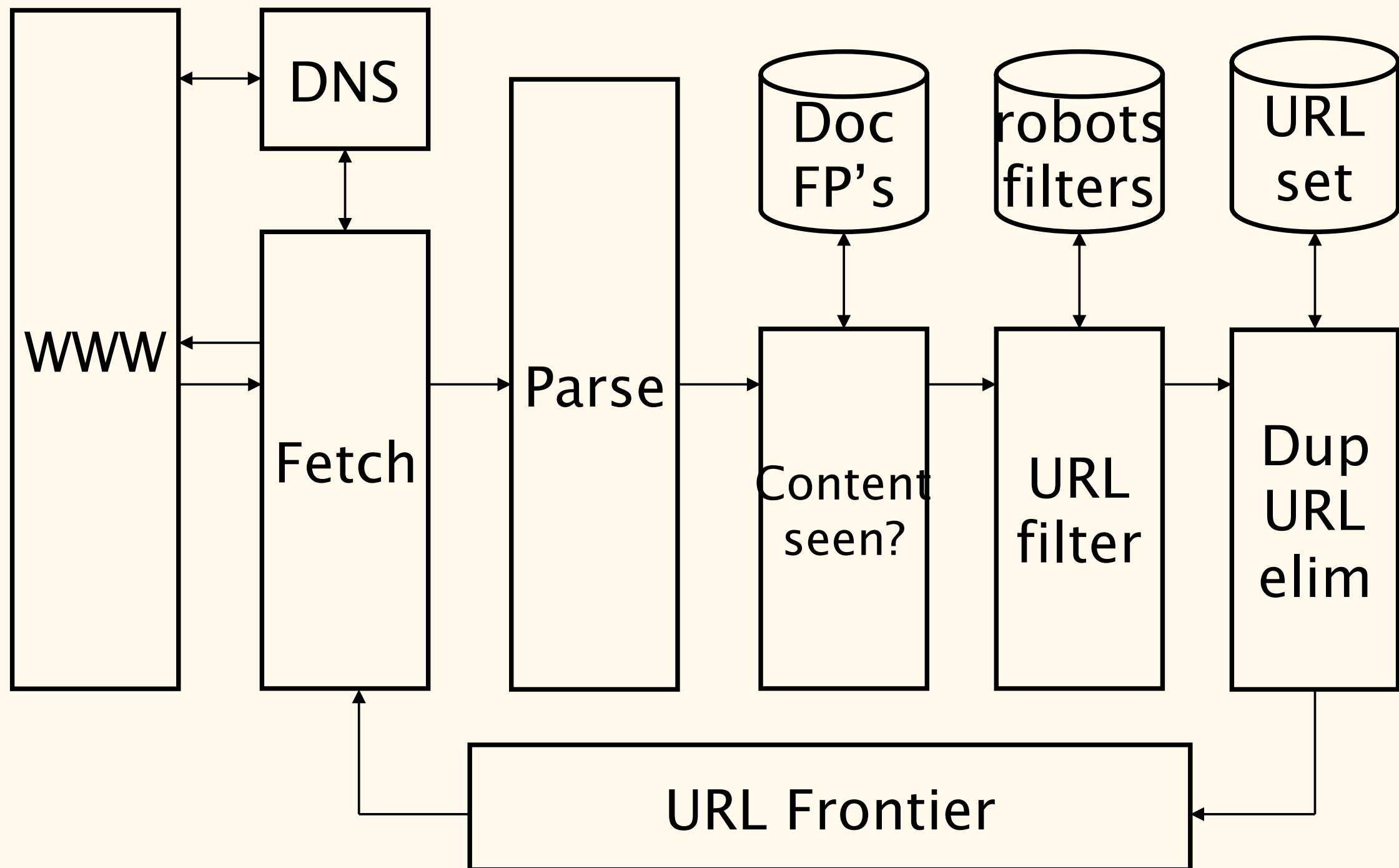
Be easily extensible to support new protocols, document types, etc.



Which one?

- Pick a URL from the frontier
- Fetch the document at the URL
- Parse the URL
 - Extract links from it to other docs (URLs)
- Check if URL has content already seen
 - If not, add to indexes
- For each extracted URL
 - Ensure it passes certain URL filter tests
 - Check if it is already in the frontier (duplicate URL elimination)

robots.txt filters, etc.



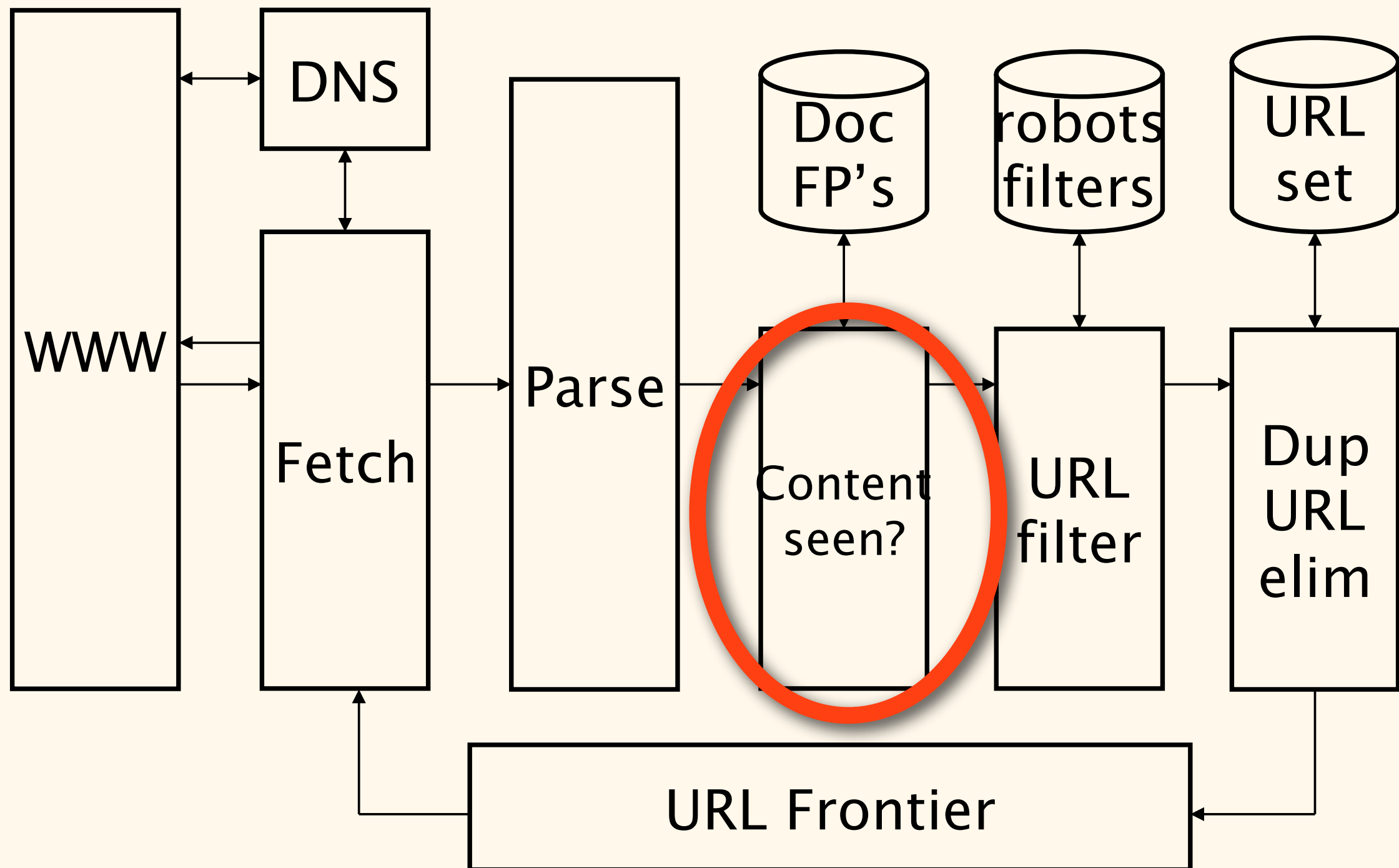
URL frontier: two main considerations

- Politeness: do not hit a web server too frequently
- Freshness: crawl some pages more often than others
 - E.g., pages (such as News sites) whose content changes often

These goals may conflict with each other.

(E.g., simple priority queue fails – many links out of a page go to its own site, creating a burst of accesses to that site.)

There are many solutions: see the book!



Duplicate documents

- The web is full of duplicated content
- Strict duplicate detection = exact match
 - Not as common
- But many, many cases of near duplicates
 - E.g., Last modified date the only difference between two copies of a page

Duplicate/Near-Duplicate Detection

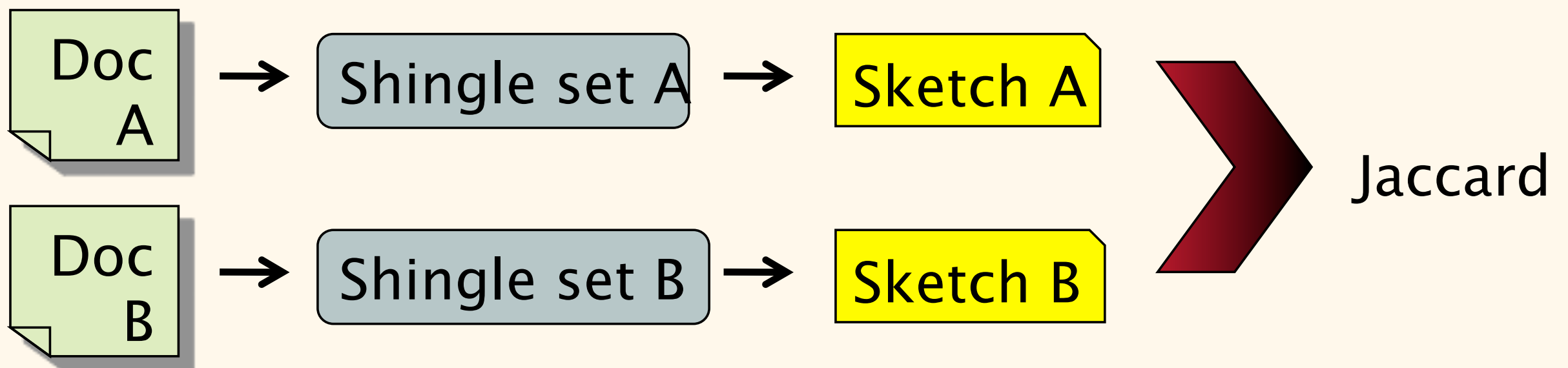
- *Duplication*: Exact match can be detected with fingerprints
- *Near-Duplication*: Approximate match
 - Overview
 - Compute syntactic similarity with an edit-distance measure
 - Use similarity threshold to detect near-duplicates
 - E.g., Similarity > 80% => Documents are “near duplicates”
 - Not transitive though sometimes used transitively

Computing Similarity

- Features:
 - Segments of a document (natural or artificial breakpoints)
 - Shingles (Word N-Grams)
 - ***a rose is a rose is a rose*** → 4-grams are
 - a_rose_is_a
 - rose_is_a_rose
 - is_a_rose_is
 - a_rose_is_a
- Similarity Measure between two docs (= sets of shingles)
 - Jaccard coefficient: (Size_of_Intersection / Size_of_Union)

Shingles + Set Intersection

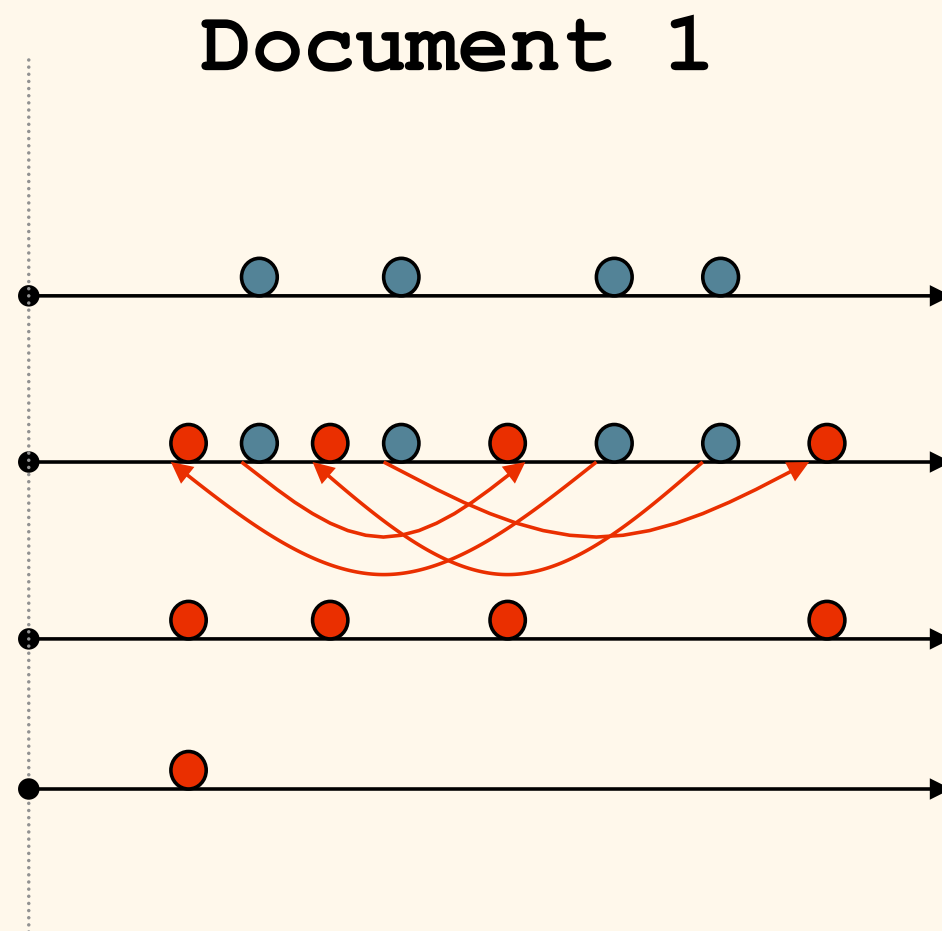
- Computing exact set intersection of shingles between all pairs of documents is expensive
- Approximate using a cleverly chosen subset of shingles from each (a *sketch*)
- Estimate ($\text{size_of_intersection} / \text{size_of_union}$) based on a short sketch



Sketch of a document

- Create a “sketch vector” (of size ~ 200) for each document
 - Documents that share $\geq t$ (say 80%) corresponding vector elements are deemed **near duplicates**
 - For doc D , $\text{sketch}_D[i]$ is as follows:
 - Let f map all shingles in the universe to $1..2^m$ (e.g., f = fingerprinting)
 - Let π_i be a *random permutation* on $1..2^m$
 - Pick $\text{MIN } \{\pi_i(f(s))\}$ over all shingles s in D

Computing Sketch[i] for Doc1

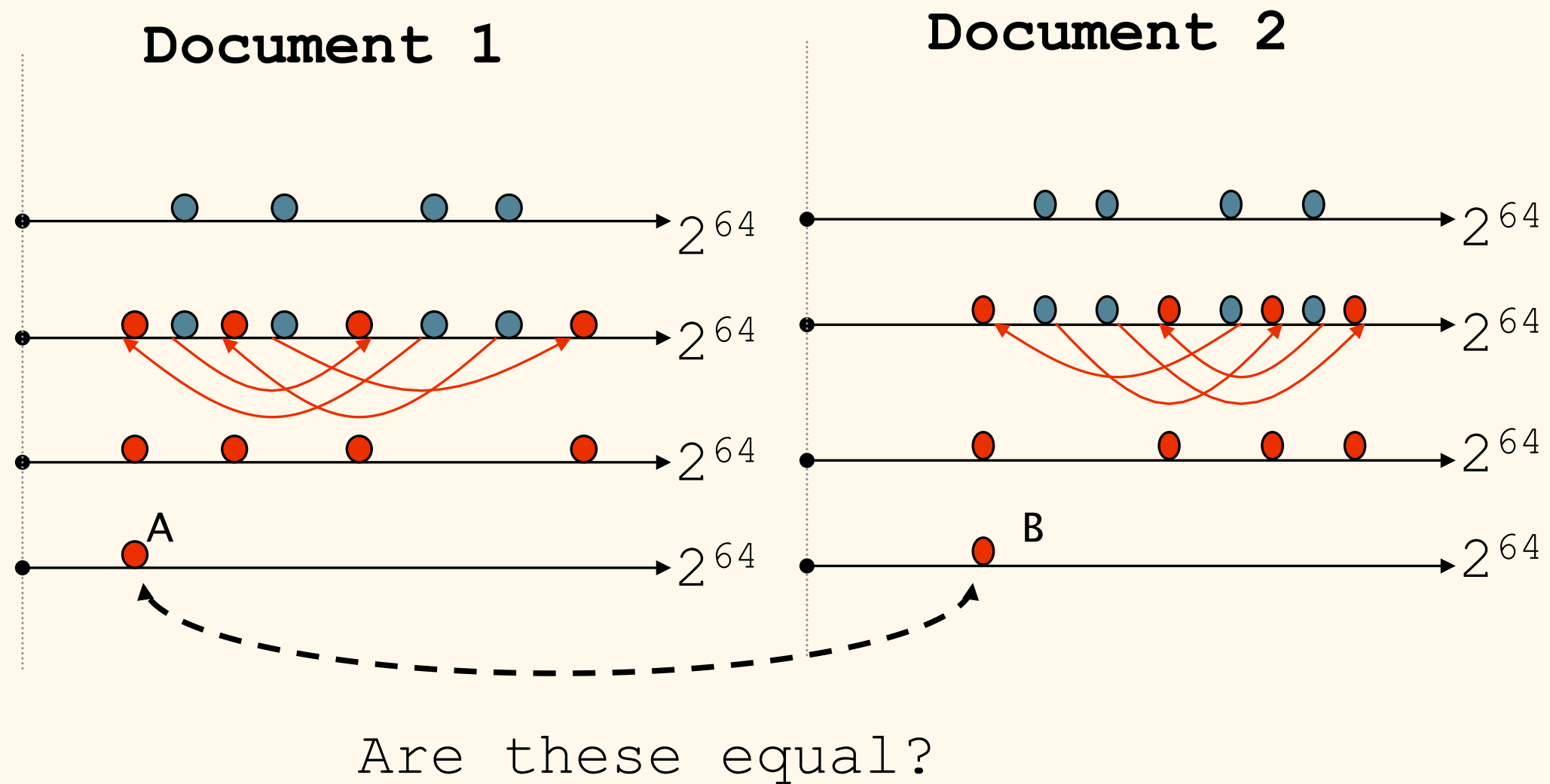


Start with 64-bit $f(\text{shingles})$

Permute on the number line
with π_i

Pick the min value

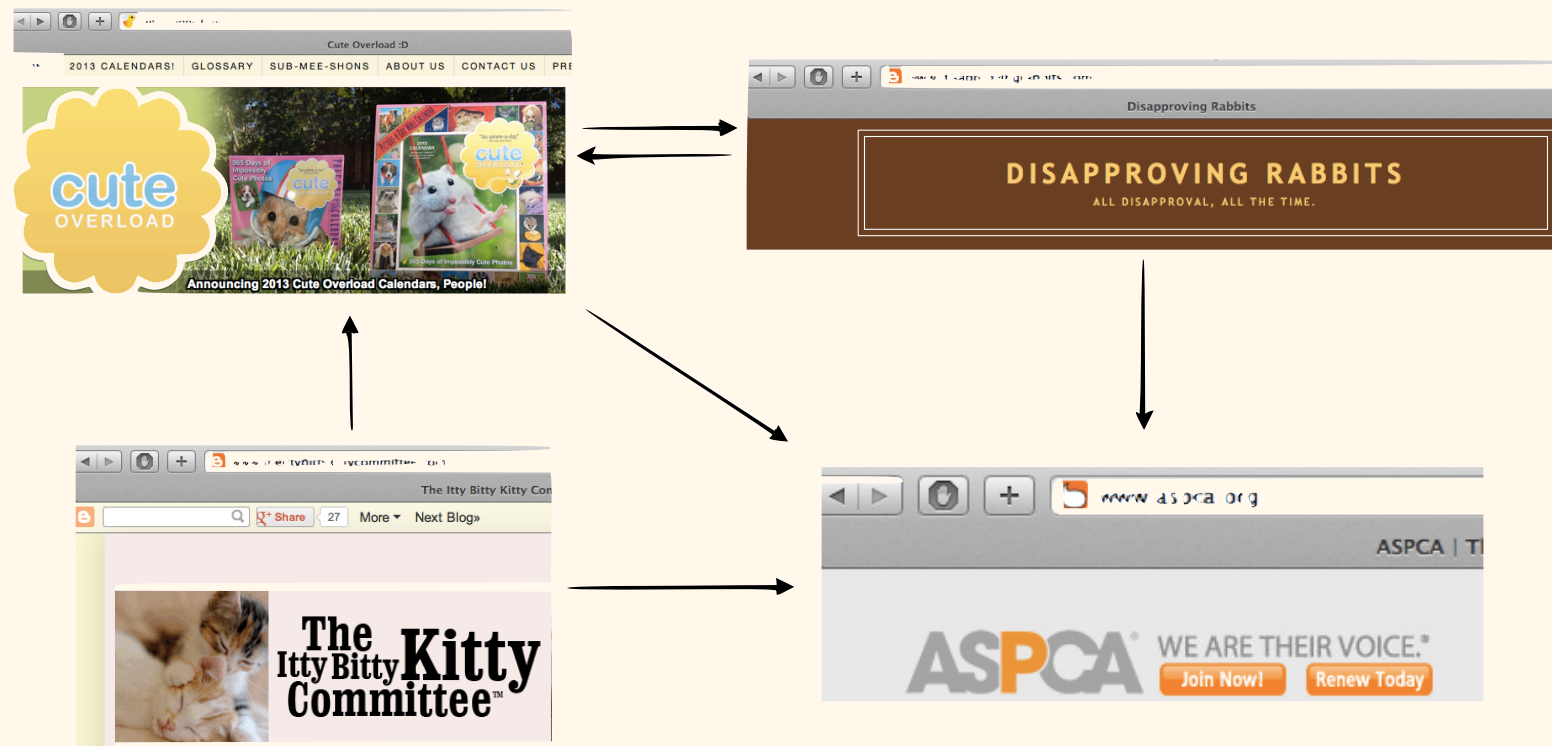
Test if $\text{Doc1.Sketch}[i] = \text{Doc2.Sketch}[i]$



Test for **200** random permutations: $\pi_1, \pi_2, \dots, \pi_{200}$

- Shingling is a *randomized algorithm*
 - Our analysis did not presume any probability model on the inputs
 - It will give us the right (wrong) answer with some probability on *any input*
- We've described how to detect near duplication in a pair of documents
- In “real life” we'll have to concurrently look at many pairs
 - See text book for details

The web can be thought of as a directed graph...



We can use link data in two main ways:

1. Inferring *authority* or *trustworthiness* of a page;
2. Identifying new pages to be added to our index.

Link analysis piece...

Challenges to all of the above: spam

Since good search results mean more traffic, more customers, etc...

... people are strongly motivated to try and “game” search engine’s indexing processes.

“As the popularity of the Web has increased, the efforts to exploit the Web for commercial, social, or political advantage have grown, making it harder for search engines to discriminate between truthful signals of content quality and deceptive attempts to game search engines’ rankings.”

Early attempts at search engine spamming involved “keyword stuffing”:

Including (many!) additional instances of key words or phrases on the page...

... but doing it in a way that was invisible to human users (hiding in a metadata field, causing them to blend in with the background, etc.).

Modern spammers have to be more clever:

Cloaking:

Returning one page for humans, and another for index crawlers

Doorway pages:

Show a page that looks legit to a web crawler, but have all the outgoing links point to commercial pages

Scraping:

Steal high-quality content from another page, and then link to commercial pages (sending “quality points” from your scraped page to the crappy spam pages)

Link buying/exchange:

What it sounds like

Comment spam:

Putting your spam in higher-ranked websites' comment sections

Adversarial Web Search

By Carlos Castillo and Brian D. Davison

Contents

1	Introduction	379
1.1	Search Engine Spam	380
1.2	Activists, Marketers, Optimizers, and Spammers	381
1.3	The Battleground for Search Engine Rankings	383
1.4	Previous Surveys and Taxonomies	384
1.5	This Survey	385
2	Overview of Search Engine Spam Detection	387
2.1	Editorial Assessment of Spam	387
2.2	Feature Extraction	390
2.3	Learning Schemes	394
2.4	Evaluation	397
2.5	Conclusions	400
3	Dealing with Content Spam and Plagiarized Content	401
3.1	Background	402
3.2	Types of Content Spamming	405
3.3	Content Spam Detection Methods	405
3.4	Malicious Mirroring and Near-Duplicates	408
3.5	Cloaking and Redirection	409
3.6	E-mail Spam Detection	413
3.7	Conclusions	413

Next up: search UI/UX.