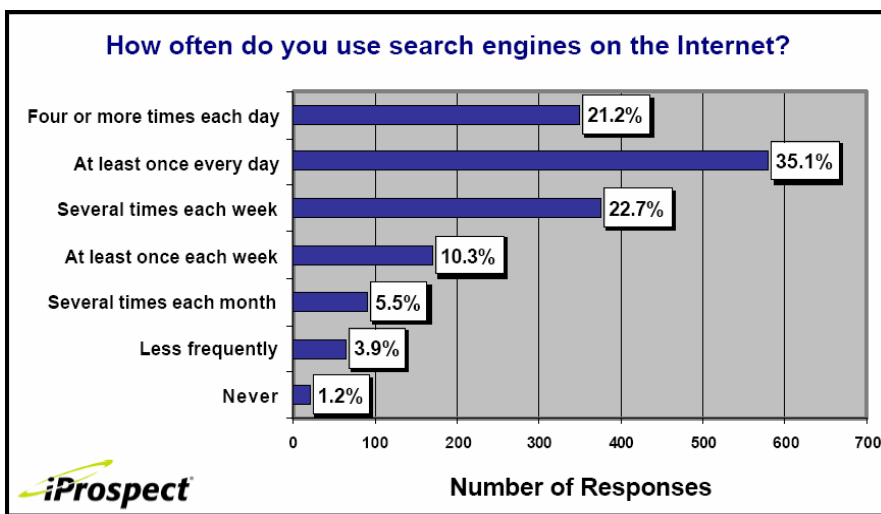**PURDUE UNIVERSITY** | Department of Computer Science

# CS47300: Web Information Search and Management

*Web Search*
Prof. Chris Clifton
14 September 2020
*Some slides courtesy*
*Manning, Raghavan, and Schütze*

Indiana Center for Database Systems ™

---

**PURDUE UNIVERSITY**
Department of Computer Science

## Usage of Web Search

**How often do you use search engines on the Internet?**

| Category | Percentage |
| --- | --- |
| Four or more times each day | 21.2% |
| At least once every day | 35.1% |
| Several times each week | 22.7% |
| At least once each week | 10.3% |
| Several times each month | 5.5% |
| Less frequently | 3.9% |
| Never | 1.2% |

**Number of Responses**

iProspect

## Without search engines the web wouldn't scale

- No incentive in creating content unless it can be easily found – other finding methods haven't kept pace (taxonomies, bookmarks, etc)
- The web is both a technology artifact and a social environment
  - "The Web has become the "new normal" in the American way of life; those who don't go online constitute an ever-shrinking minority." – [Pew Foundation report, January 2005]
- Search engines make aggregation of interest possible:
  - Create incentives for very specialized niche players
    - Economical – specialized stores, providers, etc
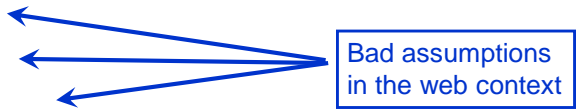    - Social – narrow interests, specialized communities, etc

## Without search engines the web wouldn't scale

- The acceptance of search interaction makes "unlimited selection" stores possible:
  - Amazon, Netflix, etc
- Search has been the best mechanism for advertising on the web, a $15+ B industry.
  - Growing very fast but entire US advertising industry $250B – huge room to grow
  - Sponsored search marketing is about $10B
  - *2020: Statista estimates search ad revenue $159B, 200B by 2024*
  - *2019: Alphabet alone had advertising revenue $142B*
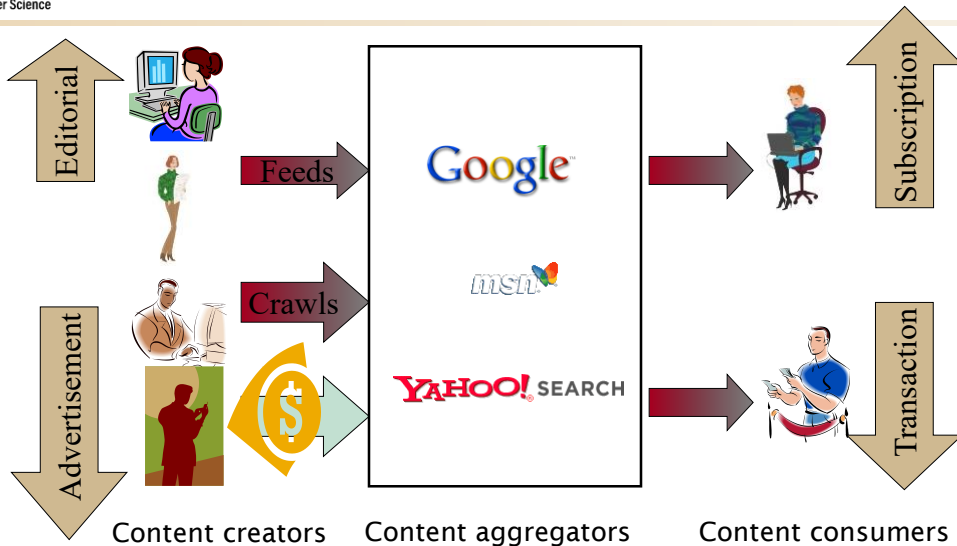
# Classic IR

## Relevance

– For each query Q and stored document D in a given corpus assume there exists relevance Score(Q, D)

   • Score is average over users U and contexts C

– Optimize Score(Q, D) as opposed to Score(Q, D, U, C)

– That is, usually:

   • Context ignored
   • Individuals ignored
   • Corpus predetermined

Bad assumptions in the web context

---

# The coarse-level dynamics

Editorial

Feeds

Crawls

Google

msn

YAHOO! SEARCH

Subscription

Transaction

Advertisement

Content creators  Content aggregators  Content consumers

## Brief (non-technical) history

Early keyword-based engines
- Altavista, Excite, Infoseek, Inktomi, ca. 1995-1997

Paid placement ranking: Goto.com (morphed into Overture.com → Yahoo!)
- Your search ranking depended on how much you paid
- Auction for keywords: ***casino*** was expensive!

## Brief (non-technical) history

1998+: Link-based ranking pioneered by Google
- Blew away all early engines:  Great user experience in search of a business model
- Meanwhile Goto/Overture's annual revenues were nearing $1 billion

Result: Google added paid-placement "ads" to the side, independent of search results
- Yahoo follows suit, acquiring Overture (for paid placement) and Inktomi (for search)

Ads vs. search results

**Department of Computer Science**

Google has maintained that ads (based on vendors bidding for keywords) do not affect vendors' rankings in search results

Search = *miele*

---

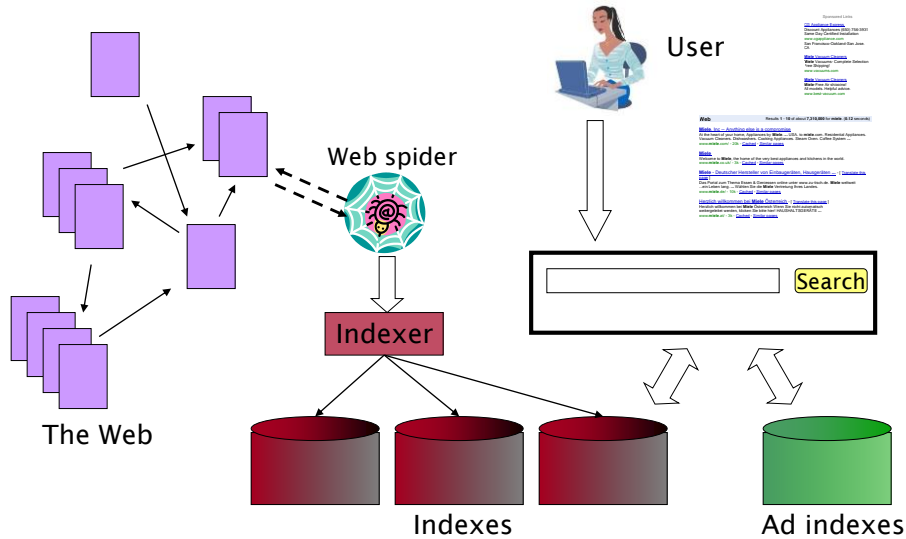Ads vs. search results

**Department of Computer Science**

Other vendors (Yahoo, MSN) have made similar statements from time to time
– Any of them can change anytime

We will focus primarily on search results independent of paid placement ads
– Although the latter is a fascinating technical subject in itself

**Web search basics**



**User Needs**

Need [Brod02, RL04]
- **<u>Informational</u>** – want to learn about something (~40% / 65%)  `P53 Cancer`
- **<u>Navigational</u>** – want to go to that page (~25% / 15%)  `United Airlines`
- **<u>Transactional</u>** – want to do something (web-mediated) (~35% / 20%)
  - Access a service  `Seattle weather`
  - Downloads  `Mars surface images`
  - Shop  `Canon S410`
- Gray areas
  - Find a good hub  `Car rental Brasil`
  - Exploratory search "see what's there"

## Web search users

- Make ill defined queries
  - Short
    - AV 2001: 2.54 terms avg, 80% < 3 words)
    - AV 1998: 2.35 terms avg, 88% < 3 words [Silv98]
  - Imprecise terms
  - Sub-optimal syntax (most queries without operator)
  - Low effort
- Wide variance in
  - Needs
  - Expectations
  - Knowledge
  - Bandwidth

- Specific behavior
  - 85% look over one result screen only
  - 78% of queries are not modified (one query/session)
  - Follow links – "the scent of information" ...
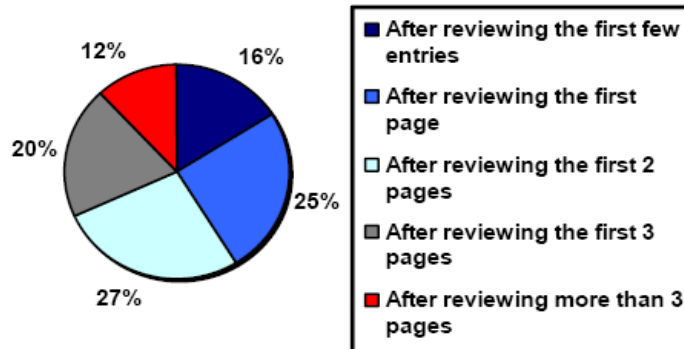
---

## Query Distribution

Legend:
- cancer
- breast cancer
- skin cancer
- prostate cancer
- lung cancer
- colon cancer
- leukemia
- lymphoma

Power law: few popular (typically broad) queries, many rare (typically more specific) queries

**How far do people look for results?**

"When you perform a search on a search engine and don't find what you are looking for, at what point do you typically either revise your search, or move on to another search engine? (Select one)"
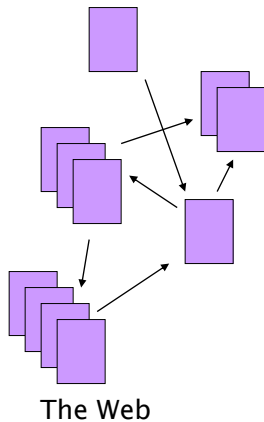
- After reviewing the first few entries — 16%
- After reviewing the first page — 25%
- After reviewing the first 2 pages — 27%
- After reviewing the first 3 pages — 20%
- After reviewing more than 3 pages — 12%

(Source: iprospect.com WhitePaper_2006_SearchEngineUserBehavior.pdf)

---

# Users' empirical evaluation of results

- Quality of pages varies widely
  - Relevance is not enough
  - Other desirable qualities (non IR!!)
    - Content: Trustworthy, new info, non-duplicates, well maintained,
    - Web readability: display correctly & fast
    - No annoyances: pop-ups, etc
- Precision vs. recall
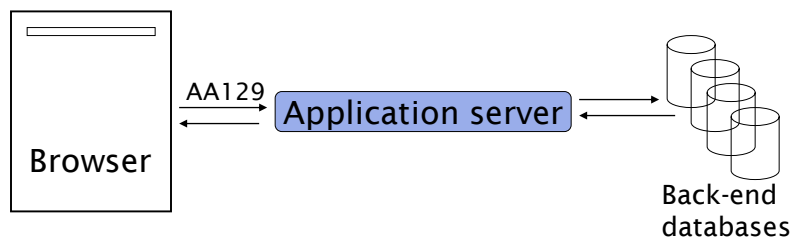  - On the web, recall seldom matters

# The Web corpus

- No design/co-ordination
- Distributed content creation, linking, democratization of publishing
- Content includes truth, lies, obsolete information, contradictions …
- Unstructured (text, html, …), semi-structured (XML, annotated photos), structured (Databases)…
- Scale much larger than previous text corpora … but corporate records are catching up.
- Content can be dynamically generated

The Web

---

# The Web: Dynamic content

- A page without a static html version
  - E.g., current status of flight AA129
  - Current availability of rooms at a hotel
- Usually, assembled at the time of a request from a browser
  - Typically, URL has a '?' character in it

Browser

AA129

Application server

Back-end databases

# Dynamic content

- Most dynamic content is ignored by web spiders
  - Many reasons including malicious spider traps
- Some dynamic content (news stories from subscriptions) are sometimes delivered as dynamic content
  - Application-specific spidering
- Spiders commonly view web pages just as Lynx (a text browser) would
- Note: even "static" pages are typically assembled on the fly (e.g., headers are common)

# Other characteristics of the Web

- Significant duplication
  - Syntactic – 30%-40% (near) duplicates [Brod97, Shiv99b, etc.]
  - Semantic – ???
- High linkage
  - More than 8 links/page in the average
- Complex graph topology
  - Not a small world; bow-tie structure [Brod00]
- Spam
  - Billions of pages