# Ranking of ads

- Goto model:
  - Rank according to how much advertiser pays
- Current model:
  - Balance auction price and relevance
  - Irrelevant ads (few click-throughs)
    - Decrease opportunities for relevant ads
    - Harm the user experience
  - Idea: Well-targeted advertising is good for everyone

# Paying for advertisements

- CPM
  - "Cost Per Mil"
  - Pay for 1000 eyeballs
  - Important for branding campaigns
- CPC
  - "Cost per Click"
  - Pay for clicking on ads
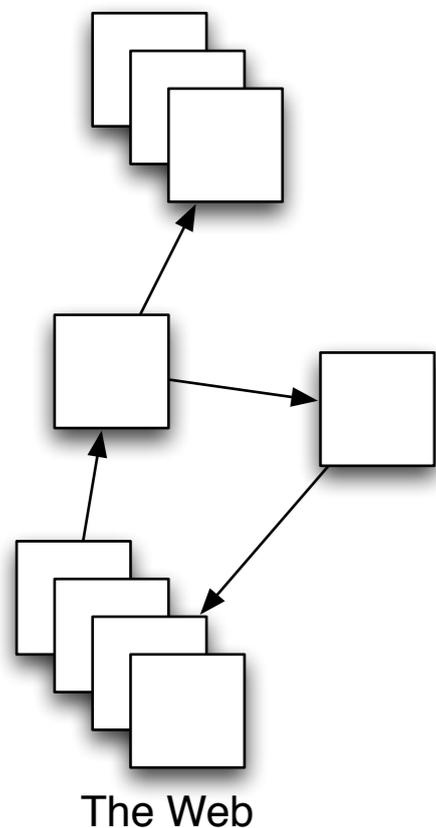  - Important for sales campaigns

# Overview

- Introduction
- Classic Information Retrieval
- Web IR
- Sponsored Search
- Web Search Basics
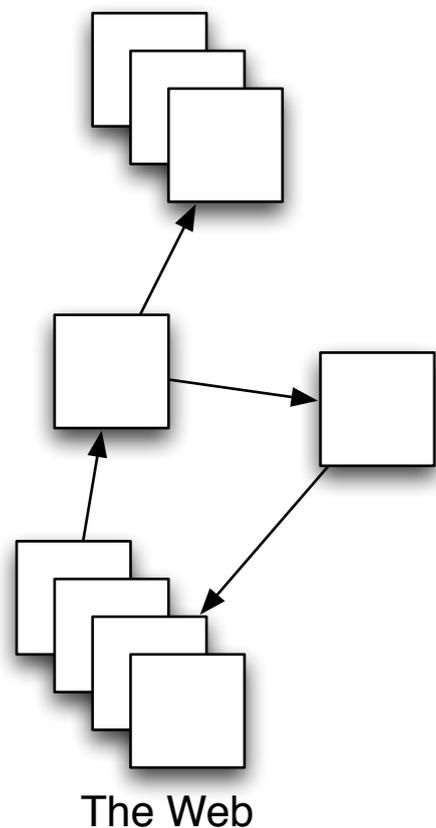  - Size of the Web
- Web Users
- Spam

# The Web Corpus

- No design/coordination
- Distributed content creation, linking
- "Democatization of publishing"
- Content includes truth, lies, contradictions, etc.
- Unstructured Data (text, html)
- Semi-Structured (XML, annotated photos)
- Structured (Databases)
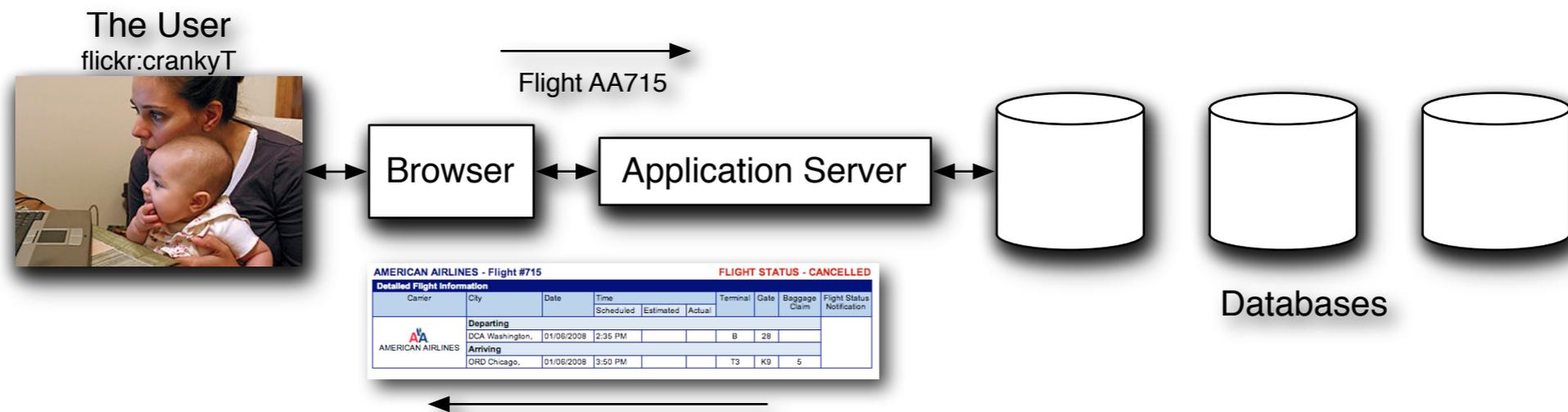- Scale is much larger than previous text copora

The Web

# The Web Corpus

- Growth - slowing from "doubling every few months", but still expanding

The Web

# Dynamic Content

- Content can by dynamically generated

  - There is no static html version

    - Flight status information, evite responses

  - Assembled on request ("?" in URL is a clue)

The User
flickr:crankyT

Flight AA715

Browser ⟷ Application Server ⟷ Databases

AMERICAN AIRLINES - Flight #715                    FLIGHT STATUS - CANCELLED

Detailed Flight Information

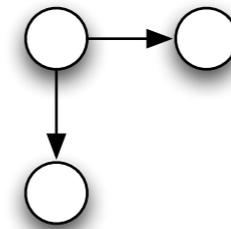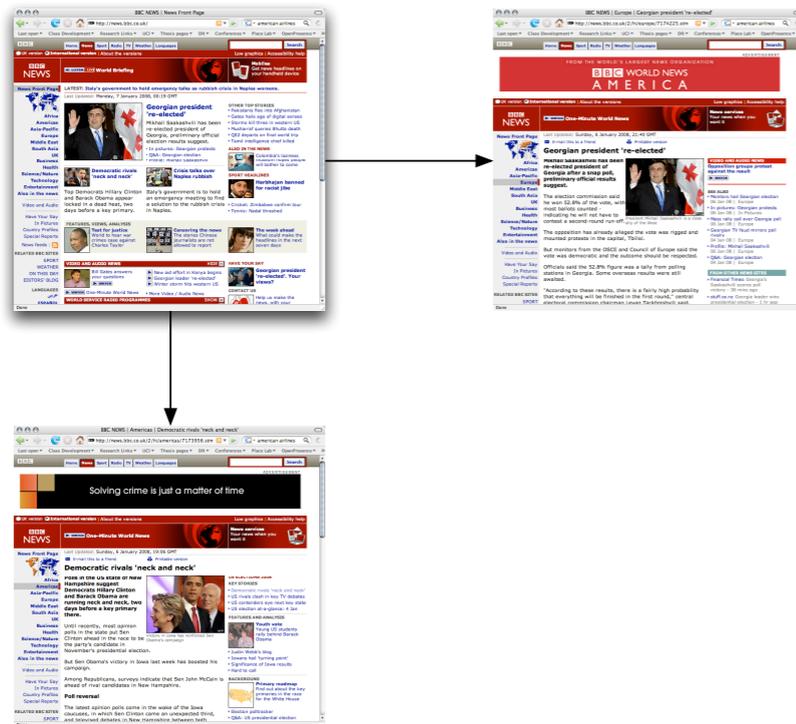| Carrier | City | Date | Time | | | Terminal | Gate | Baggage Claim | Flight Status Notification |
|---|---|---|---|---|---|---|---|---|---|
| | | | Scheduled | Estimated | Actual | | | | |
| | **Departing** | | | | | | | | |
| AMERICAN AIRLINES | DCA Washington, | 01/06/2008 | 2:35 PM | | | | B | 28 | |
| | **Arriving** | | | | | | | | |
| | ORD Chicago, | 01/06/2008 | 3:50 PM | | | T3 | K9 | 5 | |

# Dynamic Content

- Most (truly) dynamic content is ignored by web spiders
    - Too much to index
    - Static information is more important for search
    - Spider Traps look dynamic
- Actually a lot of "static" content is assembled on the fly also
    - ASP, PHP, JSP, ads, etc....

# The Web as a graph

- Web pages are nodes

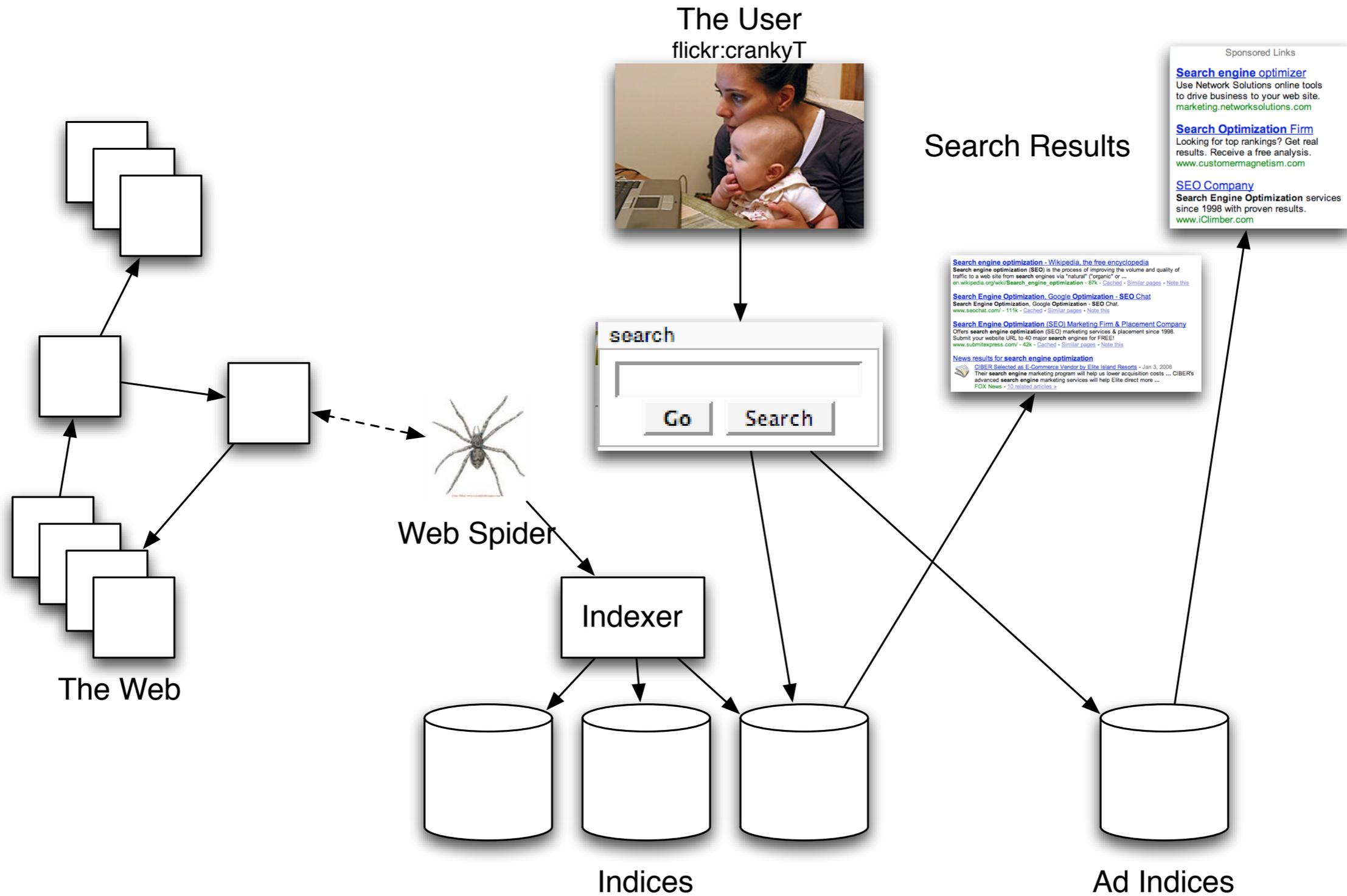- Hyperlinks are directed edges

# Characteristics of the web

- Significant Duplication
    - 30%-40% is some studies [Brod97, Shiv99]
    - www.copyscape.com
- High linkage
    - more than 8 links per page on average
- Spam
    - Billions of pages of it.

# Web Search Basics

The User
flickr:crankyT

Search Results

Sponsored Links

**Search engine** optimizer
Use Network Solutions online tools to drive business to your web site.
marketing.networksolutions.com

**Search Optimization** Firm
Looking for top rankings? Get real results. Receive a free analysis.
www.customermagnetism.com

SEO Company
**Search Engine Optimization** services since 1998 with proven results.
www.iClimber.com

The Web

Web Spider

search

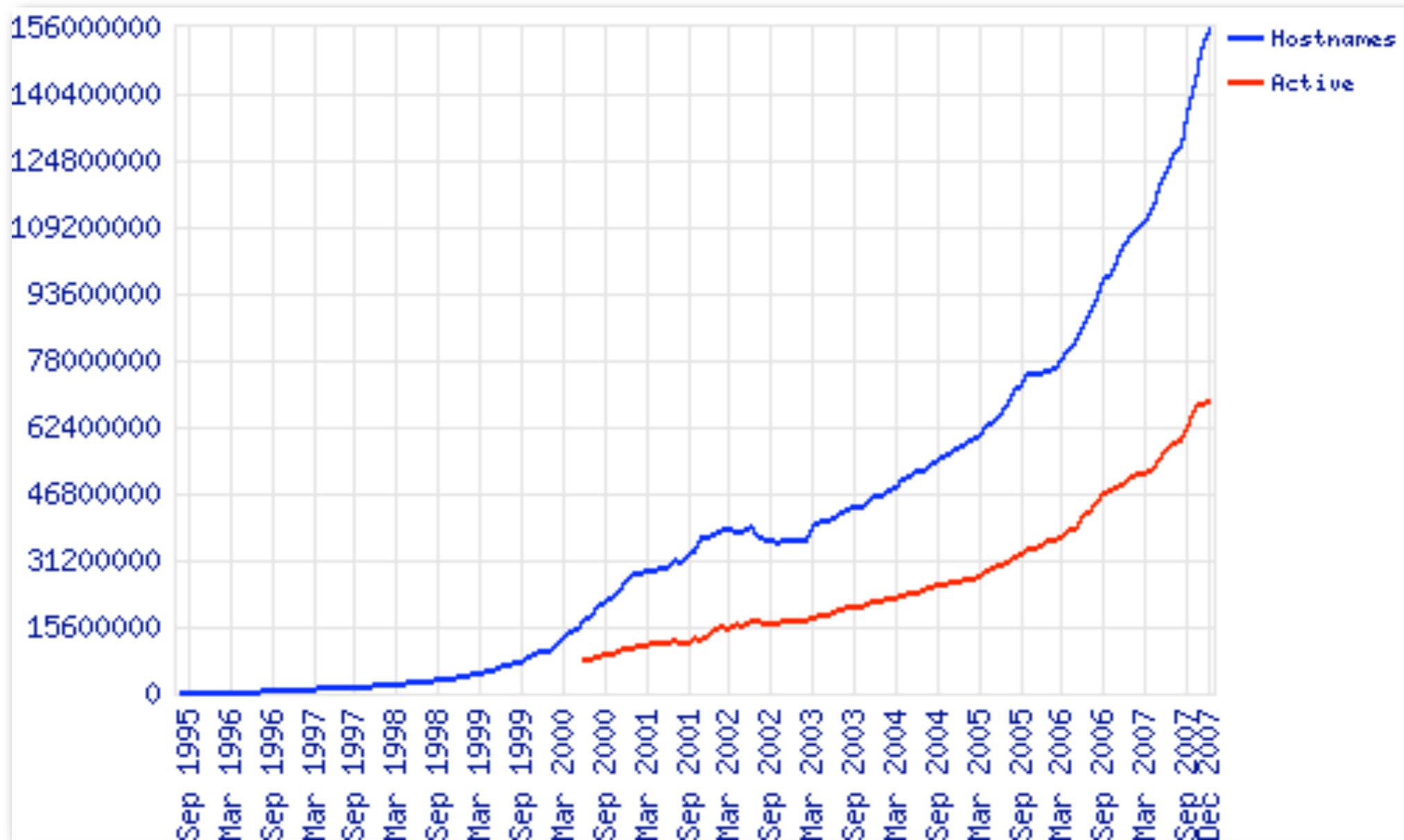Go     Search

Indexer

Indices

Ad Indices

# How big is the web?

- What is measured?

  - Number of hosts

  - Number of "static" html pages

- Number of hosts - netcraft survey

  - http://news.netcraft.com/archives/web_server_survey.html

  - Monthly report on hosts and servers

- Number of pages

  - Lots of estimates which warrant further discussion

# How big is the web?

- Netcraft Web Server Survey

# Rate of change

- [Cho00] 720k pages from 270 popular sites sample daily for 5 months in 1999
  - 40% changed weekly, 23% daily
- [Fett02] Massive study: 151M pages checked over a few months
  - Significant changes 7% weekly
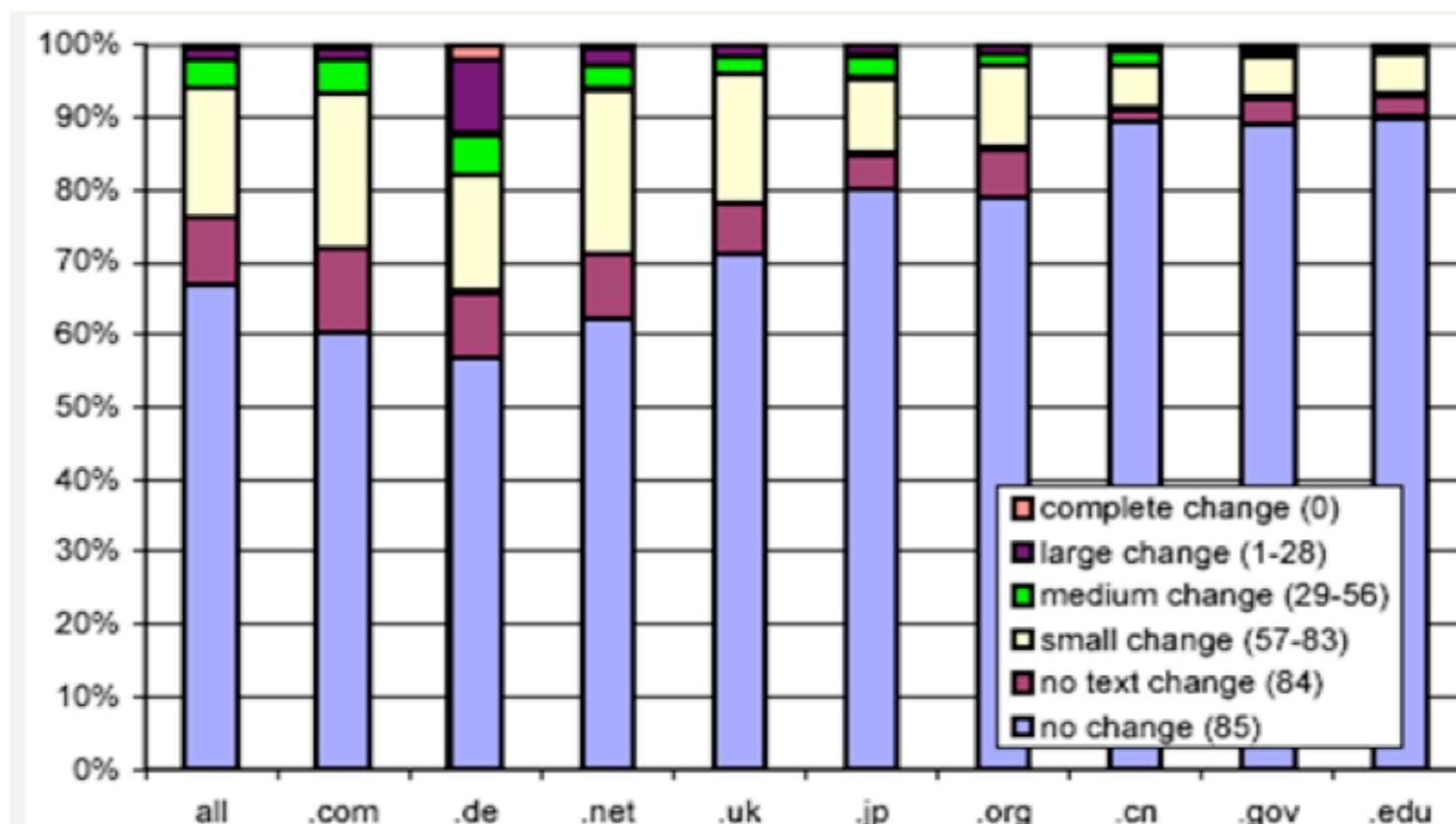  - Any change 25% weekly

# Rate of change

- [Ntul04] 154 large sites recrawled from scratch weekly
    - 8% had new pages ever week
    - 8% die
    - 5% new content
    - 25% new links per week

# Rate of change

- Fetterly et al. study in 2002

  - 150 million pages over 11 weekly crawls

  - Bucketed into 85 groups according to amount of change

# Web Evolution

- The nature of the web is change

- Not much work on studying web evolution

  - Exception is Fetterly et. al, 2003

- Some effort has been made to extrapolate from small samples using fractal models [Dill et. al. 2001]

# Overview

- Introduction
- Classic Information Retrieval
- Web IR
- Sponsored Search
- Web Search Basics
  - Size of the Web
- Web Users
- Spam

# User Search Needs in Brod02/RL04

- Informational
  - Want to learn about something (~40%/65%)
- Navigational
  - Want to go to that page (~25%/15%)
- Transactional
  - Want to do something (~35%/20%)
    - Access a service, download, shop
- Others?
  - Exploration, social, etc...

# Web Users

- Make ill defined queries
  - Short
    - Average in 2001: 2.54 terms (80% < 3 words)
    - Average in 1998: 2.35 terms (88% < 3 words) [Silv98]
  - Imprecise terms
  - Suboptimal syntax (no operators)
  - Low effort (spelling mistakes)

# Web Users

- Wide Variance in
  - Needs
  - Expectations
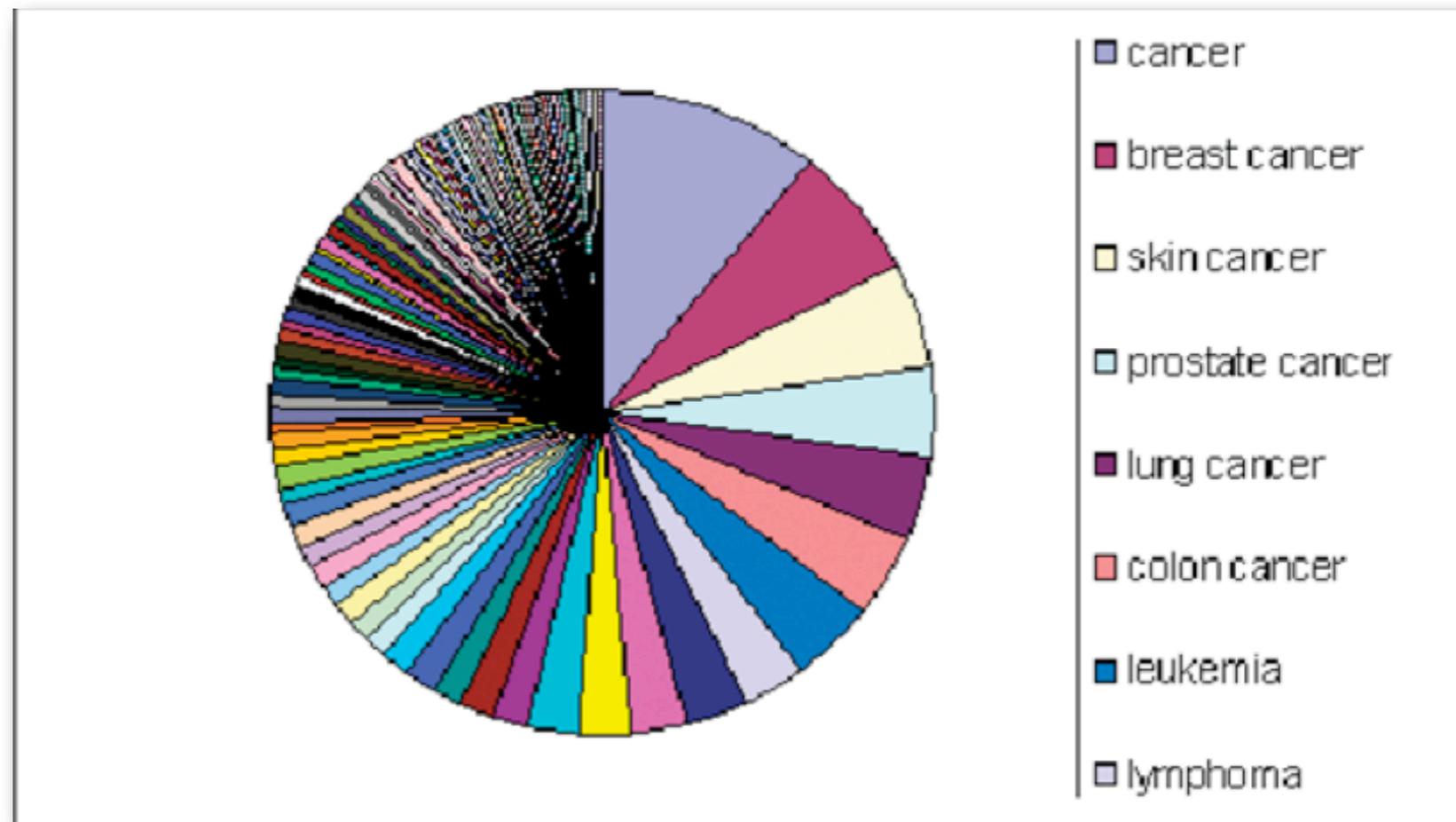  - Knowledge
  - Bandwidth

# Web Users

- Behavior
    - 85% look over one result screen only
    - 78% of queries are not modified
    - Follow links ("the scent of information")

# Power law

- Few popular broad queries

- Many rare specific queries

# Top queries

- Most are related to sex

- 2007 Who What How (Google)

| Who is... | What is... | How to... |
| --- | --- | --- |
| 1. who is god | 1. what is love | 1. how to kiss |
| 2. who is who | 2. what is autism | 2. how to draw |
| 3. who is lookup | 3. what is rss | 3. how to knit |
| 4. who is jesus | 4. what is lupus | 4. how to hack |
| 5. who is it | 5. what is sap | 5. how to dance |
| 6. who is buckethead | 6. what is bluetooth | 6. how to crochet |
| 7. who is calling | 7. what is emo | 7. how to meditate |
| 8. who is keppler | 8. what is java | 8. how to flirt |
| 9. who is this | 9. what is hpv | 9. how to levitate |
| 10. who is satan | 10. what is gout | 10. how to skateboard |

# How far do people look for results?



**If you don't find what you are looking for, at what point do you move on either to another search engine or to another search on the same engine?**

| | |
|---|---|
| After the first few entries | 22.6% |
| After the first page | 18.6% |
| After the first two pages | 25.8% |
| After the first three pages | 14.7% |
| More than three pages | 10.8% |
| The whole list, unless its dozens of pages | 7.4% |

Number of Responses (0, 50, 100, 150, 200, 250, 300, 350, 400, 450)

*iProspect*

# True Example *

The User
flickr:crankyT



| | |
|---|---|
| Task | Stop the noisy fan in the courtyard |
| Info Need | Info about EPA regulations |
| Verbal Form | What are EPA rules on noise pollution? |
| Query | EPA Sound Pollution |

search

[ Go ] [ Search ]

"To Google or to GoTo" Business Week Online 9/28/2001

# How do users evaluate search engines?

- Quality of pages
  - Classic IR relevance
  - Also important:
    - Trust
    - Duplicate elimination
    - Readability
    - Fast Access
    - No pop-ups

# How do users evaluate search engines?

- Precision is more important than recall

  - Precision:

    - How precise is a portal in locating relevant results?

  - Recall

    - How thorough is the coverage of available relevant results?

- Precision with 1 result, 10 results, 2-3 pages of results.

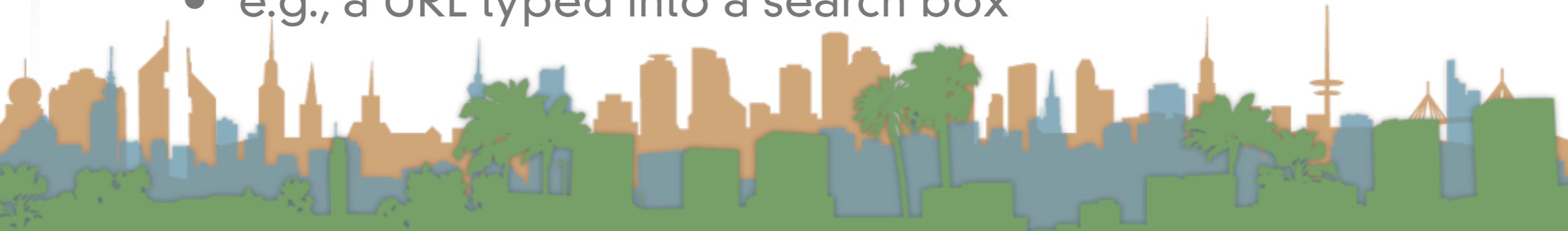- When is recall important?

# How do users evaluate search engines?

- Recall is sometimes important:
  - Googling for a new doctor
  - Googling a prospective employee
  - Googling your date

# How do users evaluate search engines?

- Good U/I
  - Simple
  - No Clutter
- Pre and post processing tools
  - Spell check ("Did you mean ....?")
  - Suggested alternative searches
  - Links to resources (maps, images, stock quotes)
- Able to deal with typical behavior
  - e.g., a URL typed into a search box

# Loyalty to a given search engine

- iProspect Survey 4/2004

**Which would you say best describes how you use search engines?**

I usually use the same search engine or directory — **56.7%**

I have several favorite search engines and use them interchangeably — **30.5%**

I use different search engines for different types of searches — **12.8%**

0    100    200    300    400    500    600    700    800    900    1000

**Number of Responses**

*iProspect*