مبانی بازیابی اطلاعات و جستجوی وب

Web Search –۱۱

# Brief (non-technical) history

- Early keyword-based engines ca. 1995-1997
- Paid search ranking: Goto ($\rightarrow$ Yahoo!)
  - Your search ranking depended on how much you paid
  - Auction for keywords

# Brief (non-technical) history

- 1998+: Link-based ranking pioneered by Google
  - Great user experience
- Result: Google added paid search "ads" to the side, independent of search results
- 2005+: Google gains search share, dominating in Europe and very strong in North America

# Without search engines, the web wouldn't work

- Without search, content is hard to find.
- → Without search, there is no incentive to create content.
- Somebody needs to pay for the web.
  - Servers, web infrastructure, content creation
  - A large part today is paid by search ads.
  - Search pays for the web.

# IR on the web vs. IR in general

- On the web, search is not just a nice feature.

  - Search is a key enabler of the web: . . .

$\rightarrow$ look at search ads

- The web is a chaotic und uncoordinated collection. $\rightarrow$ lots of duplicates – need to detect duplicates

- No control / restrictions on who can author content $\rightarrow$ lots of spam – need to detect spam

- The web is very large. $\rightarrow$ need to know how big it is

File   Edit   View   Go   Bookmarks   Yahoo!   Tools   Help

http://www.google.com/search?hl=en&q=nigritude+ultramarine&btnG=Google+Search   Go

Getting Started   Latest Headlines

Y!   Search Web   Mail   My Yahoo!   Games   Movies   Music   Answers   Personals   Sign In

pragh60@gmail.com | My Account | Sign out

**Google**

Web   Images   Groups   News   Froogle   Local   more »

nigritude ultramarine    [Search]   Advanced Search   Preferences

## Web

Results **1 - 10** of about **185,000** for **nigritude** ultramarine. (0.35 seconds)

### Anil Dash: Nigritude Ultramarine
Do me a favor: Link to this post with the phrase **Nigritude Ultramarine**. ... Just placed a link to your **Nigritude Ultramarine** article on my weblog. Cheers! ...
www.dashes.com/anil/2004/06/04/**nigritude**_ultra - 101k - Mar 1, 2006 -
Cached - Similar pages

### Nigritude Ultramarine FAQ
**Nigritude Ultramarine** FAQ - frequently asked questions about **nigritude ultramarine** and the realted SEO contest.
www.**nigritudeultramarine**s.com/ - 59k - Cached - Similar pages

### SEO contest - Wikipedia, the free encyclopedia
The **nigritude ultramarine** competition by SearchGuild is widely acclaimed as ... Comparison of search results for **nigritude ultramarine** during and after the ...
en.wikipedia.org/wiki/Nigritude_**ultramarine** - 37k - Cached - Similar pages

### Slashdot | How To Get Googled, By Hook Or By Crook
The current 3rd result showcases the "**Nigritude Ultramarine** Fighting Force" who ... When discussing **nigritude ultramarine** [slashdot.org] it is important to ...
slashdot.org/article.pl?sid=04/05/09/1840217 - 110k - Cached - Similar pages

### The Nigritude Ultramarine Search Engine Optimization Contest
It's sweeping the web -- or at least search engine optimizers -- a new contest to rank tops for the term **nigritude ultramarine** on Google.
searchenginewatch.com/sereport/article.php/3360231 - 57k - Cached - Similar pages

Sponsored Links

### Business Blogging Seminar
...g to L.A. March 16
Top bloggers reveal key techniques
www.blogbusinesssummit.com
Los Angeles, CA

### Full-Time SEO & SEM Jobs
Find companies big & small hiring full-time SEO & SEM pros right now
CareerBuilder.com

### SEO Contests
Information on SEO Contests like the **Nigritude Ultramarine** contest.
www.seo-contests.com/

### The SEO Book
**Nigritude Ultramarine** & SEO secrets Fun, free, raw, & different.
www.seobook.com

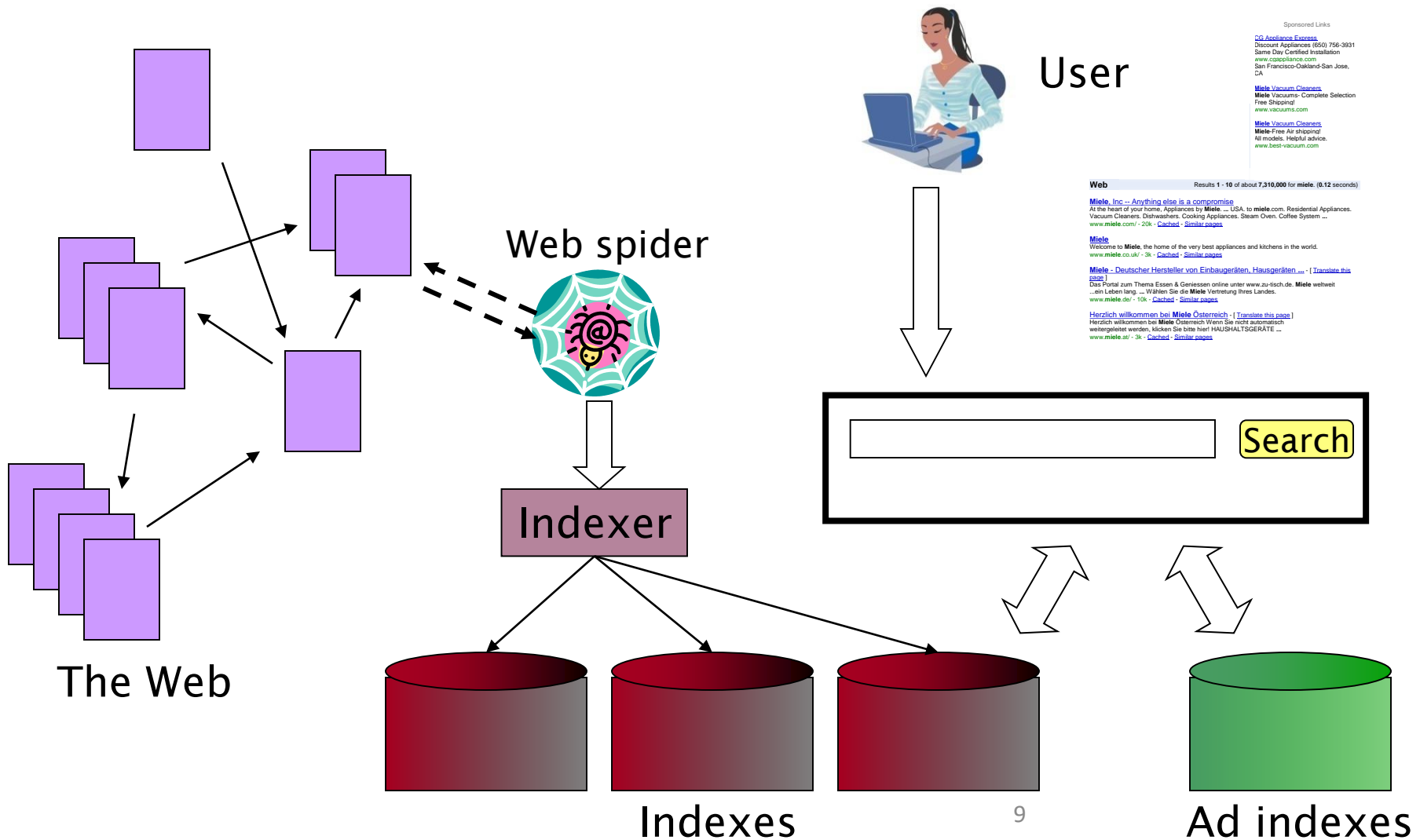Music - Dance - Electronic Overstock.com

Paid Search Ads

Algorithmic results.

7

Done

# How are ads ranked?

- First cut: according to bid price
  - Bad idea: open to abuse
  - We don't want to show nonrelevant ads.
- Instead: rank based on bid price and relevance
- Key measure of ad relevance: clickthrough rate
  - clickthrough rate = CTR = clicks per impressions
- Result: A nonrelevant ad will be ranked low.
- Other ranking factors: location, time of day, quality and loading speed of landing page
- The main ranking factor: the query

8

# Web search basics

User

Web spider

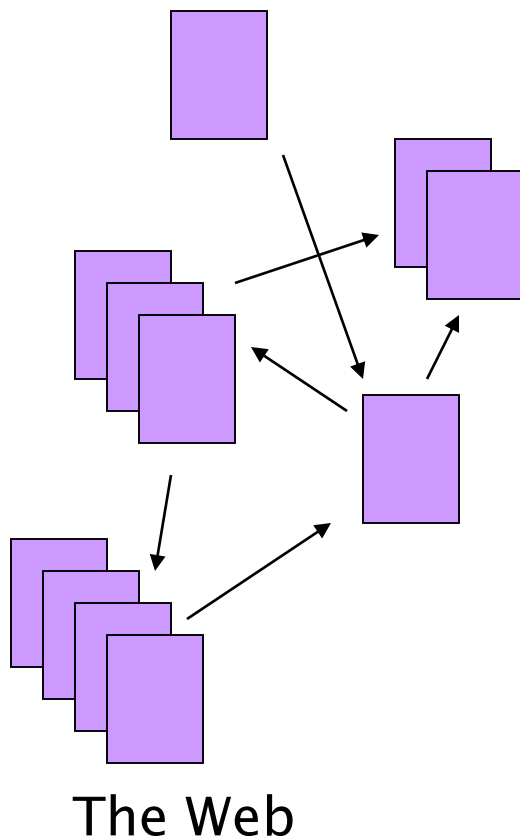Indexer

The Web

Indexes

9

Ad indexes

# Users' empirical evaluation of results

- Quality of pages varies widely
  - Relevance is not enough
  - Other desirable qualities (non IR!!)
    - Content: Trustworthy, diverse, non-duplicated, well maintained
    - Web readability: display correctly & fast
    - No annoyances: pop-ups, etc.
- Precision vs. recall
  - On the web, recall seldom matters
- What matters
  - Precision at 1? Precision above the fold?
  - Comprehensiveness – must be able to deal with obscure queries
    - Recall matters when the number of matches is very small
- <span style="color:red">User perceptions may be unscientific, but are significant over a large aggregate</span>

# Users' empirical evaluation of engines

- Relevance and validity of results
- UI – Simple, no clutter, error tolerant
- Trust – Results are objective
- Coverage of topics for polysemic queries
- Pre/Post process tools provided
  - Mitigate user errors (auto spell check, search assist,…)
  - Explicit: Search within results, more like this, refine …
  - Anticipative: related searches
- Deal with idiosyncrasies
  - Web specific vocabulary
    - Impact on stemming, spell-check, etc.
  - Web addresses typed in the search box

# The Web document collection



The Web

- 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 collections … but corporate records are catching up
- Growth – slowed down from initial "volume doubling every few months" but still expanding
- Content can be *dynamically generated*
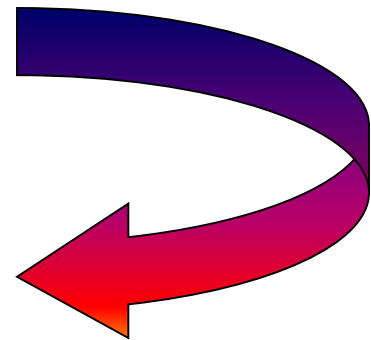
# SPAM
## (SEARCH ENGINE OPTIMIZATION)

# The trouble with paid search ads …

- It costs money (CPM, CPC, etc.).  What's the alternative?
- *Search Engine Optimization:*
    - "Tuning" your web page to rank highly in the algorithmic search results for select keywords
    - Alternative to paying for placement
    - Thus, intrinsically a marketing function
- Performed by companies, webmasters and consultants ("Search engine optimizers") for their clients
- Some perfectly legitimate, some very shady

# Simplest forms

- First generation engines relied heavily on *tf/idf*
  - The top-ranked pages for the query **Babolsar resort** were the ones containing the most **Babolsar**'s and **resort**'s

- SEOs responded with dense repetitions of chosen terms
  - e.g., **Babolsar resort Babolsar resort Babolsar resort**
  - Often, the repetitions would be in the same color as the background of the web page
    - Repeated terms got indexed by crawlers
    - But not visible to humans on browsers

Pure word density cannot be trusted as an IR signal
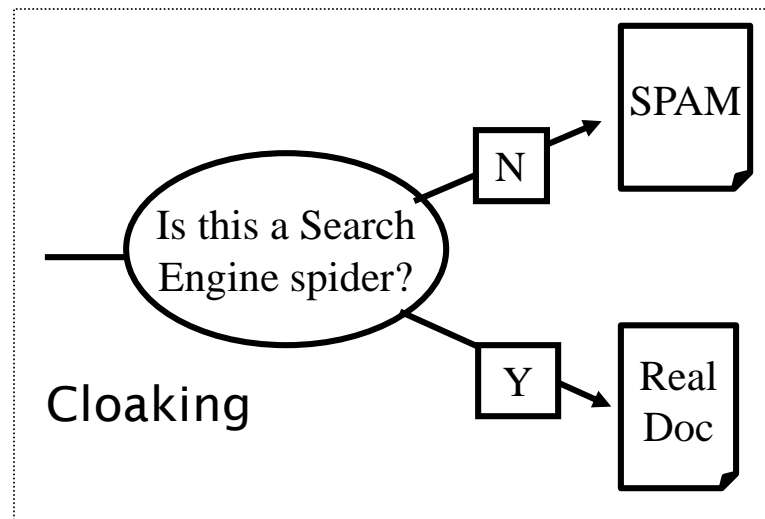
# Variants of keyword stuffing

- Misleading meta-tags, excessive repetition
- Hidden text with colors, style sheet tricks, etc.

**Meta-Tags** =
"… London hotels, hotel, holiday inn, hilton, discount, booking, reservation,mp3, …"

# Cloaking

- Serve fake content to search engine spider

# Optional: More spam techniques

- **Doorway pages**
  - Pages optimized for a single keyword that re-direct to the real target page

- **Link spamming**
  - *Domain flooding:* numerous domains that point or re-direct to a target page

- **Robots**
  - Fake query stream – rank checking programs
  - Millions of submissions

# The war against spam

- Quality signals - Prefer authoritative pages based on:
  - Votes from authors
  - Votes from users
- Policing of URL submissions
  - Anti robot test
- Limits on meta-keywords
- Robust link analysis
  - Ignore statistically implausible linkage (or text)
  - Use link analysis to detect spammers

- Spam recognition by machine learning
  - Training set based on known spam
- Family friendly filters
  - Linguistic analysis, general classification techniques, etc.
  - For images: source text analysis, etc.
- Editorial intervention
  - Blacklists
  - Top queries audited
  - Complaints addressed
  - Suspect pattern detection

# SIZE OF THE WEB

# What is the size of the web ?

- Issues
  - The web is really infinite
    - Dynamic content, e.g., calendars
  - Static web contains syntactic duplication, mostly due to mirroring (~30%)
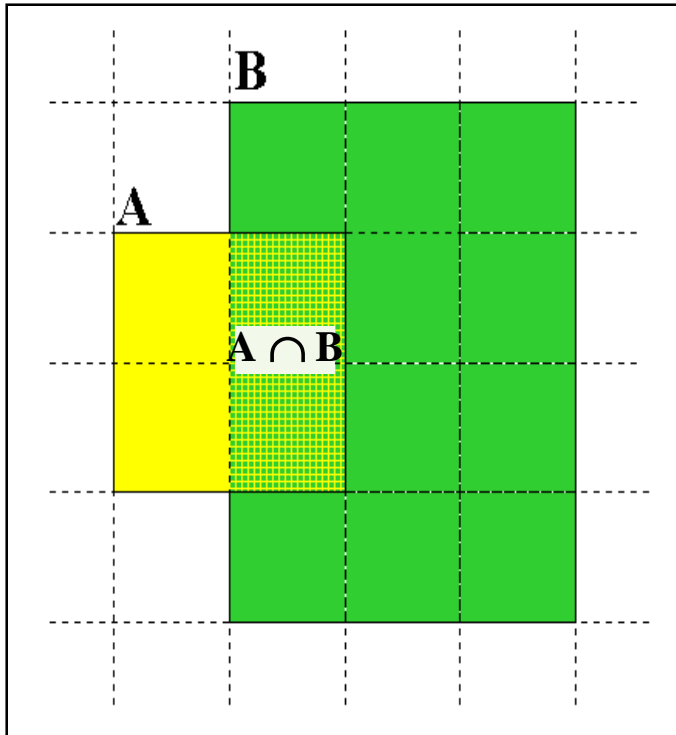  - Some servers are seldom connected
- Who cares?
  - Media, and consequently the user
  - Engine design
  - Engine crawl policy. Impact on recall.

# What can we attempt to measure?

- The relative sizes of search engines

# Relative Size from Overlap
# Given two engines A and B

Sample URLs randomly from A

Check if contained in B and vice versa

```
A ∩ B =  (1/2) * Size A
A ∩ B =  (1/6) * Size B

(1/2)*Size A = (1/6)*Size B

∴ Size A / Size B =
            (1/6)/(1/2) = 1/3
```

**Each test involves:** (i) Sampling  (ii) Checking

# Optional: Sampling URLs

■Ideal strategy: Generate a random URL and check for containment in each index.

■ Problem: Random URLs are hard to find!  Enough to generate a random URL contained in a given Engine.

■Approach 1: Generate a random URL contained in a given engine
   ■Suffices for the estimation of relative size

■Approach 2: Random walks / IP addresses
   ■In theory: might give us a true estimate of the size of the web (as opposed to just relative sizes of indexes)

■Random URLs from random queries

# منابع

- کتاب فصل نوزدهم  An introduction to information retrieval