

➤ **Web Information Retrieval**

Web Search

(SOSE 2023, 24.05.2023)

Frank Hopfgartner, Stefania Zourlidou
Institute for Web Science and Technologies

Credit for these slides

These slides have been adapted from

- Web IR (Zeyd Boukhers-WeST, SOSE 2020)

Recapitulation

- Language models
 - Query likelihood
 - Document likelihood
 - Comparison model

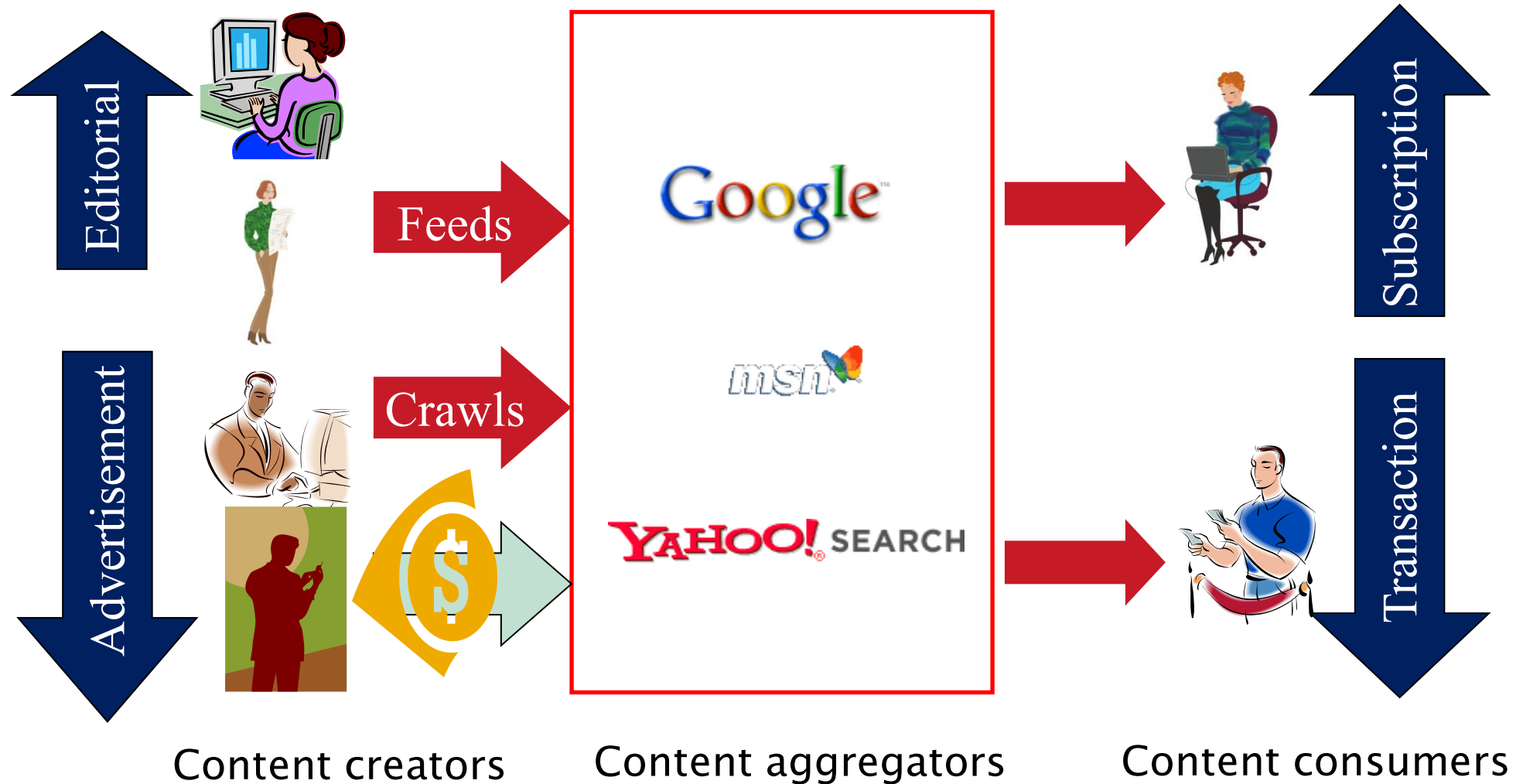
Objectives of this lecture

- Classical IR vs. Web IR
- Web IR, Web search basics
 - Ads
 - Spams
 - Duplicates
 - User Needs
 - Web Graph (anchor text)
 - Indexing

➤ 1. Classical IR

- Corpus
 - Fixed collection
 - Corpus is predetermined
- Goal
 - Retrieve documents with content relevant to user's information need
- Relevance
 - For every query q and a document d , there exists a relevance score $Score(q, d)$
 - Score is context independent
 - Score is user independent

➤ 2. Web IR



- Corpus
 - Scale much larger than in classical IR
 - Volume is expanding
 - Spam: Billions of pages
- Content
 - Truth, lies, obsolete information, contradictions
 - Unstructured, semi-structured, structured
 - Can be *dynamically generated*

- Other characteristics
 - **4.39 billion** internet users in 2019
 - Digital 2019 report
 - 1.7 billion websites in 2020
 - internetlivestats.com
 - High linkage: More than 8 links/page in the average
 - Significant duplication
 - Syntactic – 30%-40% (near) duplicates [Brod97, Shiv99b, etc.]
 - Semantic – ?

- Ads
- Spams
- Duplicates
- User Needs
- Web Graph (anchor text)
- Indexing
- Empirical Evaluation

➤ 3. Ads Web IR

Web IR: 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!)
 - 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 save Inktomi
 - 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)

About 1.220.000.000 results (0,44 seconds)

[Ad] www.qoppa.com/ ▼ +1 404-685-8733

PDF Automation Server | Tools to Streamline Processing

Rich Set of **PDF Processing** Functions for Different Environments. Try It Now! Trial Download.
Unleash the Power of **PDF**. Full Adobe Compatibility. Types: Java Developer API, SDK, Desktop
PDF Software, **PDF** Server Software.
[Contact Us](#) · [About Us](#) · [All PDF Libraries](#)

processing.org › [reference](#) › [libraries](#) › pdf ▼

PDF \ Libraries \ Processing.org

PDF Export. The **PDF** library makes it possible to write **PDF** files directly from **Processing**.
These vector graphics files can be scaled to any size and output ...

forum.processing.org › [topic](#) › [making-a-pdf-file](#) ▼

making a pdf-file - Processing Forum

Aug 20, 2013 - 11 posts - 4 authors

I saw the recodeproject and would like to know how I could make the output go to hi-quality-**PDF**
as to print it on a large scale penplotter.

forum.processing.org › [Using Processing](#) › [Library Questions](#) ▼

How to export as a PDF? - Processing 2.x and 3.x Forum

Mar 3, 2017 - ... map using tilemill and unfolding maps and now want to export/save it as a pdf.
Here's the code I've tried, however the pdf is saving as blank.

Ads

Algorithmic
result

Ads vs. search results

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

(Ad) [www.united-domains.de/](#) ▾
Domains | Die besten Adressen im Web | united-domains.de
Wunschdomain beim Spezialisten schnell und einfach suchen. Jetzt registrieren!
Zufriedenheitsgarantie. Transparente Preise. Attraktive E-Mail-Pakete.

Neue Domain-Endungen
.web, .shop, .app und viele mehr -
Die neuen Domain-Endungen sind da!

Domains registrieren
Viele Domain-Endungen einfach
und unkompliziert registrieren!

(Ad) [www.one.com/](#) ▾
Wunschdomain günstig sichern | Starten sie jetzt durch | one.com
Ihr Online-Erfolg beginnt mit dem Kauf eines Domainnamens. Alles, was Sie benötigen...

(Ad) [de.godaddy.com/domainnamen](#) ▾ 089 21094807
GoDaddy™ Domains ab 0,99 € | Kaufen Sie Ihre heute
Durchsuche die größte Domain-Datenbank und registriere ab 0,99 €! Heute Kaufen

(Ad) [www.strato.de/](#) ▾
Domain im Web reservieren | Wunschadresse inkl. E-Mail
Zahlreiche Domain-Endungen zur Auswahl. Jetzt unverwechselbar im Internet sein

[www.checkdomain.de](#) > domains > web-domain ▾ [Translate this page](#)
Web-Domain sichern - Ihre Wunschdomain preiswert ...
So sichern Sie sich eine Webdomain. Eine Web Domain ist der eigenständige Internet-Auftritt von Personen, Unternehmen oder Organisationen, um Besucher im ...

[www.domain.com](#) ▾
Website Domains Names & Hosting | Domain.com
Find and purchase your next **website domain** name and hosting without breaking the bank.
Seamlessly establish your online identify today.
[Domain Registration](#) · [Domain.com](#) | [Blog](#) · [Domain Privacy](#) · [Full service web design](#)

People also ask

What is website domain? ▾

How do I get a web domain? ▾

What is domain with example? ▾

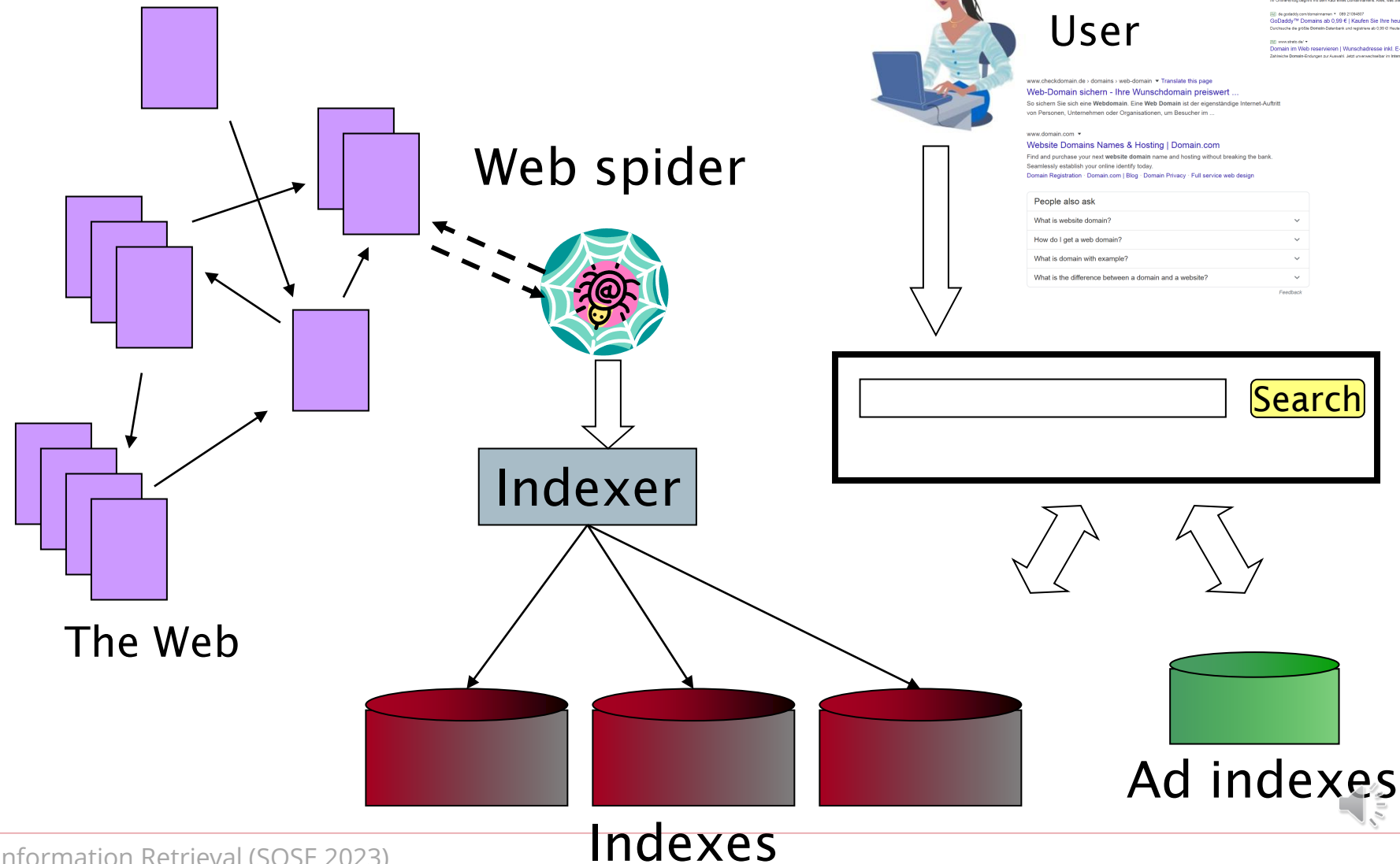
What is the difference between a domain and a website? ▾

[Feedback](#)

Search = **web domain**

- Other search engines (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



How are ads ranked?

- First cut: according to bid price à la Goto
 - Bad idea: open to abuse
 - Example: query [Buying fresh Chicken?] → ad for KFC
 - 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.
 - Even if this decreases search engine revenue short-term
 - Hope: Overall acceptance of the system and overall revenue is maximized if users get useful information
- Other ranking factors: location, time of day, quality and loading speed of landing page
- The main ranking factor: the query

Google's second price auction

advertiser	bid	CTR	ad rank	rank	paid
A	\$4.00	0.01	0.04	4	(minimum)
B	\$3.00	0.03	0.09	2	\$2.68
C	\$2.00	0.06	0.12	1	\$1.51
D	\$1.00	0.08	0.08	3	\$0.51

- bid: maximum bid for a click by advertiser
- CTR: click-through rate: when an ad is displayed, what percentage of time do users click on it? CTR is a measure of relevance.
- ad rank: $\text{bid} \times \text{CTR}$: this trades off (i) how much money the advertiser is willing to pay against (ii) how relevant the ad is
- rank: rank in auction
- paid: second price auction price paid by advertiser

Google's second price auction

advertiser	bid	CTR	ad rank	rank	paid
A	\$4.00	0.01	0.04	4	(minimum)
B	\$3.00	0.03	0.09	2	\$2.68
C	\$2.00	0.06	0.12	1	\$1.51
D	\$1.00	0.08	0.08	3	\$0.51

Second price auction: The advertiser pays the minimum amount necessary to maintain their position in the auction (plus 1 cent)

- $\text{price}_1 = \text{bid}_2 \times \text{CTR}_2 / \text{CTR}_1$
- $p_1 = \text{bid}_2 \times \text{CTR}_2 / \text{CTR}_1 = 3.00 \times 0.03 / 0.06 = 1.50$
- $p_2 = \text{bid}_3 \times \text{CTR}_3 / \text{CTR}_2 = 1.00 \times 0.08 / 0.03 = 2.67$
- $p_3 = \text{bid}_4 \times \text{CTR}_4 / \text{CTR}_3 = 4.00 \times 0.01 / 0.08 = 0.50$

Keywords with high bids

- According to <https://www.wordstream.com/articles/most-expensive-keywords>

Insurance	\$54.91
Loans	\$44.28
Mortgage	\$47.12
Attorney	\$47.07
Credit	\$36.06
Lawyer	\$42.51
Donate	\$42.02
Degree	\$40.61
Hosting	\$31.91
Claim	\$45.51
Conference Call	\$42.05
Trading	\$33.19
Software	\$35.29

Search ads: a win-win-win?

- The search engine company gets revenue every time somebody clicks on an ad
- The user only clicks on an ad if they are interested in the ad
 - Search engines punish misleading and nonrelevant ads
 - As a result, users are often satisfied with what they find after clicking on an ad
- The advertiser finds new customers in a cost-effective way

- Why is web search potentially more attractive for advertisers than TV spots, newspaper ads or radio spots?
- The advertiser pays for all this. How can the advertiser be cheated?
- Any way this could be bad for the user?
- Any way this could be bad for the search engine?

Not a win-win-win: keyword arbitrage

- Buy a keyword on Google
- Then redirect traffic to a third party that is paying much more than you are paying Google
 - E.g., redirect to a page full of ads
- This rarely makes sense for the user
- Ad spammers keep inventing new tricks
- The search engines need time to catch up with them

➤ 4. Spam

The trouble with paid placement

- It costs money. What's the alternative?
- *Search Engine Optimization*
 - "Tuning" your web page to rank highly in the 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

- First generation engines relied heavily on *tf/idf*
 - The top-ranked pages for the query **maui resort** were the ones containing the most **maui's** and **resort's**
- SEOs responded with dense repetitions of chosen terms
 - e.g., **maui resort maui resort maui 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

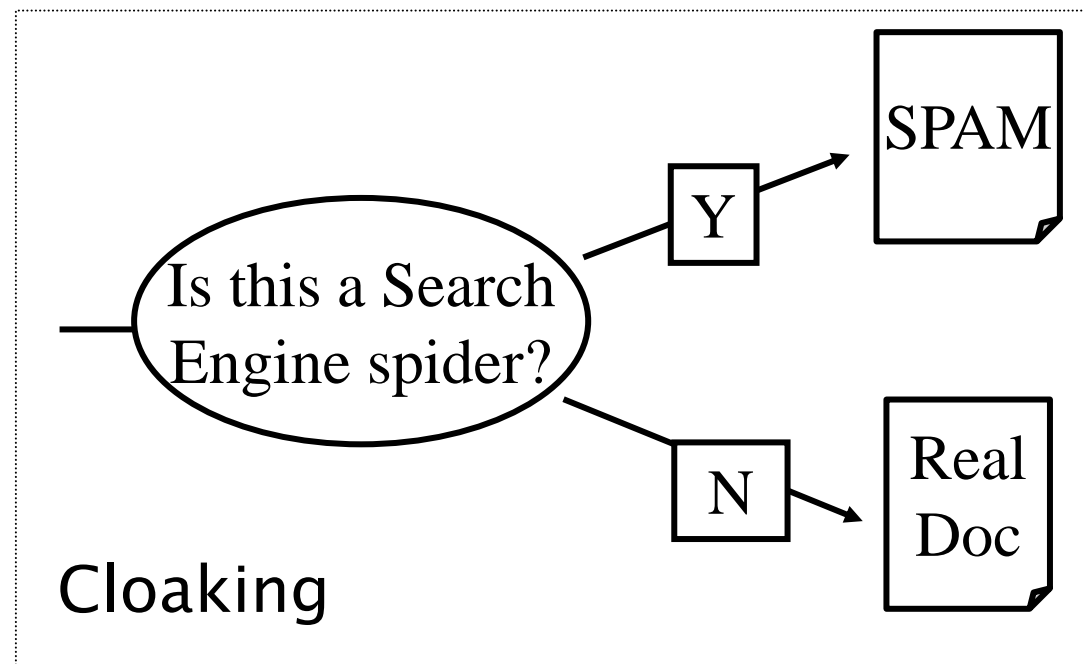
- Misleading meta-tags, excessive repetition
- Hidden text with colors, style sheet tricks, etc.

Meta-Tags =

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

Cloaking

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



Search engine optimization (Spam)

- Motives
 - Commercial, political, religious, lobbies
 - Promotion funded by advertising budget
- Operators
 - Contractors (Search Engine Optimizers) for lobbies, companies
 - Web masters
 - Hosting services
- Forums
 - E.g., Web master world (www.webmasterworld.com)
 - Search engine specific tricks

The spam industry

Search Engine Cloaker
Need more search engine listings?
fantomas

Web Guide
Our hand-picked directory of the best business links on the web.
Cloaking
Category Path
[Home](#) > [Guide Topics](#) > [Technology](#) > [Internet](#) > [Search Technology](#) > [Search Engines](#) > [Search Engine Placement](#) > [Cloaking](#)

OUTSMART
Free Domain Forwarding - Domain Cloaking - DNS Forwarding
Web site is cloaked when the web address of a web site is hidden from viewers in their browser window.
For example your user would type in www.yourname.com into their browser window. They are then automatically redirected to your web site:
(<http://www.someisp.com/~users/yourname/yoursite.html>) or any where you like.
However your users would continue to www.yourname.com as they browsed.
Cloaking Services: Included Branded Email Services 5 Mail boxes mailboxename@yourDomain.com \$49/Year

PhantomLine™ — the ultimate stealth
News Best Keywords! SE
Understanding Cloaking
al: Cloaking and Stealth Technology
Page 2 | Page 3 | Page 4 | Page 5
g, stealth or phantom page technology constitutes the sophisticated and efficient approach towards search engine optimization. A mystique surrounding cloaking or stealth tech

The war against spam

- Quality signals - Prefer authoritative pages based on
 - Votes from authors (linkage signals)
 - Votes from users (usage signals)
- Limits on meta-keywords
- Robust link analysis
 - Ignore statistically implausible linkage (or text)
 - Use link analysis to detect spammers (guilt by association)
- Spam recognition by machine learning
 - Training set based on known spam
- Family friendly filters
 - Linguistic analysis, general classification techniques, etc.
 - For images: flesh tone detectors, source text analysis, etc.
- Editorial intervention
 - Blacklists
 - Top queries audited
 - Complaints addressed
 - Suspect pattern detection

More on spam

- Web search engines have policies on SEO practices they tolerate/block
 - <http://help.yahoo.com/help/us/ysearch/index.html>
 - <http://www.google.com/intl/en/webmasters/>
- Adversarial IR: the unending (technical) battle between SEO's and web search engines
- Research <http://airweb.cse.lehigh.edu/>

➤ 5. Duplicates – WEB IR

Duplicate detection

- The web is full of duplicated content
- More so than many other collections
- Exact duplicates
 - Easy to eliminate
 - E.g., use hash/fingerprint
- Near-duplicates
 - Abundant on the web
 - Difficult to eliminate
- For the user, it's annoying to get a search result with near-identical documents
- Marginal relevance is zero: even a highly relevant document becomes nonrelevant if it appears below a (near-)duplicate
- We need to eliminate near-duplicates

Near-duplicates: example

The image shows a side-by-side comparison of two web pages in a browser window. The left page is the official Wikipedia article for Michael Jackson, while the right page is a parody titled 'wapedia'. Both pages contain identical text, illustrating near-duplicate content.

Left Page (Wikipedia):

- Navigation: Main page, Contents, Featured content, Current events, Random article
- Search: Go, Search
- Interaction: About Wikipedia, Community portal, Recent changes, Contact Wikipedia
- Article Title: Michael Jackson
- Text: From Wikipedia, the free encyclopedia. For other persons named Michael Jackson, see Michael Jackson (disambiguation).
- Text: Michael Joseph Jackson (August 29, 1958 – June 25, 2009) was an American recording artist, entertainer and businessman. The seventh child of the Jackson family, he made his debut as an entertainer in 1968 as a member of The
- Image: Michael Jackson in a black suit and sunglasses.

Right Page (wapedia):

- Title: wapedia.
- Section: Wiki: Michael Jackson (1/6)
- Text: For other persons named Michael Jackson, see Michael Jackson (disambiguation).
- Text: Michael Joseph Jackson (August 29, 1958 - June 25, 2009) was an American recording artist, entertainer and businessman. The seventh child of the Jackson family, he made his debut as an entertainer in 1968 as a member of The Jackson 5. He then began a solo

Detecting near-duplicates

- Compute similarity with an edit-distance measure
- We want “syntactic” (as opposed to semantic) similarity
 - True semantic similarity (similarity in content) is too difficult to compute
- We do not consider documents near-duplicates if they have the same content, but express it with different words
- Use similarity threshold θ to make the call “is/isn’t a near-duplicate”
- E.g., two documents are near-duplicates if similarity
$$> \theta = 80\%.$$

Represent each document as set of shingles

- A shingle is simply a word n-gram
- Shingles are used as features to measure syntactic similarity of documents
- For example, for $n = 3$, “a rose is a rose is a rose” would be represented as this set of shingles
 - { a-rose-is, rose-is-a, is-a-rose }
- We define the similarity of two documents as the Jaccard coefficient of their shingle sets

Recall: Jaccard coefficient

- A commonly used measure of overlap of two sets
- Let A and B be two sets
- Jaccard coefficient

$$\text{JACCARD}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- $(A \neq \emptyset \text{ or } B \neq \emptyset)$
- $\text{JACCARD}(A, A) = 1$
- $\text{JACCARD}(A, B) = 0$ if $A \cap B = \emptyset$
- A and B don't have to be the same size
- Always assigns a number between 0 and 1

Jaccard coefficient: example

- Three documents
 - $d1$: "Jack London traveled to Oakland"
 - $d2$: "Jack London traveled to the city of Oakland"
 - $d3$: "Jack traveled from Oakland to London"
- Based on shingles of size 2 (2-grams or bigrams), what are the Jaccard coefficients $J(d1, d2)$ and $J(d1, d3)$?
 - $J(d1, d2) = 3/8 = 0.375$
 - $J(d1, d3) = 0$
 - Note: very sensitive to dissimilarity

➤ 6. User need – WEB IR

User needs

- **Informational** – want to learn about something (~40% / 65%) Low hemoglobin
- **Navigational** – want to go to that page (~25% / 15%) United Airlines
- **Transactional** – 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”
 - Need [Brod02, RL04]

Answering “the need behind the query”

- Semantic analysis
 - Query language determination
 - Auto filtering
 - Different ranking (if query in Japanese do not return English)
 - Hard & soft (partial) matches
 - Personalities (triggered on names)
 - Cities (travel info, maps)
 - Medical info (triggered on names and/or results)
 - Stock quotes, news (triggered on stock symbol)
 - Company info
 - Etc.
 - Natural Language reformulation
 - Integration of Search and Text Analysis

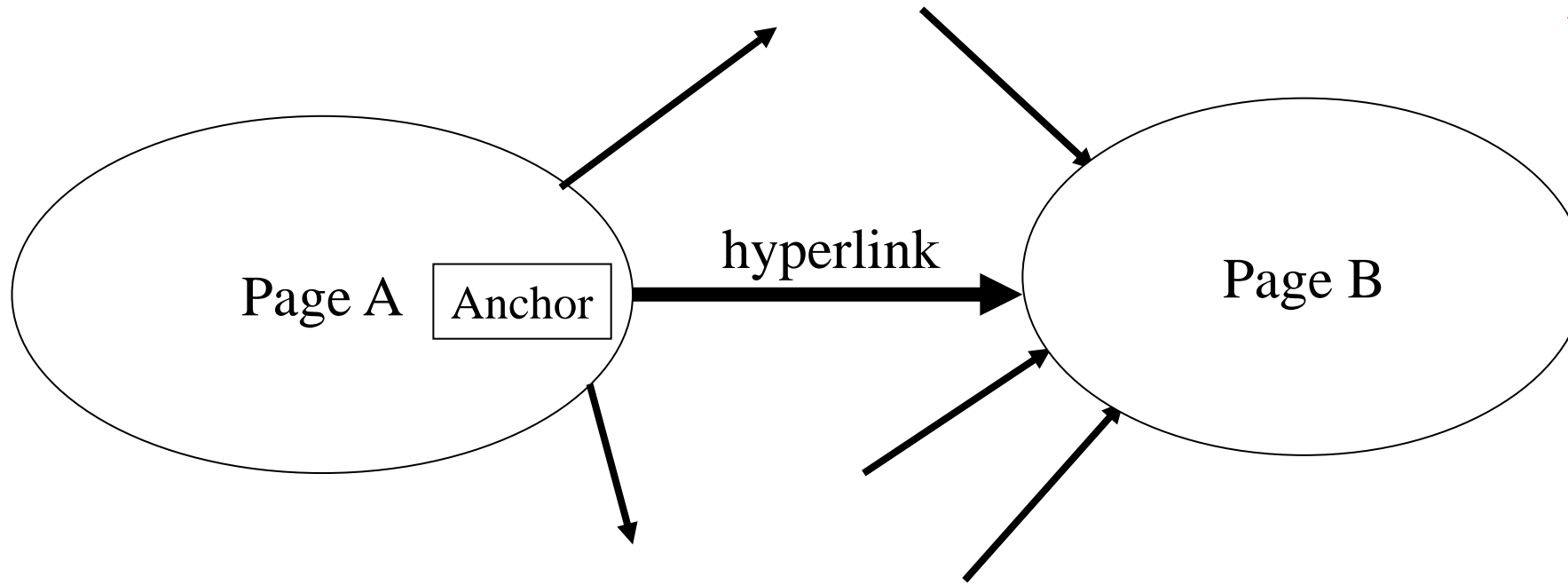
Answering “the need behind the query”: context

- Context determination
 - spatial (user location/target location)
 - query stream (previous queries)
 - personal (user profile)
 - explicit (user choice of a vertical search)
 - implicit (use Google from France, use google.fr)

- Context use
 - Result restriction
 - Kill inappropriate results
 - Ranking modulation
 - Use a “rough” generic ranking, but personalize later

➤ 7. Web graph

The Web as a directed graph

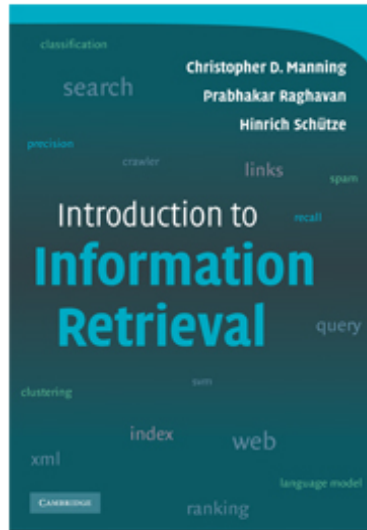


Hypothesis 1: A hyperlink between pages denotes a conferral of authority (quality signal)

Hypothesis 2: The text in the anchor of the hyperlink on page A describes the target page B

Assumption 1: reputed sites

Introduction to Information Retrieval



This is the companion website for the following book.

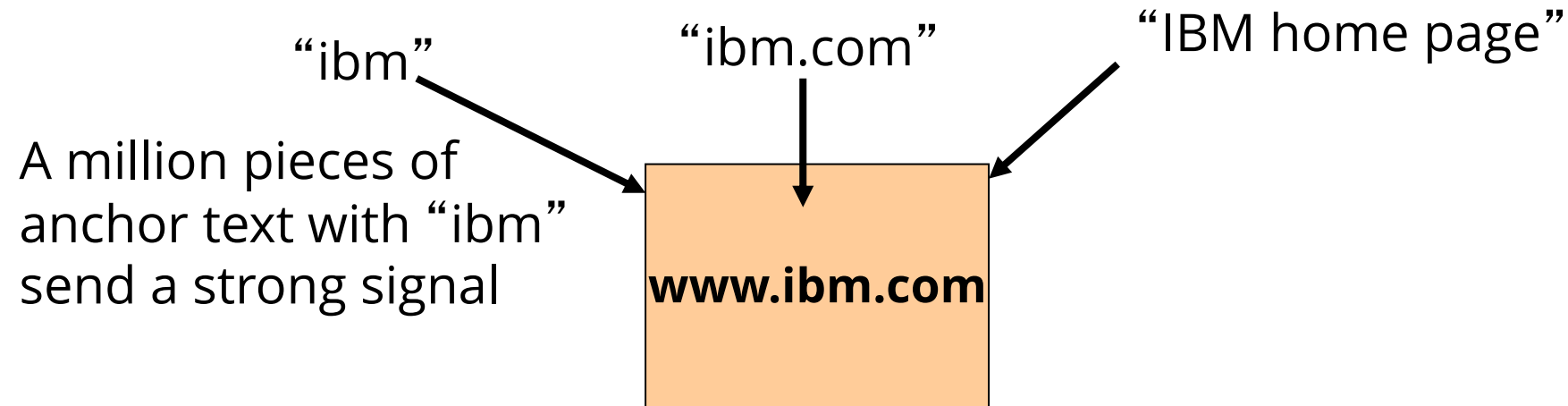
[Christopher D. Manning](#), [Prabhakar Raghavan](#) and [Hinrich Schütze](#), *Introduction to Information Retrieval*

You can order this book at [CUP](#), at your local bookstore or on the internet. The best search

The book aims to provide a modern approach to information retrieval from a computer science perspective. It is available at the [University of Cambridge](#) and at the [University of Stuttgart](#).

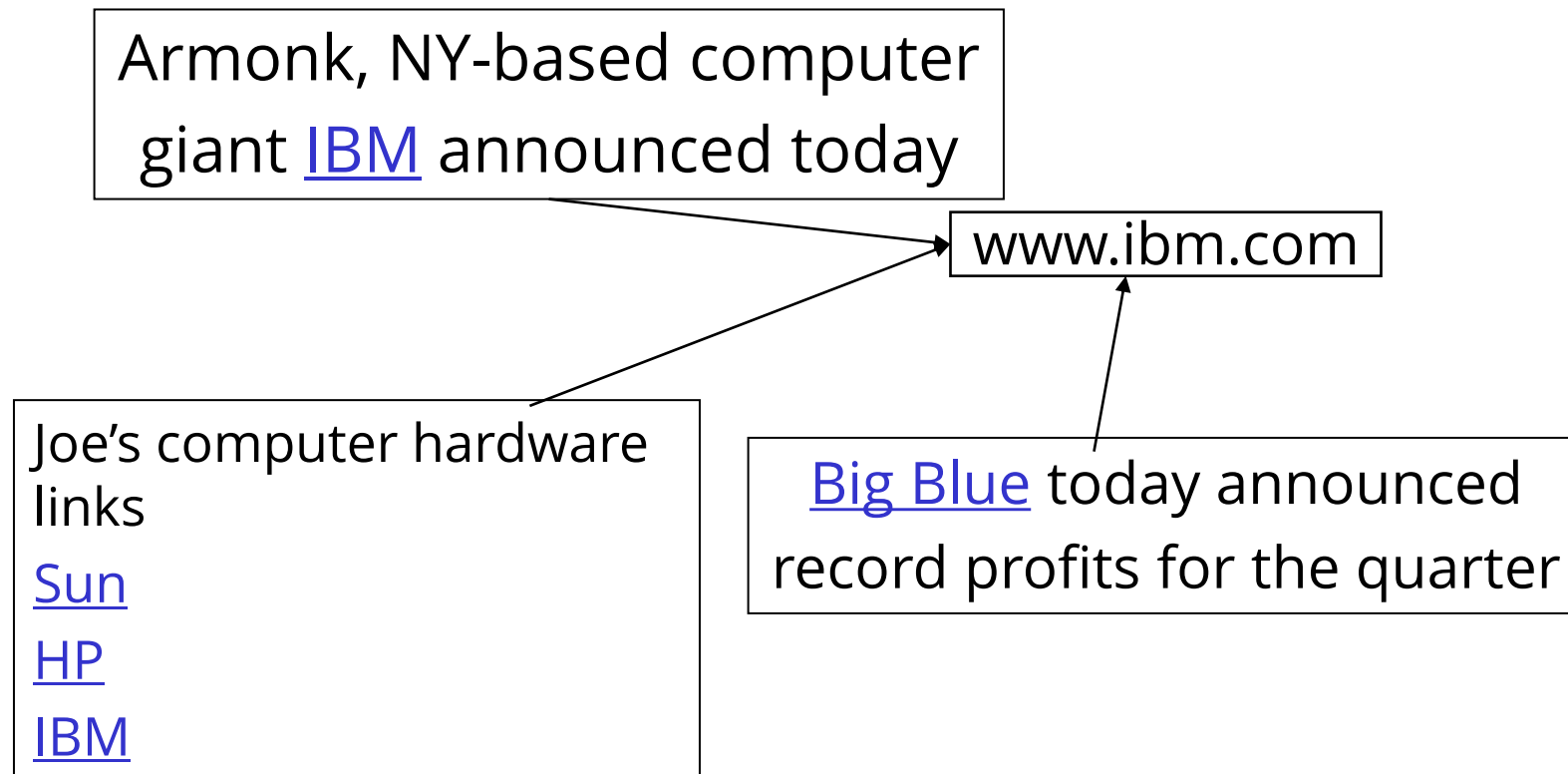
We'd be pleased to get feedback about how this book works out as a textbook, what is missing, and what is not. Please send your comments to: [informationretrieval \(at\) yahoogroups \(dot\) com](mailto:informationretrieval(at)yahoogroups(dot)com)

- For **ibm** how to distinguish between
 - IBM' s home page (mostly graphical)
 - IBM' s copyright page (high term freq. for 'ibm')
 - Rival' s spam page (arbitrarily high term freq.)



Indexing anchor text

- When indexing a document D , include (with some weight) anchor text from links pointing to D



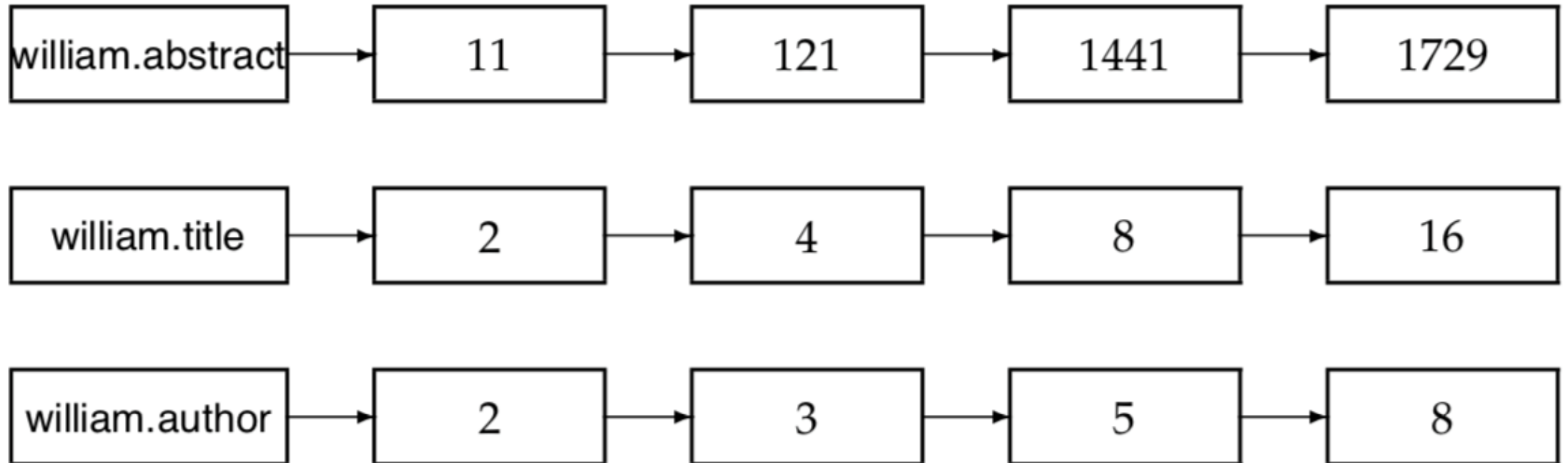
- Thus: Anchor text is often a better description of a page's content than the page itself
- Anchor text can be weighted more highly than document text

➤ 8. Indexing

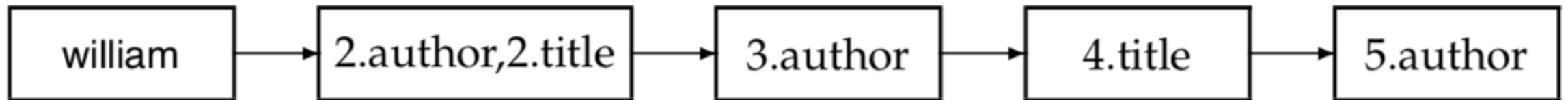
- Web pages with different zones
 - Title, metadata, author
- How to do weighted ranking

“find documents with *merchant* in the title and *william* in the author list and the phrase *gentle rain* in the body”

- zones are encoded as extensions of dictionary entry



- zone is encoded in the postings rather than dictionary



- Consider a set of documents each of which has l zones

$$g_1, \dots, g_l \in [0,1] \quad \sum_{i=1}^l g_i = 1$$

- s is the Boolean score denoting a match (or absence thereof) between *query* and the zone

$$\sum_{i=1}^l g_i s_i$$

- Web is large
- Term-partitioned index
- Document-partitioned index
- Map-reduce phenomena for the computation

➤ 9. User centric evaluation measures

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
- User perceptions may be unscientific, but are significant over a large aggregate

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 first, the last, the best and the worst ...”

➤ 10. Summary

Summary

- Web search basic
 - Ads
 - Spams
 - Duplicates
 - User Needs
 - Web Graph (anchor text)

[1] <https://olat.vcrp.de/auth/RepositoryEntry/4071063853>

[2] <https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze,
Introduction to Information Retrieval, Cambridge University Press. 2008

- ▶ Chapter 4 (Distributed indexing)
- ▶ Chapter 6 (Parametric and zone indexes)
- ▶ Chapter 19 (Web search basics)