

# ➤ Web Retrieval

## Search on the Web

Frank Hopfgartner  
Institute for Web Science and Technologies

# Intended Learning Outcomes

**At the end of the lecture, you will be able to:**

- Describe the structure of the Web
- Outline personalisation approaches
- Understand how ads are embedded on search result lists
- Describe spamming approaches to improve search ranking

# Outline

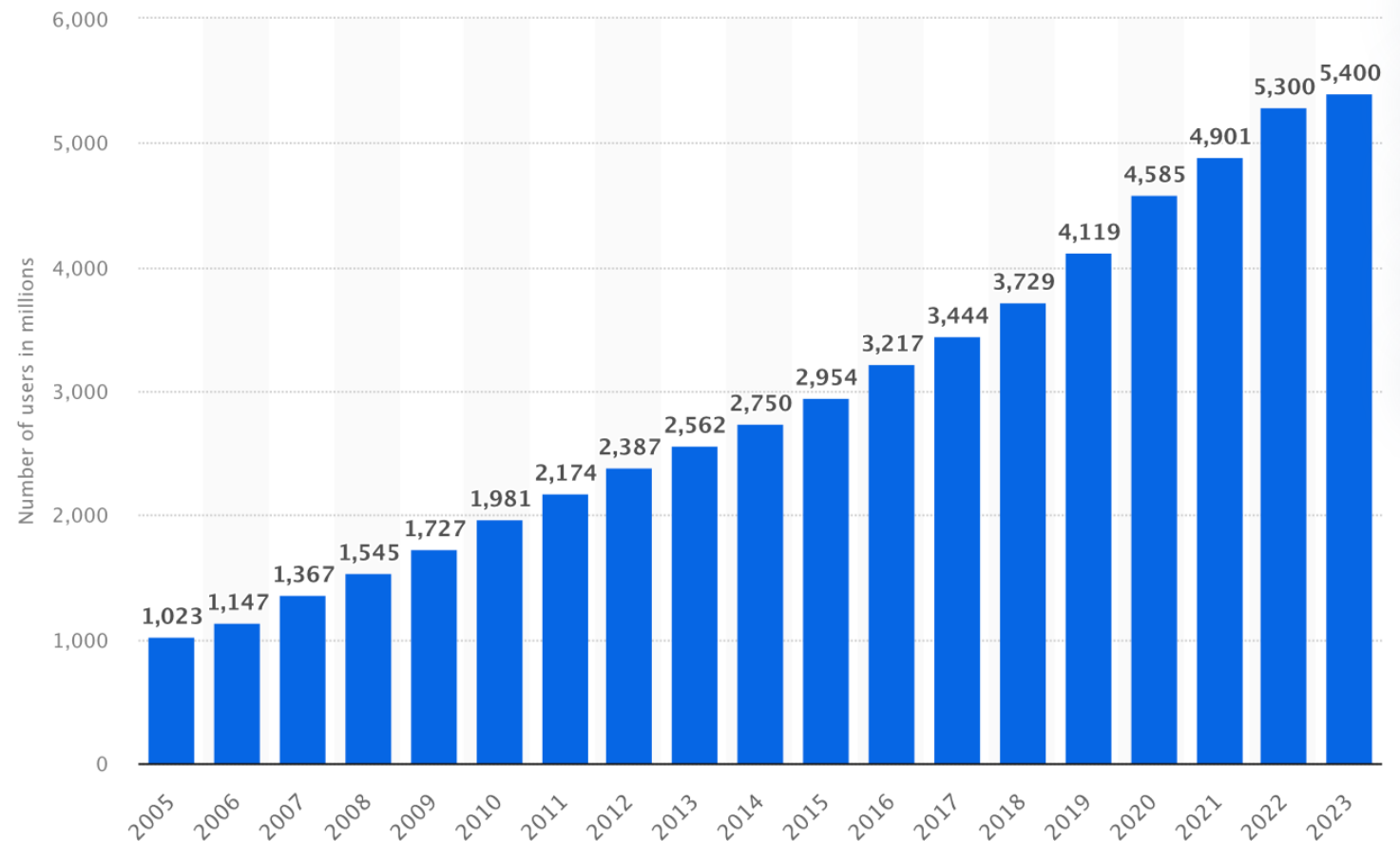
- Motivation
- Web Graph
- Personalisation
- Ads
- Spamming

# ➤ 1. Motivation

- 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

## Number of internet users worldwide from 2005-2023

- Billions of information needs to satisfy
- Relevance is context-dependent

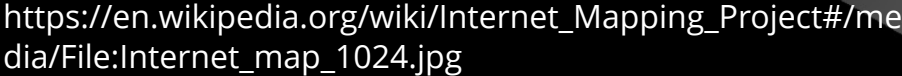


- Editorial content
- User-generated content
- AI-generated content
  
- Truth, lies, obsolete information, contradictions
- Unstructured, semi-structured, structured
- Can be dynamically generated



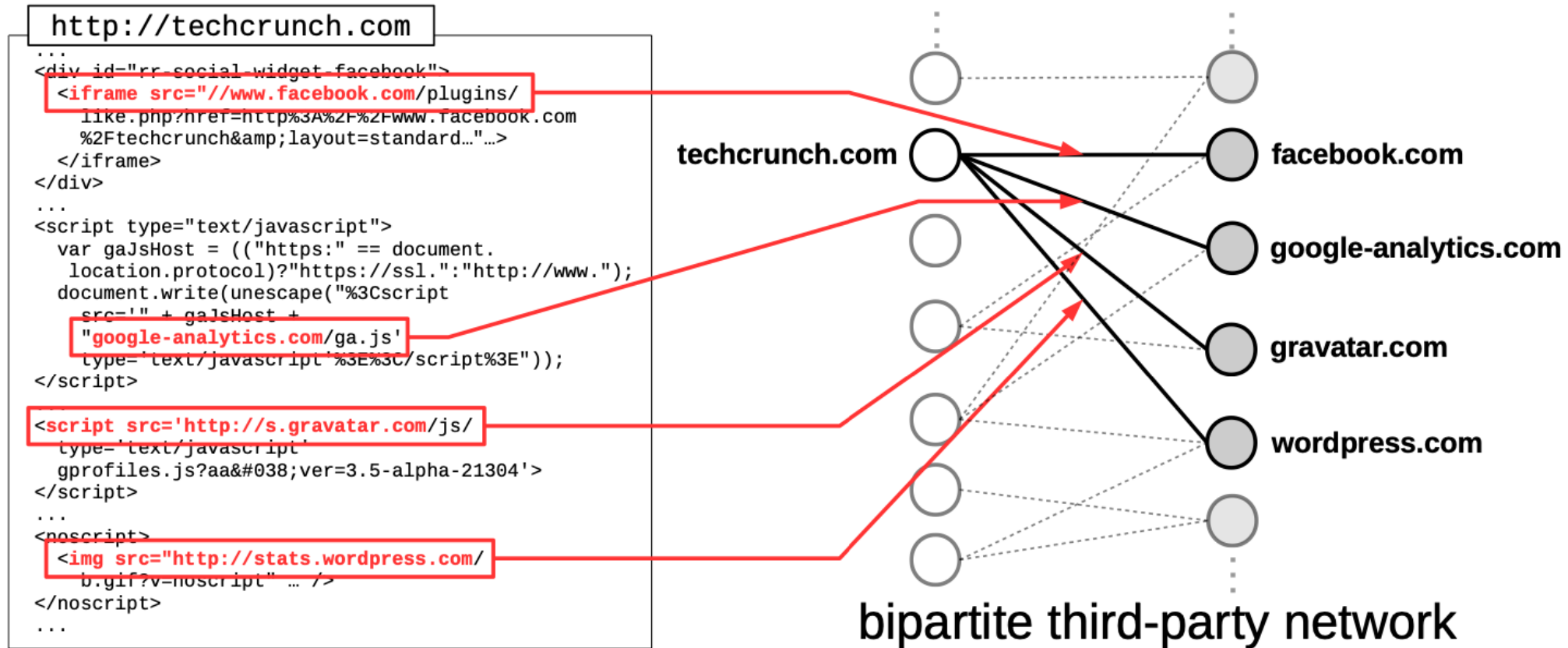
- HTML formatting
- Hyperlinks between web pages

Much more  
about this in  
Week 7

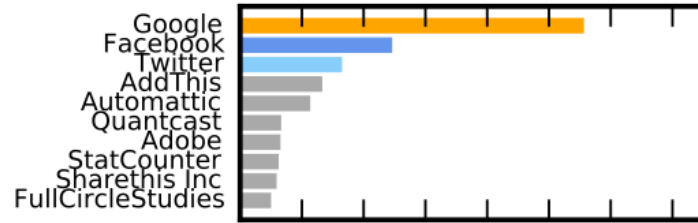




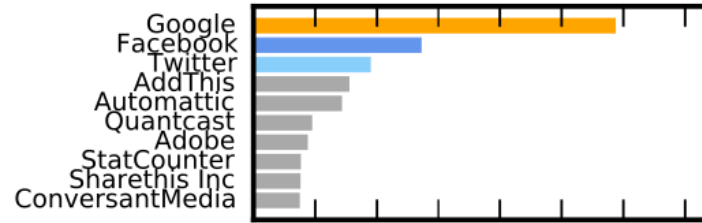
# Web Tracking is happening



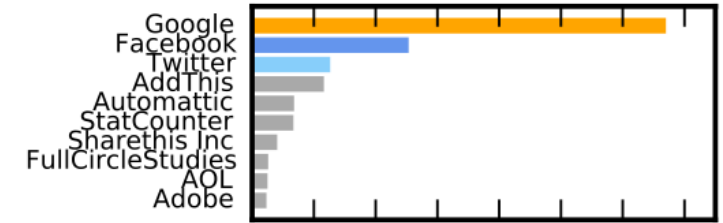
overall



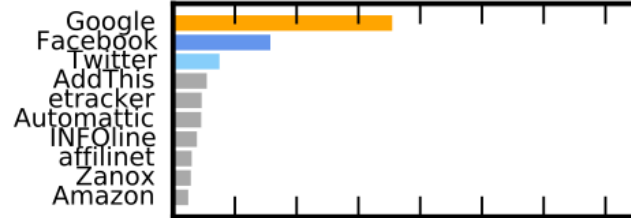
.com



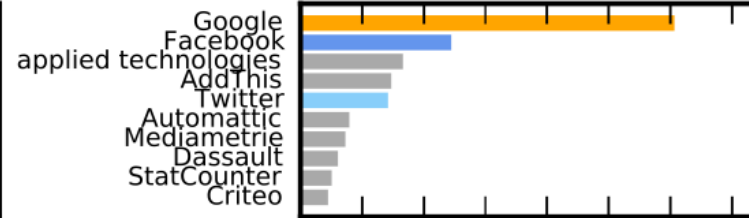
Ireland (.ie)



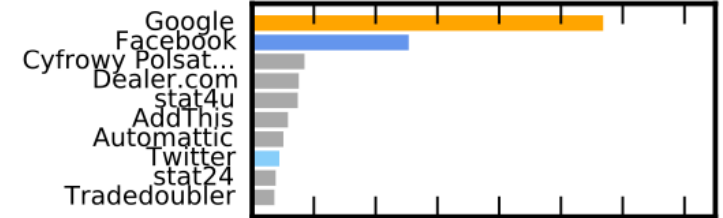
Germany (.de)



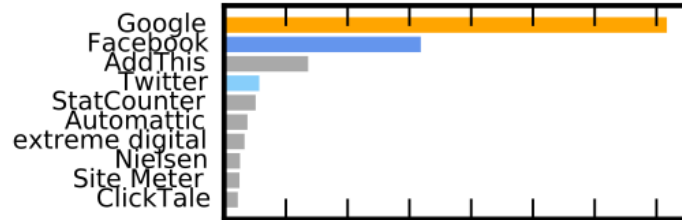
France (.fr)



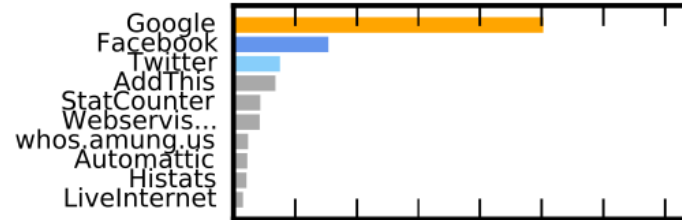
Poland (.pl)



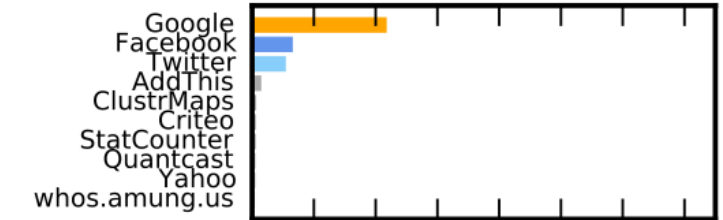
Israel (.il)



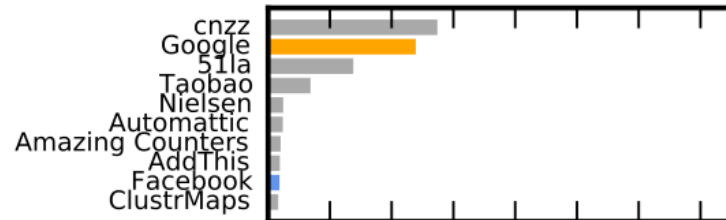
Turkey (.tr)



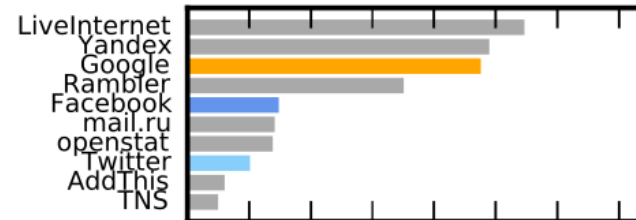
Korea (.kr)



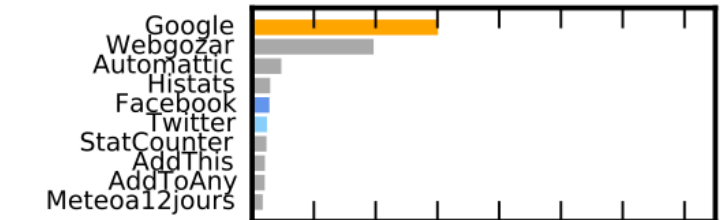
China (.cn)



Russia (.ru)



Iran (.ir)



0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7

rank share

0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7

rank share

0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7

rank share

# Evolution of Web Search

- Content (1st Generation)
- Links (2nd Generation)
- Personalisation (3rd Generation)

## 1st generation Web search

- Early 1990s
- Examples: Lycos, Altavista, AlltheWeb,...
- Ranking signals
  - Term frequency (TF)
  - Inverse document frequency (IDF)
  - TF-IDF

## 2nd generation Web search

- Take the link structure of the Web into account
- Second half of 1990s
- Examples: Google (PageRank), Ask! (HITS)
- Ranking signals
  - Website popularity

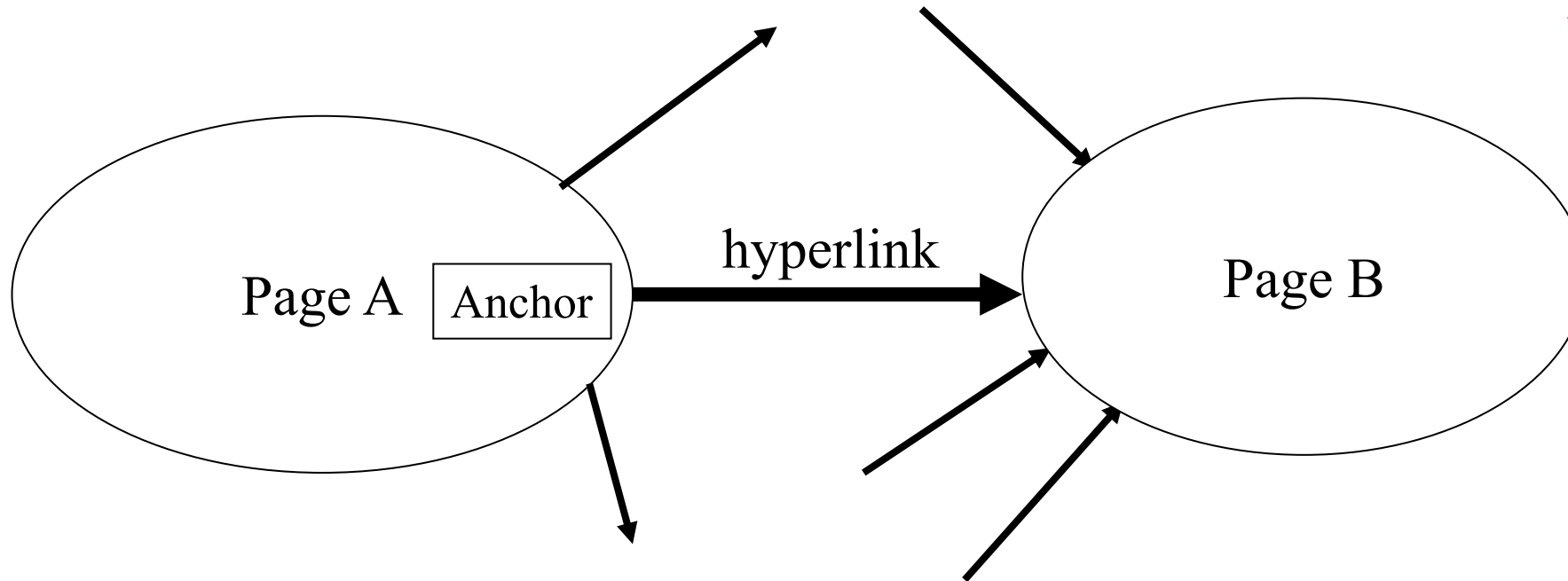
## 3rd generation Web search

- Provide search results tailored to the individual user
- 2004: Google introduces personalised search
- Ranking signals
  - Users' relevance feedback
  - Context

## ➤ 2. Web graph



# The Web as a directed graph

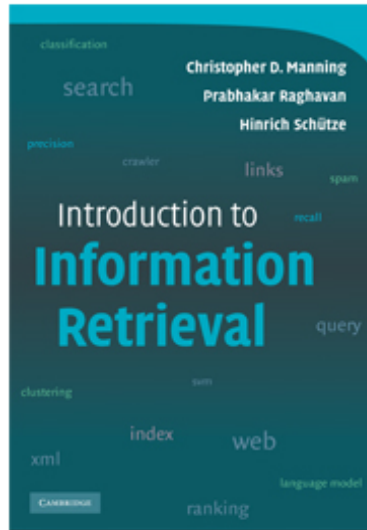


**Hypothesis 1:** A hyperlink between pages denotes a conferral of authority (quality signal)

**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*

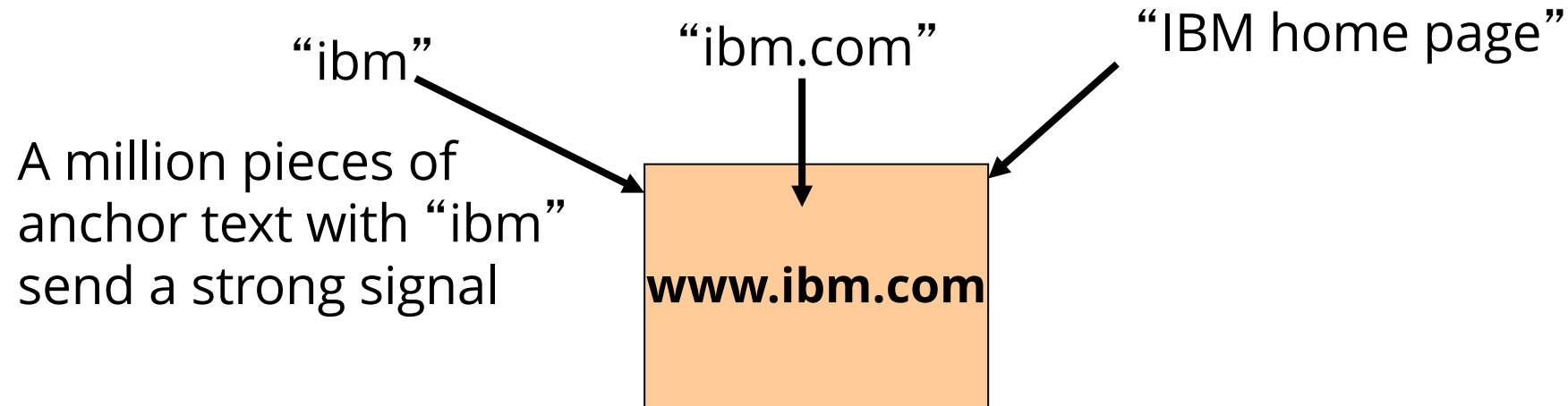
You can order this book at [CUP](#), at your local bookstore or on the internet. The best search

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

We'd be pleased to get feedback about how this book works out as a textbook, what is missing, and what you think. Please send comments to: [informationretrieval@yahoo.com](mailto:informationretrieval@yahoo.com)

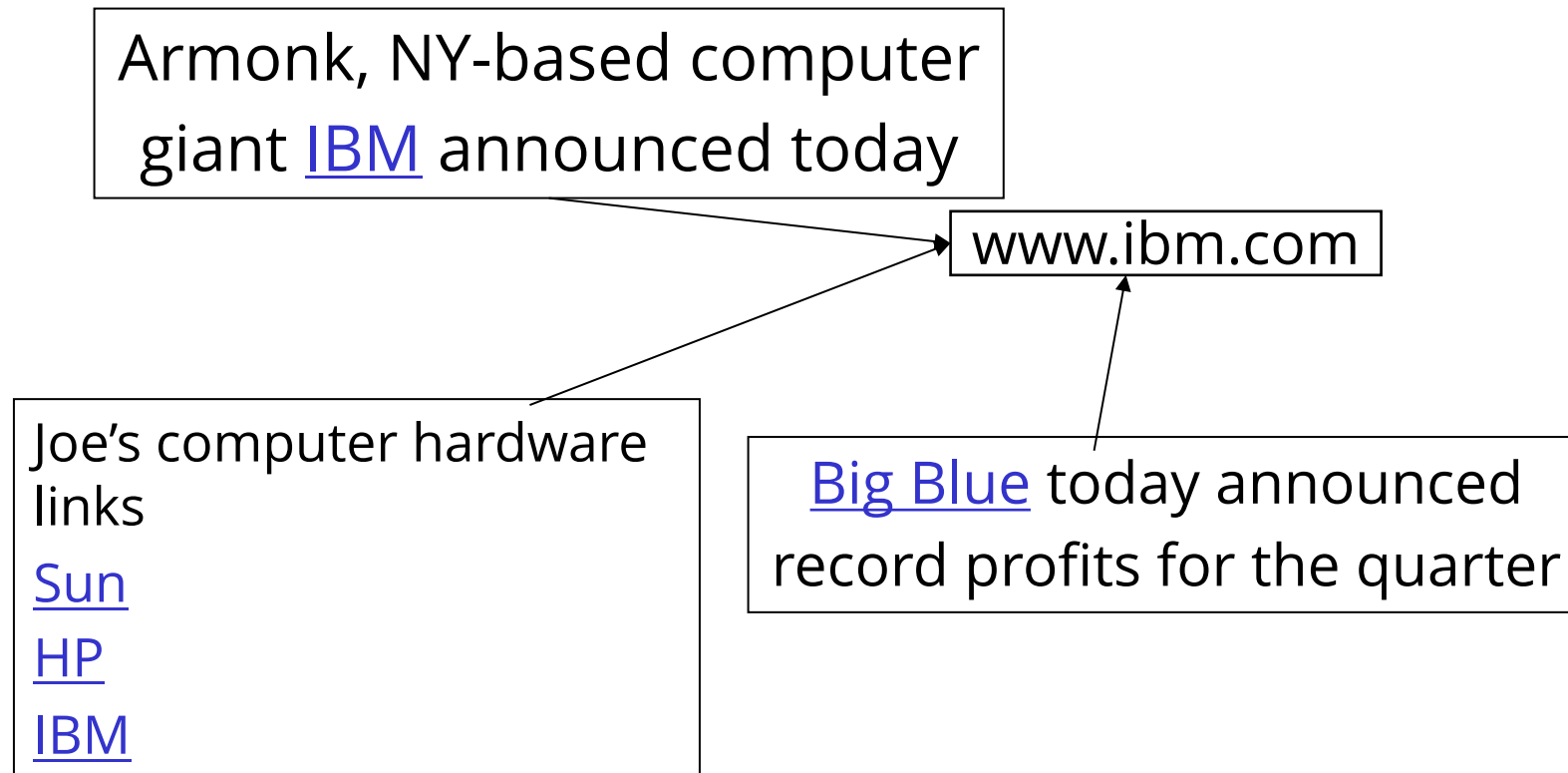
# Anchor text

- 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  $D$



# Indexing anchor text

- 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

- Initiative by Bing, Google, Yahoo!, Yandex, ...
- De-facto standard vocabulary for structured data on the Web
- Can be used to describe the meaning of websites

## ➤ 3. Personalisation



# Pros and Cons

- Saves time by reducing number of results to inspect
- Better decision making by filtering out inferior information
- Filter bubble (as much a personal decision as an algorithmic restriction)
- Users are products (using search history for advertisement)

## Personal details

- Information about the users
- Ranking signals
  - Language
    - Language preferences can be used to filter out results
  - Demographics
    - Usually predicted
    - Results selected by other users from similar cohorts can be ranked higher

## Social networks

- Information about a user's social network
- Ranking signals
  - Social network connections
    - Results selected by friends for similar searches could be given more weight
    - Web pages shared by friends could be given more weight

## Prior activities (query logs)

- Information about the queries submitted by the user and other users in the past
- Ranking signals
  - Query suggestion
    - Other users entered queries A and B in the same session -> B might be a good suggestion for a user entering Query A
  - Spellign correction
    - Immediately after query X other query Y -> Y might be the correct version of query X

## Context

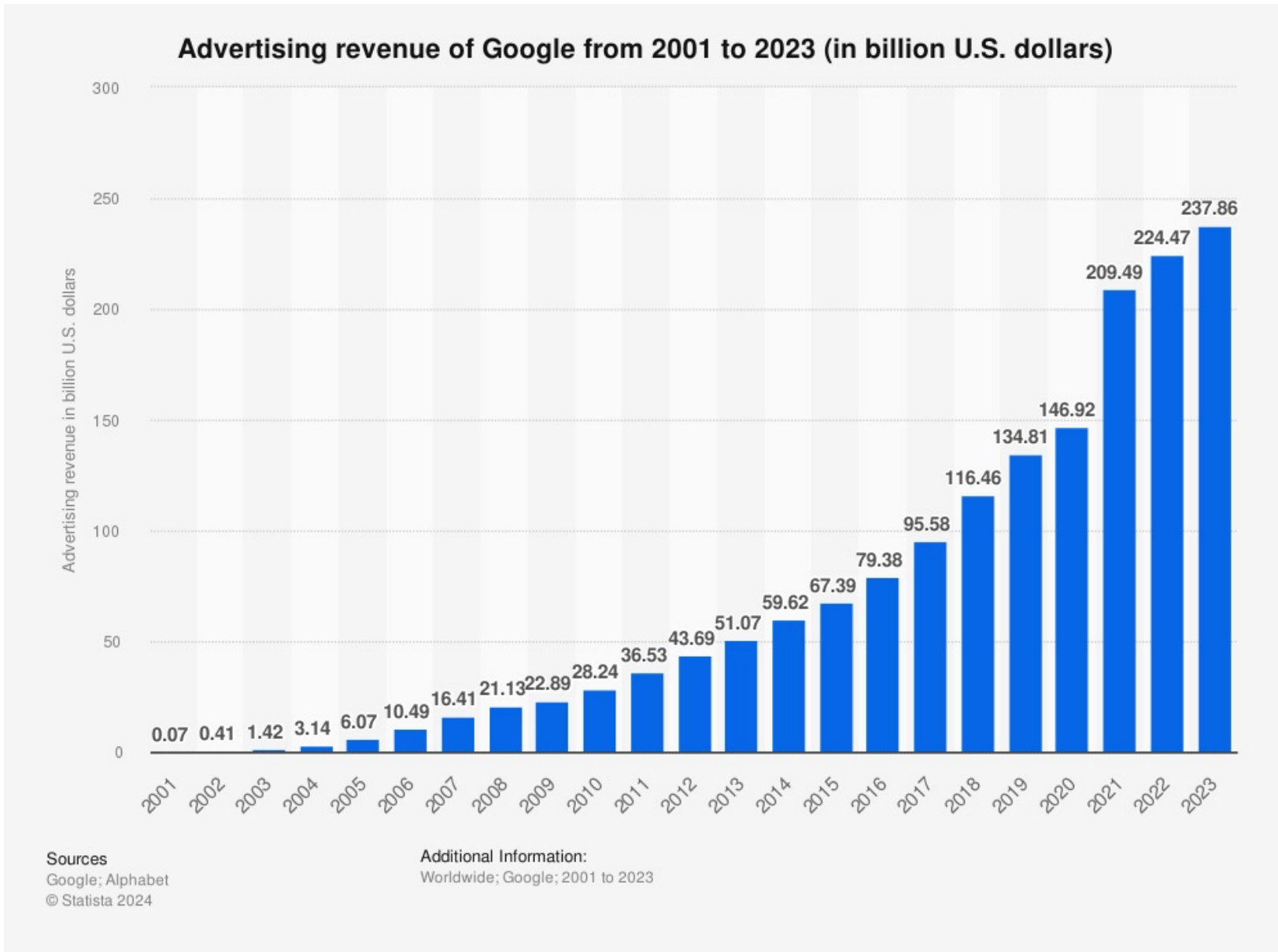
- Information about the context in which the search is performed
- Ranking signals
  - Location
    - Used to prioritise locally relevant results
    - Essential for mobile search
  - Date & Time
    - Seasonal influences, home vs. work, ....

- Learning the optimal combination of all ranking signals
- Goal: to do this continuously and automatically using machine learning
  - Predict for each query-result pair whether the result is relevant for that user's query at this specific time

## ➤ 4. Ads



# Advertising revenue of Google from 2001-2023



# Google Ads



PDF processing



[All](#)

[Images](#)

[News](#)

[Videos](#)

[Shopping](#)

[More](#)

[Settings](#)

[Tools](#)

About 1.220.000.000 results (0,44 seconds)

[\(Ad\)](#) [www.qoppa.com/](http://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](http://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](http://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](http://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**

Search = **web domain**

(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? ▼

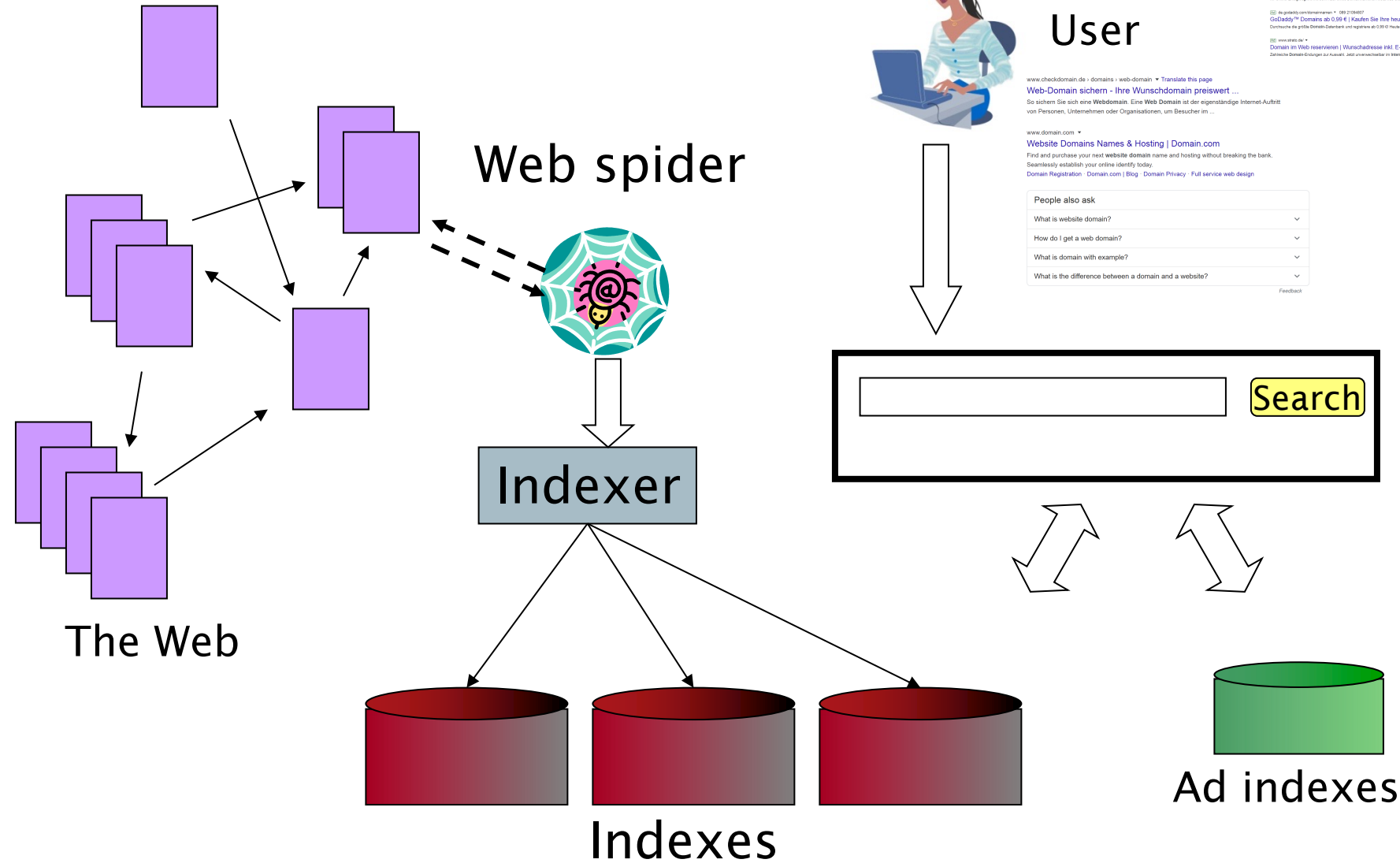
What is the difference between a domain and a website? ▼

[Feedback](#)

# Ads vs. search results

- 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

## ➤ 5. 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

# Simplest forms

- 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

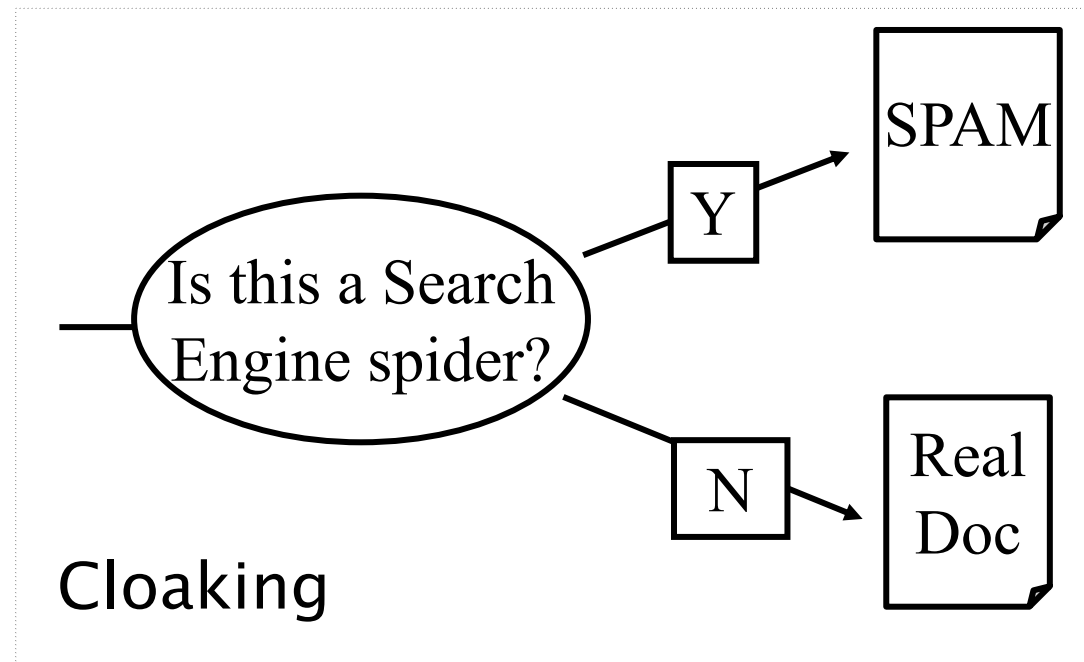
- Misleading meta-tags, excessive repetition
- Hidden text with colors, style sheet tricks, etc.

## 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](http://www.webmasterworld.com) )
    - Search engine specific tricks

# The spam industry

## 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](#)

Search Engine  
*Cloaker*

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](http://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](http://www.yourname.com) as they browsed.

Cloaking Services: Included Branded Email Services 5  
Mail boxes [mailboxname@yourDomain.com](mailto:mailboxname@yourDomain.com) \$49/Year



antomLine™ — the ultimate stealth

## Understanding Cloaking

al: Cloaking and Stealth Technology  
[Page 2](#) | [Page 3](#) | [Page 4](#) | [Page 5](#)

g, stealth or phantom page technology constitutes the  
h sophisticated and efficient approach towards search engine  
on. 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/>

## ➤ 6. Summary

# Intended Learning Outcomes

## You are now able to:

- Describe the structure of the Web
- Outline personalisation approaches
- Understand how ads are embedded on search result lists
- Describe spamming approaches to improve search ranking