

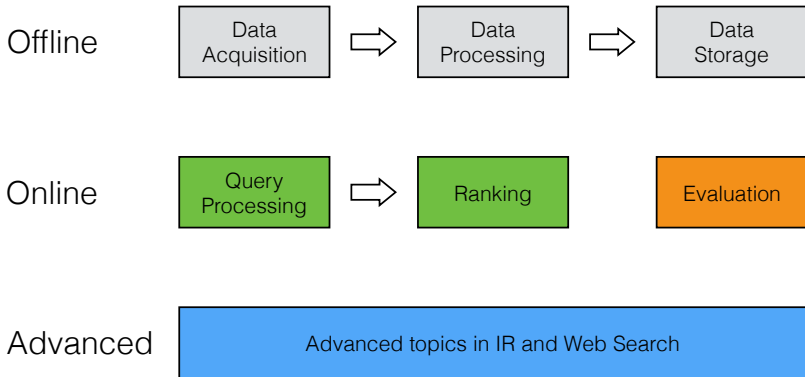
Information Retrieval

Data Acquisition

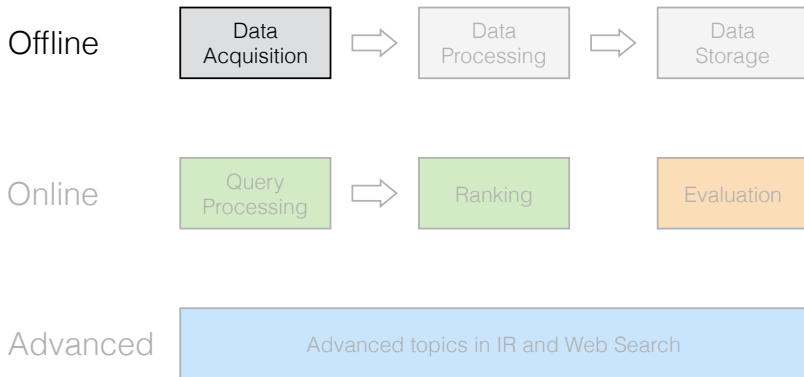
Ilya Markov
i.markov@uva.nl

University of Amsterdam

Course overview



This lecture



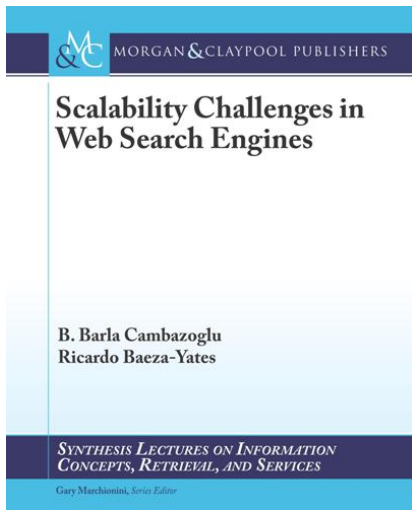
Data acquisition methods

- Downloading
- Feeds
- API
- Metadata harvesting
- **Crawling**

Outline

- 1 Crawling
- 2 Practical considerations
- 3 Duplicate detection
- 4 Spam
- 5 Summary

Scalability challenges in web search engines



Outline

- 1 **Crawling**
 - Basic architecture
 - Politeness
 - Extending the Web repository
 - Refreshing the Web repository
- 2 Practical considerations
- 3 Duplicate detection
- 4 Spam
- 5 Summary

Outline

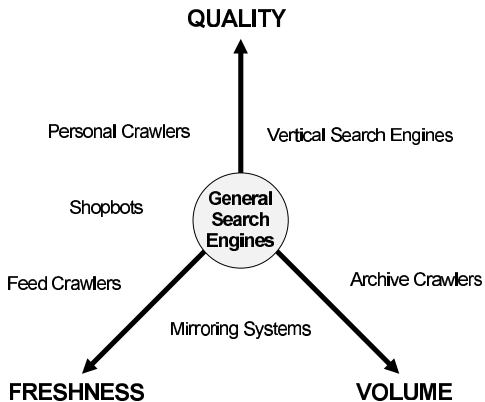
- 1 **Crawling**
 - Basic architecture
 - Politeness
 - Extending the Web repository
 - Refreshing the Web repository

Crawling

<https://staff.fnwi.uva.nl/i.markov/>

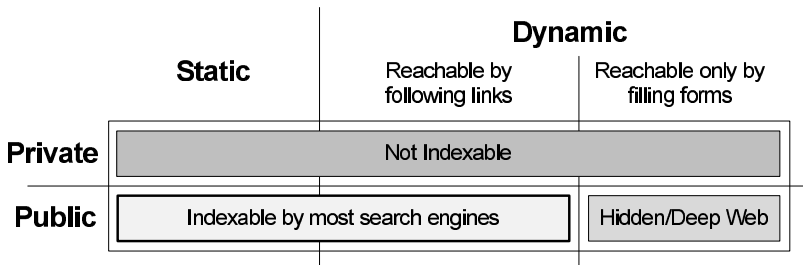


Taxonomy of crawlers



Baeza-Yates and Ribeiro-Neto, "Modern Information Retrieval"

Taxonomy of pages



Baeza-Yates and Ribeiro-Neto, "Modern Information Retrieval"

Outline

- 1 **Crawling**
 - Basic architecture
 - **Politeness**
 - Extending the Web repository
 - Refreshing the Web repository

How to be polite?

A Web crawler must. . .

- ① identify itself
- ② obey the robots exclusion protocol
- ③ keep a low bandwidth usage in a given web site

Robot identification

- Fill the `user-agent` field in the HTTP request
- Include the word “crawler”, “robot”, “bot”, etc.

Robot exclusion protocol

- Server-wise exclusion

`robots.txt`

- Page-wise exclusion

```
<META NAME="ROBOTS" CONTENT="NOINDEX, NOFOLLOW">
```

- Cache exclusion

```
<META NAME="ROBOTS" CONTENT="NOARCHIVE">
```

Bandwidth usage

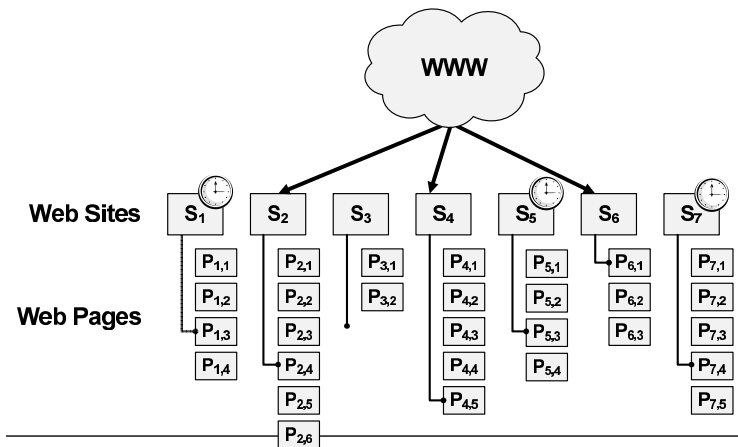
- Empirical thresholds
- Adaptive politeness policy, e.g., $10 \times t$
- `crawl-delay:45` (in seconds in `robots.txt`)

Simple crawling thread implementation

```
procedure CRAWLERTHREAD(frontier)
  while not frontier.done() do
    website ← frontier.nextSite()
    url ← website.nextURL()
    if website.permitsCrawl(url) then
      text ← retrieveURL(url)
      storeDocument(url, text)
      for each url in parse(text) do
        frontier.addURL(url)
      end for
    end if
    frontier.releaseSite(website)
  end while
end procedure
```

Croft et al., "Search Engines, Information Retrieval in Practice"

Simple crawling thread implementation

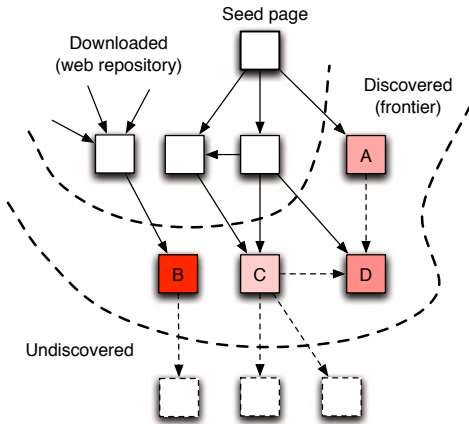


Baeza-Yates and Ribeiro-Neto, "Modern Information Retrieval"

Outline

- 1 **Crawling**
 - Basic architecture
 - Politeness
 - **Extending the Web repository**
 - Refreshing the Web repository

Extending the Web repository



- Random ordering
- Breadth-first
- In-degree
- Potential impact on search quality

B. Cambazoglu and R. Baeza-Yates, "Scalability Challenges in Web Search Engines"

Extending the Web repository

What data structures could be used to implement the frontier?

- Random ordering
- Breadth-first
- In-degree
- Potential impact on search quality

Outline

- 1 **Crawling**
 - Basic architecture
 - Politeness
 - Extending the Web repository
 - Refreshing the Web repository

Refreshing the Web repository

- ① Identify changes in web content
- ② Measure these changes
- ③ Predict changes
- ④ Select pages to update

Identifying changes

Client request: HEAD /csinfo/people.html HTTP/1.1
Host: www.cs.umass.edu

HTTP/1.1 200 OK

Date: Thu, 03 Apr 2008 05:17:54 GMT

Server: Apache/2.0.52 (CentOS)

Last-Modified: Fri, 04 Jan 2008 15:28:39 GMT

Server response: ETag: "239c33-2576-2a2837c0"

Accept-Ranges: bytes

Content-Length: 9590

Connection: close

Content-Type: text/html; charset=ISO-8859-1

Croft et al., "Search Engines, Information Retrieval in Practice"

Measuring changes

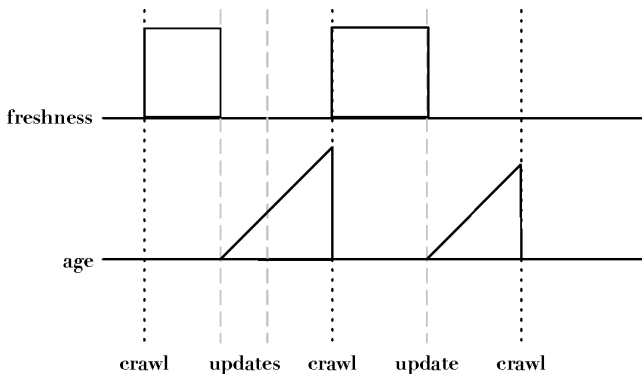
- Freshness

$$F_p(t) = \begin{cases} 1 & \text{if } p \text{ has not changed since last crawl} \\ 0 & \text{otherwise} \end{cases}$$

- Age

$$A_p(t) = \int_0^t P(p \text{ changed at time } x)(t - x)dx$$

Measuring changes



Croft et al., "Search Engines, Information Retrieval in Practice"

Predicting changes

- Age

$$A_p(t) = \int_0^t P(p \text{ changed at time } x)(t - x)dx$$

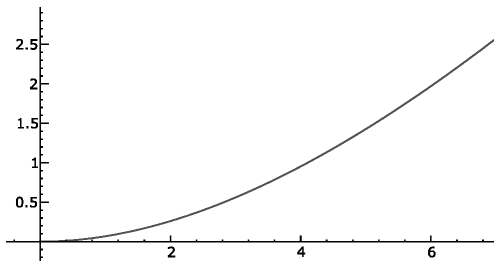
- Time between two changes follows an exponential distribution

$$A_p(t) = \int_0^t \lambda e^{-\lambda x}(t - x)dx$$

- λ – average number of changes

Example

- One update a week, $\lambda = 1/7$
- In the end of the week, the expected age is approx. 2.6



Croft et al., "Search Engines, Information Retrieval in Practice"

Selecting pages to update

- Optimizing freshness or age?
- Optimizing freshness
 - Pages with low change frequency do not need to be refreshed often
 - Pages with high change frequency are never fresh
 - Refresh pages with medium change frequency
- Optimizing age
 - The older a page gets, the more it costs not to crawl it

Selecting pages to update

	A	B	C	D
PageRank	0.0003	0.0007	0.0002	0.0001
Average daily click count	47	332	2	1974
Last download time	2 hours ago	1 day ago	8 days ago	6 hours ago
Estimated update frequency	daily	never	minutely	yearly

What data structures could be used to implement refreshing?

B. Cambazoglu and R. Baeza-Yates, "Scalability Challenges in Web Search Engines"

Practical implementation

Maintain several queues

- A queue for news sites that is refreshed several times a day
- A daily or weekly queue for popular or relevant sites
- A large queue for the rest of the Web

Outline

- 1 Crawling
- 2 **Practical considerations**
 - Storage and data structures
 - Distributed crawling
 - Factors affecting crawling performance
 - Deep web
- 3 Duplicate detection
- 4 Spam
- 5 Summary

Outline

- 2 Practical considerations
 - Storage and data structures
 - Distributed crawling
 - Factors affecting crawling performance
 - Deep web

Storing documents

- Support random and bulk access
- Flat file system
- A page is accessed through an identifier, e.g., a hash value of its URL

Storing document info

- Location on disk, size, history of HTTP status codes, download times, etc.
- Catalog-like data structures, e.g., B+ tree

Distributed storage

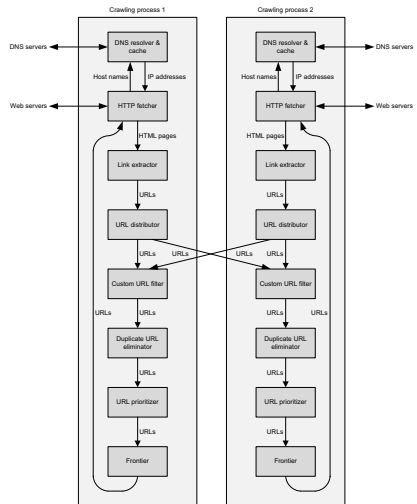
- Pages are distributed over storage nodes
- Uniformly or based on hash value ranges
- The mapping between pages and storage nodes is maintained in an index
- Distribution based on hash value ranges → smaller index
- Uniform distribution → simplified redistribution

Outline

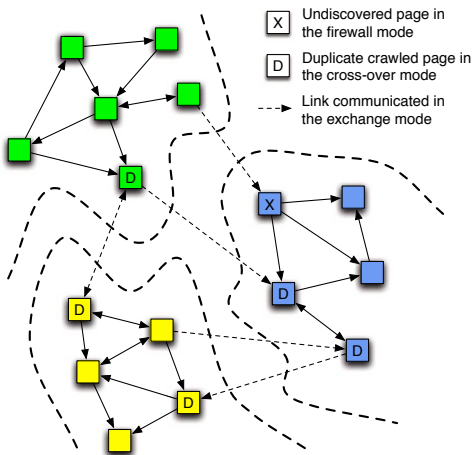
- 2 Practical considerations
 - Storage and data structures
 - Distributed crawling**
 - Factors affecting crawling performance
 - Deep web

Distributed crawling

- Synchronize URLs between threads
- Minimize communication overhead

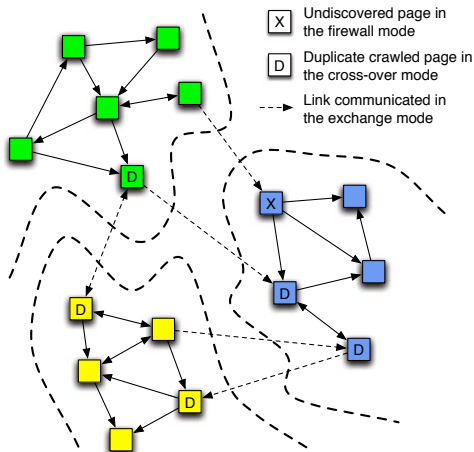


Firewall and cross-over modes



B. Cambazoglu and R. Baeza-Yates, "Scalability Challenges in Web Search Engines"

Exchange mode



- ① Partition pages between crawling nodes based on domains
- ② Exchange only non-local links
- ③ Send URLs in small batches

B. Cambazoglu and R. Baeza-Yates, "Scalability Challenges in Web Search Engines"

Outline

- 2 Practical considerations
 - Storage and data structures
 - Distributed crawling
 - Factors affecting crawling performance
 - Deep web

Factors affecting crawling performance

- Delay attack
- Spider trap
- Link farm
- Website mirroring
- Soft 404 error

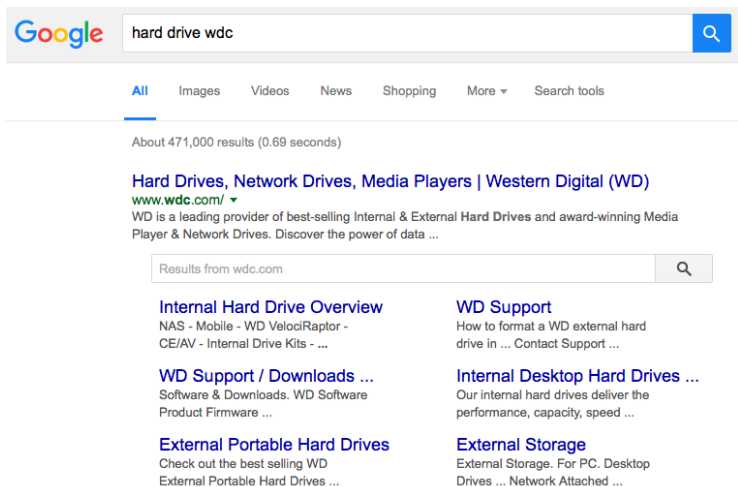
Outline


- ② Practical considerations
 - Storage and data structures
 - Distributed crawling
 - Factors affecting crawling performance
 - Deep web

Deep web

- No incoming links
- Password-protected
- Dynamic content
- Web forms

Deep web example




Google 

[All](#) [Images](#) [Videos](#) [News](#) [Shopping](#) [More ▾](#) [Search tools](#)


About 471,000 results (0.69 seconds)

Hard Drives, Network Drives, Media Players | Western Digital (WD)
www.wdc.com/ ▾
WD is a leading provider of best-selling Internal & External Hard Drives and award-winning Media Player & Network Drives. Discover the power of data ...

Results from wdc.com 

Internal Hard Drive Overview NAS - Mobile - WD VelociRaptor - CE/AV - Internal Drive Kits - ...	WD Support How to format a WD external hard drive in ... Contact Support ...
WD Support / Downloads ... Software & Downloads. WD Software Product Firmware ...	Internal Desktop Hard Drives ... Our internal hard drives deliver the performance, capacity, speed ...
External Portable Hard Drives Check out the best selling WD External Portable Hard Drives ...	External Storage External Storage. For PC. Desktop Drives ... Network Attached ...

Deep web example

hard drive site:wdc.com  

[All](#) [Shopping](#) [Images](#) [News](#) [Videos](#) [More](#) [Search tools](#)

About 15,000 results (0.42 seconds)

Hard Drives, Network Drives, Media Players | Western Digital (WD)

www.wdc.com/ • Western Digital •

WD is a leading provider of best-selling Internal & External Hard Drives and award-winning Media Player & Network Drives. Discover the power of data ...

Internal Hard Drive Overview - Western Digital

www.wdc.com/en/products/internal/ • Western Digital •

Our internal hard drives deliver the performance, capacity, speed and reliability that you expect from the best selling hard drive manufacturer in the world.

External Portable Hard Drives | Western Digital (WD)

www.wdc.com/en/products/external/portable/ • Western Digital •

Check out the best selling WD External Portable Hard Drives. Discover which type of data storage is right for you with an easy side-by-side comparison.

Internal Mobile Hard Drives Overview - Western Digital

www.wdc.com/en/products/internal/mobile/ • Western Digital •

Our internal hard drives deliver the performance, capacity, speed and reliability that you expect from the best selling hard drive and storage manufacturer in the ...

WD Elements - Portable Hard Drives

www.wdc.com/EN/PRODUCTS/products.aspx?id=470 • Western Digital •

WD Elements Portable Hard Drives. ... Passport Ultra Metal Edition · My Passport Ultra · WD Elements Portable · Compare All. For Mac, Desktop Drives for Mac.



The screenshot shows the top portion of the WD website. At the top left is the WD logo with the tagline "a Western Digital brand". To the right is a "Blog" link with a plus icon. Below the logo is a navigation menu with four items: "External Storage", "Internal Hard Drives", "Personal Cloud", and "Network Attached Storage". Underneath the menu is a breadcrumb trail that reads "WD Home /".

About 2,590 results (0.40 seconds)

[Internal Hard Drive Overview](#)

www.wdc.com/en/products/internal/

Our internal **hard drives** deliver the performance, capacity, speed and reliability that you expect from the world.

[External Portable Hard Drives | Western Digital \(WD\)](#)

www.wdc.com/en/products/external/portable/

Check out the best selling WD External Portable Hard Drives. Discover which type of data storage

[Internal Mobile Hard Drives Overview](#)

www.wdc.com/en/products/internal/mobile/

Our internal **hard drives** deliver the performance, capacity, speed and reliability that you expect from the manufacturer in the ...

[External Hard Drives | Western Digital \(WD\)](#)

www.wdc.com/en/products/external/desktop/

Check out the best selling WD External Hard Drives. Discover which type of data storage is rig

[WD Gold - Datacenter Hard Drives | Western Digital](#)

www.wdc.com/en/products/products.aspx?id=1670

WD Gold **hard drives** employ advanced technology to deliver among the best in reliability, capa

Outline

- 1 Crawling
- 2 Practical considerations
- 3 Duplicate detection**
- 4 Spam
- 5 Summary

Duplicates

- Mirroring
- Different URLs, same content
- URL modifiers, e.g., `/dir/page.html&jssid=09A89732`

Duplicate detection: Hashing

- ① Compute a hash value of a document
- ② Directly compare documents with the same hash value

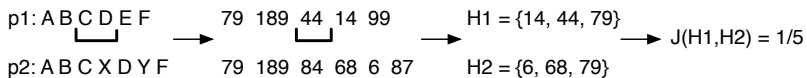
Near duplicate detection: Shingling

- ① Extract all unique (consecutive) sequences of n words from a document (n -grams/shingles)
- ② Compute hash values for the extracted shingles
- ③ For two documents compute the Jaccard coefficient

$$J(d_1, d_2) = \frac{|H(d_1) \cap H(d_2)|}{|H(d_1) \cup H(d_2)|}$$

- ④ Documents are near-duplicates if $J(d_1, d_2) > \textit{threshold}$
- ⑤ In practice, hashes are sorted and only the top- k hashes are considered to calculate $J(d_1, d_2)$

Near duplicate detection: Shingling



B. Cambazoglu and R. Baeza-Yates, "Scalability Challenges in Web Search Engines"

Outline

- 1 Crawling
- 2 Practical considerations
- 3 Duplicate detection
- 4 Spam**
- 5 Summary

Spam

- Cloacking, redirection spam
- Link spam
- Content spam
- Link spam

Cloaking, redirection spam

Serving crawlers and users with different content

Link spam

- Link farms
 - Buy a large number of domains
 - Create a large number of sites
 - Link them to each other
- Blog spam, comment spam, wiki spam
- Hidden links

Content spam

- Keyword stuffing
 - Raise the keyword count, variety, and density
- Hidden or invisible text
 - Unrelated text using same color as the background
 - Using tiny font size
 - Hiding within HTML code
- Doorway pages
 - Designed to rank highly within search results
 - Redirect users to spam content
- Scraper sites
 - “Scrape” other sources of content
 - Create “content” for a website

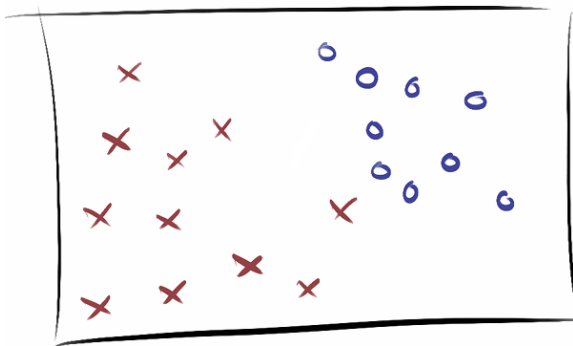
<https://en.wikipedia.org/wiki/Spamdexing>

Click spam

- Reasons
 - Make competitors pay for clicks
 - Compromise publishers
- Implementation
 - Many computers in different geographic locations
 - Use Trojan code on personal computers
 - Hit inflation attack

https://en.wikipedia.org/wiki/Click_fraud

Identifying spam



Picture taken from http://smerity.com/media/talks/ml_for_your_robotic_army/template.html

Outline

- 1 Crawling
- 2 Practical considerations
- 3 Duplicate detection
- 4 Spam
- 5 Summary**

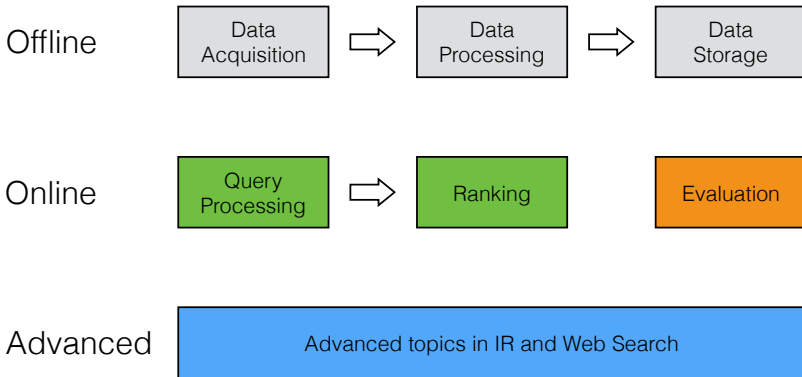
Data acquisition summary

- 1 **Crawling**
 - Basic architecture
 - Politeness
 - Extending the Web repository
 - Refreshing the Web repository
- 2 **Practical considerations**
 - Storage and data structures
 - Distributed crawling
 - Factors affecting crawling performance
 - Deep web
- 3 **Duplicate detection**
- 4 **Spam**
- 5 **Summary**

Materials

- Croft et al., Chapters 3.1–3.4, 9.1.5
- Baeza-Yates and Ribeiro-Neto, Chapter 12
- B. Barla Cambazoglu and Ricardo Baeza-Yates
Scalability Challenges in Web Search Engines
Morgan & Claypool Publishers, 2017

Course overview



Next lecture

