

CES Data Scientist

Web Crawling





Discovering new URLs

Identifying duplicates

Crawling architecture

Recrawling URLs?

Modern Crawling

Conclusion

13 June 2016





- The World Wide Web: **largest information repository**
- By nature (mostly) **public** and **freely accessible** information
- Numerous applications require **navigating** the Web and **retrieving Web content** in bulk:
 - **Web search engines** (e.g., Google, Bing), to create an index of Web content
 - **Web archiving** (e.g., Internet Archive), to preserve parts of the Web for future generations
 - **Shopping aggregators** (e.g., Google Shopping, Kelkoo), to build a database of products and their prices on different sites
 - etc.





https : // www.example.com / path/to/doc
scheme hostname path

scheme: way the resource can be accessed; generally **http** or **https**

hostname: **domain name** of a host; hostname of a website may start with **www.**, but not a rule

path: **logical path** of the document

- Unique identification of a resource (page, document, etc.) on the Web
- Empty path: root of the Web server
- Relative URLs with respect to a **context** (e.g., the URL above):

/titi `https://www.example.com/titi`

tata `https://www.example.com/path/to/tata`





- **crawlers, (Web) spiders, (Web) robots**: autonomous user agents that retrieve pages from the Web
- Basics of crawling:
 1. Start from a given URL or set of URLs
 2. Retrieve and process the corresponding page
 3. Discover new URLs (see further)
 4. Repeat on each found URL
- No real termination condition (the Web is **virtually infinite!**)
- **Graph-traversing** problem





A

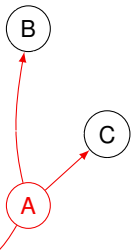
13 June 2016

6 / 39



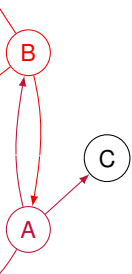
Licence de droits d'usage





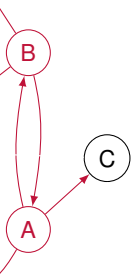
13 June 2016





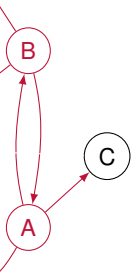
13 June 2016





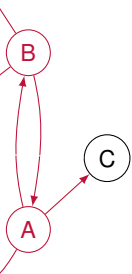
13 June 2016





13 June 2016





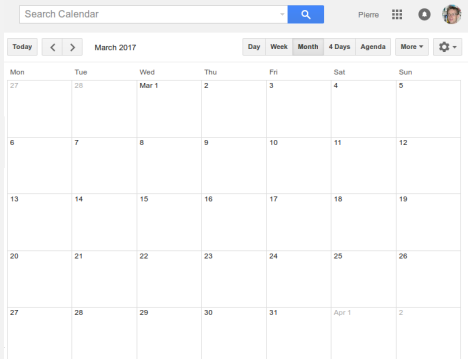
13 June 2016



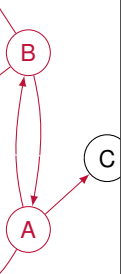


Characteristics

- Very susceptible to robot traps



- Ensures better temporal consistency [Spaniol et al., 2009]



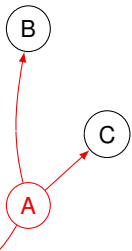


A

13 June 2016

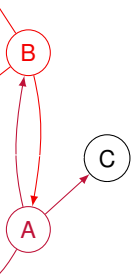


Licence de droits d'usage



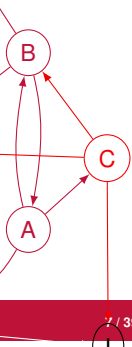
13 June 2016

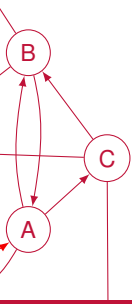


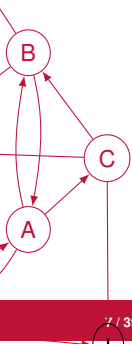


13 June 2016





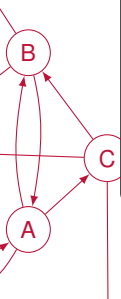






Characteristics

- **Better overall** for exploring a large part of the Web
- Risk of crawling **too superficially** given Web sites
- In practice, **pragmatic** strategy:
 - breadth-first crawl at the level of the whole Web
 - **limited-depth** depth-first crawl within a given Web site





from HTML pages:

- hyperlinks `...`
 - media `` `<embed src="...">`
`<object data="...">`
 - frames `<frame src="...">` `<iframe src="...">`
 - JavaScript links `window.open("...")`
 - etc.
- Other hyperlinked content (e.g., PDF files)
 - Non-hyperlinked URLs that appear anywhere on the Web (in HTML text, text files, etc.): use regular expressions to extract them
 - Sitemaps [sitemaps.org, 2008]





Web-scale

- The Web is infinite! Avoid robot traps by putting depth or page number **limits** on each Web server
- Focus on **important** pages [Abiteboul et al., 2003]
- Web servers under a list of **DNS domains**: easy filtering of URLs
- A given topic: **focused crawling** techniques [Chakrabarti et al., 1999, Diligenti et al., 2000] based on classifiers of Web page content and predictors of the interest of a link.
- The national Web (cf. **public deposit**, national libraries): what is this? [Abiteboul et al., 2002]
- A given Web site: what is a Web site? [Senellart, 2005]



Used to be plenty of free APIs to Web search engines... not the case any more

- **Paid-for** Web search APIs:

Yahoo! BOSS 0.80 USD per 1,000 queries (uses Bing's index)

<https://developer.yahoo.com/boss/search/>

Google Custom Search Engine 100 free queries per day, 5 USD per further 1,000 queries, up to 10,000 queries per day

<https://developers.google.com/custom-search/>

Bing Search API free for 5,000 queries per **month**; \approx 20 USD per further 5,000 queries

<https://datamarket.azure.com/dataset/5BA839F1-12CE-4CCE-BF57-A49D98D29A44>

- Anything else?





Discovering new URLs

Identifying duplicates

Crawling architecture

Recrawling URLs?

Modern Crawling

Conclusion

13 June 2016





A **hash function** is a deterministic mathematical function transforming objects (numbers, character strings, binary...) into fixed-size, seemingly random, numbers. The more random the transformation is, the better.

Example

Java hash function for the `String` class:

$$\sum_{i=0}^{n-1} s_i \times 31^{n-i-1} \bmod 2^{32}$$

where s_i is the (Unicode) code of character i of a string s .





Problem

Identifying duplicates or near-duplicates on the Web to prevent multiple indexing

trivial duplicates: same resource at the same **canonized** URL:

`http://example.com:80/toto`

`http://example.com/titi/../toto`

exact duplicates: identification by **hashing**

near-duplicates: (timestamps, tip of the day, etc.) more complex!





Edit distance. Count the **minimum number of basic modifications** (additions or deletions of characters or words, etc.) to obtain a document from another one. Good measure of similarity, and can be computed in $O(mn)$ where m and n are the size of the documents. But: **does not scale** to a large collection of documents (unreasonable to compute the edit distance for every pair!).

Shingles. Idea: two documents similar if they mostly share the same **succession of k -grams** (succession of tokens of length k).

Example

I like to watch the sun set with my friend.

My friend and I like to watch the sun set.

$S = \{i \text{ like, like to, my friend, set with, sun set, the sun, to watch, watch the, with my}\}$

$T = \{and \text{ i, friend and, i like, like to, my friend, sun set, the sun, to watch, watch the}\}$





Similarity: **Jaccard coefficient** on the set of shingles:

$$J(S, T) = \frac{|S \cap T|}{|S \cup T|}$$

- Still **costly to compute!** But can be approximated as follows:
 1. Choose N **different hash functions**
 2. For each hash function h_i and each set of shingles $S_k = \{s_{k1} \dots s_{kn}\}$, store $\phi_{ik} = \min_j h_i(s_{kj})$
 3. Approximate $J(S_k, S_l)$ as the **proportion** of ϕ_{ik} and ϕ_{il} that are equal
- Possibly to repeat in a hierarchical way with **super-shingles** (we are only interested in **very** similar documents)





Discovering new URLs

Identifying duplicates

Crawling architecture

Recrawling URLs?

Modern Crawling

Conclusion

13 June 2016





Standard for robot exclusion: **robots.txt** at the root of a Web server [Koster, 1994].

```
User-agent: *  
Allow: /searchhistory/  
Disallow: /search
```

- Per-page exclusion.

```
<meta name="ROBOTS" content="NOINDEX,NOFOLLOW">
```

- Per-link exclusion.

```
<a href="toto.html" rel="nofollow">Toto</a>
```

- Avoid **Denial Of Service** (DOS), wait ≈ 1 s between two repeated requests to the same Web server



General principles:

- It is to access or keep access to a “system for automated data processing” *in a fraudulent manner* is punished of two years of prison and 60,000 euros fine (Code pénal 323-1, modified by law 2015-912 on “Renseignement”)
- to disrupt the functioning of a “system for automated data processing” is punished of five years of prison and 150,000 euros fine, extended to seven years and 300,000 euros when the system is a public one containing personal information (Code pénal 323-2, modified by law 2015-912 on “Renseignement”)
- A Web site hosted in a different country may invoke completely different legal principles, under a different jurisdiction
- Crawling content can be considered accessing and keeping access to a “system for automated data processing” (Cour d’appel de Paris, 5 February 2014, “Bluetouff case”)





robots.txt files can be taken as a receivable way to specify what can be crawled (Cour d'appel de Paris, 26 January 2011, Google vs SAIF)

- Frequent requests to a Web site can be considered as a way to disrupt the functioning of a “system for automated data processing” (Cour d'appel de Bordeaux, 15 November 2011, Cédric M. vs C-Discount), but only if it reaches abusive levels and can be shown to have cause disruption
- Web content is subject to “droit d’auteur” (Code de la propriété intellectuelle, Première partie, Livre Ier) and cannot generally be broadcast by third-parties; only transient copies are allowed (CJEU, 5 June 2014, PRCA vs NLA)
- Web content containing personal data is even more sensitive (Loi “Informatique et Libertés”): personal data should be collected for a specific purpose, and should be kept updated

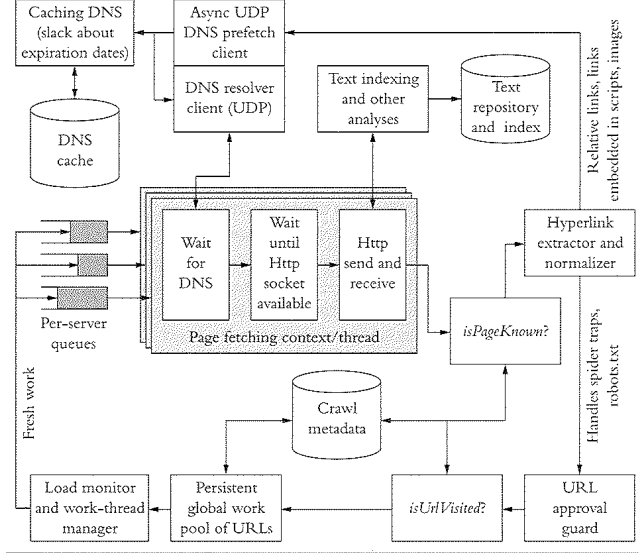




Network delays, waits between requests:

- **Per-server queue** of URLs
- Parallel processing of requests to different hosts:
 - **multi-threaded** programming
 - **asynchronous** inputs and outputs (`select`, classes from `java.util.concurrent`): less overhead
- Many engineering tricks (HTTP keep-alive, DNS pre-caching, etc.) to optimize a crawler and reduce connection overheads







Discovering new URLs

Identifying duplicates

Crawling architecture

Recrawling URLs?

Modern Crawling

Conclusion

13 June 2016





- Content on the Web **changes**
- Different **change rates**:
 - online newspaper main page: every hour or so
 - published article: virtually no change
- **Continuous** crawling, and identification of change rates for **adaptive** crawling: how to know the **time of last modification** of a Web page?





1. Check HTTP timestamp.
2. Check content timestamp.
3. Compare a hash of the page with a stored hash.
4. Non-significant differences (ads, fortunes, request timestamp):
 - only hash text content, or “useful” text content;
 - compare distribution of n -grams (shingling);
 - or even compute edit distance with previous version.

Adapting strategy to each different archived website?





g new URLs

Identifying duplicates

Crawling architecture

Recrawling URLs?

Modern Crawling

Modern Web Sites

Social Networking Sites

Conclusion

13 June 2016





g new URLs

Identifying duplicates

Crawling architecture

Recrawling URLs?

Modern Crawling

Modern Web Sites

Social Networking Sites

Conclusion

13 June 2016





- Some modern Web sites only work when cookies are activated (**session cookies**), or when **JavaScript code** is interpreted
- Regular Web crawlers (**wget**, **Heritrix**, **Apache Nutch**) do not usually perform any cookie management and do not interpret JavaScript code
- Crawling of some Websites therefore require more **advanced tools**





Web scraping frameworks such as **scrapy** (Python) or **WWW::Mechanize** (Perl) simulate a Web browser interaction and cookie management (but no JS interpretation)

Headless browsers such as **htmlunit** simulate a Web browser, including simple JavaScript processing

Browser instrumentors such as **Selenium** allow full instrumentation of a regular Web browser (Chrome, Firefox, Internet Explorer)

OXPath: a **full-fledged navigation and extraction language** for complex Web sites [Sellers et al., 2011]





g new URLs

Identifying duplicates

Crawling architecture

Recrawling URLs?

Modern Crawling

Modern Web Sites

Social Networking Sites

Conclusion

13 June 2016



- 1 google.com
- 2 facebook.com
- 3 yahoo.com
- 4 yahoo.com
- 5 baidu.com
- 6 wikipedia.org
- 7 live.com
- 8 twitter.com
- 9 qq.com
- 10 amazon.com
- 11 blogspot.com
- 12 linkedin.com
- 13 google.co.in
- 14 taobao.com
- 15 sina.com.cn
- 16 yahoo.co.jp
- 17 msn.com
- 18 wordpress.com
- 19 google.com.hk
- 20 t.co
- 21 google.de
- 22 ebay.com
- 23 google.co.jp
- 24 googleusercontent.com
- 25 google.co.uk
- 26 yandex.ru
- 27 163.com
- 28 weibo.com

(Alexa)

13 June 2016



1 google.com
2 facebook.com
3 yahoo.com
4 yahoo.com
5 baidu.com
6 wikipedia.org
7 live.com
8 twitter.com
9 qq.com
10 amazon.com
11 blogspot.com
12 linkedin.com
13 google.co.in
14 taobao.com
15 sina.com.cn
16 yahoo.co.jp
17 msn.com
18 wordpress.com
19 google.com.hk
20 t.co
21 google.de
22 ebay.com
23 google.co.jp
24 googleusercontent.com
25 google.co.uk
26 yandex.ru
27 163.com
28 weibo.com

(Alexa)

Social networking sites



1 google.com
2 facebook.com
3 wikipedia.org
4 baidu.com
5 wikipedia.org
6 wikipedia.org
7 live.com
8 twitter.com
9 qq.com
10 amazon.com
11 blogspot.com
12 linkedin.com
13 google.co.in
14 taobao.com
15 sina.com.cn
16 yahoo.co.jp
17 msn.com
18 wordpress.com
19 google.com.hk
20 t.co
21 google.de
22 ebay.com
23 google.co.jp
24 googleusercontent.com
25 google.co.uk
26 yandex.ru
27 163.com
28 weibo.com

(Alexa)

Social networking sites

Sites with social networking features (friends, user-shared content, user profiles, etc.)





Huge numbers of users (2012):

Facebook 900 million

QQ 540 million

W. Live 330 million

Weibo 310 million

Google+ 170 million

Twitter 140 million

LinkedIn 100 million





Huge numbers of users
(2012):

Facebook 900 million

QQ 540 million

W. Live 330 million

Weibo 310 million

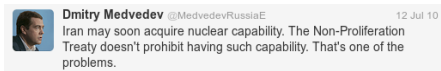
Google+ 170 million

Twitter 140 million

LinkedIn 100 million

Huge volume of shared data:
250 million tweets per day on Twitter
(3,000 per second on average!) . . .

. . . including statements by heads of
states, revelations of political activists, etc.





Theoretically possible to crawl social networking sites using a regular Web crawler

- Sometimes not possible:

`https://www.facebook.com/robots.txt`

- Often **very inefficient**, considering politeness constraints

- Better solution: Use provided social networking APIs

`https://dev.twitter.com/docs/api/1.1`

`https://developers.facebook.com/docs/graph-api/reference/v2.1/`

`https://developer.linkedin.com/apis`

`https://developers.google.com/youtube/v3/`

- Also possible to buy access to the data, directly from the social network or from brokers such as `http://gnip.com/`





- Most social networking Web sites (and some other kinds of Web sites) provide **APIs** to effectively access their content
- Usually a **RESTful** API, occasionally SOAP-based
- Usually require a **token** identifying the application using the API, sometimes a cryptographic signature as well
- May access the API as an authenticated user of the social network, or as an **external party**
- APIs seriously limit the **rate of requests**:
`https://dev.twitter.com/docs/api/1.1/get/search/tweets`





- Mode of interaction with a **Web service**
- Follow the KISS (**Keep it Simple, Stupid**) principle
- Each request to the service is a **simple HTTP GET method**
- Base URL is the **URL of the service**
- Parameters of the service are sent as **HTTP parameters** (in the URL)
- **HTTP response code** indicates success or failure
- Response contains **structured output**, usually as JSON or XML
- **No side effect**, each request independent of previous ones





- Two main APIs:
 - **REST APIs**, including search, getting information about a user, a list, followers, etc. <https://dev.twitter.com/docs/api/1.1>
 - **Streaming API**, providing real-time result
- **Very limited history** available
- Search can be on **keywords**, **language**, **geolocation** (for a small portion of tweets)





- Often useful to combine results from **different social networks**
- Numerous libraries facilitating SN API accesses (twipy, Facebook4J, FourSquare VP C++ API...) **incompatible with each other**... Some efforts at generic APIs (OneAll, APIBlender [Gouriten and Senellart, 2012])
- **Example use case:** No API to get all check-ins from FourSquare, but a number of check-ins are available on Twitter; given results of Twitter Search/Streaming, use FourSquare API to get information about check-in locations.





Discovering new URLs

Identifying duplicates

Crawling architecture

Recrawling URLs?

Modern Crawling

Conclusion

13 June 2016



What you should remember

- Crawling as a **graph-browsing** problem.
- **Shingling** for identifying duplicates.
- Numerous **engineering issues** in building a Web-scale crawler.
- Crawling modern Web content is **not as easy** as launching a traditional Web crawler
- Often critical to **focus the crawl** towards content of interest
- Ideally: a traditional large-scale crawler that knows **when to delegate** to more specialized crawling mechanisms (tools querying social networking APIs, deep Web crawlers, JS-aware crawlers, etc.)
- Huge variety of tools, techniques, suitable for different needs





- Wget, a simple yet effective Web spider (free software)
- Heritrix, a Web-scale highly configurable Web crawler, used by the Internet Archive (free software)
- HTML Parser, TagSoup: Java libraries for parsing real-world Web pages

To go further

- A good textbook [Chakrabarti, 2003]
- Main references:
 - HTML 4.01 recommendation [W3C, 1999]
 - HTTP/1.1 RFC [IETF, 1999]



Serge Abiteboul, Grégory Cobena, Julien Masanès, and Gerald Sedrati. A first experience in archiving the French Web. In *Proc. ECDL*, Roma, Italie, September 2002.

Serge Abiteboul, Mihai Preda, and Gregory Cobena. Adaptive on-line page importance computation. In *Proc. WWW*, May 2003.

Soumen Chakrabarti. *Mining the Web: Discovering Knowledge from Hypertext Data*. Morgan Kaufmann, San Fransisco, USA, 2003.

Soumen Chakrabarti, Martin van den Berg, and Byron Dom. Focused crawling: A new approach to topic-specific Web resource discovery. *Computer Networks*, 31(11–16):1623–1640, 1999.

Michelangelo Diligenti, Frans Coetzee, Steve Lawrence, C. Lee Giles, and Marco Gori. Focused crawling using context graphs. In *Proc. VLDB*, Cairo, Egypt, September 2000.

Georges Gouriten and Pierre Senellart. API Blender: A uniform interface to social platform APIs. In *Proc. WWW*, Lyon, France, April 2012. Developer track.

ietf. Request For Comments 2616. Hypertext transfer protocol—HTTP/1.1. <http://www.ietf.org/rfc/rfc2616.txt>, June 1999.

Martijn Koster. A standard for robot exclusion. <http://www.robotstxt.org/orig.html>, June 1994.

Andrew Sellers, Tim Furche, Georg Gottlob, Giovanni Grasso, and Christian Schallhart. Exploring the Web with OXPath. In *LWDM*, 2011.

Pierre Senellart. Identifying Websites with flow simulation. In *Proc. ICWE*, pages 124–129, Sydney, Australia, July 2005.

sitemaps.org. Sitemaps XML format. <http://www.sitemaps.org/protocol.php>, February 2008.

Marc Spaniol, Dimitar Denev, Arturas Mazeika, Pierre Senellart, and Gerhard Weikum. Data quality in Web archiving. In *Proc. WICOW*, pages 19–26, Madrid, Spain, April 2009.

W3C. HTML 4.01 specification, September 1999. <http://www.w3.org/TR/REC-html40/>.



Par le téléchargement ou la consultation de ce document, l'utilisateur accepte la licence d'utilisation qui y est attachée, telle que détaillée dans les dispositions suivantes, et s'engage à la respecter intégralement.

La licence confère à l'utilisateur un droit d'usage sur le document consulté ou téléchargé, totalement ou en partie, dans les conditions définies ci-après et à l'exclusion expresse de toute utilisation commerciale.

Le droit d'usage défini par la licence autorise un usage à destination de tout public qui comprend :

- le droit de reproduire tout ou partie du document sur support informatique ou papier,
- le droit de diffuser tout ou partie du document au public sur support papier ou informatique, y compris par la mise à la disposition du public sur un réseau numérique,
- le droit de modifier la forme ou la présentation du document,
- le droit d'intégrer tout ou partie du document dans un document composite et de le diffuser dans ce nouveau document, à condition que :
 - L'auteur soit informé.

Les mentions relatives à la source du document et/ou à son auteur doivent être conservées dans leur intégralité.

Le droit d'usage défini par la licence est personnel et non exclusif.

Tout autre usage que ceux prévus par la licence est soumis à autorisation préalable et expresse de l'auteur : sitepedago@telecom-paristech.fr

