# Archiving Web Content

Jean-Christophe Peyssard

HAL Id: **cel-02130558**

**https://shs.hal.science/cel-02130558**

Submitted on 15 May 2019

# ARCHIVING WEB CONTENT

> DHIB. 2019

Jean-Christophe Peyssard

Head of Digital Humanities

Institut français du Proche-Orient (Ifpo)

jc.peyssard@ifporient.org

Research Ingineer from CNRS at the French Institute of the Near East, Head of Digital Humanities

- Ifpo : http://www.ifporient.org/jean-christophe-peyssard/
- Linkedin : http://fr.linkedin.com/pub/jean-christophe-peyssard/22/705/782/
- ORCID : 0000-0002-8503-2217
- Google Scholar : https://scholar.google.com/citations?user=32cZHPsAAAAJ
- HAL : https://cv.archives-ouvertes.fr/jcpeyssard



Institut français du Proche-Orient (Ifpo) MEAE / CNRS – UMIFRE 6 / USR 3135

# Ifpo : Institut français du Proche-Orient / French Institute of the Near East



**http://www.ifporient.org/**

27 IFRE

34 pays

150 chercheurs
350 doctorants

800 manifestations scientifiques par an

+ 45 000 documents scientifiques dans le moteur de recherche

180 partenariats locaux

Un réseau unique de 27 centres de recherche français, répartis sur tous les continents, abordant toutes les sciences humaines et sociales, en contact avec les institutions de recherche locales et françaises.

http://www.ifre.fr/

# French research infrastructures for Scholarly Communication & Digital Humanities



https://www.huma-num.fr/
https://www.openedition.org/
http://persee.fr/
https://www.ccsd.cnrs.fr/

https://shakk.hypotheses.org/

# Three main references for this workshop



- Musiani, F, Paloque-Bergès, C., Schafer, V., & G. Thierry, B. (2019). *Qu'est-ce qu'une archive du web ?* Retrieved from http://books.openedition.org/oep/8713 ; DOI : 10.4000/books.oep.8713

- Brügger, N. (2018). *The Archived Web*. Retrieved from http://www.worldcat.org/oclc/1057466800

- Nielsen, J. (2016). *Using web archives in research: an introduction*. Retrieved from http://www.worldcat.org/oclc/960018046

- "The first-ever website (info.cern.ch) was published on August 6, 1991 by British physicist Tim Berners-Lee while at CERN, in Switzerland"

- "There are over 1.5 billion websites on the World Wide Web today. Of these, less than 200 million are active" (http://www.internetlivestats.com)

- "In the April 2019 survey we received responses from 1,445,266,139 sites across 233,886,577 unique domain names and 8,613,630 web-facing computers. **This reflects a loss of 16.8 million sites** [from previous March 2019 survey]" https://news.netcraft.com/archives/2019/04/22/april-2019-web-server-survey.html)

- "The Indexed Web contains at least 5.22 billion pages" (Tuesday, 30 April, 2019, https://www.worldwidewebsize.com/).

- "The average life span of a Web page is only 44 days, and 44 percent of the Web sites found in 1998 could not be found in 1999. […] As ubiquitous as the Web seems to be, it is also ephemeral, and much of today's Web will have disappeared by tomorrow." (Lyman, 2002 p. 38)

- "40% of the material on the Internet disappears within a year, while another 40% has been changed, which is why today we can only expect to find 20% of the material that was on the Internet one year ago." (Brügger, 2005 p. 15)

- "We now know that Web pages only last about 100 days on average before they change or disappear." (Kahle, 2015)

- In 2013, the average life span of a URL is 9.3 years (Musiani at al., 2019, https://books.openedition.org/oep/8743)

- In 2019, according to the Wayback Machine team the life span of a Web page is 92 days

**The New York Times**

## *In Supreme Court Opinions, Web Links to Nowhere*

By Adam Liptak

Sept. 23, 2013

WASHINGTON — Supreme Court opinions have come down with a bad case of link rot. According to a new study, 49 percent of the hyperlinks in Supreme Court decisions no longer work.

This can sometimes be amusing. A link in a 2011 Supreme Court opinion about violent video games by Justice Samuel A. Alito Jr. now leads to a mischievous error message.

"Aren't you glad you didn't cite to this Web page?" it asks. "If you had, like Justice Alito did, the original content would have long since disappeared and someone else might have come along and purchased the domain in order to make a comment about the transience of linked information in the Internet age."



**Conservatives**

### Conservative party deletes archive of speeches from internet

**Decade's worth of records is erased, including PM's speech praising internet for making more information available**

Randeep Ramesh and Alex Hern
Wed 13 Nov 2013 15.40 GMT

▲ A speech in which David Cameron said the internet would help people hold politicians to account was among those deleted. Photograph: Barcroft Media

The Conservatives have removed a decade of speeches from their website and from the main internet library - including one in which David Cameron claimed that being able to search the web would democratise politics by making "more information available to more people".

https://www.theguardian.com/politics/2013/nov/13/conservative-party-archive-speeches-internet

https://www.nytimes.com/2013/09/24/us/politics/in-supreme-court-opinions-clicks-that-lead-nowhere.html

# A problem for research and scholarly communication

- Error 404, Broken links, Link rot, Reference rot, Infosuicide, digital ruins, content drift, zombie media*,..*

- Shut down & take down, mergers and acquisitions:

    On March 18, 2019, it was revealed that MySpace lost all of their user content from 2016 and earlier in "a server migration gone wrong". It was widely reported that over 50 million songs and 12 years worth of content was permanently lost, and there was no backup (https://en.wikipedia.org/wiki/Myspace)

- History :

    - Yougoslavia (.yu) breakup (now Serbia and Montenegro, .rs and .me)
    - Czechoslovakia (.cs) dissolution (now Czech Republic and Slovakia, .cz and .sk)

- **Reference rot**, a combination of:

    - **Content decay**: The content of the linked resource may change over time and, as a result, the degree to which that content remains representative of the content that was intended to be linked to may decrease over time.

- **Link rot**: The linked resource may disappear altogether. (Thoughts on Referencing, Linking, Reference Rot http://mementoweb.org/missing-link/)

- The integrity of research is at risk ! (James G. Neal, http://library.ifla.org/id/eprint/907)

# Questions

- What is your interest in Web archiving?

- Did you ever experience a Web content loss?

- What is your experience in Web archiving?

# E-corpus 2009-2016



**https://web.archive.org/web/20180103235239/http://www.e-corpus.org/**

- To maintain our digital cultural heritage
- To stabilize and preserve web materials as a research object
- To be able to document and illustrate claims based on analyses of web materials (whether the web itself is the research object or a source of knowledge about other research objects).

Nielsen, J. (2016). *Using web archives in research: an introduction*. Retrieved from http://www.worldcat.org/oclc/960018046 p. 7

- Web archiving is the process of collecting portions of the World Wide Web, preserving the collections in an archival format, and then serving the archives for access and use. (IIPC Web site, http://netpreserve.org/web-archiving/)

- "Web archiving is the process of gathering up data that has been recorded on the World Wide Web, storing it, ensuring the data is preserved in an archive, and making the collected data available for future research." (Niu, J. (2012). An Overview of Web Archiving. *D-Lib Magazine*, *18* (3/4). https://doi.org/10.1045/march2012-niu1)

## A short chronology of Web archives

**1537** Legal deposit in France (1619 Spain, 1710 United Kingdom)
**1989** The World Wide Web was invented by Tim Berners-Lee
**1996** Internet Archive is founded by Brewster Kahle
**1996** Kulturarw3 in Sweden for archiving the .se top level domain name
**1998** Google is launched
**2001** The Wayback Machine (Internet Archive)
**2003** UNESCO Charter on the Preservation of Digital Heritage
**2003** the International Internet Preservation Consortium is formally chartered at the National Library of France with 12 participating institution
**2005** Youtube is launched
**2006** [The] Facebook and Twitter are launched
**2006** National Library of France is in charge of collecting and preserving the "French Internet" (new French Copyright Law) alongside with the National Audiovisual Institute (Ina)
**2013** EU Web Archive

# Different strategies & methods for Web archiving

- Macro archiving
- Micro archiving
- Thematic or selective archiving
- Bulk or snapshot harvesting, or broad crawls
- "Exhaustivity" of a National domain name archiving (.se in Sweden)
- Event and 'real-time' institutional archiving (after 2015 terrorist attacks and Notre Dame fire in 2019 France)
- Shared archiving among institutions (in France btwn BnF and Ina)
- …

- Finnish Web Archive (since 2006)
  https://www.kansalliskirjasto.fi/en/collections-and-content-online#finnish-web-archive

  The contents of the archive can be only accessed from special legal deposit workstations that are available in selected libraries within Finland (including The National Library of Finland).

- Portugal (since 2008) accessible in Open Access at Arquivo.pt

- The Wayback Machine (since 2006) accessible in Open Access at https://archive.org/web/

As much as for other kind of archives, one must know the history of a Web archive and how it was constructed to better understand it and use it in research work. What you see is a reconstruction, not a copy of the site

- "What is harvested is both a point in time (the time of harvesting) and a period of time (the period up to the time of harvesting)." (Brügger, 2008 p. 158)

- "On the one hand the archive does not look like the internet as it actually was in the past (we have lost something), but on the other hand the archive might look like the internet as it never was in the past (we get something different)." (Brügger, 2001 p. 6)

Web archiving projects often needs to gather diverse and multiples expertises and skills : archivists and librarians, researchers, legal officers, IT and computer specialists and… users and stakeholders

As for any other kind of archives one must act lawfully and ethically when archiving and using Web archives:

- The materials in the web archives are protected by copyright law as they were on the live web
- There is "tensions around the archival principles of preserving the public record vs the individual's expectation of the right to be forgotten" (http://netpreserve.org/ga2018/programme/abstracts/#paper21)
- The processing of personal data is submitted to laws and even more to the research project ethics
- New laws to take into account ex. General Data Protection Regulation (GDPR, https://gdpr-info.eu/)

**RESAW**
*A Research infrastructure for the Study of Archived Web materials*

Home    About RESAW    Events    Projects    Participants    Web Archives    Resources    Forum    Contact

## Home

Recherche

**RESAW**
A Twitter list by @resaw_eu
Members tweet about...

**UK Web Archive**
@UKWebArchive
Web Archiving Roundup: April, 2019
webarchivingrt.wordpress.com/2019/04/29/web… via @WebArch_RT

Web Archiving Roundup: April, 2019
The Library of Congress Digital Collections Development Coordinator positioncloses soon (May 1st)! Archive-It is hosting a training webinar on Web Archiving Systems API

*Stories and News*

*Peter Webster, Web Historian, UK*
A. S. Byatt's faded church

*Ian Milligan, Web Historian, Canada*
New Grant: "Continuing Education to Advance Web Archiving"

http://resaw.eu/

**Web90 – Patrimoine, Mémoires et Histoire du Web dans les années 1990**
"Technological progress has merely provided us with more efficient means for going backwards." — Aldous Huxley, Ends and Means

Le projet ANR Web90    Recherches en cours    TTOW Symposium    Veille    L'équipe    Crédits

Publié le **20 août 2018** par **Valérie Schafer**

**jahendler**
@jahendler

Tim Berners-Lee and I had a paper rejected because reviewer said we didn't really understand how the web worked.

https://web90.hypotheses.org/

**Welcome to the Archives Unleashed Project**

**The Archives Unleashed Project**

Home
  Welcome
  Contact Us

**About the Project**

**Archives Unleashed Toolkit**

**Archives Unleashed Cloud**

**Archives Unleashed Notebooks**

**Warclight**

**Get Involved**

**Events**

**Publications**

Welcome to the Archives Unleashed Project

**The Archives Unleashed Project**

### Welcome

Archives Unleashed aims to make petabytes of historical internet content accessible to scholars and others interested in researching the recent past. Supported by a grant from the Andrew W. Mellon Foundation, we are developing web archive search and data analysis tools to enable scholars, librarians and archivists to access, share, and investigate recent history since the early days of the World Wide Web.

https://archivesunleashed.org/

Why GitHub? Enterprise Explore Marketplace Pricing    Search    Sign in  Sign up

internetarchive / heritrix3                                   Watch 173   Star 1,454   Fork 601

<> Code   ⚠ Issues 19   ⑂ Pull requests 2   ⊞ Projects 0   📖 Wiki   Insights

Heritrix is the Internet Archive's open-source, extensible, web-scale, archival-quality web crawler project.   https://heritrix.readthedocs.io/

java   webcrawling   warc   heritrix

Heritrix is the Internet Archive's open-source, extensible, web-scale, archival-quality web crawler project

| commons | [maven-release-plugin] prepare for next development iteration | 9 days ago |
|---|---|---|
| contrib | [maven-release-plugin] prepare for next development iteration | 9 days ago |
| dist | [maven-release-plugin] prepare for next development iteration | 9 days ago |
| docs | Note feature only applies to forthcoming 3.3 release | 6 months ago |
| engine | [maven-release-plugin] prepare for next development iteration | 9 days ago |
| modules | [maven-release-plugin] prepare for next development iteration | 9 days ago |
| .gitignore | Format HTTP request lines in API guide | 10 months ago |
| .travis.yml | allow travis-ci jdk7 failures because contrib... | a month ago |
| CHANGELOG.md | Update changelog for 3.4.0-20190418 | 9 days ago |

https://github.com/internetarchive/heritrix3

- WARC file format = Web ARChive archive format
- ARC was accepted as an international standard in 2009 (ISO 28500:2009)
- WARC is now recognised by most national library systems as the standard to follow for web archival

**WARC file**

**WARC record**

Text header

Content block

[image/jpeg binary data]

```
WARC/1.0
WARC-Type: resource
WARC-Target-URI: file://var/www/htdoc/images/logoc.jpg
WARC-Date: 2006-09-19T17:20:24Z
WARC-Record-ID: <urn:uuid:92283950-ef2f-4d72-b224-f54c6ec90bb0>
Content-Type: image/jpeg
WARC-Payload-Digest: sha1:CCHXETFVJD2MUZY6ND6SS7ZENMWF7KQ2
WARC-Block-Digest: sha1:CCHXETFVJD2MUZY6ND6SS7ZENMWF7KQ2
Content-Length: 1662
```

...etc.

https://en.wikipedia.org/wiki/Web_ARChive
https://wiki.archivematica.org/Significant_characteristics_of_websites
https://wiki.archivematica.org/File:WARCdiagram.png

- Robots.txt
- Captcha (ie Completely Automated Public Turing-test to tell Computers and Humans Apart)
- User interaction needed
- Password protected content
- Technologies and dynamic content : Flash, java scripts,…
- Distant content
- Temporal inconsistencies
- Bot traps
- …

- Define a strategy
- Use a log and write throughout the life of the project
- You may need to use additional methods and tools

- Screen capture and screen recording
- Link crawling
- On demand archiving

- We will practice:
  - The Wayback machine
  - Web recorder
  - Archive IT

The structure of URLs on the World Wide Web (www):

protocol://subdomain.domain.top-level domain/path/page/

Ex:

https://dhibeirut.wordpress.com/archive/dhi-b-2017/

https://en.wikipedia.org/wiki/URL
http://dac.au.dk/forskning/forskningsprogrammer/ p. 51

# Internet Archive's Wayback Machine



- Launched in 2001
- 357 billion archived [Web pages](#) so far
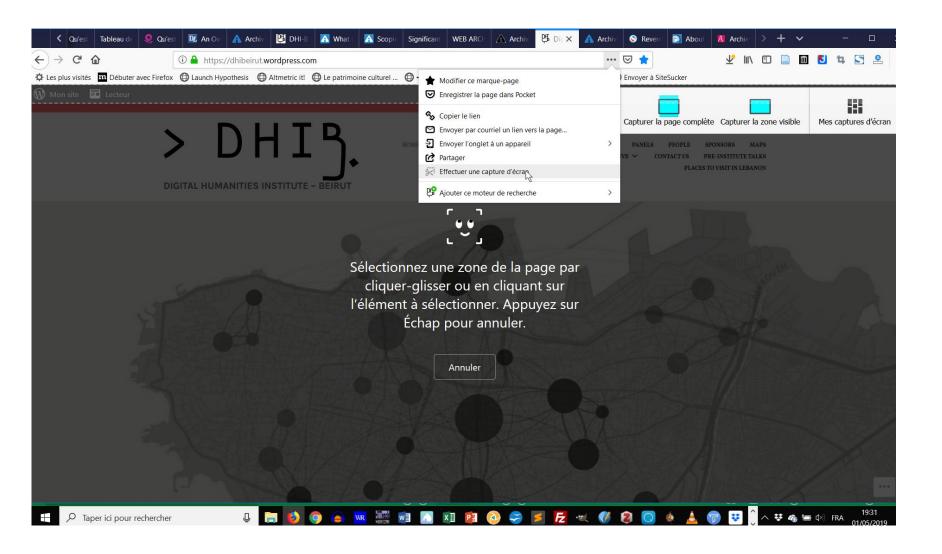- Archived content going back to 1996

https://archive.org/
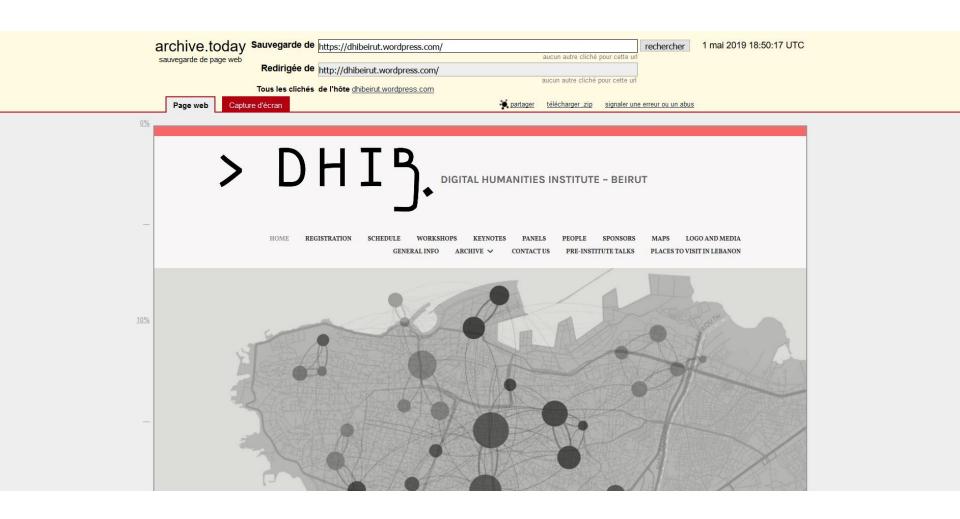
Structure of URLs in the Internet Archive's Wayback Machine:

Wayback Machine URL/collection/ time shown as yyyymmddhhmmss/URL

Ex:

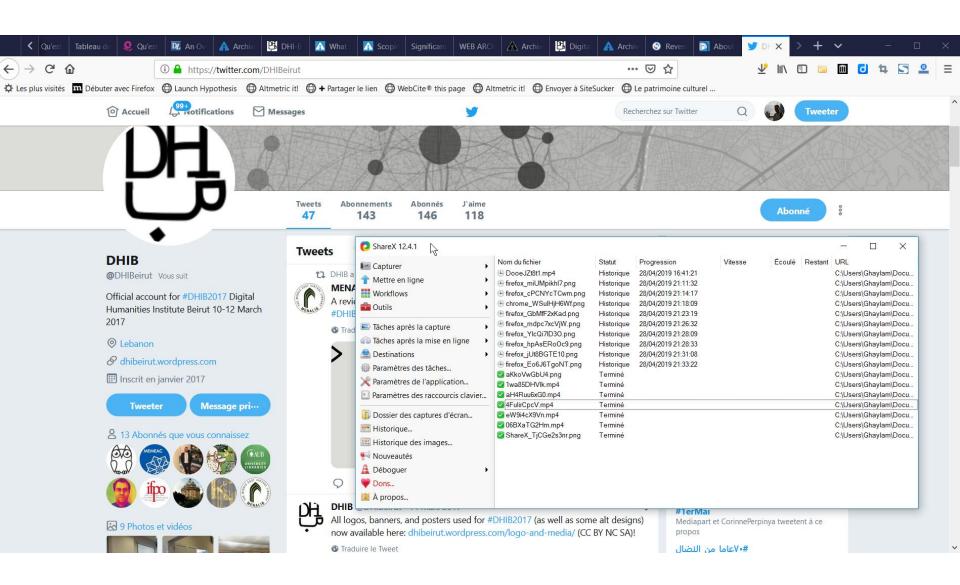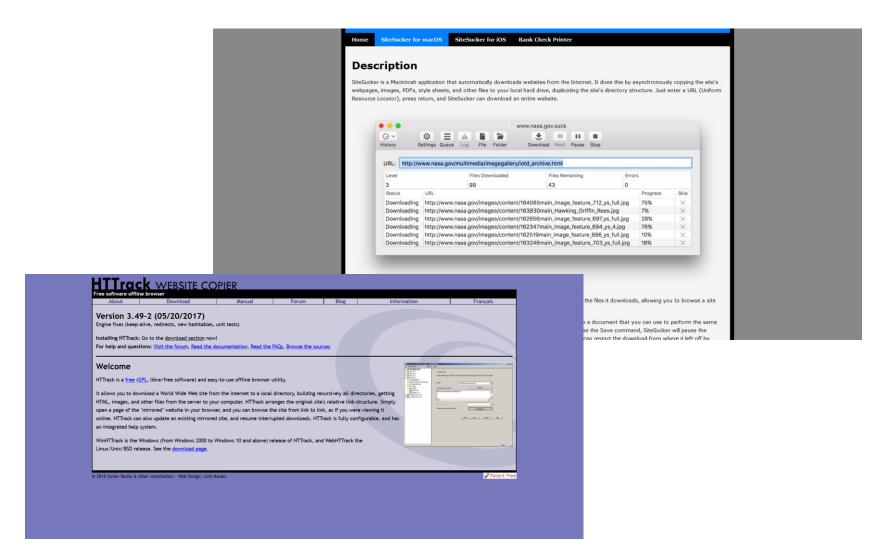https://web.archive.org/web/19980118071158/http://www.aub.edu.lb/

http://archive.fo/ZAFsE

# Citation and cached version with Webcite



http://www.webcitation.org/783EYmwuV
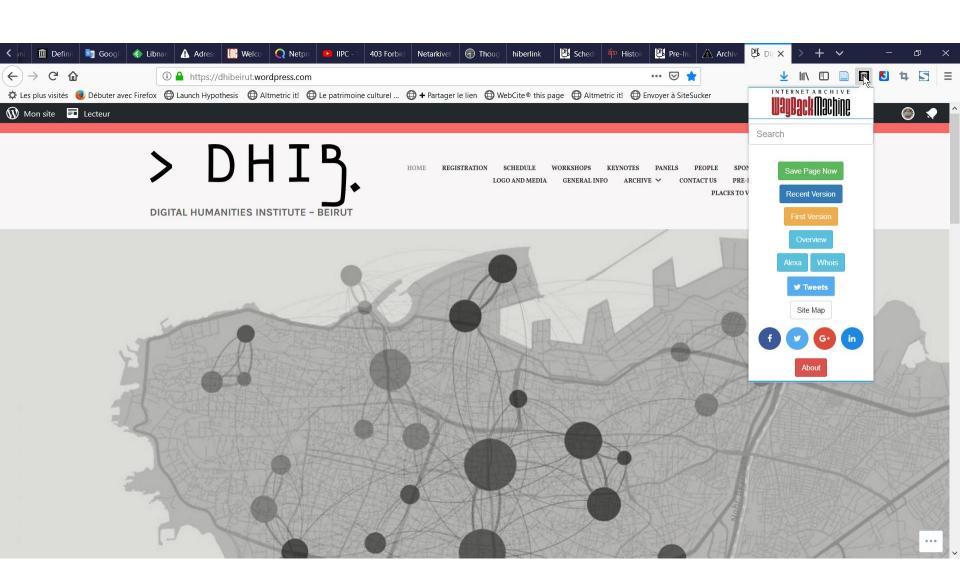
https://github.com/ShareX/ShareX

http://ricks-apps.com/
https://www.httrack.com/

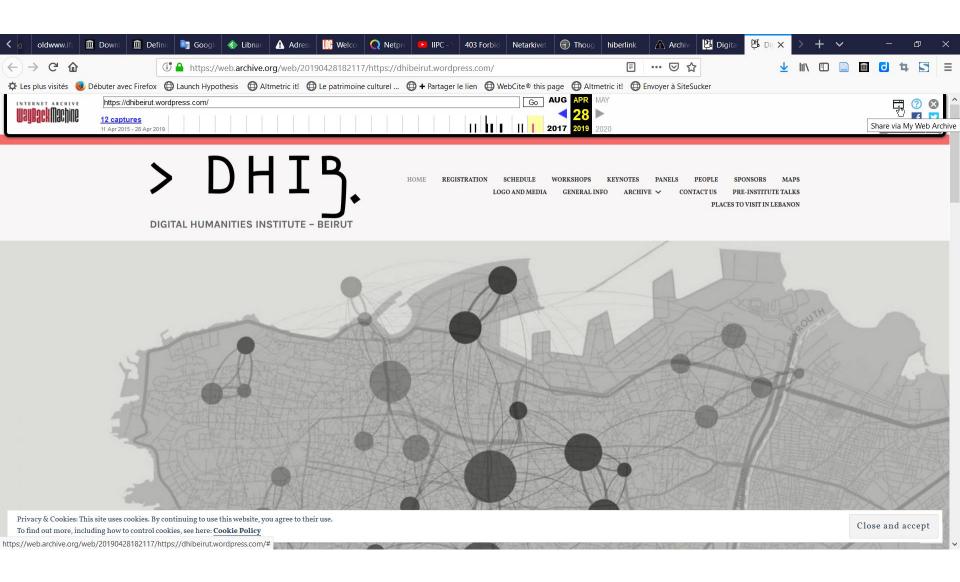# Internet Archive / Wayback Machine extension for browsers



https://chrome.google.com/webstore/detail/wayback-machine/fpnmgdkabkmnadcjpehmlllkndpkmiak?hl=fr

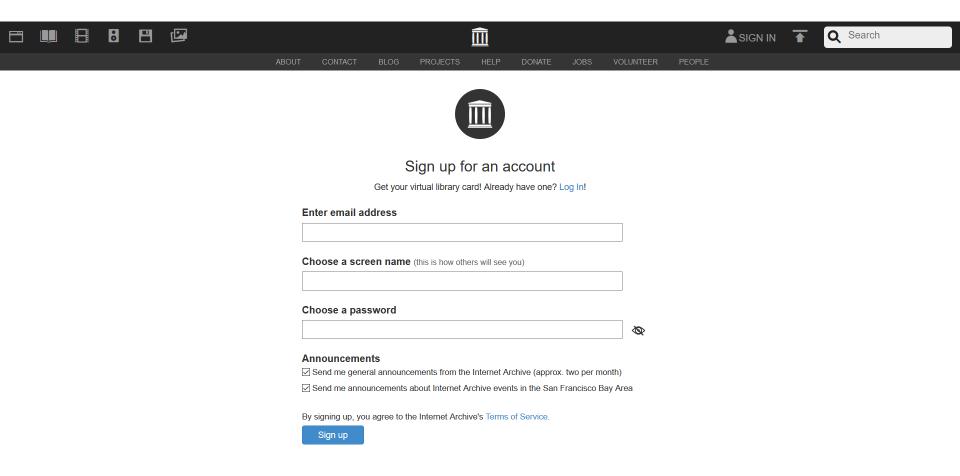https://addons.mozilla.org/en-US/firefox/user/12373129/

# Internet Archive / Wayback Machine extension for browsers

https://web.archive.org/web/20190428182117/https://dhibeirut.wordpress.com/

# Your own collection of archived Web pages in Internet Archive

ABOUT    CONTACT    BLOG    PROJECTS    HELP    DONATE    JOBS    VOLUNTEER    PEOPLE

## Sign up for an account

Get your virtual library card! Already have one? Log In!

**Enter email address**

**Choose a screen name** (this is how others will see you)

**Choose a password**

**Announcements**

☑ Send me general announcements from the Internet Archive (approx. two per month)

☑ Send me announcements about Internet Archive events in the San Francisco Bay Area

By signing up, you agree to the Internet Archive's Terms of Service.

Sign up
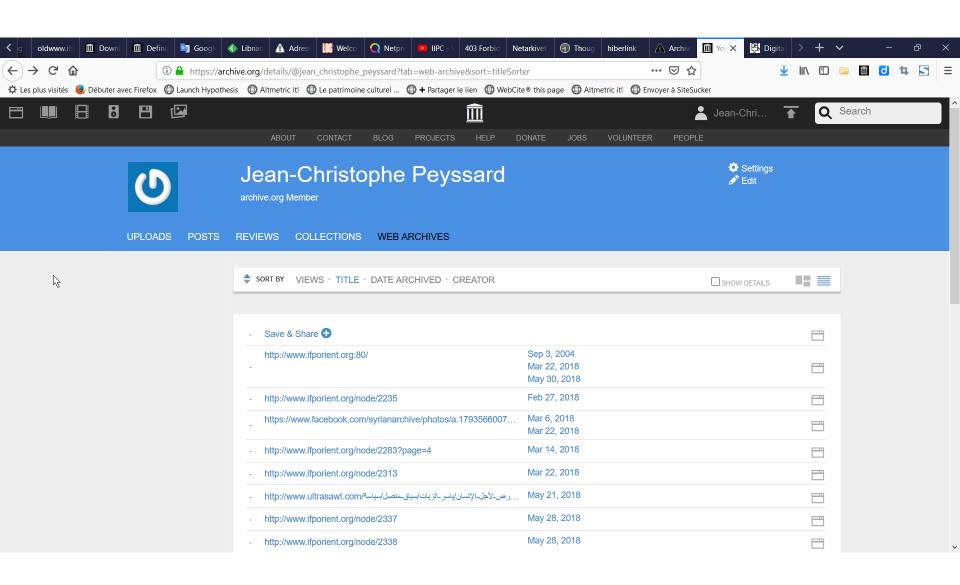
https://archive.org/account/signup

# Your own collection of archived Web pages in Internet Archive

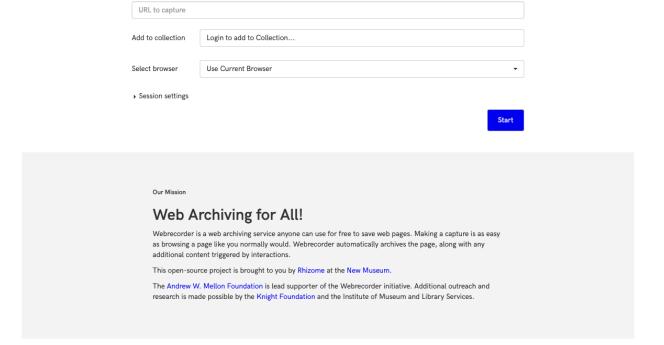# Your own collection of archived Web pages in Internet Archive



## Jean-Christophe Peyssard
archive.org Member

⚙ Settings
✎ Edit

UPLOADS  POSTS  REVIEWS  COLLECTIONS  **WEB ARCHIVES**

SORT BY · VIEWS · **TITLE** · DATE ARCHIVED · CREATOR

☐ SHOW DETAILS

| | | |
|---|---|---|
| – Save & Share ⊕ | | |
| – http://www.ifporient.org:80/ | Sep 3, 2004 Mar 22, 2018 May 30, 2018 | |
| – http://www.ifporient.org/node/2235 | Feb 27, 2018 | |
| – https://www.facebook.com/syrianarchive/photos/a.1793566007... | Mar 6, 2018 Mar 22, 2018 | |
| – http://www.ifporient.org/node/2283?page=4 | Mar 14, 2018 | |
| – http://www.ifporient.org/node/2313 | Mar 22, 2018 | |
| – ...رض-لأجل-الإنسان/ياسر-الزيات/سياق-متصل/سياسة/http://www.ultrasawt.com | May 21, 2018 | |
| – http://www.ifporient.org/node/2337 | May 28, 2018 | |
| – http://www.ifporient.org/node/2338 | May 28, 2018 | |

https://github.com/webrecorder/webrecorder-player

# Archive-It: a service for collecting and accessing Web archives



https://archive-it.org/

# Archive-It: project public collection



https://archive-it.org/home/dhib

# Archive-It: the archving back office user interface



https://partner.archive-it.org/1562

> DHIB.

**DIGITAL HUMANITIES INSTITUTE – BEIRUT**

HOME     REGISTRATION     SCHEDULE     WORKSHOPS     KEYNOTES
PANELS     PEOPLE     SPONSORS     MAPS     LOGO AND MEDIA
GENERAL INFO     ARCHIVE ⌄     CONTACT US     PRE-INSTITUTE TALKS
PLACES TO VISIT IN LEBANON

Privacy & Cookies: This site uses cookies. By continuing to use this website, you agree to their use.
To find out more, including how to control cookies, see here: **Cookie Policy**

Close and accept

https://wayback.archive-it.org/12102/20190501122014/https://dhibeirut.wordpress.com/archive/dhi-b-2017/

# Bibliography

Home > People > Jean-Christophe Peyssard > Library > Web archiving

📁 Library

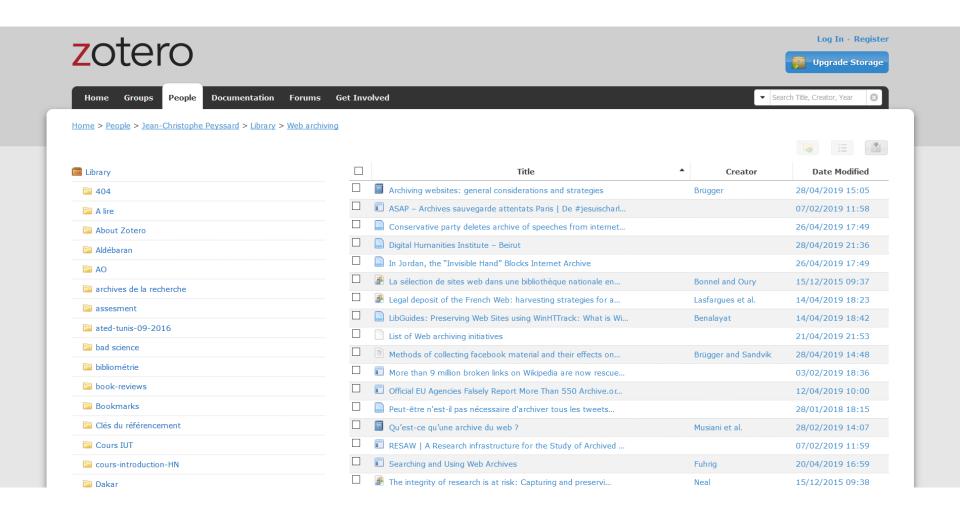| | Title | Creator | Date Modified |
|---|---|---|---|
| | Archiving websites: general considerations and strategies | Brügger | 28/04/2019 15:05 |
| | ASAP – Archives sauvegarde attentats Paris | De #jesuischarl... | | 07/02/2019 11:58 |
| | Conservative party deletes archive of speeches from internet... | | 26/04/2019 17:49 |
| | Digital Humanities Institute – Beirut | | 28/04/2019 21:36 |
| | In Jordan, the "Invisible Hand" Blocks Internet Archive | | 26/04/2019 17:49 |
| | La sélection de sites web dans une bibliothèque nationale en... | Bonnel and Oury | 15/12/2015 09:37 |
| | Legal deposit of the French Web: harvesting strategies for a... | Lasfargues et al. | 14/04/2019 18:23 |
| | LibGuides: Preserving Web Sites using WinHTTrack: What is Wi... | Benalayat | 14/04/2019 18:42 |
| | List of Web archiving initiatives | | 21/04/2019 21:53 |
| | Methods of collecting facebook material and their effects on... | Brügger and Sandvik | 28/04/2019 14:48 |
| | More than 9 million broken links on Wikipedia are now rescue... | | 03/02/2019 18:36 |
| | Official EU Agencies Falsely Report More Than 550 Archive.or... | | 12/04/2019 10:00 |
| | Peut-être n'est-il pas nécessaire d'archiver tous les tweets... | | 28/01/2018 18:15 |
| | Qu'est-ce qu'une archive du web ? | Musiani et al. | 28/02/2019 14:07 |
| | RESAW | A Research infrastructure for the Study of Archived ... | | 07/02/2019 11:59 |
| | Searching and Using Web Archives | Fuhrig | 20/04/2019 16:59 |
| | The integrity of research is at risk: Capturing and preservi... | Neal | 15/12/2015 09:38 |

Library collection folders:
- 📁 404
- 📁 A lire
- 📁 About Zotero
- 📁 Aldébaran
- 📁 AO
- 📁 archives de la recherche
- 📁 assesment
- 📁 ated-tunis-09-2016
- 📁 bad science
- 📁 bibliométrie
- 📁 book-reviews
- 📁 Bookmarks
- 📁 Clés du référencement
- 📁 Cours IUT
- 📁 cours-introduction-HN
- 📁 Dakar

https://www.zotero.org/peyssard/items/collectionKey/CNGX47Z3

# Thank you / Merci / شكراً

> DHI♭ 2019

**jc.peyssard@ifporient.org**