

D	I	G
H	U	M
L	A	B

# Web Archiving

## Part 2: EXISTING WEB ARCHIVES

Asger Harlung

Ulrich Karstoft Have



# Module 2: Existing Web Collections

- Introduction to web archives
- The Danish Netarkivet
- Internet Archive
- Library of Congress
- Other (US) web archives



Important characteristics:

Strategies

Access

Search

Documentation





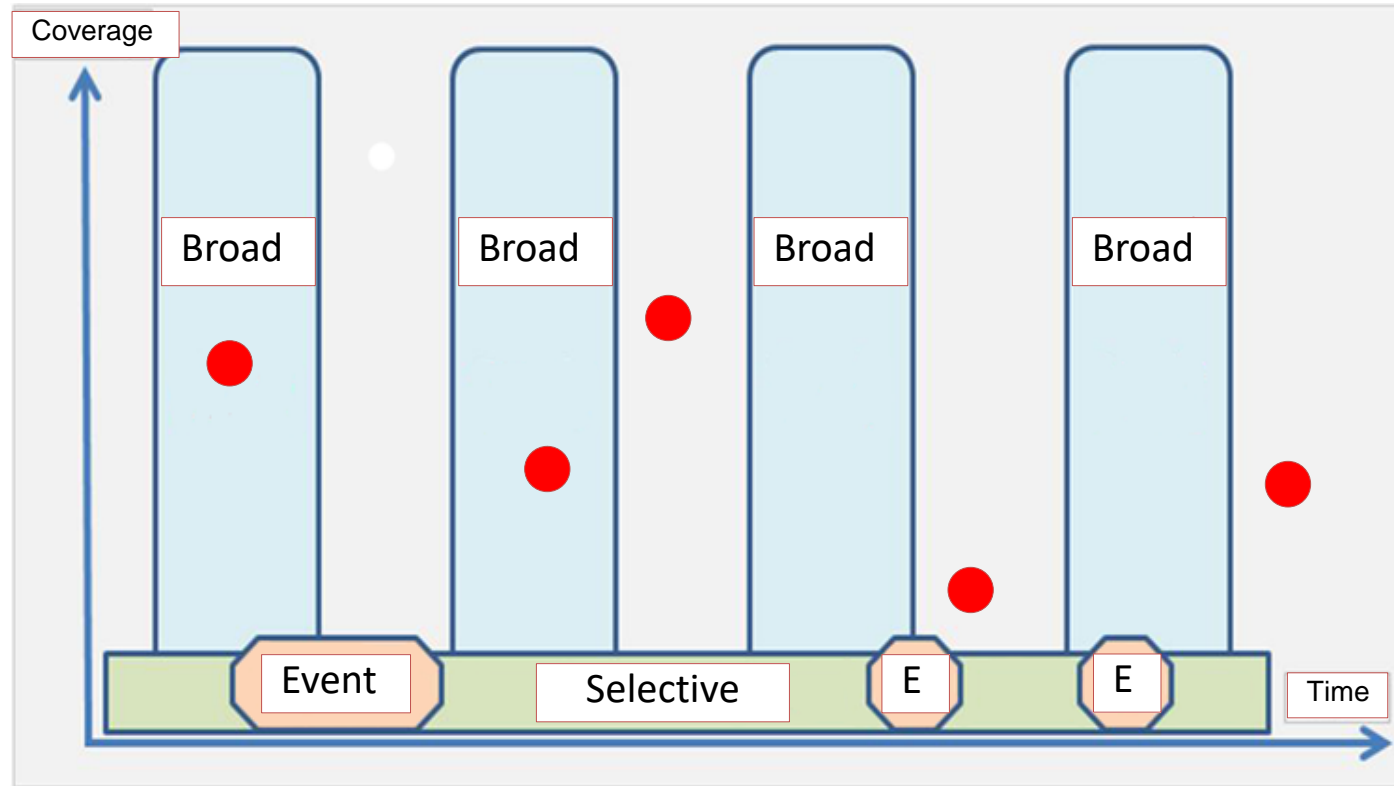
- Netarkivet is run by the State and University Library (Aarhus) and the Royal Library (National Library of Denmark, Copenhagen).
- The Danish part of the Internet is defined as cultural heritage in the Legal Deposit Act (Act no. 1439 of 22.12.2004), effective from June 1<sup>st</sup>, 2005
- The "Danish part of the Internet" = all Internet content in Danish or meant for Danes → the top level domain .dk and danica (e.g. sites in Danish or addressing Danes on other domains such as .com, .eu, .nu, etc.)
- .dk domain names: 607.000 in July 2005, 960.000 in January 2013
- Dead .dk domains from July 2005 to January 2013: 741.838
- 2011: Roughly 222 TB; 6 m objects, most common file types are html, jpeg, gif and png
- 2013: Most common file types are html, jpeg, pdf and mp4 (video)
- 2014: On July 27 the data in Netarkivet amounted to 501 TB
- 2015: On November 15 the data comprised 654 TB

# Netarkivet

2005 →

Strategies:

- Broad/bulk
- Selective
- Event
- Special



From <http://netarkivet.dk/om-netarkivet>

## The Internet Archive:

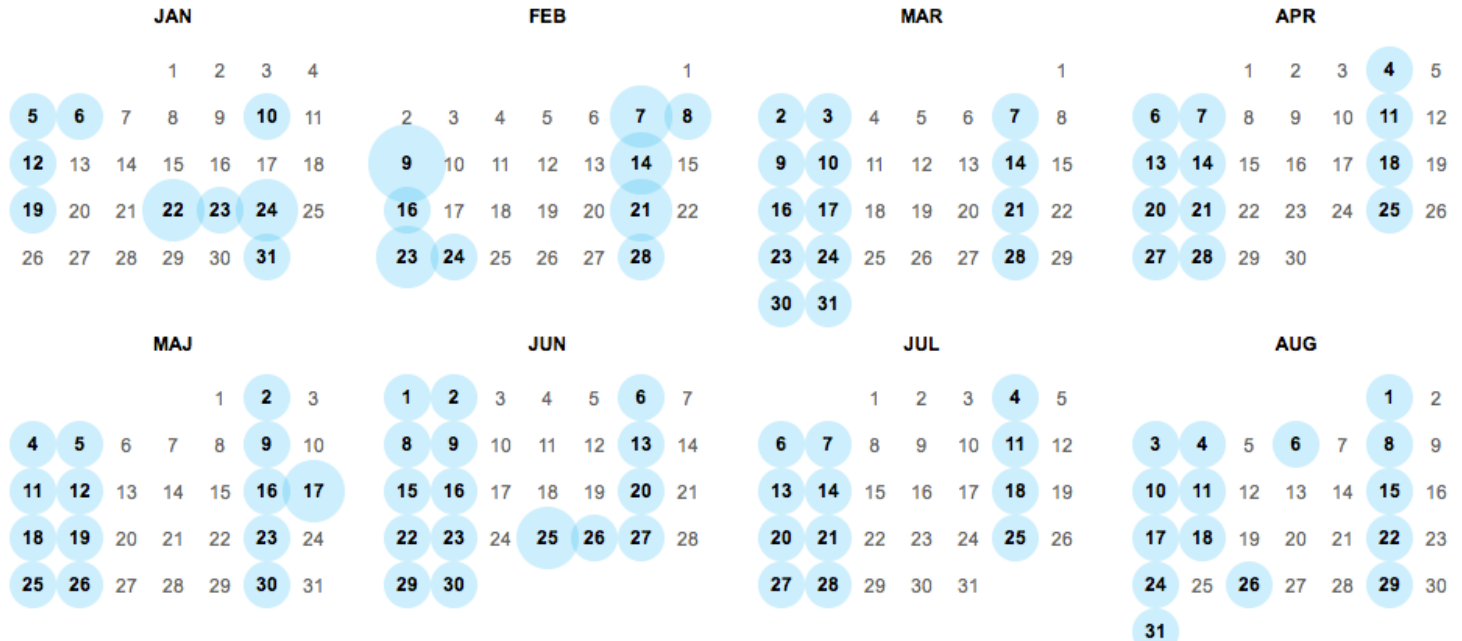
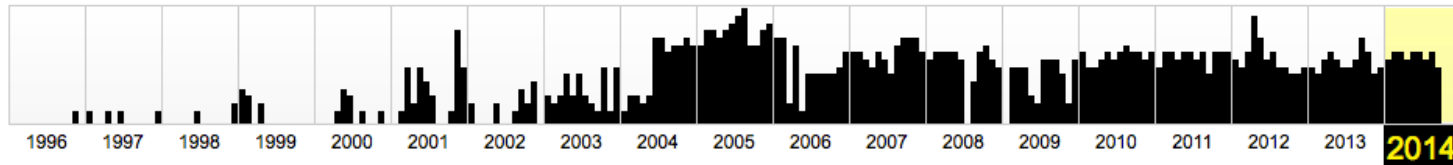
- american non-profit
- from 1996
- not based on national legislation
- in general based on cumulative archiving, following hyperlinks from what was already archived
- the worlds largest collection of archived web
- more than 491 billion web pages, collects app. 1 billion pages per week
- quality is erratic — often only top level(s)
- heterogenous collection, no overall strategy, including donations...



<http://jp.dk>

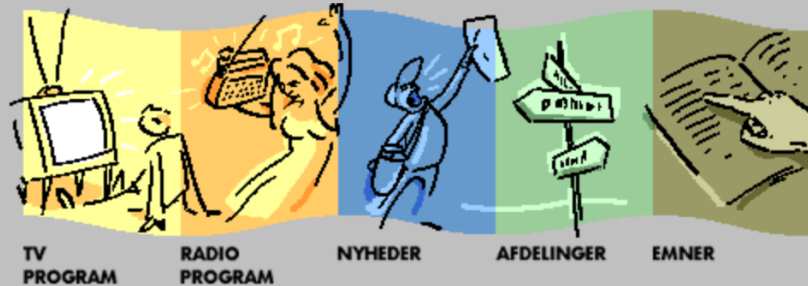
Saved **2.204 times** between [november 8, 1996](#) and [september 19, 2014](#).

**PLEASE DONATE TODAY.** Your generosity preserves knowledge for future generations. Thank you.



# DR ONLINE

[Om DR Online](#) - [English version](#) - [Tekst-version](#)



Så er tiden endelig kommet, hvor DR begynder egentlige radiosendinger over Internettet. De sidste to måneder af 1996 kører vi et forsøg med nyheder og magasinstof sendt via [RealAudio](#).



Dokumentarprogrammet 'Grevinden på tredje' om Erna Hamilton blev modtaget med begejstring af seere og presse. Programmet genudsendes søndag 27.10. kl. 14.55. Læs manuskriptet, anmeldelser og instruktørens artikel fra Politiken [på DR Online](#).



Rapporten - DR1's dybdeborende journalistiske program - har sin egen hjemmeside, hvor du kan finde mere information om programmets [afsløringer](#).



Archive-It — the Internet Archive's subscription web archiving service



- A number of collections from their partners, including event collections
- Full-text searchable
- Archive-It Research Services (ARS) — provides access to data sets extracted from collections (metadata, link graphs, named entities, other data).
- <https://archive-it.org/>



## Library of Congress web archive:

- from 2000
- curated, topic based and selective collections
- harvested by the Internet Archive (not Archive-It)
- 763TB
- <https://loc.gov/websites/collections>



LIBRARY OF CONGRESS





### Brazilian Presidential Election 2010 Web Archive

Collection Period: September 2010 to January 2011. The Brazilian elections were held October 3, 2010. The president, Luiz Inácio Lula da Silva, was not allowed by law to ...

[View 48 Items](#)



### Burma/Myanmar General Election 2010 Web Archive

Collection Period: May 2010 to January 2011. The November 7, 2010 Burma/Myanmar elections provided for under the 2008 constitution were watched by the world as critical to a ...

[View 34 Items](#)



### Crisis in Darfur 2006 Web Archive

Collection Period: February 20, 2006 to December 5, 2006. Described by the United Nations as the world's worst humanitarian crisis, the current conflict in Sudan's western region of ...

[View 220 Items](#)



### Egypt 2008 Web Archive

Collection Period: April 12, 2008 to May 31, 2008. In 2008 Egypt witnessed a remarkable experience rich in political and democratic practices, and the sites captured back then ...

[View 12 Items](#)



### Indian General Election 2009 Web Archive

Collection Period: April 9, 2009 to June 29, 2009. This collection of websites documents the 2009 Indian General Elections. As the world's most populated democratic nation with the largest ...

[View 59 Items](#)



### Indonesian General Election 2009 Web Archive

Collection Period: March 10, 2009 to October 21, 2009. This collection of websites documents the Indonesian general elections held in 2009. The collection covers three separate rounds of elections ...

[View 80 Items](#)



### Iraq War 2003 Web Archive

Collection Period: March 13, 2003 to June 16, 2003 and December 31, 2004 to July 21, 2009. On March 20, 2003, the United States began military action against ...

[View 232 Items](#)



### Laotian General Election 2011 Web Archive

Collection Period: June 2011 to November 2011. Laos, one of the world's few remaining communist regimes, holds general elections for the National Assembly every five years. Parliamentary Election ...

[View 17 Items](#)



# Other Web Archives

## IIPC Member Archives

<http://netpreserve.org/resources/member-archives>

List of Web archiving initiatives,

[https://en.wikipedia.org/wiki/List\\_of\\_Web\\_archiving\\_initiatives](https://en.wikipedia.org/wiki/List_of_Web_archiving_initiatives)

Truman, G. (2016). WebArchiving Environmental Scan. Harvard Library Report. <http://nrs.harvard.edu/urn-3:HUL.InstRepos:25658314>

