



Crash course webarchivering

NDE Webinar Webarchivering

Jesse de Vos, juni 2020



**netwerk
digitaal
erfgoed**



Wat gaan we behandelen?

- Wat is webarchivering
- Het waarom van webarchivering
- Vier fases in webarchivering
- Quickstart guide for het archiveren van websites

Je kunt vandaag nog beginnen met het archiveren van het web!



**netwerk
digitaal
erfgoed**

Wat is webarchivering?

“**Web archivering** is het proces van het verzamelen van onderdelen van het Wereld Wijde Web, het conserveren van deze collecties in een archiveerbaar formaat, en deze archieven aanbieden voor toegang en gebruik.”



Waarom webarchivering?

Wettelijke redenen:

Het vastleggen van overheidsinformatie of zakelijke communicatie voor eventueel gebruik bij geschillen (resource: [Richtlijn webarchivering NA](#))

Cultuur-historische redenen:

Het webcontent als drager van creativiteit, communicatie, entertainment, technologische ontwikkeling, etc.



Waarom webarchivering?

It's Time to Archive the Internet Archive

Publishers are suing the Internet Archive for its emergency library, putting the whole project in danger.

By [Samantha Cole](#) and [Jason Koebler](#)

Jun 9 2020, 2:00pm [Share](#) [Tweet](#) [Snap](#)



HEADQUARTERS OF INTERNET ARCHIVE, LOCATED IN RICHMOND DISTRICT, SAN FRANCISCO, CALIFORNIA.
CREDIT: [GIRL2K/WIKIMEDIA COMMONS](#)

Related Articles

You Can Now Access 1.4 Million Books for Free Thanks to the Internet Archive

MADDIE BENDER



Viacom Forced Internet Archive to Remove Hundreds of Hours of MTV Broadcasts

SAMANTHA COLE



Please Don't Sue LeVar Burton for Reading Soothing Stories to Scared Children

SAMANTHA COLE



Het web: een korte geschiedenis

Eerste NL website -
1992 Nikhef



Google - 1998

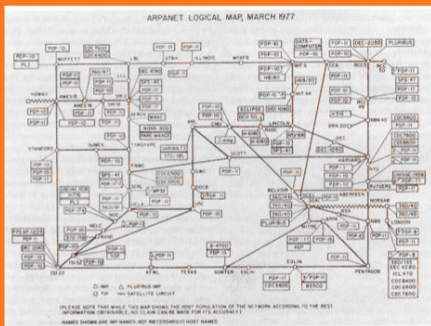
“QWERTYIOP”

Eerste e-mail - 1971

```
WWW - World Wide Web at CERN
Help [1] About this program, and the World-Wide Web
XXXXXXXX [2] XXXXXX information
CERN [3] CERN Information, ftp access to file server www1.cern.ch [4]
HEP [5] High Energy Physics
WWW [6] The home page of WWW on info.cern.ch
XXXXXXXX [7] World Wide Web Lookup of mail addresses
[End]
I-7, Quit, or Help: 2
```



ARPANET - 1969



Eerste website - 1991
CERN

```
CERN
CERN Welcome

The European Laboratory for Particle Physics, located near Geneva[1] in
Switzerland[2] and France[3]. Also the birthplace of the World-Wide
Web[4].

This is the CERN laboratory main server. The support team provides a set of
Services[5] to the physics experiments and the lab. For questions and
suggestions, see WWW Support Contacts[6] at CERN

About the Laboratory[7] - Hot News[8] - Activities[9] - About Physics[10] -
Other Subjects[11] - Search[12]

About the Laboratory

Help[13] and General information[14], divisions, groups and
activities[15] (structure), Scientific committees[16]

Directories[17] (phone & mail, services & people), Scientific
Information Service[18] (library, archives or Alice), Preprint[19] Server

I-45, Back, Up, <RETURN> for more, Quit, or Help: █
```



De Digitale Stad - 1995



Het web: een korte geschiedenis



Wikipedia - 2001

MySpace - 2003



YouTube - 2005



Blogger - 1999



Hyves - 2004



Facebook - 2004



"just setting up my twttr"

Twitter - 2006



Het web: de statistieken

4.4 miljard gebruikers

3.5 miljard op social media

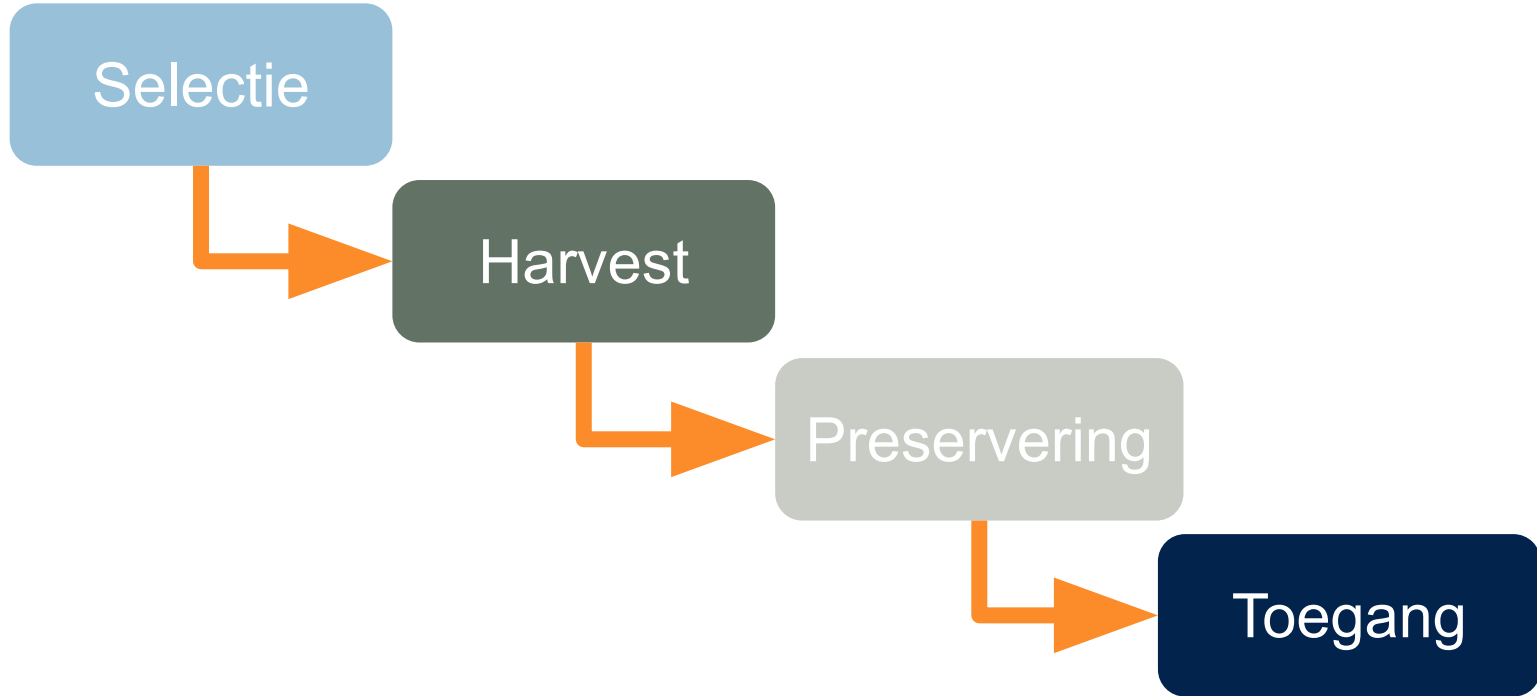
52% van het webverkeer is
door bots

6 miljoen .nl domeinen

1.7 miljard websites
(200 miljoen actief)



Vier fases in webarchivering



Vier fases in webarchivering



Fase 1: Selectie

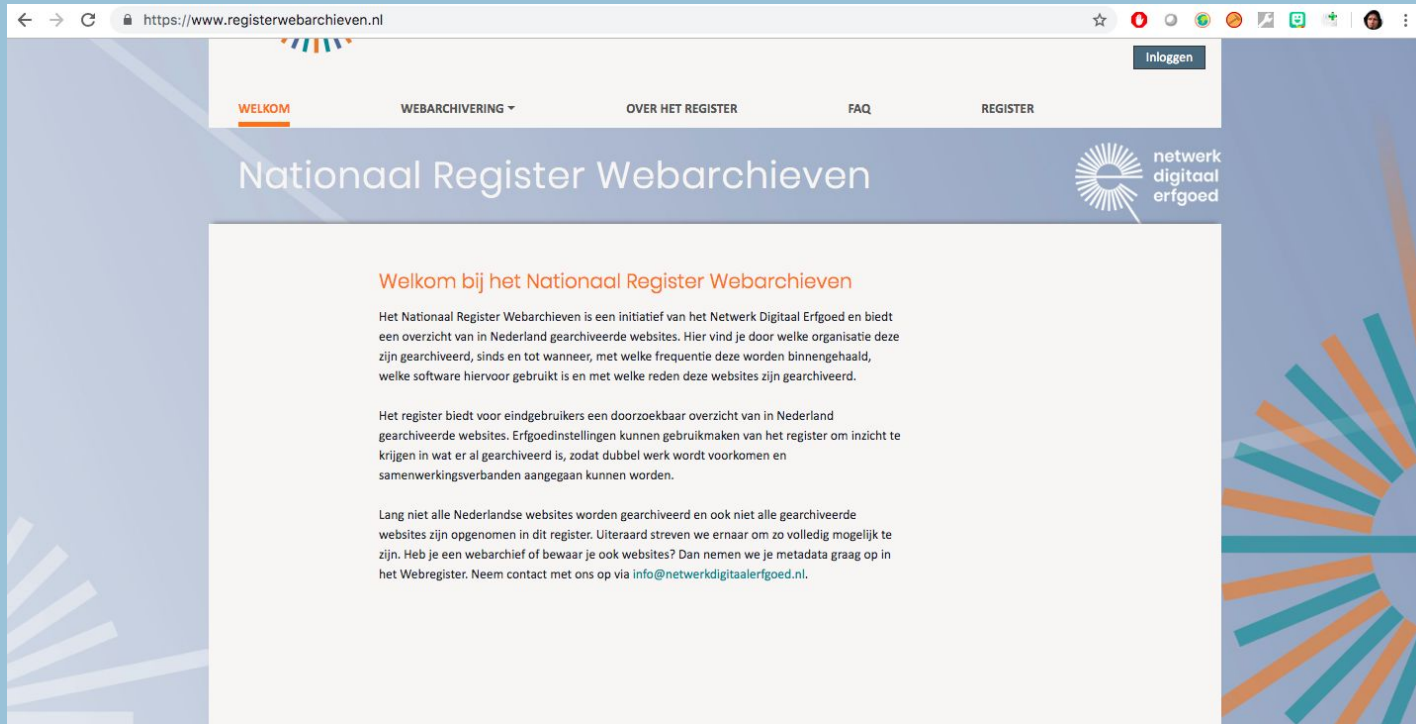
- Kan gekoppeld zijn aan wettelijke taak
- Kan gebaseerd zijn op afwegingen over de cultuur-historische waarde van webcontent
- Denk na over de doelgroep
- Afstemming met andere webarchieven in Nederland kan via het Nationaal Register Webarchieven



Fase 1: Selectie



Fase 1: Selectie



← → ↻ https://www.registerwebarchieven.nl ☆ 🔒 🌐 📄 📧 📞

Inloggen

WELKOM WEBARCHIVERING OVER HET REGISTER FAQ REGISTER

Nationaal Register Webarchieven

netwerk digitaal erfgoed

Welkom bij het Nationaal Register Webarchieven

Het Nationaal Register Webarchieven is een initiatief van het Netwerk Digitaal Erfgoed en biedt een overzicht van in Nederland gearcheverde websites. Hier vind je door welke organisatie deze zijn gearcheveerd, sinds en tot wanneer, met welke frequentie deze worden binnengehaald, welke software hiervoor gebruikt is en met welke reden deze websites zijn gearcheveerd.

Het register biedt voor eindgebruikers een doorzoekbaar overzicht van in Nederland gearcheverde websites. Erfgoedinstellingen kunnen gebruikmaken van het register om inzicht te krijgen in wat er al gearcheveerd is, zodat dubbel werk wordt voorkomen en samenwerkingsverbanden aangegaan kunnen worden.

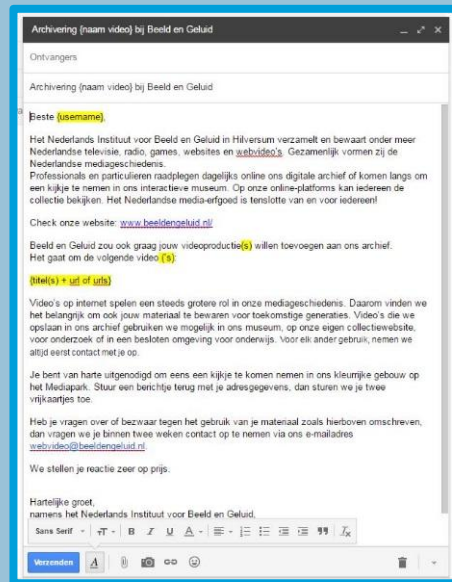
Lang niet alle Nederlandse websites worden gearcheveerd en ook niet alle gearcheverde websites zijn opgenomen in dit register. Uiteraard streven we ernaar om zo volledig mogelijk te zijn. Heb je een webarchief of bewaar je ook websites? Dan nemen we je metadata graag op in het Webregister. Neem contact met ons op via info@netwerkdigitaal erfgoed.nl.



Fase 1: Juridisch

- Opt-out, maar is niet juridisch houdbaar
- Meerdere rechthebbenden
- Ethische zaken m.n. Privacy

Maar: hoop op verruiming van regelgeving voor crawlen van webcontent.



Fase 2: Harvest

Ook wel 'crawling' of 'capture'

- Het gebruik van een tool ("crawler" of "spider") om systematisch door een website te browsen
- Configureer de crawler, definieer de scope
- De crawler download code (html, css), afbeeldingen, documenten, en andere bestanden
- Uitsluitend toegestaan met toestemming van rechthebbenden (opt-out is nu gangbaar)



Fase 2: Harvest

Domain crawl
Een beetje van alles



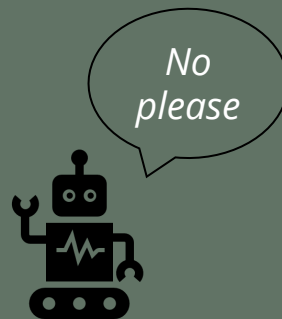
Select crawl
Alles van een beetje



Fase 2: Harvest

Seed URLs of **URIs**: startpunt(en) voor webcrawlers; de crawler volgt de links vanuit deze eerste URL of set van URLs, een zgn. “Seed List”

Robots.txt: een bestand besloten in een website die een crawler instrueert bepaalde content niet vast te leggen



Fase 2: Harvest

Crawl diepte of “**Hops**”: aantal links vanaf de Seed URL die de crawler crawlt



Crawl-Frequentie: hoe vaak dezelfde webcontent wordt gecrawld



Fase 2: Harvest

Crawler tools, o.a.:

- wget
- Heritrix
- Brozzler
- Web Curator Tool (WCT)

Crawler diensten, o.a.:

- NetarchiveSuite
- Archive-It
- Archiefweb



Dynamische webcontent

Mogelijkheden voor het vastleggen van dynamische webcontent:

- Conifer (was Webrecorder): door mens aangestuurde tool ipv automatische crawler.
- Server-side webarchivering: kopieer de bestanden op de server en maak een virtuele server.

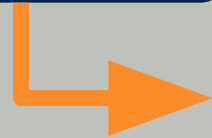


Fase 3: Preservering

Gedownloadde files worden gechecked op kwaliteit, geconverteerd naar een stabiel bestandstype indien noodzakelijk en er wordt op termijn zorg voor gedragen.



Preservering



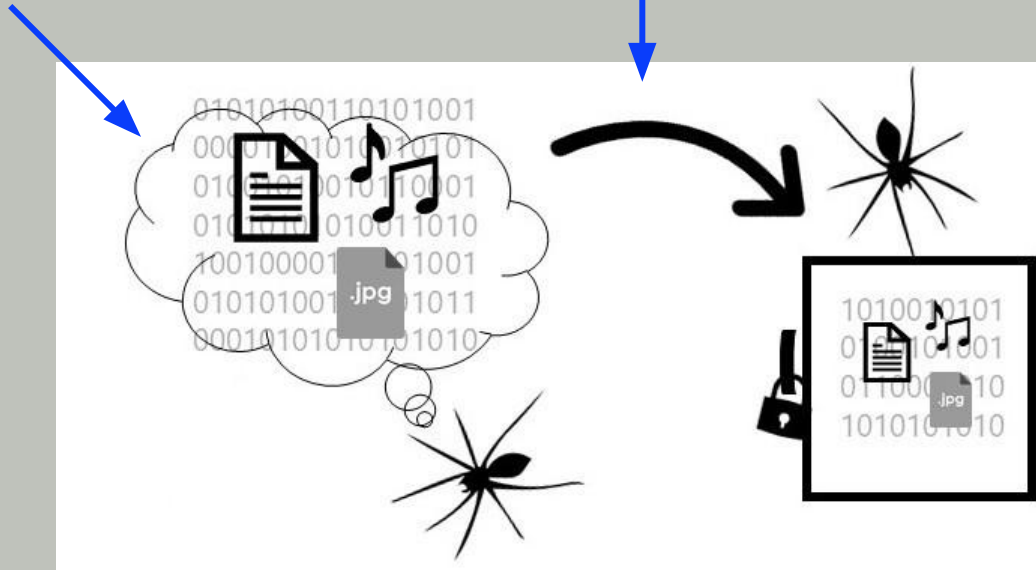
Fase 3: Kwaliteitscontrole

Gecrawlde webcontent
SIP

Kwaliteitscontrole!
Blogpost met tips

Webarchief
duurzaam
opgeslagen

AIP



Fase 3: Kwaliteitscontrole



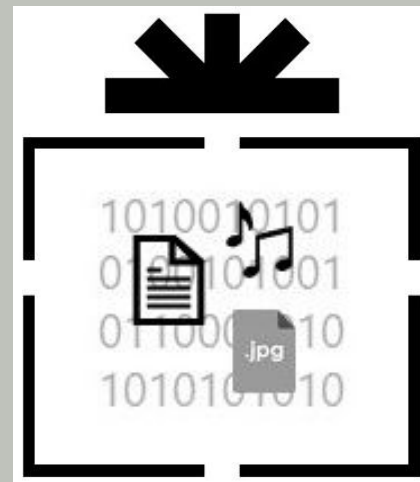
Fouten in de crawl? Onderzoek de broncode van de online website (dmv 'inspecteren' in Chrome Browser). Voeg ontbrekende URLs toe aan de 'seedlist' in de instellingen van de crawler, of pas het aantal 'hops' aan en laat de crawler opnieuw draaien.



Fase 3: Preservering

WARC (Web ARChive)

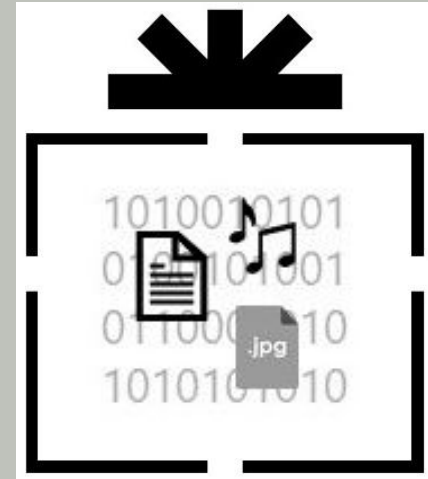
- Een WARC wordt aangemaakt door alle standaard crawling tools
- WARC is een wrapper voor gearchiveerde webobjecten ontwikkeld door de IIPC.
- ISO 28500:2017 (formerly ISO 28500:2009)
- Een WARC-bestand kan instromen in een digitaal preserveringssysteem.
- Voor het WARC-formaat was er het ARC (.arc) formaat



Fase 3: Preservering

WARC (Web ARChive)

- Een bestand dat verschillende bestandstypen bevat die verkregen zijn door het crawlen
- Behoudt de relatie tussen webpagina's en gerelateerde content
- Bevat beperkte metadata in de header
- Gebruikt specifieke toegangstools

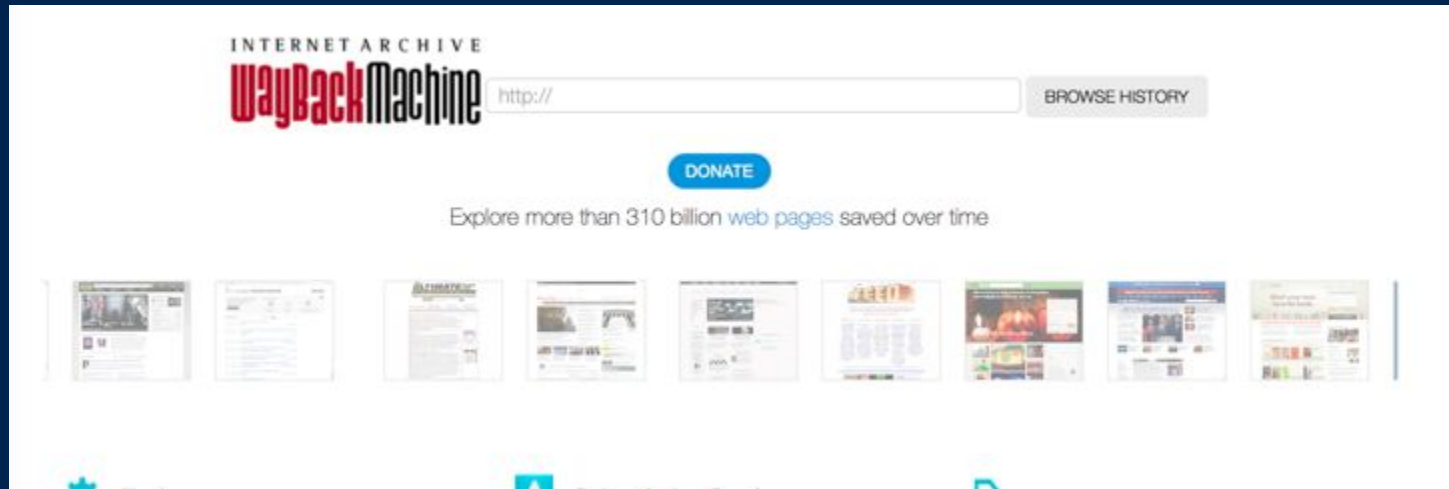


Fase 4: Toegang

- Het on-site of online beschikbaarstellen van het webarchief
- Uitsluitend mogelijk met toestemming van rechthebbenden



Fase 4: Toegang



Fase 4: Toegang

Tools en diensten voor beschikbaarstelling

- Webrecorder Player
- Web Archieven Dashboard (Archiefweb)
- Web Curator Tool (WCT)
- Archivelt
- Archives Unleashed Toolset (voor onderzoekers)

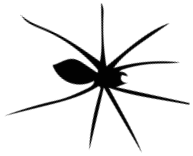


Fase 4: Toegang

Browseremulatie voor toegang



Vier fases in webarchivering



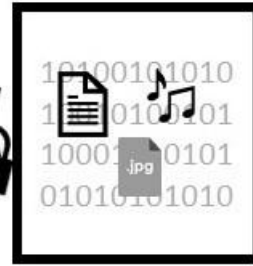
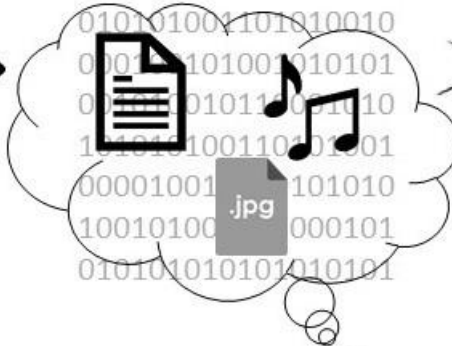
Crawler



Code en bestanden
nodig voor reproductie
van website



Playback Tool
(bijv.
Wayback)



Webarchivering quickstart guide

- Ga naar <https://conifer.rhizome.org/>
- Crawl de gewenste website
- Download de WARC vanuit je Conifer archief
- Check de kwaliteit van de WARC
- Neem de WARC op in je digitale archief (e-depot)
- Gebruik Webrecorder Player om de gearchiveerde WARC weer te openen

Of doe een suggestie aan een ander Nederlands webarchief via het [Nationaal Register Webarchieven](#)



Bedankt!

Jesse de Vos

jdvos@beeldengeluid.nl

