

Tobias Steinke

# Das deutsche Internet archivieren? Zwischen selektivem Ansatz und .de-Domain-Crawl

# Sammelauftrag

- Gesetz über die Deutsche Nationalbibliothek von 2006: Sammelauftrag auch für „alle Darstellungen in öffentlichen Netzen“
- Sammlung von Netzpublikationen über Ablieferung: E-Books, elektronische Zeitschriften, E-Paper, Hochschulprüfungsarbeiten, Musikdateien, Hörbücher, Digitalisate
- Erschließung (Katalog), Zugriff in den Lesesälen, Langzeitarchivierung

## Das deutsche Internet?

- Gedruckte Publikationen: Deutscher Erscheinungsort im Impressum
- Webseiten: Impressumspflicht nur für kommerzielle Seiten in Deutschland
- .de-Domain nicht zwingend nötig für Webseiten deutscher Herkunft
- Deutsche Sprache und Serverstandort sind keine hinreichenden Merkmale
- Jede Sammlung nur Annäherung, keine Vollständigkeit

## Webharvesting an der DNB

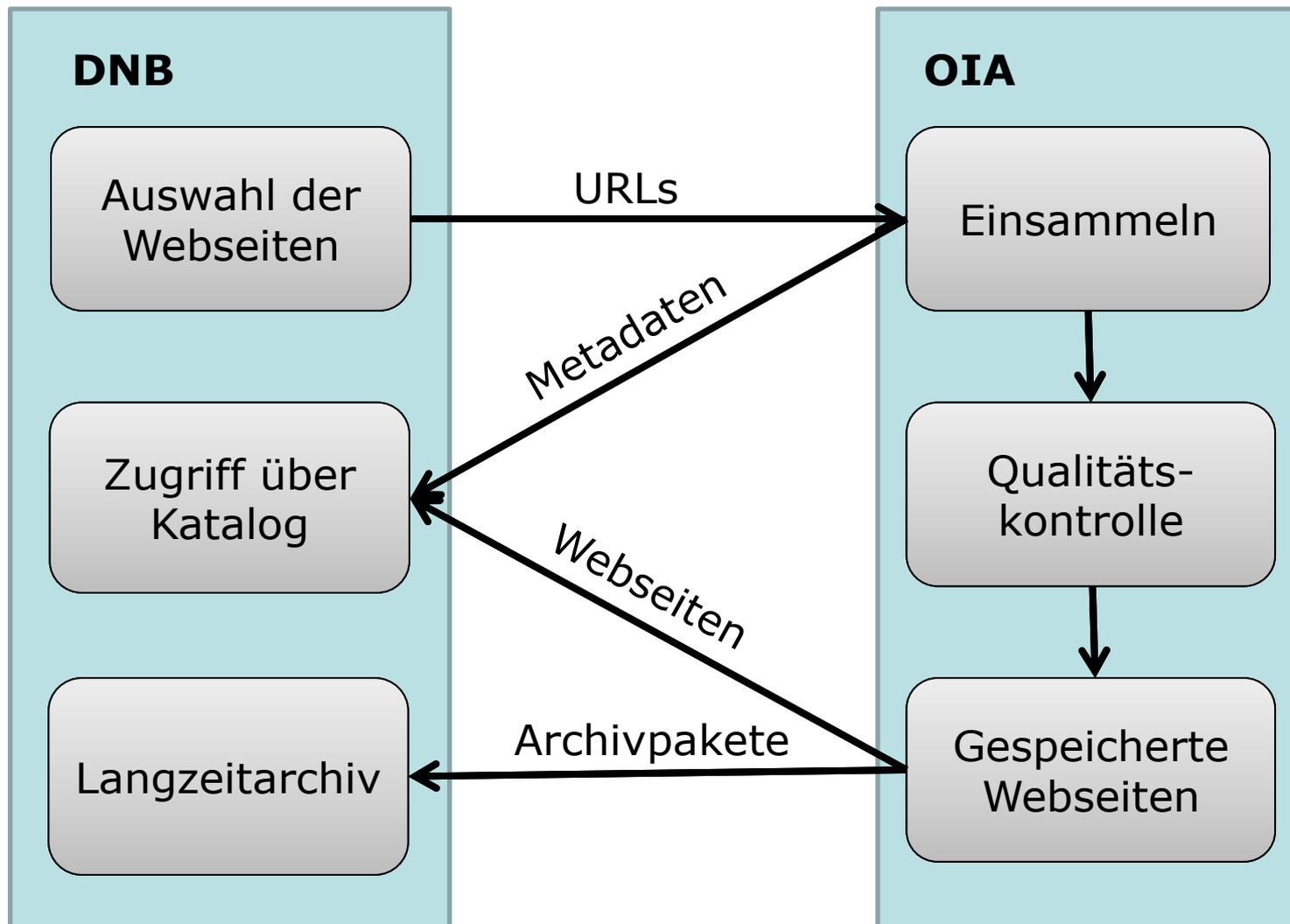
- Event-Crawls zur Bundestagswahl 2004 und EU-Ratspräsidentschaft 2007 mit European Archive
- Mitgliedschaft beim International Internet Preservation Consortium (IIPC) seit 2007
- Seit 2010 Aufbau eines Webarchivierungs-Geschäftsgangs
- 2014: Experimenteller .de-Domain-Crawl

## Aufbau eines Webarchivierungs- Geschäftsgangs

- Internes Projekt zum Vergleich der Optionen für alle Schritte im Geschäftsgang (Auswahl, Sammlung, Erschließung, Bereitstellung, Archivierung)
- Aufwandsabschätzung für Sammlungsentwicklung, Erfahrungen aus vergleichbaren Institutionen
- Ergebnis: Selektives Crawling nach thematischer Auswahl und zu Events; Sammlung, Speicherung, Qualitätssicherung und Bereitstellung über Dienstleister; Langzeitarchivierung bei DNB; .de-Crawl als Ergänzung

## Selektiver Workflow

- Europaweite Ausschreibung 2011 für Dienstleister:  
Vergabe an deutsche Firma oia ([www.oia-duesseldorf.de](http://www.oia-duesseldorf.de))
- Auswahl von Seiten durch DNB, Erschließung mit Titel und Kategorie
- Sammlung mit eigener Harvester-Software bei oia, manuelle Qualitätskontrolle
- Metadaten pro Site automatisch in DNB-Katalog
- Zugriff im DNB-Lesesaal per Katalog und Volltextsuche
- Archivierung als WARC-Dateien in DNB-Langzeitarchiv



The screenshot displays the OWA Client interface. The main window title is 'System zur Archivierung von Webseiten'. On the left, a tree view shows the project structure under 'Offline Web Archiv', with 'Bundesamt für Strahlenschutz, BfS' selected. The central pane shows the 'Eigenschaften' (Properties) dialog for this site, with the 'Allgemein' (General) tab active. The 'Allgemein' tab contains the following information:

- URL:** 172.16.2.132
- Name:** Bundesamt für Strahlenschutz, BfS
- URLs:** http://www.bfs.de/, http://www.bfs.de/de/bfs
- Options:**  Auto Ressourcen Management,  Tiefengrenze (0),  Linkverfolgung projektübergreifend

At the bottom of the dialog are buttons for 'XML Export', 'XML Import', 'Abbrechen', and 'übernehmen'. The status bar at the bottom left indicates 'Projekt: 434'. On the right side of the main window, a table shows storage statistics:

Volumen	AIU
6 GB	149.941
13 GB	67.900
52 GB	32.031
21 GB	35.869
2 GB	82.041

DEUTSCHE NATIONAL BIBLIOTHEK

LEIPZIG  
FRANKFURT AM MAIN

Übersicht Volltextsuche

Startseite → Behörden und Institutionen des Bundes → B → Bundeskartellamt Filter

### Bundeskartellamt

09.06.2014 00:03:05  
27.03.2014 18:31:49  
07.07.2013 00:01:10  
07.04.2013 00:03:24  
06.01.2013 00:00:51  
07.10.2012 00:08:44  
10.07.2012 10:08:44

Impressum Hilfe powered by oia

100%

DNB http://dnb.oia-dwa.de/show.aspx DNB - Webarchive

DNB Archivierte Netzressource vom 07.10.2012 www.bundeskartellamt.de Datensatz im Katalog

Home | Sitemap | RSS | Suche | Kontakt | Impressum | English | Français | Drucken

Bundeskartellamt Offene Märkte | Fairer Wettbewerb

Suchbegriff eingeben

**Aktuelle Meldungen**

- 05.10.2012 Aktuelle Entscheidungen der Vergabekammern
- 05.10.2012 Liste der Neuerwerbungen der Bibliothek
- 05.10.2012 Aktuelle Entscheidung: Klinikum Worms / Hochstift Worms
- 02.10.2012 Das Bundeskartellamt stellt ein: studentische Aushilfskraft
- 02.10.2012 Aktuelle Liste der laufenden Zusammenschlussvorhaben (02.10.2012)
- 02.10.2012 Aktuelle Liste der Hauptprüfverfahren (02.10.2012)

**Pressemeldungen**

- 01.10.2012 Unternehmensverflechtungen auf dem Prüfstand - Bundeskartellamt veröffentlicht Abschlussbericht der Sektoruntersuchung Walzasphalt
- 27.09.2012 Einleitung der Sektoruntersuchung Raffinerien und Mineralölgroßhandel
- 21.09.2012 OLG Düsseldorf weist Beschwerde von ConocoPhillips gegen

**Weitere Meldungen**

- Hinweise auf Kartellverstöße
- Gemeinsames Energie-Monitoring 2012 der Bundesnetzagentur und des Bundeskartellamtes
- Häufige Fragen zum Thema Kraftstoffpreise

16. Internationale Kartellkonferenz Berlin, 20.-22. März 2013

Impressum Hilfe powered by oia

100%

## Stand zu selektivem Harvesting

- Thematische Kategorien: Behörden und Institutionen des Bundes, Interessenverbände, Kultureinrichtungen, Parteien, Politiker, Religionsgemeinschaften, Sozialversicherung, Sportverbände
- Event-Crawls: Z. B. 100. Geburtstag Willy Brandt, Berlinale 2013, Bundestagswahl 2013, Grimme Online Award 2013, Hochwasser 2013, Olympia 2014
- Insgesamt derzeit ca. 700 Sites, geplante Steigerung bis Ende 2015 auf ca. 4000 Sites
- Thematische Ausweitung geplant in Kooperationen mit Aggregatoren wie Academic Linkshare und Projekten wie Alexandria

## **.de-Domain-Crawl**

- Top-Level-Domain-Crawl als Ergänzung zum selektiven
- Großer Umfang: ca. 16 Millionen registrierte Domains für .de, Crawl der BnF für .fr mit 33 TB
- Ausschreibung 2013, Vergabe an französische Firma Internet Memory Research ([www.internetmemory.org](http://www.internetmemory.org))
- Experiment über Machbarkeit und Umfang: Maximal 100 TB, Dateien mit maximal 10 MB
- Zugriff per Volltextsuche im Lesesaal
- Durchführung erste Hälfte 2014, bei Erfolg alle 2 Jahre

# Vielen Dank!

- Deutsche Nationalbibliothek: [www.dnb.de](http://www.dnb.de)