

Tobias Steinke

# Erfahrungsbericht der Deutschen Nationalbibliothek zum ersten Crawl der .de-Domain

# Inhaltsverzeichnis

## **1. Webarchivierung der DNB**

- 1. Sammelauftrag**
- 2. Selektives Webharvesting**

## **2. DE-Crawl**

- 1. Beauftragung**
- 2. Durchführung**
- 3. Bereitstellung**
- 4. Erkenntnisse**

# Sammelauftrag der Deutschen Nationalbibliothek

- 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

## Selektives Webharvesting

- Workflow zum regelmäßigen Einsammeln von Momentaufnahmen
- Auswahl von Seiten durch DNB, Erschließung mit Titel und Kategorie
- Sammlung mit eigener Harvester-Software bei Dienstleister 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

# Selektives Webharvesting

- Thematische Kategorien, z. B. Behörden und Institutionen des Bundes, Interessenverbände, Kultureinrichtungen
- Event-Crawls, z. B. Berlinale, Bundestagswahl, Olympia
- Seit 2012: Ca. 1.700 Sites, ca. 8.300 Crawls
- Regelmäßige Crawls zweimal pro Jahr, Events variierend

## DE-Crawl: Beauftragung

- Ergänzung zum selektiven Harvesting
- Top-Level-Domain-Crawls seit Jahren üblich in anderen Ländern, jedoch kleinere Größenordnung
- Optionaler Teil der Ausschreibung 2011, keine überzeugenden Angebote
- Erneute separate Ausschreibung 2013: Einmaliger experimenteller Crawl
- Internet Memory Foundation / Research (IMR)

## DE-Crawl: Durchführung

- Keine Liste aller registrierten DE-Domains von DENIC
- IMR hat Liste aus früherem Projekt als Ausgangspunkt
- Erster Crawl ohne Speicherung zur Ausweitung der Liste von Start-URLs
- Sammlung mit eigenen Crawlern auf IMR-Servern
- Einschränkungen: Keine Videos, 5.000 Dateien pro Site, 10 MB pro Datei, Gesamtgröße 100 TB
- Berücksichtigung von robots.txt

## DE-Crawl: Ergebnis

- Crawl im Juni 2014, Dauer 33 Tage
- 2,6 Millionen Start-URLs, 6 Millionen gesammelte Sites, 2,4 Milliarden Dateien
- 120 TB gesammelt (nicht genau 100 TB möglich)
- Bekannte nicht gesammelte Sites: 150 tausend
- Geschätzter Datenumfang ohne Einschränkungen: Mindestens 200 TB
- Abgleich mit Internet Archive: 62% auch dort



## DE-Crawl: Bereitstellung

- Eigene Oberfläche von IMR: Per URL und Volltextsuche
- Exklusiver Aufruf aus Lesesälen über Link in Portal und Eintrag im Katalog
- Anpassungen an DNB-Vorgaben und komplette Erneuerung der Oberfläche führte zu Verzögerungen
- Offiziell verfügbar seit 29.02.2016

The screenshot shows a web browser window with the URL `http://collections.europarchive.org/dnb/search/?v`. The page header includes the DNB logo and navigation links for 'Startseite' and 'Wir über uns'. A search bar at the top contains the text 'München' and a 'Finden' button. Below the search bar, there is a section for 'Erweiterte Suche' (Advanced Search) with various filters: 'Alle diese Wörter' (München), 'Keines dieser Wörter', 'Genau diese Phrase', 'Format' (Alle Formate), 'Zeitraum' (Von TT-MM-JJJJ bis TT-MM-JJJJ), and a 'Finden' button. The search results section shows '1-10 von 66231968' results. The first two results are:
 

- München Oben**: Original-URL: `http://eurogast.de/unbenutzt/frame-oben/muenchenob.htm`, Archiviert am: 2014-06-28 17:47:27. Action: Alle Zeitschnitte zu
- Bavaria**: Original-URL: `http://ting-chris.de/trips/europe/germany/bavaria/bavaria/bavaria_5.html`, Archiviert am: 2014-06-15 10:23:47. Action: Alle Zeitschnitte zu

 The third result is partially visible: **München, Landschaft, Architektur, Tania Reh**. The browser's status bar at the bottom indicates a zoom level of 100%.

The screenshot shows a web browser window displaying the website <http://www.br.de/fernsehen/index.html>. The page features a navigation bar with links for 'NACHRICHTEN', 'RADIO', 'FERNSEHEN', 'THEMEN', and 'MEDIATHEK'. Below this, there are several content blocks:

- BAYERISCHES FERNSEHEN**: A section with a blue header and logo, containing a description of the channel's focus on regional life and a link to 'Alle Sendungen'.
- BR-ALPHA**: A section with a white header and logo, describing it as 'Klassisches Bildungsfernsehen' and providing a link to 'Alle Sendungen'.
- BAYERNTEXT**: A section with a blue header and logo, described as 'Videotext aus Bayern', with a link to 'Alle Sendungen'.
- Wetter**: A weather forecast section showing conditions for 'Heute', 'Morgen', and 'Sonntag' with temperature ranges and icons for sun, clouds, and rain.
- Meldungen**: A news section with several headlines, including 'SPD verzichtet auf Posten in der EU-Kommission' and 'US-Präsident Obama ist zu gezielten Angriffen im Irak bereit'.

The browser's address bar shows the URL <http://collections.europarchive.org/dnb/20140610> and the page title is '[ARCHIVED CONTENT] DN...'. The browser interface includes standard navigation buttons and a search bar.

## DE-Crawl: Erkenntnisse

- Sehr gute Quantität trotz einschränkender Vorgaben
- Teilweise enttäuschende Qualität (fehlende CSS)
- Hoher Ressourcenaufwand
- Ständig steigender Umfang erhöht den zu erwartenden Ressourcenaufwand in Zukunft
- Bisher keine Entscheidung und Planung zu weiterem DE-Crawl