

# Webharvesting als Dienstleistung

Renate Hannemann | nestor for newbies | 15.05.2018

---

## BSZ:

- seit 2004 Dienstleistung für Webharvesting mit Eigenentwicklung SWBcontent:
  - 3 LBs: Stuttgart, Karlsruhe, Saarbrücken
  - Archive: Landesarchiv BW, Dt. Literaturarchiv, Kommunalarchive
  - 6 Installationen mit Webharvesting-Komponente
- rd. 1.700 Webpräsenzen mit rd. 12.000 Snapshots
- Größtmögliche Vollständigkeit, hohe Spiegelungstiefe, -genauigkeit und Authentizität
- selektives Harvesting insbes. institutioneller Websites und Event-Harvesting

### SWBcontent (2004 – 2017/2018) :

- Oberfläche in deutscher Sprache
- bibliograph. Metadaten in SWBcontent oder Übernahme aus SWB
- Parametrisierung, Crawling und Qualitätskontrolle (QA)
- HTTrack, ab 2012 Heritrix3
- Recherche, Präsentation in SWBcontent / Viewer
- zuverlässiges Tool mit überschaubarem Komfort bei Reporting, Logging
- keine automatischen Crawls (Scheduling)
- Erschließung und Crawling autonom durch Kunden

# Webharvesting als Dienstleistung dezentrales Crawling

The collage shows several library search interfaces:

- SaarDok**: Search interface for Saarland libraries with options for simple, extended, or advanced search.
- deutsches literatur archiv marbach**: Search interface for the German Literature Archive in Marbach.
- Literatur im Netz**: Search interface for digital literature resources.
- SWBregio**: Search interface for the Digital Website Archive (Digitales Webseitenarchiv).
- BOA**: Search interface for the Baden-Württemberg Online-Archiv.

### Konfiguration des Heritrix3 Crawls mittels UI

▼ Minimalkonfiguration

- URIs:  Zurücksetzen
- include URIs/SURTs as prefix:  Zurücksetzen
- Tiefe einer Spiegelung - MaxHops: 20 Zurücksetzen
- robots.txt: robots.txt immer berücksichtigen ▼
- User Agent String: Mozilla/5.0 Zurücksetzen
- include URIs by RegEx:  Zurücksetzen
- exclude URIs/SURTs as prefix:  Zurücksetzen
- exclude URIs by RegEx:  Zurücksetzen
- Tiefe einer Spiegelung - MaxPathDepth: 50 Zurücksetzen
- Obergrenze für den Download: GByte 0 Zurücksetzen MByte 0 Zurücksetzen
- Zeitlimit für den Download: Tage 0 Zurücksetzen Stunden 0 Zurücksetzen Minuten 0 Zurücksetzen

Weiter

### Alternativ - Heritrix3 Crawl XML File (cxml) - Upload

XML Heritrix Job Control Template XML-File  
Pfad des hochzuladenden cxml-Files?  
Datei auswählen Keine ausgewählt MDS ▼

UPLOAD

BOA: Baden-Württembergisches Online-Archiv 2008 - © Bibliothekservice-Zentrum Baden-Württemberg

## Aufwände BSZ und Erfahrungen:

- Administration und Entwicklung Applikation, Viewer
- Administration Präsentationsdaten im BSZ
- Technische Infrastruktur, Speicherinfrastruktur lokal im BSZ und Landesspeicher (LSDF am KIT Karlsruhe)
- Support: Schulungen, Dokumentation, Parametrisierung Spezialfälle, Problem- und Fehleranalyse
- hohes Kommunikationsaufkommen mit Kunden, zeitintensive Fehlerbehandlung
- Datenqualität: schwankend und nicht beeinflussbar

## Entwicklung im BSZ seit 2016:

- ! Ablösung SWBcontent
- ? Entwicklung eigenes Tool?
- ? Nutzung Fremdsystem? Evaluation
- ! Nutzung Dienstleister Archive-It, vollständige Migration der Altdaten
- ! neues Servicemodell

## Archive-It

- seit 2006 Angebot des Internet Archive in S.F.
- rd. 500 Partner weltweit
- komfortables Workflow-Tool:  
Erschließung, Parametrisierung, Jobmanagement,  
Zugriffsschutz, QA, Schnittstellen
- Heritrix3 + Eigenentwicklungen (Umbra, Brozzler)
- Individualisierbare Präsentation der archivierten  
Ressourcen als eigenständige Archive im Web
- Volltextrecherche (Websites und pdf)
- Oberfläche, Dokumentation, Support in englischer  
Sprache



Explore >> Saarländische Universitäts- und Landesbibliothek (SULB) >> Saarland – frei zugängliche Ressourcen



## Saarland – frei zugängliche Ressourcen

Collected by: [Saarländische Universitäts- und Landesbibliothek \(SULB\)](#)

Archived since: Jan, 2018

**Description:** In dieser Sammlung werden Webauftritte, die im Saarland erscheinen oder sich inhaltlich mit dem Land, seinen Orten und Personen beschäftigen, gesammelt, erschlossen und der Öffentlichkeit zur Verfügung gestellt. - In this collection, websites that are published in the Saarland or deal with the content of the region, its places and people are collected, indexed and made available to the public.

### Narrow Your Results

**Group** Sort By: **Count** | (A-Z)

- Bildung (11)
- Geschichte (16)
- Gesellschaft (10)
- Kreise, Städte und Gemeinden (24)
- Kunst und Kultur (42)

More ▾

**Creator** Sort By: **Count** | (A-Z)

- Universität des Saarlandes (3)
- Landesjugendring Saar (3)
- Gemeinde Kirkel (2)
- Gesellschaft zur Medienförderung Saarland - Saarland-Medien <Saarbrücken> (2)
- Industrie- und Handelskammer des Saarlandes (2)

More ▾

**Language** Sort By: **Count** | (A-Z)

- Ger (170)

**Collector**

Sites for this collection are listed below. Narrow your results at left, or enter a search query below to find a site, specific URL or to search the text of archived webpages.

Enter search terms here

Search

Clear

Sites

Search Page Text

Page 1 of 2 (170 Total Results)

Next Page ▶

Sort By: **Title (A-Z)** | Title (Z-A) | URL (A-Z) | URL (Z-A)

**Title:** Ahnen im Saarland

**URL:** <http://ahnen-im-saarland.de/>

Captured once on Apr 4, 2018

**Group:** Geschichte

**Language:** ger

**Rights:** public

**ZDB:** 2860484-3

**Identifier:** <http://ahnen-im-saarland.de/>

**Collector:** Saarländische Universitäts- und Landesbibliothek Saarbrücken (SULB)

**Title:** Die Bosener Gruppe

**URL:** <http://bosenergruppe.saar.de/>

No Captures were found for this URL

**Group:** Sprache und Literatur

**Creator:** Die Bosener Gruppe



### „SWBregio“ seit 2017




- Zentraler Crawlingservice durch BSZ
- Konsortiallösung für derz. 20 Kommunal- und Kreisarchive
- Auftragscrawling nach schriftlicher Beauftragung
- rd. 320 Websites
- elastisches Preismodell nach Nutzungsintensität
- Archiv: Spiegelungserlaubnis; Beauftragung zu Crawling und Präsentation der Ressource
- BSZ: Einrichten; Metadaten in Archive-It; Parametrisierung; Tests; produktiver Betrieb und QA
- Recherche, Präsentation in Archive-It
- Archiv hat keinen Zugriff auf das Workflow-Tool



### Aufwände BSZ und Erfahrungen:

- Mehraufwand: Erschließung, Parametrisierung, Crawling (Test und Produktion); automatisierte Datenübernahme für LZA in B-W
- Es entfallen: Kommunikation, Support, Schulung, Dokumentation
- Es entfallen: Administration Applikation + Kommunikations- und Präsentationsumgebung, lokale Datenverwaltung, lokale Speicherinfrastruktur
- Inanspruchnahme Helpdesk, Dokumentation, Schulungen, Webinare etc. von Archive-It durch BSZ
- Datenqualität: homogen und sehr gut

### „SaarDok“ 2018

- Zentraler Crawlingservice im BSZ
- Auftragscrawling nach schriftlicher Beauftragung
- rd. 165 Websites
- SULB Saarbrücken: Metadatenerfassung ZDB/SWB; Spiegelungserlaubnis; Beauftragung zu Crawling und Präsentation der Ressource; Erstabnahme Tests; QA für Produktion
- Recherche SWB, Präsentation Archive-It
- BSZ: Einrichten, Parametrisierung, Tests, produktiver Betrieb
- SULB: eingeschränkter Zugriff auf Workflow-Tool
- Aufwände erwartet wie Modell I + Schulung und Support QA

-  zentrale Webharvesting-Dienstleistung hat die Erwartungen erfüllt
-  personelle und technische Ressourcen eingespart bzw. effizienter eingesetzt
-  bessere Qualität in Crawling und Präsentation, zufriedene Kunden

-  leistungsfähiges, nutzerfreundliches Workflow-Tool
-  Inanspruchnahme Fremddienstleistung

Vielen Dank für Ihre Aufmerksamkeit!

Für Fragen wenden Sie sich gerne an  
[renate.hannemann@bsz-bw.de](mailto:renate.hannemann@bsz-bw.de)