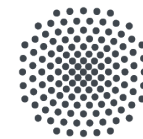

Digitale Objekte in der Webarchivierung aus textanalytischer Perspektive

Claus-Michael Schlesinger, Mona Ulrich, André Blessing
Science Data Center for Literature
Deutsches Literaturarchiv Marbach
Universität Stuttgart

deutsches
literatur
archiv **marbach**

H L R I S
High-Performance Computing Center | Stuttgart



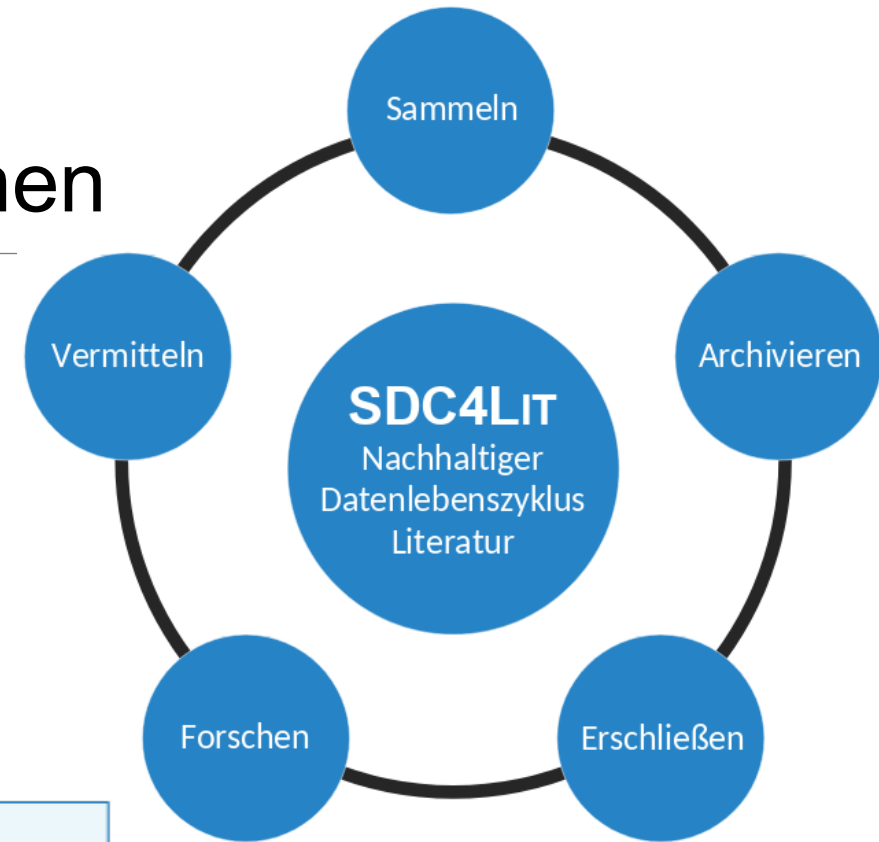
Universität Stuttgart



Baden-Württemberg
MINISTERIUM FÜR WISSENSCHAFT, FORSCHUNG UND KUNST

Literatur im Netz archivieren und erforschen

Deutsches Literaturarchiv Marbach (DLA)
Höchstleistungsrechenzentrum Universität Stuttgart (HLRS)
Institut für Maschinelle Sprachverarbeitung Universität Stuttgart (IMS)
Institut für Literaturwissenschaft Universität Stuttgart (ILW)



 SDC4LIT PLATTFORM

EXTERNE REPOSITORIEN

PROJEKT-REPOSITORIEN

ANALYSEUMGEBUNG

STORAGE

STORAGE

Das Internet ist kein Speichermedium

Literatur im Netz 2010/2022: Wieviele Seiten sind noch online?

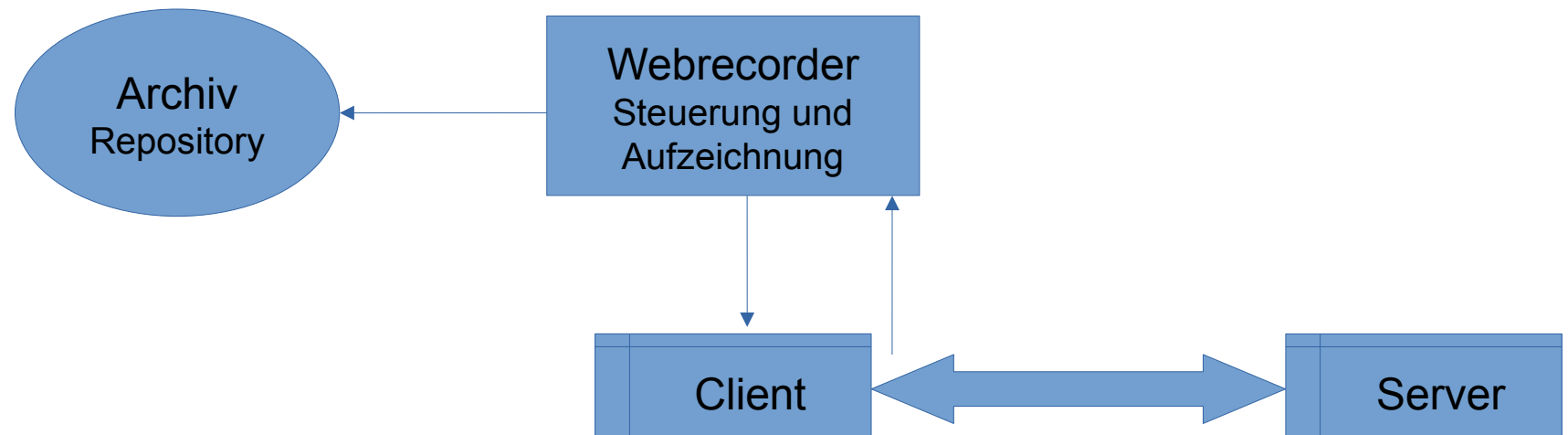
- Prüfkorpus: 134 Webseiten, davon 71 literarische Werke
- Seiten und Werke online 2022 (live web)
 - Webseiten: 79 von 134 (59 %)
 - Werke: 47 von 71 (66 %)

Webseiten und Werke
in Datenbanken und Archiven

Catalogue	Sites	Works
swb	22	22
iza	3	3
the_next	9	8
ia	107	59
kallias	33	28
cell	36	35
adel	45	42
elmcip	55	33

Literatur im Netz archivieren

- Live Web
 - Seiten können nicht stabil referenziert werden
 - Seiten verschwinden
- Archiv
 - stabile Referenz
 - Erschließung
 - Erhaltung
- Bestand DLA
 - 330 Blogs
 - 60 Netzliteraturwerke
 - 90 literarische Zeitschriften
 - WARC-Format
 - regelmäßige vollständige Crawls



URL: <http://www.aaleskorte.de/>

Titel: Die Aaleskorte der Ölig

Verfasser: Frank Klötgen, Dirk Günther

Jahr: 1998

WARC-Format

- ISO 28500:2017
- enthält alle durch den Client angefragten Ressourcen und die Kommunikation (requests und responses) zwischen Client und Server

```
1 WARC/1.0
2 WARC-Type: response
3 WARC-Target-URI: http://www.aaleskorte.de/anfang.htm
4 WARC-Date: 2015-07-02T07:25:47Z
5 WARC-Payload-Digest: sha1:IJLVDQYXWECXQ5AXT6GXEQ25NCSNDCJP
6 WARC-IP-Address: 82.165.104.8
7 WARC-Record-ID: <urn:uuid:796b9674-0a37-4f2a-8207-d48a5c4b9824>
8 Content-Type: application/http; msgtype=response
9 Content-Length: 1744
10
11 HTTP/1.1 200 OK
12 Date: Thu, 02 Jul 2015 07:25:47 GMT
13 Server: Apache
14 Last-Modified: Mon, 27 Oct 2003 17:52:49 GMT
15 ETag: "f9bb4026-5de-3caaed488fe40"
16 Accept-Ranges: bytes
17 Content-Length: 1502
18 Connection: close
19 Content-Type: text/html
20
21 <!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 3.2 Final//EN">
22
23 <HTML>
24 <HEAD>
25 |   <TITLE>Aaleskorte der Ölig</TITLE>
26 </HEAD>
27
28 <BODY bgcolor="black" text="#cccccc" link="red" alink="red"
29 vlink="red">
30 <table width="500"><TR><td>
31 <a href="index.htm">Stop!</A><P>
32 Zunächst müssen Sie Ihr Drehbuch für die "Aaleskorte der Ölig"
33 zusammenstellen. Keine Angst, diese Mühe kostet Sie ganze zwanzig
34 Mausklicks - dann startet der Film.<P>
35 <HR width="500" align="left">
36 Auf den folgenden Seiten werden Sie sehen:<BR>
37 die Portraits der Protagonisten, die für die jeweilige Szene anwählbar
38 sind, sowie eine kurze Regieanweisung.<P>
39 </TD></TR></TABLE>
40
41 <table width="500"><TR><td width="100">



Web Archive Graph Visualization

<https://linkgate.bibalex.org/>  
<https://netpreserveblog.wordpress.com/2020/04/23/linkgate-update/>



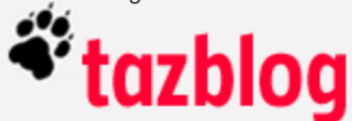
<https://web.archive.org>

# WARC-Objekte analysieren: Literarische Blogs



Die Anzeige dieser Ressource stammt aus "Literatur im Netz" . Sie wurde am **04.12.2007** archiviert.

zaehler szmtag



[--> zur aktuellen taz](#)

[--> Blogübersicht](#)

## AutorIn



Bio

## Archiv

December 2007

November 2007

October 2007

## Auf der Borderline nachts um halb eins.

### Gabor Steingart

Monday, 03.12.2007 von lottmann

Judith las gerade einen Artikel in der Frankfurter Allgemeinen Sonntagszeitung über Ludwig Erhard und das Deutsche Wirtschaftswunder, während ich mir nochmal Gabor Steingarts China-Artikel vornahm, für den er gerade den Helmut-Schmidt-Preis für herausragende journalistische Leistungen bekommen hatte. Der Mann wr ja eigentlich ein Wirtschaftler. Als Chefredakteur mußte er natürlich ein Allrounder sein, und das war er auch. Er leitete das Hauptstadtbüro in Berlin, und zuletzt berichtete er von New York ...

Mehr...»

Geschrieben in [Allgemein](#) | [1 Kommentar](#) »

### Stefan Aust

Sunday, 02.12.2007 von lottmann

„Das Papier ist nicht das Problem, aber die Tintenpatronen sind so arschteuer“, meinte Judith, als sie von der Möglichkeit erfuhr, das 'Borderline'-Buch bald herunterladen zu können. „Bei 500 Seiten brauchst du mehrere Patronen. Du mußt schon stinkereich sein, um das bezahlen zu können.“ Sie sagte immer 'stinkereich', niemals 'reich'. Ich selbst drückte mich nie so aus, von mir konnte sie es nicht haben. Ich sagte auch nicht 'arschteuer', sondern 'preislich ...

Joachim Lottmann: Auf der Borderline nachts um halb eins, Crawl vom 12.4.2007, URN: urn:nbn:de:bsz:mar1-dd001-49f6aea4-3d5a-44d4-b29f-38e487e166ae1, Screenshot ([http://literatur-im-netz.dla-marbach.de/js\\_pview/downloads/frei/49f6aea4-3d5a-44d4-b29f-38e487e166ae/0/index.html](http://literatur-im-netz.dla-marbach.de/js_pview/downloads/frei/49f6aea4-3d5a-44d4-b29f-38e487e166ae/0/index.html))

# WARC-Objekte analysieren: Objektstruktur Crawls

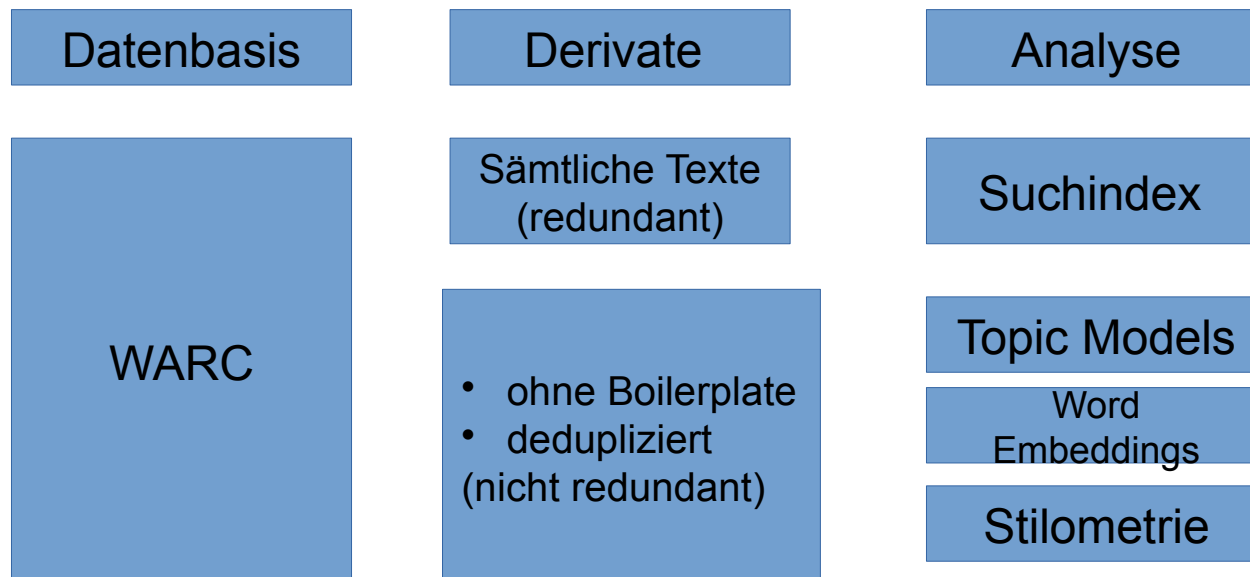
---

- Der Server liefert die Texte in verschiedenen Ansichten aus, z.B. Übersicht, Einzelbeitrag, Kategorienansicht usw.
- Ein vollständiger Crawl speichert alle Ansichten, die der Server ausliefert.
  
- => Einzeltexte sind im Crawl mehrfach redundant enthalten
- => Redundanz ist abhängig von serverseitigen Einstellungen und kann deshalb nicht kontrolliert werden
  
- Analyse: Die Redundanz der Daten ist ein Problem
- Replay: Die Redundanz der Daten ist ein Feature!



# WARC-Objekte analysieren: Preprocessing und Derivate

Textstatistische Methoden wie Word Embeddings, Topic Modeling oder Most-Frequent-Words-Ansätze erfordern nicht-redundante Daten.



Schwierigkeiten bei der Einzelobjektanalyse:

- WARC-Crawls enthalten oft Daten außerhalb einer Seite
- Entfernung Boilerplate nicht exakt
- Deduplizierung nur eingeschränkt automatisierbar

# Strukturierte Daten via API

Viele Content Management Systeme und Plattformen bieten Schnittstellen für den Download von strukturierten Daten an.



twitter.com



tumblr.com

```
{'contributors': None,
 'truncated': True,
 'text': 'katzenfutter\nkatzenfutter und rettich\n\nrettich\nrettich und brokkoli\n\nkatzenfutter\nkatzenfutter und brokkoli\n\nkatzen... https://t.co/cJk5n0vRsv',
 'is_quote_status': False,
 'in_reply_to_status_id': None,
 'id': 1447516333000859653,
 'favorite_count': 2,
 'source': 'gomringador',
 'retweeted': False,
 'coordinates': None,
 'entities': {'symbols': [],
 'user_mentions': [],
 'hashtags': [],
 'urls': [{'url': 'https://t.co/cJk5n0vRsv',
 'indices': [117, 140],
 'expanded_url': 'https://twitter.com/i/web/status/1447516333000859653'}
```

Twitterbot @gomringador von Kathrin Passig,  
Ausschnitt aus dem Datensatz für einen Tweet, Tweet-ID: 1447516333000859653, Download 11.10.2021.

# Ergänzende Archivierung via API

---

Struktur eines Datenpakets für @gomringador

```
.
├── screenshots
│ ├── 2019-10-15-000130_762x684_scrot.png
│ ├── Screenshot 2022-05-17 at 20-30-33 gomringador (gomringador) Twitter.png
│ ├── Screenshot 2022-05-17 at 20-31-22 gomringador on Twitter.png
│ ├── Screenshot 2022-05-17 at 20-31-46 gomringador on Twitter.png
│ └── Screenshot 2022-05-17 at 20-32-11 gomringador (gomringador) Twitter.png
├── sourceCodes
│ ├── source.v1.zip
│ └── source.v2.zip
├── structuredData
│ └── structuredData.zip
├── webArchive
│ └── gomringador_20220519.wacz
```

# Strukturierte Daten herunterladen

---

"Dschungel. Anderswelt" von Alban Nikolai Herbst

- 14995 Einträge (Mo 27. Jun 19:43:12 CEST 2022)
- 97 Kategorien
- 9875 Tags
- CMS: Wordpress, self-hosted
- API: Wordpress REST-API v2
- API Documentation: <https://developer.wordpress.org/rest-api/>
- API URL: <https://dschungel-anderswelt.de/wp-json/wp/v2/>
- Hosting Provider: Hetzner
- **Download without authentication: yes**

# Strukturierte Daten herunterladen

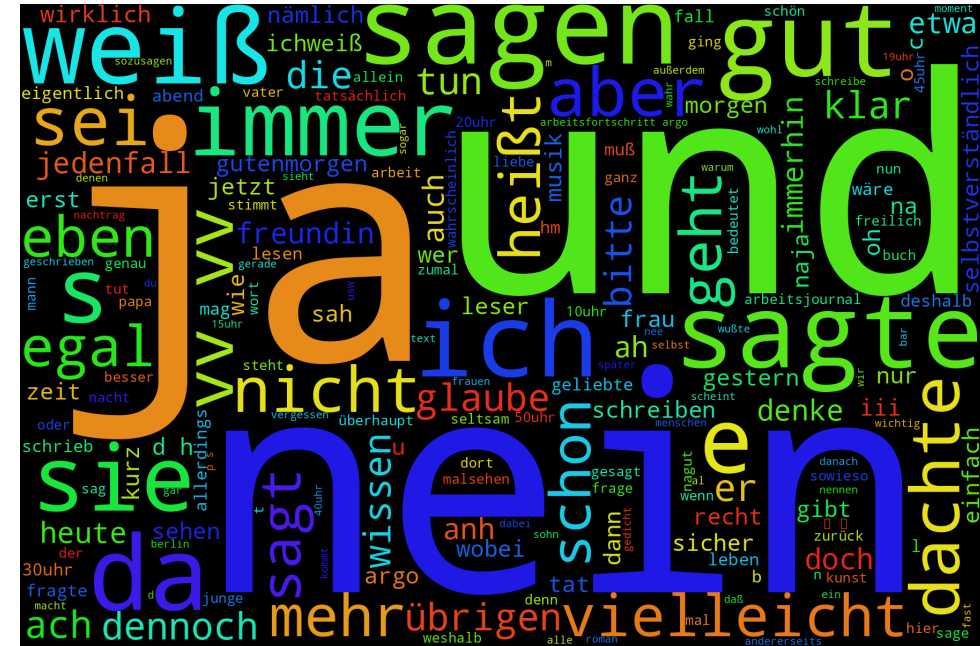
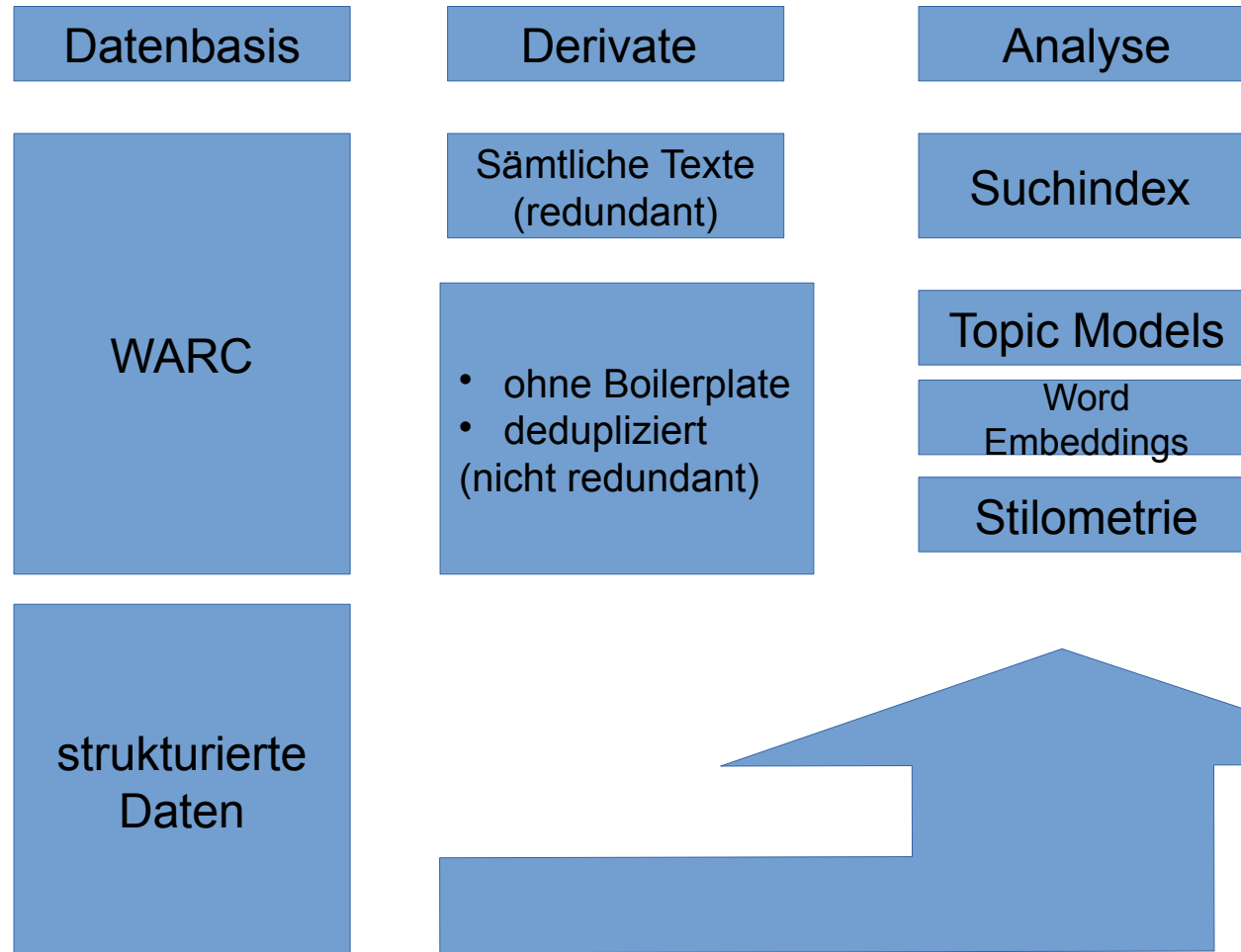
---

"Techniktagebuch" von Kathrin Passig u.a.

- ca. 7000 Einträge (Stand 11/2021)
- API: Tumblr
- **Download without authentication: yes**

|       | Webcrawl | API-Crawl   |
|-------|----------|-------------|
| Größe | 9.7 GB   | <1 MB       |
| Dauer | ~25 h    | 30 Sekunden |

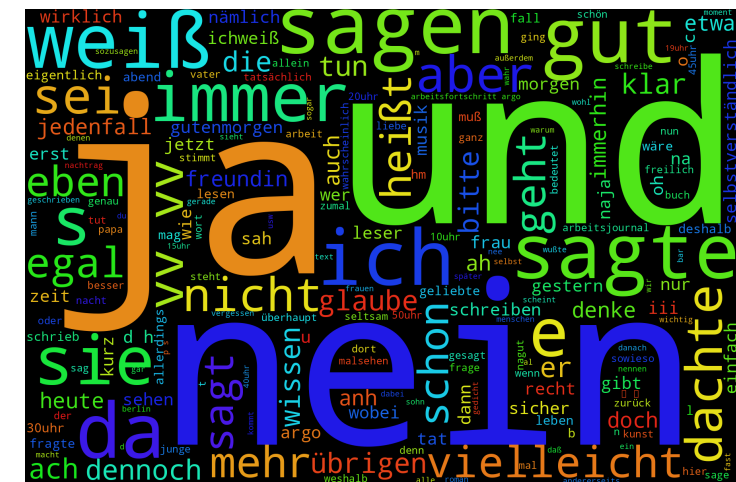
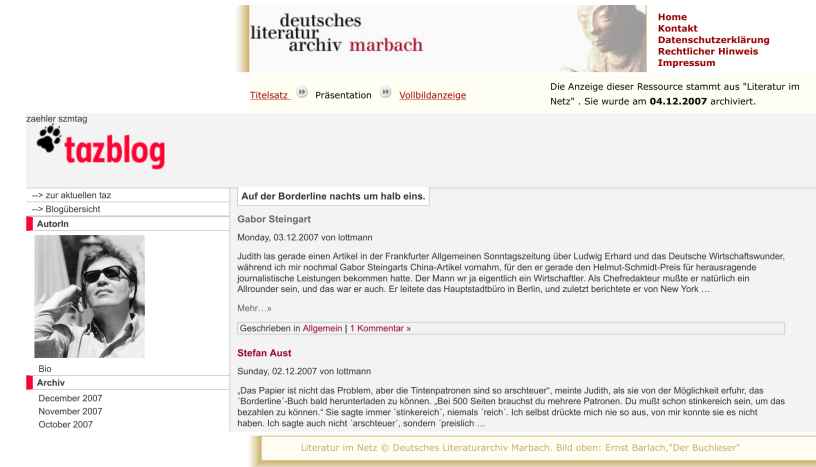
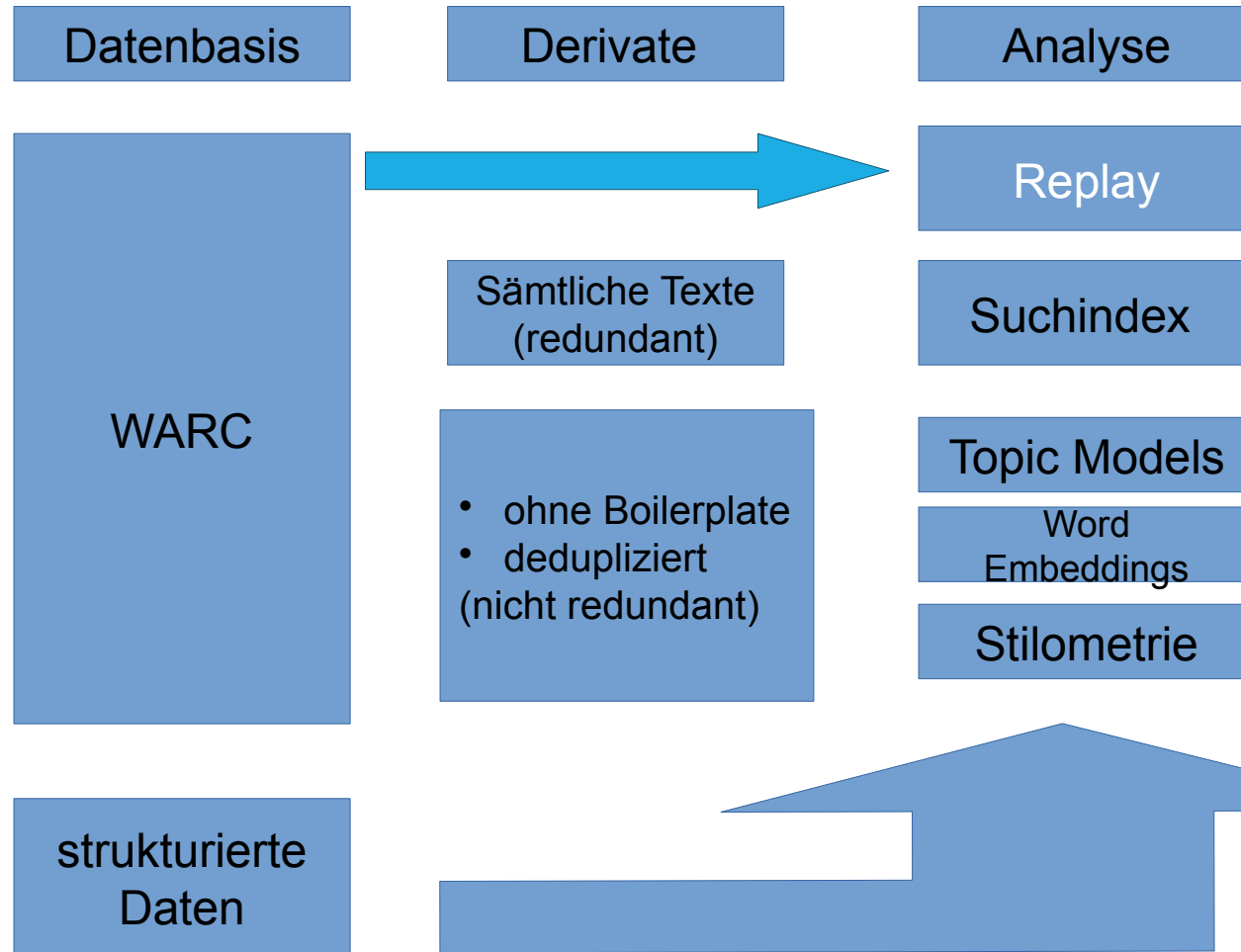
# Literatur im Netz analysieren



Wordclouds!

Wordcloud "Dschungel. Anderswelt" von Alban Nikolai Herbst, Datenbasis: Vollständiger API-Crawl vom 28.6.2022.

# Literatur im Netz analysieren



---

*Danke für Ihre Aufmerksamkeit!*

Science Data Center for Literature (SDC4Lit) <<https://sdc4lit.de>>

Claus-Michael Schlesinger <[cms@ilw.uni-stuttgart.de](mailto:cms@ilw.uni-stuttgart.de)>

Mona Ulrich <[mona.ulrich@dla-marbach.de](mailto:mona.ulrich@dla-marbach.de)>

Andre Blessing <[andre.blessing@ims.uni-stuttgart.de](mailto:andre.blessing@ims.uni-stuttgart.de)>

License: CC-BY 4.0 International

DOI: 10.5281/zenodo.7107005

---