

# Vernetzung von Daten im Deutschen Textarchiv

Susanne Haaf, Matthias Boenig, Christian Thomas,  
Alexander Geyken, Bryan Jurish, Frank Wiegand

Berlin-Brandenburgische Akademie der Wissenschaften/Deutsches Textarchiv

## Einleitung

In den vergangenen acht Jahren wurde im Rahmen des Projekts Deutsches Textarchiv<sup>1</sup> ein historisches Korpus für die Entwicklung der neuhochdeutschen Sprache aufgebaut und interessierten Nutzern in hoher Qualität, unter freien Lizenzen und mit vielfältigen Recherche- und Nachnutzungsmöglichkeiten zur Verfügung gestellt. Die DTA-Plattform zieht mittlerweile Nutzer aus ganz unterschiedlichen Disziplinen und mit vielfältigen Fragestellungen an. Sie analysieren die DTA-Korpora in Hinblick auf sprachhistorische, begriffsgeschichtliche, medienhistorische oder wissenschaftsgeschichtliche Phänomene, oder sie nutzen die Plattform, um ihre eigenen Korpora oder Editionen zu publizieren und weiter zu bearbeiten. Für diese Nutzungsszenarien bedarf es sowohl zuverlässiger Forschungsdaten und Analysetools als auch konkreter Richtlinien und Workflows für die Datenaufbereitung und -integration. Wesentliche Faktoren dabei sind zum einen sehr hohe Anforderungen an die Qualität der angebotenen Daten sowie deren Interoperabilität. So wird es möglich, Daten aus ganz verschiedenen Projektkontexten zusammenzuführen und miteinander zu vernetzen, wie im Folgenden gezeigt wird.

## Vernetzung von Daten über das DTA

Das DTA besteht aus einem Kernkorpus aus ~1600 Druckpublikationen (17.–19. Jh.) sowie mehreren Erweiterungskorpora (DTAE) auf Grundlage kuratierter Ressourcen aus externen Projekten und dem CLARIN-D-Netzwerk. Die Texterfassung für das Kernkorpus erfolgte mit sehr hoher Zuverlässigkeit im Double-Keying-Verfahren, die aufgenommenen Textquellen in DTAE wurden gemäß den DTA-Richtlinien und dem im DTA entwickelten DTA-Basisformat (DTABf)<sup>2</sup> aufbereitet und in das Korpus integriert. Diese konsequente Vereinheitlichung von Daten unterschiedlicher Herkunft stellt deren Interoperabilität sicher und erlaubt somit deren Vernetzung.

---

<sup>1</sup> Cf. die Projektwebseite [www.deutschestextarchiv.de](http://www.deutschestextarchiv.de).

<sup>2</sup> Format und Dokumentation unter [www.deutschestextarchiv.de/doku/basisformat](http://www.deutschestextarchiv.de/doku/basisformat).

Beispielsweise unternahm das genuin linguistische Projekt „Volltextdigitalisierung der Staats- und Gelehrte[n] Zeitung des Hamburgischen Unpartheyischen Correspondenten und ihrer Vorläufer (1712–1848)“<sup>3</sup> an der Universität Paderborn die Digitalisierung von 204 Zeitungsausgaben entsprechend dem Workflow und den Richtlinien des DTA; auf ähnliche Weise wurde die „Neue Rheinische Zeitung“ (301 Ausgaben) durch die Marx-Engels-Gesamtausgabe an der BBAW digitalisiert. Das am IDS Mannheim produzierte „Mannheimer Korpus Historischer Zeitungen“<sup>4</sup> (652 Ausgaben) konnte aus dem TEI-Format I5<sup>5</sup> fast vollautomatisch in das DTABf überführt werden. Die Daten dieser drei zunächst unabhängigen Projekte werden derzeit im DTA zu einem Spezialkorpus historischer Zeitungen zusammengeführt.

Aus dem Bereich der historischen Wissenschaften wurde im Rahmen des Projekts „AEDit Frühe Neuzeit“<sup>6</sup> ein Korpus von 337 Funeralschriften aus der Universitätsbibliothek Breslau entsprechend den DTA-Richtlinien erfasst und in das DTA integriert. Ergänzt wird es durch 110 Funeralgedichte aus der Staatsbibliothek zu Berlin, die dort unter bibliothekarischen Gesichtspunkten im Rahmen einer OCR-Studie<sup>7</sup> digitalisiert und im Anschluss an die Qualitätsrichtlinien des DTA angepasst wurden.

Das DTA wurde auf diese Weise durch zwei Spezialkorpora um neue Textsorten erweitert. Dabei wurden Daten unterschiedlicher Herkunft miteinander und mit der bestehenden DTA-Infrastruktur (Korpora, Analyse-Tools, Recherchemöglichkeiten) verknüpft, was je nach Format, Aufbereitungsstand und Qualität der Quelldaten mit unterschiedlich hohem Aufwand verbunden war. Voraussetzung für diese Aufgabe waren eindeutige und gut dokumentierte Richtlinien für die Texterfassung und -annotation.

## **Nutzungsmöglichkeiten historischer Daten im DTA**

Das DTA bietet vielfältige Möglichkeiten zur Recherche. Hierfür durchlaufen sämtliche Daten eine automatische linguistische Analyse (inkl. Tokenisierung, Lemmatisierung, morphologischer Analyse und orthographischer Normierung).<sup>8</sup> Mithilfe der linguistischen Suchanfragesprache DDC sind sodann komplexe Suchanfragen möglich (Abb. 1).

---

<sup>3</sup> Näheres zum Projekt unter:

<http://kw1.uni-paderborn.de/institute-einrichtungen/institut-fuer-germanistik-und-vergleichende-literaturwissenschaft/germanistik/personal/schuster/projekte/volltextdigitalisierung-der-staats-und-gelehrten-zeitung-des-hamburgischen-unpartheyischen-correspondenten-und-ihrer-vorlaeufer-1712-1851>.

<sup>4</sup> Näheres zum Projekt unter: <https://repos.ids-mannheim.de/mkhz-beschreibung.html>

<sup>5</sup> Cf. Lungen/Sperberg-McQueen 2012.

<sup>6</sup> Näheres zum Projekt unter:

<http://www.hab.de/de/home/wissenschaft/forschungsprofil-und-projekte/aedit-fruehe-neuzeit-archiv--editions--und-distributionsplattform-fuer-werke-der-fruehen-neuzeit.html>.

<sup>7</sup> Cf. Federbusch/Polzin 2013.

<sup>8</sup> Zur ling. Analyse im DTA vgl. auch <http://www.deutschestextarchiv.de/doku/software>.

1: [anonym_relation_1609:15]	... einziehen/ die vornehmsten Türcken aber auff die	<b>Galleren</b>	<b>schmieden/</b> auch den Lateinischen Bischoff vbel tractiren ...
2: [albertinus_landtstoertzer01_1615:397]	... gehalten/ aber alle jhre beste vnd gelehrteste	<b>Bücher</b>	<b>schmieden</b> sie auff jhrer Bibliothec an eysenen Ketten ...
3: [haarer_bawrenkrieg_1625:85]	... helfen/ denselben Jöckeln ließ er an ein	<b>Ketten</b>	<b>schmieden/</b> vnd bey einem Feuer lebendig/ ...
4: [kentz_handwerksboden_1629:78]	... sagt obgemelter Autor insonderheit von den Zeug- vnd	<b>Waffen</b>	<b>Schmiden/</b> daß selbige vor andern sind genennet ...
5: [gottfried_historia_1631:239]	... noch im Meerhafen gefangen genommen/ vnd in	<b>Eisen</b>	<b>schmieden</b> lassen/ vnd damit sie desto grössere ...
6: [gottfried_historia_1631:296]	.../ in Eyl auß Gold vnd Silber etliche	<b>Rüstung</b>	<b>schmiden/</b> vnd die Statt allenthalben mit starcker ...
7: [wartmann_germania0203_1650:51]	... dem Königreich Neapels nicht vnterhalten/ vnd keine	<b>Practicken</b>	<b>schmiden.</b>
8: [wartmann_germania04_1652:291]	... vmb diejenigen Geistlichen Gewissen/ die am	<b>Frieden</b>	<b>schmieden</b> od' hindern solten/ ob sie nicht ...
9: [wartmann_germania05_1653:118]	.../ welche kein andere Arbeit vornehmen/ als	<b>Schlußreden</b>	<b>schmieden/</b> vnd spintisieren/ sonderlich aber vber ...
10: [wartmann_germania05_1653:258]	... sondern man müste ohne Odem schöpfen an diesem	<b>Eysen</b>	<b>schmieden.</b>
11: [greflinger_krieg_1657:82]	... göldnen Frieden/ Lieb aber mitlerzeit viel grimme	<b>Waffen</b>	<b>schmieden/</b> und unter Holckens Hand bey zehen ...
12: [glauber_opera01_1658:376]	... so weich machen/ daß mans auff einem	<b>Amboß</b>	<b>schmieden</b> kan.
13: [buchholtz_herkules01_1659:393]	... aus tausend Markomiren könne man nicht einen einzigen	<b>Herkules</b>	<b>schmieden/</b> welches ich zu dem Ende andeute ...

Abb. 1: Abfrage: "\$p=NN schmieden with \$p=VVINF" ('Normales Nomen, gefolgt von einer Form des Worts "schmieden", dessen Wortart als infinite Verbform bestimmt wurde')

Des Weiteren können die DTA-Korpora im Zusammenhang und Vergleich mit den synchronen Korpora des DWDS<sup>9</sup> ausgewertet werden (Abb. 2).

The screenshot shows the DWDS (Das Wortauskunftssystem zur deutschen Sprache in Geschichte und Gegenwart) search interface. The search term 'schmieden' is entered in the search bar. The results are displayed in a list format, showing the first two results. The interface includes a navigation menu at the top, a search bar, and a list of results with their respective text snippets and source information.

**Wortinformationen zu „schmieden“ ...**

**Korpusbelege (DWDS-Kernkorpus)**

schmieden

**Korpus:** DWDS-Kernkorpus

**Start:** 1900 **Ende:** 1999

**Textklassen:**  Belletristik  Wissenschaft  Gebrauchsliteratur  Zeitung

**Sortierung:** Datum absteigend **Anzahl Treffer pro Seite:** 50

1-50 von 398 Treffern (556 insgesamt)

1: Moers, Walter: Die 13 1/2 Leben des Käpt'n Blaubär, Frankfurt a. M.: Eichborn 1999, S. 60  
Ich konnte eine Rede schwingen, einen Toast ausbringen, einen Schwur schwören (und wieder brechen), einen Fluch ausstoßen, einen Monolog deklamieren, einen Vers **schmieden**, ein Kompliment schleimen, Stuß reden und Unverständliches lallen.

2: Moers, Walter: Die 13 1/2 Leben des Käpt'n Blaubär, Frankfurt a. M.: Eichborn 1999, S. 153  
Professor Nachtigaller erwähnte in einem Vortrag beiläufig »Die Zyklopenkrone«, das epische Sagenwerk von

**Belege in Korpora**

**Referenzkorpora**

- DWDS-Kernkorpus (556)
- DWDS-Kernkorpus 21 (61)
- Deutsches Textarchiv (1970)

**Zeitungskorpora**

- Berliner Zeitung (978)
- Tagesspiegel (842)
- Die Zeit (2908)

**Spezialkorpora**

- Blogs (300)
- Polytechnisches Journal (789)
- Filmuntertitel

Abb. 2: Suche nach "schmieden" auf <http://zwei.dwds.de>.

<sup>9</sup> Digitales Wörterbuch der Deutschen Sprache, Website (beta): <http://zwei.dwds.de>.

Die Analyseergebnisse lassen sich mittels Wortverlaufskurven visualisieren, welche die chronologische Verteilung der Treffer nach Gattungsbereichen darstellen (Abb. 3).

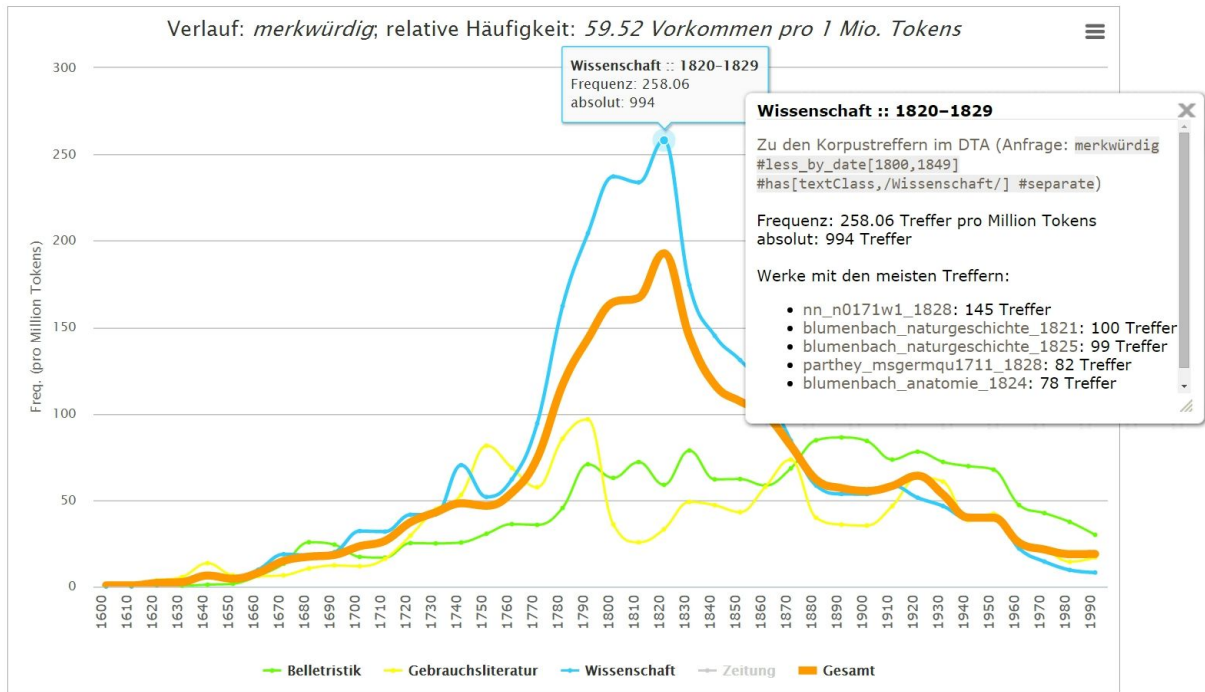
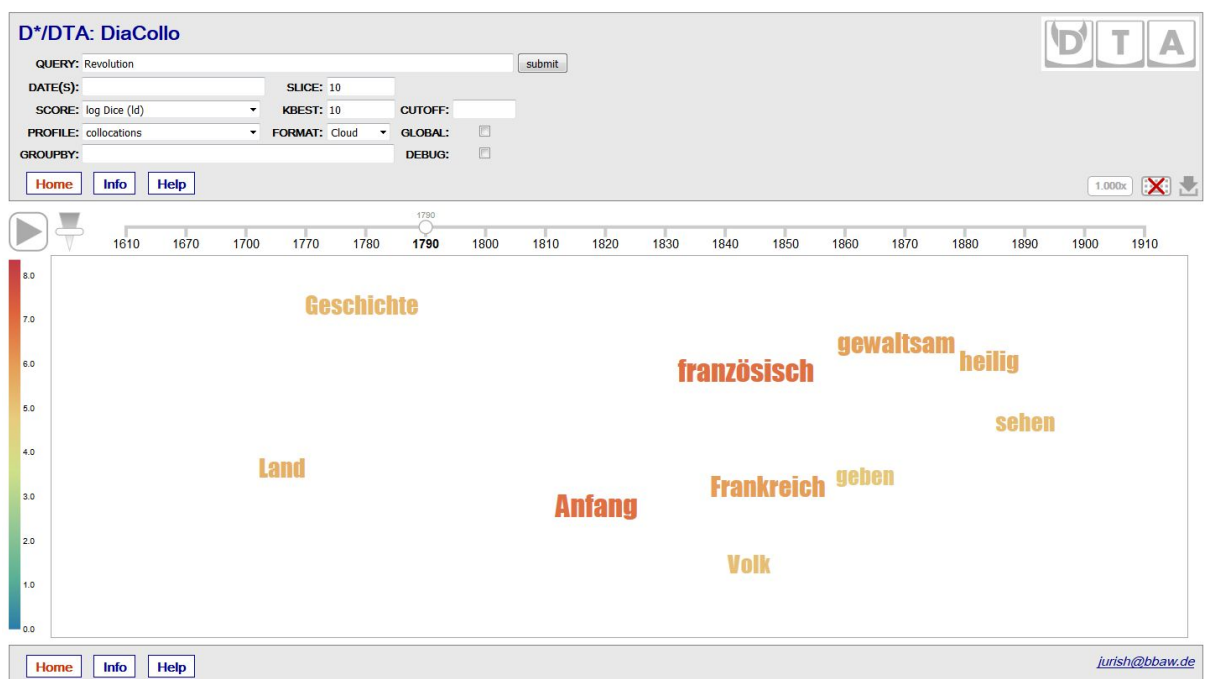


Abb. 3: Wortverlaufskurve DTA+DWDS (1600–2000),

<http://www.deutschestextarchiv.de/search/plot/?query=merkw%C3%BCrdig&mode=extended;norm=date%2Bclass&smooth=spline&single=0&grand=1&slice=10&prune=0.01&window=2&wbase=3&logavg=0&logscale=0&xrange=1600%3A2000&totals=0>

Das Tool DiaCollo erlaubt schließlich die Betrachtung von Termini im Zusammenhang mit den für sie typischen Kollokationen im zeitlichen Verlauf. Hierbei sind auch vergleichende Analysen von Termini möglich (Abb. 4 & 5).



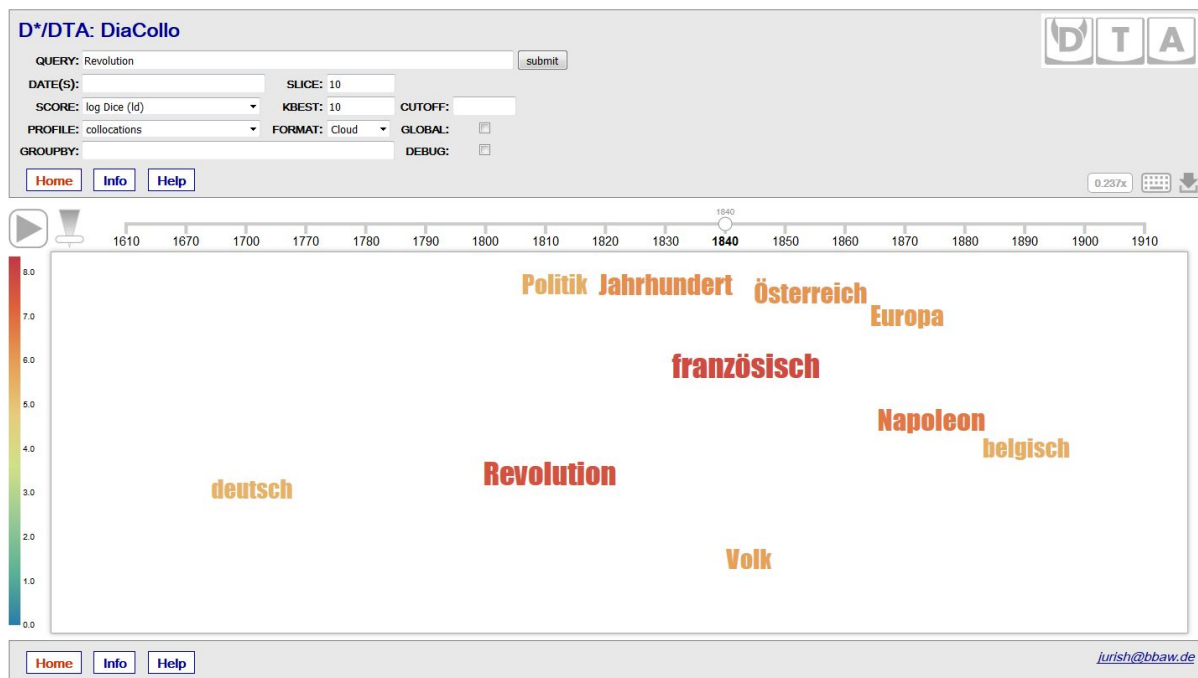


Abb. 4 & 5: DiaCollo-Ansicht "Revolution" mit seinen häufigsten Kollokationen, Zeitschnitte 1790 und 1840

## Integration von Daten in die CLARIN-Infrastruktur

Über das Zentrum Sprache der BBAW ist das DTA Partner im Verbundprojekt CLARIN-D.<sup>10</sup> Diese Anbindung an die CLARIN-Infrastruktur eröffnet Nutzern ausgehend vom DTA weitere Recherche- und Analysemöglichkeiten.

Hierfür stellt das DTA sämtliche Texte auch in den CLARIN-kompatiblen Formaten CMDI (Metadaten) und TCF (Textdaten inkl. linguistischer Analyse) bereit. Die Metadaten können beispielsweise durch das CLARIN-übergreifende Virtual Language Observatory<sup>11</sup> geharvestet werden. Die Textdaten sind über CLARINs Federated Content Search<sup>12</sup> zugänglich und somit im Vergleich mit den Korpora anderer CLARIN-Zentren durchsuchbar. Die TCF-Daten können in der CLARIN-Umgebung WebLicht<sup>13</sup> weiter analysiert werden. Ein TEI-to-TCF-Converter ermöglicht die Konvertierung von TEI-Daten in das TCF-Format und vice versa.

<sup>10</sup> Cf. <http://clarin.bbaw.de>.

<sup>11</sup> Cf. <https://vlo.clarin.eu>.

<sup>12</sup> Cf. <http://weblicht.sfs.uni-tuebingen.de/Aggregator/>.

<sup>13</sup> Cf. [http://weblicht.sfs.uni-tuebingen.de/weblichtwiki/index.php/Main\\_Page](http://weblicht.sfs.uni-tuebingen.de/weblichtwiki/index.php/Main_Page).

## Bibliographie

*Federbusch/Polzin 2013*: Maria Federbusch und Christian Polzin: Volltext via OCR – Möglichkeiten und Grenzen. Testszenarien zu den Funeralschriften der Staatsbibliothek zu Berlin – Preußischer Kulturbesitz. Mit einem Erfahrungsbericht von Thomas Stäcker aus dem Projekt „Helmstedter Drucke Online“ der Herzog August Bibliothek Wolfenbüttel. Berlin 2013 (=Beiträge aus der Staatsbibliothek zu Berlin – Preußischer Kulturbesitz 43).  
[http://staatsbibliothek-berlin.de/fileadmin/user\\_upload/zentrale\\_Seiten/historische\\_drucke/pdf/SBB\\_OCR\\_STUDIE\\_WEBVERSION\\_Final.pdf](http://staatsbibliothek-berlin.de/fileadmin/user_upload/zentrale_Seiten/historische_drucke/pdf/SBB_OCR_STUDIE_WEBVERSION_Final.pdf)

*Lüngen/Sperberg-McQueen 2012*: Harald Lüngen und C. M. Sperberg-McQueen: A TEI P5 Document Grammar for the IDS Text Model. In: Journal of the Text Encoding Initiative 3 (2012): TEI and Linguistics. <https://jtei.revues.org/508>

Für Publikationen zum DTA, zu den Richtlinien und Tools sowie zu angrenzenden Projekten cf. [www.deutschestextarchiv.de/doku/publikationen](http://www.deutschestextarchiv.de/doku/publikationen).