

# Deutsches Referenzkorpus zur internetbasierten Kommunikation: Fragen der Standardisierung und Datenerhebung

2. DTA-/CLARIN-D-Konferenz und CLARIN-D-Workshop

TEXTKORPORA IN INFRASTRUKTUREN FÜR DIE GEISTES- UND  
SOZIALWISSENSCHAFTEN

Michael Beißwenger (Dortmund)  
Lothar Lemnitzer (Berlin)

## Internetbasierte Kommunikation (IBK) *engl. Computer-Mediated Communication (CMC)*

Im Fokus der Erforschung internetbasierter Kommunikation stehen nicht beliebige Webtexte, sondern die sprachlichen Äußerungen in dialogischen Webgenres wie z.B. Online-Foren, Chats, Instant Messaging und Wiki-Diskussionen, in Twitter-Postings, in Kommentaren und Diskussionen in Weblogs, auf Videoplattformen (*YouTube*) und auf den Profelseiten „sozialer Netzwerke“ (*Facebook, MySpace*) sowie in multi-modalen Kommunikationsumgebungen wie *Skype*, MMORPGs und „virtuellen Welten“ (*SecondLife* u.a.).

## Aus der Perspektive der Linguistik (Medienlinguistik etc.):

Für Forschungszwecke allgemein und frei zugängliche Korpora bieten mehrere Vorteile:

- (1) Der/die Forscher/in spart Zeit und muss Datensets/Korpora nicht selbst erheben und für Analysezwecke aufbereiten, bevor er/sie sich den eigentlich interessanten Fragen zuwenden kann;
- (2) Untersuchungen zu bestimmten Phänomenen sind replizierbar und vergleichbar

## Im Bereich der empirischen Erforschung und lexikographischen Beschreibung der deutschen Gegenwartssprache:

- Existierende Korpusansammlungen zur deutschen Gegenwartssprache decken IBK nicht ab. Ein wichtiger Kommunikationsbereich ist somit in den verfügbaren Korpora nicht erfasst.
  - ⇒ Kann man die deutsche Gegenwartssprache zum Stand 2014 beschreiben, wenn man den Kommunikationsbereich „IBK“ (aufgrund fehlender Abdeckung in Korpora) einfach ausspart?

# Warum sind bisher so wenige IBK-Korpora allgemein verfügbar?

- ... weil die **rechtlichen Rahmenbedingungen** für die Erhebung, Speicherung, Bearbeitung und Bereitstellung von IBK-Daten für Forschungszwecke bislang noch größtenteils unklar sind.
- ... weil IBK eine „moving target“ ist: schnelle Veränderlichkeit von Kommunikationstechnologien und Genres  
⇒ hohe Anforderungen an **Metadaten / Dokumentation**
- ... weil sprachtechnologische Verfahren für die **linguistische Basisannotation** von schriftlichen Korpusdaten (Wortatenannotation, Lemmatisierung, syntaktische Strukturen) mit den Besonderheiten von IBK-Schriftlichkeit bislang nicht sinnvoll umgehen können.
- ... weil bislang **keine Standards für die Annotation der strukturellen Besonderheiten von IBK-Genres existieren** (Threads in Foren / sozialen Netzwerken, Chat-Logfiles, Twitter-Timelines, Wikipedia-Diskussionen, multi-modale Genres usw.) ⇒ Wer ein IBK-Korpus aufbaut, kann i.d.R. keine existierenden Schemata übernehmen, sondern muss selbst Lösungen entwickeln.

⇒ **Der Aufbau von IBK-Korpora ist ein eigener methodischer Teilbereich der Korpuslinguistik.**



## Projekt: *Deutsches Referenzkorpus zur internetbasierten Kommunikation (DeRiK)*

**Ziel:** Aufbau einer Zusatzkomponente zu den Korpora des Projekts „Digitales Wörterbuch der deutschen Sprache“ (DWDS), die Sprachdaten aus den wichtigsten Genres internetbasierter Kommunikation umfasst und gemeinsam mit den bereits existierenden Korpusressourcen abgefragt werden kann.

( M. Beißwenger, A. Geyken, L. Lemnitzer, A. Storrer )



⇒ Schließung der „IBK-Lücke“ in den DWDS-Korpora zur deutschen Gegenwartssprache.

Vgl. Beißwenger et al. (2013), Beißwenger/Lemnitzer (2013).



http://www.dwds.de

Ressourcen ▼ Erschließung ▼ Projekt ▼ Aktuelles ▼

Troll



DWDS Standardsicht

+Ressourcen



DWDS-Wörterbuch

Troll

mask., -s/-es, -e

Aussprache: ▶

gespenstisches Wesen des Volksaberglaubens, besonders in der nordischen Mythologie, Unhold

Die Trolle sind die Anwälte der Tiere. Sie suchen den heim, der Tiere quält – Jahnn in Dt. Erzähler 2,89

Dazu:

- Trollblume

Version: 0.4.22 | Quelle: WDG | Artikeltyp: Vollartikel

Etymologisches Wörterbuch (nach Pfeifer)

Troll

**Troll** m. (in der nordischen Mythologie) 'dämonisches Wesen, Kobold', anord. *troll*, *tröll*, schwed. *troll*, dän. *trol/d*. Die im 18. Jh. ins Dt. entlehnte Bezeichnung trifft auf einheimisches (in den Mundarten bewahrtes) *Troll* 'grober, ungeschlachter Kerl' (15. Jh.), auch 'Dämon, Unhold' (15. bis 17. Jh.). Vielleicht hervorgegangen aus germ. \**truz/a-* und verwandt mit *trollen* (s. d.).

Version: 1.0.65

DWDS-Wortprofil 3.0

Abfragewort: Troll

Vergleichswort:

Substantiv

logDice



36

Überblick zu 'Troll'

aussehen wie böser Dämonen Elfen  
erzählte von Feen Fjorde Geister Gnome  
Hexen Koblode Land Reich Riesen schwarze  
tolle versteinerten Welt Zwerge

'Troll' hat Attribut

Koordination mit 'Troll'

'Troll' ist Aktivsubjekt von

Version: 3.0

Einstellungen

Kernkorpus 20

Treffer: 133, davon anzeigbar: 98

Filter

KWiC	Datum ↓	Datum ↑	Zufällig	Links	Rechts
1	[1999]				erliegende , Chemluth und dem Troll machte das nichts aus , abe
2	[1999]				Dann gab er dem Troll die Fackel und verschwand ir
3	[1999]				t eine Abkürzung « , sagte der Troll , als hätte er meine Gedanke
4	[1999]				r Drache schlief , trampelte der Troll auf ihm herum , als sei er eir
5	[1999]				schrie der Troll atemlos während unseres nä
6	[1999]				anderen nach « , hechelte der Troll .
7	[1999]				rief der Troll .
8	[1999]				Der Troll stieg an mir hoch .
9	[1999]				Chemluth und der Troll sahen durch das Loch zu mir
10	[1999]				verkehrt gemacht « , sagte der Troll .
11	[1999]				Chemluth und der Troll waren verschwunden , und ic
12	[1999]				Der Troll zog Chemluth hoch , und bev
13	[1999]				rief der Troll .
14	[1999]				» Ja , ich weiß « , sagte der Troll kleinlaut .

DIE ZEIT

Treffer: 648

Filter

KWiC	Datum ↓	Datum ↑	Zufällig	Links	Rechts
1	[2009]				Ob Trolle überhaupt vom Influenzaviru
2	[2009]				Harry und Weasley hatten den Troll zur Strecke gebracht , indem
3	[2009]				Weil sie dort Hexen , Trolle und Teufel vermuteten .
4	[2009]				cht die Figur des Paladins den Troll besiegt : » Ich war das « ; ni
5	[2009]				lf und Sauron stehen könnten : Trolle kämpfen gegen Krieger , Kr
6	[2009]				Internet Unter Trollen Wo Jungen zu Männern ui
7	[2009]				Unter Trollen
8	[2009]				ütteln den Kopf , bevor sie sich trollen .
9	[2009]				auf die Elfen zu achten und die Trolle , die seien in diesem Jahr e
10	[2008]				rgischen Einläufen ( von Lothar Trolle ) auf eine Länge von fünf Sti
11	[2008]				ndenklar , da wird es für einen Troll ein Leichtes gewesen sein ,
12	[2008]				entstanden sein , als ein böser Troll die Stadt Farum begraben wc
13	[2008]				en Universitätscampus voll mit Trollen , sprechenden Tieren und
14	[2008]				vor Heinrich Böll und Thaddäus Troll getan hatten , aus Sicht des

DeRiK

Treffer: ???

Filter

KWiC	Datum ↓	Datum ↑	Zufällig	Links	Rechts
1					
2					
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					

... under construction ...

DeRiK ist konzipiert als

- ein **zeitlich gestaffeltes Korpus**, bei dem die Datenerhebung nicht nur einmalig, sondern mehrfach in regelmäßigen Abständen erfolgen soll, wodurch es möglich wird, auch sprachlichen Wandel *innerhalb* der internetbasierten Kommunikation darzustellen;
- ein spezialisiertes **Referenzkorpus** zur internetbasierten Kommunikation, das der Fachcommunity als Basis für korpusgestützte Untersuchungen und für die Vermittlung der sprachlichen Besonderheiten von IBK-Genres in der Lehre zur Verfügung gestellt werden soll;

## Zusammensetzung des DWDS-Kernkorpus:

⇒ Ideal: **Ausgewogenheit** des Korpus nach Zeitabschnitten und Textsortenbereichen:

⇒ Für jede Dekade (ab 1900):

Anteil Zeitungstexte:	~27%
Anteil Werke der Literatur:	~28%
Anteil Fachtexte:	~23%
Anteil Gebrauchstexte:	~21%

Vgl. Geyken (2007), Klein/Geyken (2010)



## Offene Fragen:

- ⇒ Was sind relevante „IBK-Sorten“ (vergleichbar zu Textsortenbereichen)?
- ⇒ Wie lässt sich ein *ausgewogenes Verhältnis* festlegen?

## Idee:

Kopplung des Schlüssels für die Auswahl und Zusammensetzung der Korpusdaten an die jährlich durchgeführte **ARD/ZDF-Onlinestudie**.



## Idealer Schlüssel:

- Ermittlung der meistgenutzten Kommunikationsanwendungen im Internet im Untersuchungszeitraum
- Gewichtung der Präferenzen unterschiedlicher Altersgruppen nach dem Grad ihrer ‚Online-Affinität‘
- Neuberechnung des Schlüssels für die Datenaufnahme auf Basis der jeweils aktuellen Ausgabe der Studie

## Offene Fragen:

- ⇒ Was sind relevante „IBK-Sorten“ (vergleichbar zu Textsortenbereichen)?
- ⇒ Wie lässt sich ein *ausgewogenes Verhältnis* festlegen?

## Pragmatischer Schlüssel:



Angesichts der unklaren Rechtslage zunächst nur Aufnahme von Daten aus solchen Anwendungen, bei denen die Nutzung juristisch unbedenklich ist (= Anwendungen, bei denen alle Inhalte explizit für eine Nutzung, Bearbeitung und Wiederbereitstellung durch Dritte lizenziert sind).

## Idealer Schlüssel:

- Ermittlung der meistgenutzten Kommunikationsanwendungen im Internet im Untersuchungszeitraum
- Gewichtung der Präferenzen unterschiedlicher Altersgruppen nach dem Grad ihrer ‚Online-Affinität‘
- Neuberechnung des Schlüssels für die Datenaufnahme auf Basis der jeweils aktuellen Ausgabe der Studie

DeRiK ist konzipiert als

- ein **zeitlich gestaffeltes Korpus**, bei dem die Datenerhebung nicht nur einmalig, sondern mehrfach in regelmäßigen Abständen erfolgen soll, wodurch es möglich wird, auch sprachlichen Wandel *innerhalb* der internetbasierten Kommunikation darzustellen;
- ein **Referenzkorpus** zur internetbasierten Kommunikation, das der Fachcommunity als Basis für korpusgestützte Untersuchungen und für die Vermittlung der sprachlichen Besonderheiten von IBK-Genres in der Lehre zur Verfügung gestellt wird;
- ein **annotiertes Korpus**, das neben einer linguistischen Basisannotation auch Annotationen zu charakteristischen sprachlichen und strukturellen Besonderheiten bei der Sprachverwendung im Netz umfasst.

- Für die Annotation von Primärdaten aus Genres internetbasierter Kommunikation gibt es derzeit noch **keine Standards**.
- Linguistische und texttechnologische Modelle, die für die Annotation von Textkorpora etabliert sind (Linguistische Basisannotation, Strukturbeschreibung von Texten), decken die strukturellen und linguistischen Besonderheiten internetbasierter Kommunikation nicht vollständig ab (Beißwenger et al. 2012).
- Sprachtechnologische Verfahren für die automatische linguistische Annotation von Textkorpora führen bei IBK-Daten zu nicht akzeptablen Ergebnissen (Phänomene „nichtstandardisierter“ Schriftlichkeit; vgl. Bartz et al. 2013).

## Desiderate:

- 1) Anpassung texttechnologischer Standards für die Repräsentation von IBK-Genres und -Korpora (Strukturannotation, Metadaten)  
⇒ *Interoperabilität*
- 2) Anpassung von Tagsets und Verfahren/Werkzeugen für die linguistische Annotation ⇒ *Suche über linguistischen Annotationen*



## Computer-Mediated Communication SIG

### Contents

- [Context](#)
- [Scope and Tasks](#)
- [Convener](#)
- [Wiki space and mailing lists](#)
- [Activities and cooperations](#)

### Context

In the past three decades, computer networks and especially the internet have enabled *computer-mediated communication*, henceforth "CMC". Even though there's been significant progress in the sciences as well as in the field of natural language processing, there's still a need for communication and their structural and linguistic peculiarities. Being a broadly acknowledged within the field of digital humanities will allow for the development of standards for different languages

### Scope and Tasks

This special interest group is elaborating on suggestions for adapting existing TEI Schemas. The focus of the group's work is on (but not limited to) tasks such as

- modelling user contributions (*posts*) to written CMC dialogues (e.g. chat logs)
- modelling CMC document structures ("*CMC macrostructures*" – e.g. forum threads)
- annotating linguistic features within user posts ("*CMC microstructures*" – e.g. emoticons)
- representing linked data and media objects connected with/mentioned in CMC
- metadata schemata for the description of CMC resources;
- developing perspectives for the representation of discourse in a variety of modalities from written, spoken and non-verbal modes of communication

### Convener

[Michael Beißwenger](#), TU Dortmund University

### Wiki space and mailing lists

For exchange on the issues and tasks listed above, the SIG uses the talk pages in the TEI wiki and a mailing list.

## Special Interest Group im Rahmen der TEI:

seit 2013: Erarbeitung eines Entwurfs zu einem **TEI-Standard für die Annotation von IBK-Genres**

(unter Berücksichtigung sowohl schriftlicher als auch multimodaler Genres)

... auf der Grundlage eines spezialisierten TEI-Schemas, das 2011-12 im Rahmen des DeRiK-Projekts entwickelt wurde (Beißwenger et al. 2012)

... und das in 2013-14 von SIG-Mitgliedern am IDS und an der Universität Clermont-Ferrand für weitere Typen von IBK-Korpora (Deutsch und Französisch) angepasst wurde (Margaretha/Lüngen (forthc.); Chanier et al. (forthc.))

Beteiligte aus Korpusprojekten zu versch. Sprachen (CoMeRe, **DeRiK**, SoNaR, Web2Corpus\_it, Dortmunder Chat-Korpus, Mannheimer Wikipedia-Korpus, ...)

## Community Shared Task

zur automatischen linguistischen Annotation von IBK-Daten (2015), initiiert von Mitgliedern des Empirikom-Netzwerks und unterstützt durch die GSCL

## Fokus:

- Tokenisierung
- Part-of-Speech-Annotation

## Grundlage:

- Handannotierte Trainings- und Evaluationsdaten („Goldstandard“)

## Ziel:

- Anpassung automatischer Verfahren an den Umgang mit IBK-Schriftlichkeit (in einem kompetitiven Szenario)

Tag	Kategorie	Beispiele
<b>I. Tags für IBK-spezifische Phänomene:</b>		
EMO ASC	Emoticon, als Zeichenfolge dargestellt (Typ „ASCII“)	:-) :-( ^ O.O
EMO IMG	Emoticon, als Grafik-Ikon dargestellt (Typ „Image“)	😊 🍌 😊
AKW	Aktionswort	*lach*, freu, grübel, *lol*
HST	Hashtag	Kreta war super! <a href="#">#urlaub</a>
ADR	Adressierung	<a href="#">@lothar</a> : Wie isset so?
URL	Uniform Resource Locator	<a href="http://www.tu-dortmund.de">http://www.tu-dortmund.de</a>
EML	E-Mail-Adresse	<a href="mailto:peterklein@web.de">peterklein@web.de</a>
<b>II. Tags für Phänomene der konzeptionellen Mündlichkeit:</b>		
VV PPER	Tags für die häufigsten Bildungsmuster kontraktierter Formen (APPRART ist in STTS bereits vorhanden)	schreibste, machste
APPR ART		vorm, überm, fürn
VM PPER		willste, darfst, musste
VA PPER		haste, biste, isses
KOUS PPER		wenns, weils, obse
PPER PPER		ichs, dus, ers
ADV ART		son, sone
PTK IFG	Intensitäts-, Fokus- oder Gradpartikel	<u>sehr</u> schön, <u>höchst</u> eigenartig, <u>nur</u> sie, <u>voll</u> geil
PTK MA	Modal- oder Abtönungspartikel	Das ist <u>ja</u> / <u>vielleicht</u> doof. Ist das <u>denn</u> richtig so? Das war <u>halt</u> echt nicht einfach.
DM	Diskursmarker	<i>prototypisch: weil, obwohl, nur, also als Einheiten mit projektivem Potenzial im Vorvorfeld von V2-Sätzen</i>
ONO	Onomatopoetikon	boing, miau, zisch

STTS-Tagset mit Erweiterungen für IBK und für gesprochene Sprache (Beißwenger, Bartz, Storrer, Westpfahl; forthc.)

## a) Aufgabenbereich „Repräsentation“:

- Ausarbeitung des in der TEI-SIG diskutierten Schemas (in Kooperation mit den KollegInnen aus F, IT, NL und aus dem IDS) und Anpassung für die Repräsentation des Dortmunder Chat-Korpus (→ CLARIN-D-Kurationsprojekt)
- Ziel: Bis Ende 2015 Dokumentation eines vollständigen Entwurfs, der in verschiedenen Korpora zum Deutschen und zum Französischen umgesetzt wurde und Formulierung eines Vorschlags für die Erweiterung der TEI-Guidelines um entsprechende Modelle für den Bereich IBK

## b) Korpusaufbau:

- Anwendung von Verfahren zum automatischen Crawling von CC-lizenzierten IBK-Ressourcen (Blogkommentare, soziale Netzwerke, Foren, ...) inkl. Lizenzerkennung und Spracherkennung
- In 2015: Erhebung eines ersten Teilkorpus für den Zeitraum 2010-2015
- Integration des „Dortmunder Chat-Korpus“ als Teilkorpus für den Zeitraum vor 2010 (Kurationsprojekt im Rahmen von CLARIN)

## c) Aufgabenbereich „Linguistische Basisannotation“:

- eigene Vorarbeiten zur Anpassung sprachtechnologischer Verfahren an IBK-Schriftlichkeit an der BBAW
- Mitarbeit im GSCL-AK „Social Media / Internetbasierte Kommunikation“
- 2015: Linguistische Basisannotation für das Dortmunder Chat-Korpus (unter Nutzung aktueller Ansätze, die z. B. im Rahmen der GSCL-Shared-Task entwickelt werden; s. auch die Ansätze, die beim KONVENS-Workshop zu IBK im Okt. 2014 diskutiert wurden)
- Nutzung von Verfahren aus dem BMBF-Verbundprojekt „Korpus-basierte linguistische Recherche und Analyse mit Hilfe von Data-Mining“ (KobRA, Leitung: Angelika Storrer)



<http://www.kobra.tu-dortmund.de>



## DeRiK:

- Beißwenger, Michael; Ermakova, Maria; Geyken, Alexander; Lemnitzer, Lothar; Storrer, Angelika (2013): **DeRiK: A German Reference Corpus of Computer-Mediated Communication**. In: Literary and Linguistic Computing (DOI: 10.1093/lc/fqt038).  
<http://llc.oxfordjournals.org/cgi/reprint/fqt038?ijkey=GXFixqNNy0uW7cO&keytype=ref>
- Beißwenger, Michael; Lemnitzer, Lothar (2013): **Aufbau eines Referenzkorpus zur deutschsprachigen internetbasierten Kommunikation als Zusatzkomponente für die Korpora im Projekt „Digitales Wörterbuch der deutschen Sprache“ (DWDS)**. In: Journal for Language Technology and Computational Linguistics 26 (2), 1-22. [http://www.jlcl.org/2013\\_Heft2/1BeiLem.pdf](http://www.jlcl.org/2013_Heft2/1BeiLem.pdf)

## DeRiK-TEI-Schema:

- Beißwenger, Michael; Ermakova, Maria; Geyken, Alexander; Lemnitzer, Lothar; Storrer, Angelika (2012): **A TEI Schema for the Representation of Computer-Mediated Communication**. In: Journal of the Text Encoding Initiative (jTEI), Issue 3, November 2012 (DOI: 10.4000/jtei.476).  
<http://jtei.revues.org/476>

## Projekte / Netzwerke:

- DFG-Netzwerk Empirikom: <http://www.empirikom.net>
- DWDS: <http://www.dwds.de>
- Dortmunder Chat-Korpus: <http://www.chatkorpus.tu-dortmund.de>
- TEI-SIG „Computer-Mediated Communication“: <http://www.tei-c.org/Activities/SIG/CMC/>
- BMBF-Projekt *KobRA*: <http://www.kobra.tu-dortmund.de>
- GSCL-AK „Social Media / Internetbasierte Kommunikation“: <http://gscl.org/ak-ibk.html>



## Wissenschaftliches Netzwerk: Empirische Erforschung internetbasierter Kommunikation

### Aktuelle Aktivitäten und Initiativen

- 7. Arbeitstagung des Empirikom-Netzwerks: "**Social Media Corpora for the eHumanities: Standards, Challenges, and Perspectives**" (TU Dortmund, 20./21.2.2014). [Tagungsprogramm](#)
- Neuer Arbeitskreis der Gesellschaft für Sprachtechnologie und Computerlinguistik (GSCL): "[Social Media / Internetbasierte Kommunikation](#)"
- Special Topic Panel im Rahmen der TEI-Konferenz 2013 in Rom: "[Computer-Mediated Communication in TEI: What Lies Ahead](#)"
- Empirikom-Workshop im Rahmen der GSCL-Konferenz 2013 in Darmstadt: "[Verarbeitung und Annotation von Sprachdaten aus Genres internetbasierter Kommunikation](#)"
- Neue Special Interest Group "[Computer-Mediated Communication](#)" als Teil der [Text Encoding Initiative \(TEI\)](#)
- Empirikom-Projekt: [Vorbereitung einer Shared Task zur automatischen linguistischen Annotation von IBK-Daten](#)
- DGFS-AG "Modellierung nichtstandardisierter Schriftlichkeit" (März 2013): Die Dokumentation der Vorträge [ist online](#).

### Über das Netzwerk

[ [Description of the network in English](#) ]

Im Internet und speziell in den sozialen Netzwerken des Web 2.0 entstehen neue Formen der internetbasierten Kommunikation, die interdisziplinär erforscht werden. Das wissenschaftliche Netzwerk "Empirische Erforschung internetbasierter Kommunikation" vereint [fünfzehn Forscherinnen und Forscher aus zwölf verschiedenen Hochschulen und Forschungseinrichtungen](#), die an den [theoretischen und methodologischen Grundlagen der datengestützten Analyse internetbasierter Kommunikation](#) arbeiten und dafür Techniken und Methoden aus Korpuslinguistik, Computerlinguistik und Informatik nutzen. Aufgrund des digitalen Ausgangsformats sind Datensammlungen zur internetbasierten Kommunikation zwar zunächst einfach zu erheben; es fehlen aber bisher Standards sowie Annotations- und Analysekatoren, um die **sprachlichen und interaktionalen Besonderheiten in neuen Kommunikationsformen wie z.B. E-Mail, Instant Messaging, Chats, Twitter, Weblogs, Skype sowie Diskussionen in Foren, Wikis und sozialen Netzwerken** zu erfassen. Außerdem müssen existierende Verfahren zur automatischen Aufbereitung und Verarbeitung von Sprachdaten, die häufig für standardsprachliche Schrifttexte entwickelt sind, an die sprachlichen Besonderheiten von internetbasierter Schriftlichkeit angepasst werden

**Ziel des Netzwerks**, das 2010-2014 durch die [Deutsche Forschungsgemeinschaft \(DFG\)](#) gefördert wird, ist es, Kompetenzen aus germanistischer Sprachwissenschaft, Computerlinguistik, Informatik und Psychologie zu bündeln, um anhand einer Reihe [konkreter Forschungsfragen Lösungsansätze zu offenen Fragen beim Aufbau, bei der Aufbereitung und bei der Analyse von Korpora internetbasierter Kommunikation](#) zu entwickeln und **Vorschläge für Standards zur Modellierung und Annotation ihrer sprachlichen und strukturellen Besonderheiten** zu erarbeiten. Der Fokus liegt auf dem Deutschen; sprachübergreifende Fragen werden aber in engem Austausch mit Projekten zu anderen Sprachen und aus dem europäischen Ausland diskutiert.

Die Ergebnisse der Arbeit im Netzwerk werden in Publikationen und online dokumentiert.

## Dortmunder Chat-Korpus

<http://www.chatkorpus.tu-dortmund.de>

Ergebnis eines  
Lehrstuhlprojekts an der TU  
Dortmund (2002-2008)  
(A. Storrer / M. Beißwenger)

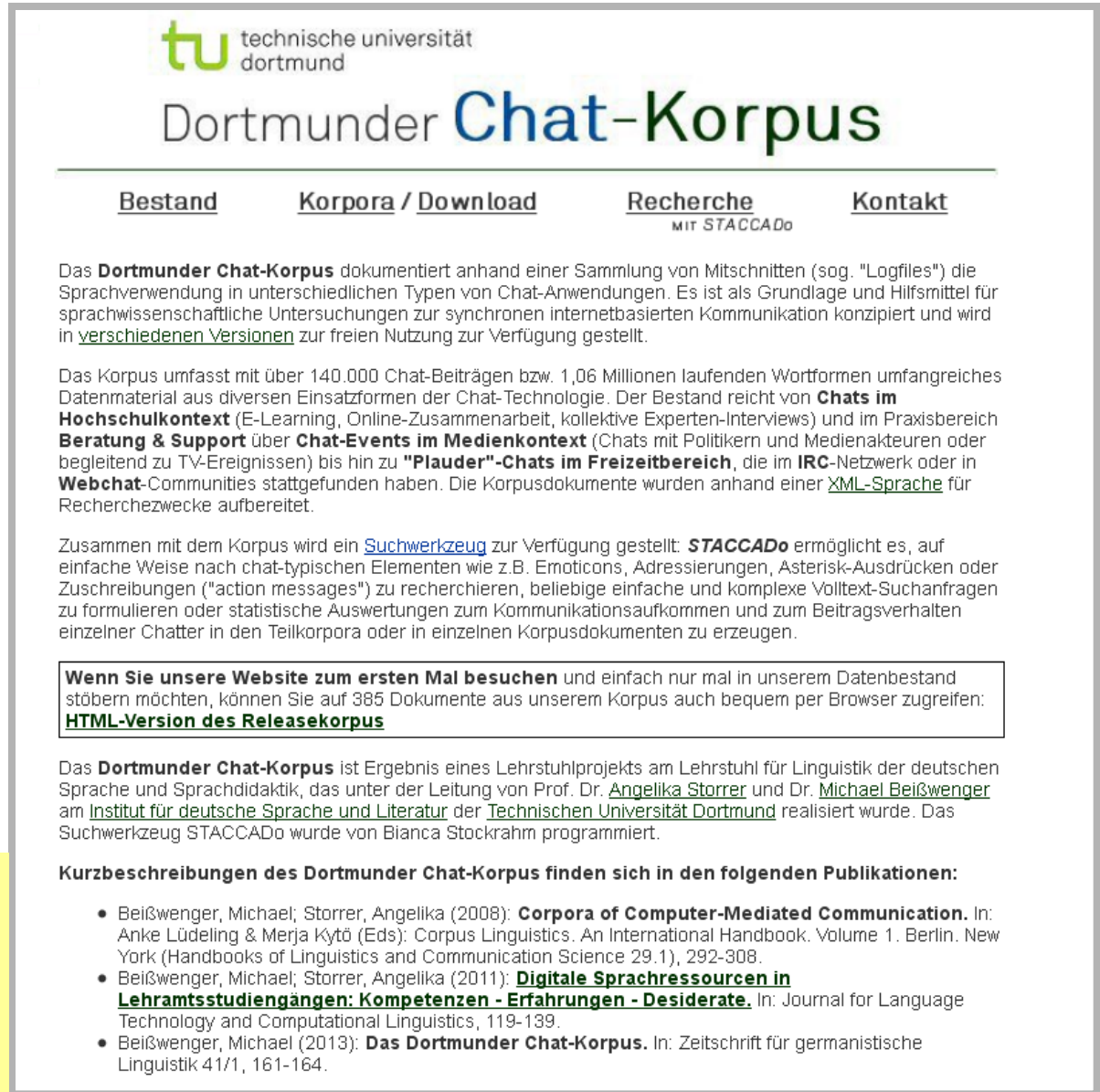


CLARIN-D-Kurationsprojekt  
(F-AG 1 „Deutsche Philologie“, 2015):

### **ChatCorpus2CLARIN:**

Integration des Dortmunder  
Chat-Korpus in die CLARIN-D-  
Korpusinfrastrukturen am  
Institut für deutsche Sprache  
(IDS) und an der Berlin-  
Brandenburgischen Akademie  
der Wissenschaften (BBAW)

⇒ Anpassung des aktuellen  
Standes des Schemas aus  
der TEI-SIG + Integration  
der Ressource u.a. in  
DeRiK/DWDS



The screenshot shows the homepage of the Dortmund Chat-Korpus project. At the top left is the TU Dortmund logo. The main title is 'Dortmunder Chat-Korpus'. Below the title are four navigation links: 'Bestand', 'Korpora / Download', 'Recherche MIT STACCADo', and 'Kontakt'. The main content area contains three paragraphs of text. The first paragraph describes the project's goal of documenting chat logs. The second paragraph details the scope of the corpus, including over 140,000 chat contributions and various contexts like E-Learning and media events. The third paragraph mentions a search tool (STACCADo) and a list of publications. A box highlights that the website is accessible via a browser. The footer of the screenshot lists several publications related to the corpus.

tu technische universität  
dortmund

## Dortmunder Chat-Korpus

Bestand    Korpora / Download    Recherche  
MIT STACCADo    Kontakt

Das **Dortmunder Chat-Korpus** dokumentiert anhand einer Sammlung von Mitschnitten (sog. "Logfiles") die Sprachverwendung in unterschiedlichen Typen von Chat-Anwendungen. Es ist als Grundlage und Hilfsmittel für sprachwissenschaftliche Untersuchungen zur synchronen internetbasierten Kommunikation konzipiert und wird in verschiedenen Versionen zur freien Nutzung zur Verfügung gestellt.

Das Korpus umfasst mit über 140.000 Chat-Beiträgen bzw. 1,06 Millionen laufenden Wortformen umfangreiches Datenmaterial aus diversen Einsatzformen der Chat-Technologie. Der Bestand reicht von **Chats im Hochschulkontext** (E-Learning, Online-Zusammenarbeit, kollektive Experten-Interviews) und im Praxisbereich **Beratung & Support über Chat-Events im Medienkontext** (Chats mit Politikern und Medienakteuren oder begleitend zu TV-Ereignissen) bis hin zu **"Plauder"-Chats im Freizeitbereich**, die im IRC-Netzwerk oder in **Webchat-Communities** stattgefunden haben. Die Korpusdokumente wurden anhand einer XML-Sprache für Recherchezwecke aufbereitet.

Zusammen mit dem Korpus wird ein Suchwerkzeug zur Verfügung gestellt: **STACCADo** ermöglicht es, auf einfache Weise nach chat-typischen Elementen wie z.B. Emoticons, Adressierungen, Asterisk-Ausdrücken oder Zuschreibungen ("action messages") zu recherchieren, beliebige einfache und komplexe Volltext-Suchanfragen zu formulieren oder statistische Auswertungen zum Kommunikationsaufkommen und zum Beitragsverhalten einzelner Chatter in den Teilkorpora oder in einzelnen Korpusdokumenten zu erzeugen.

**Wenn Sie unsere Website zum ersten Mal besuchen** und einfach nur mal in unserem Datenbestand stöbern möchten, können Sie auf 385 Dokumente aus unserem Korpus auch bequem per Browser zugreifen: **HTML-Version des Releasekorpus**

Das **Dortmunder Chat-Korpus** ist Ergebnis eines Lehrstuhlprojekts am Lehrstuhl für Linguistik der deutschen Sprache und Sprachdidaktik, das unter der Leitung von Prof. Dr. Angelika Storrer und Dr. Michael Beißwenger am Institut für deutsche Sprache und Literatur der Technischen Universität Dortmund realisiert wurde. Das Suchwerkzeug STACCADo wurde von Bianca Stockrahm programmiert.

**Kurzbeschreibungen des Dortmunder Chat-Korpus finden sich in den folgenden Publikationen:**

- Beißwenger, Michael; Storrer, Angelika (2008): **Corpora of Computer-Mediated Communication**. In: Anke Lüdeling & Merja Kytö (Eds): Corpus Linguistics. An International Handbook. Volume 1. Berlin. New York (Handbooks of Linguistics and Communication Science 29.1), 292-308.
- Beißwenger, Michael; Storrer, Angelika (2011): **Digitale Sprachressourcen in Lehramtsstudiengängen: Kompetenzen - Erfahrungen - Desiderate**. In: Journal for Language Technology and Computational Linguistics, 119-139.
- Beißwenger, Michael (2013): **Das Dortmunder Chat-Korpus**. In: Zeitschrift für germanistische Linguistik 41/1, 161-164.