

# An extended tag set for annotating parts of speech in CMC corpora

Thomas Bartz<sup>1</sup>, Michael Beißwenger<sup>1</sup>,  
Eric Ehrhardt<sup>2</sup>, Angelika Storrer<sup>2</sup>

1)  technische universität  
dortmund

2)  UNIVERSITÄT  
MANNHEIM



International Research Days:  
Social Media and CMC Corpora for the eHumanities



Journées Internationales de recherche

«Médias sociaux et corpus de communication médiée  
par les réseaux. Annotation, analyse, données libres»

23-24 octobre 2015

# Part-of-speech tagging for CMC corpora

## Without a part-of-speech (PoS) annotation:

- only very limited querying options;
- no basis for advanced processing steps which require a useful linguistic preprocessing (e.g., parse trees).

## The Problem:

Part-of-speech taggers (NLP tools in general) do not perform very well on written CMC discourse:

- new elements which don't fit into any established PoS category (emojicons, addressings, action words, hashtags);
- speedwriing phenomena (typos, omission of characters, norm-deviating use of whitespace);
- colloquial (*Wazzup?*) and creative spellings (*nyce2meetU*)

# The problem

What's up,  
Deutschland?

A: Rufst an wenn du  
köpenick bist!

B: Ja

B: Wir sehn uns ja gleich

A: Jo

B: Ersatzverkejr

B: Ich hab keine ahnubg  
wo der hinfähr-.-

'I have no idea, where it  
*[that one]* is going to-.-'

# The problem

What's up,  
Deutschland?

A: Rufst an wenn du  
köpenick bist!

B: Ja

B: Wir sehn uns ja gleich

A: Jo

B: Ersatzverkejr

B: Ich hab keine ahnubg  
wo der hinfähr-.-

'I have no idea, where it  
[that one] is going to-.-'

## WebLicht

<https://weblicht.sfs.uni-tuebingen.de/weblicht/>

TreeTagger

Annotation Layers: <

language = de

Simple view

text

sentences

Table view

tokens

POStags

lemmas

token ID

tokens

POStags

lemmas

t1

Ich

t2

hab

t3

keine

t4

ahnubg

t5

wo

t6

der

t7

hinfähr-.-

WebLicht-  
Toolchain:

SFS: To TCF Converter

Document Type

Language: German

IMS: Tokenizer

Sentences

Tokens

IMS: TreeTagger

Part of Speech: STTS Ta

Lemmas

# The problem

What's up,  
Deutschland?

A: Rufst an wenn du  
köpenick bist!

B: Ja

B: Wir sehn uns ja gleich

A: Jo

B: Ersatzverkejr

B: Ich hab keine ahnubg  
wo der hinfähr-.-

'I have no idea, where it  
[that one] is going to-.-'

## WebLicht

<https://weblicht.sfs.uni-tuebingen.de/weblicht/>

TreeTagger

Annotation Layers: <

language = de

Simple view

- text
- sentences
- Table view
  - tokens
  - POStags
  - lemmas

token ID	tokens	POStags	lemmas
t1	Ich	PPER	ich
t2	hab	VAIMP	haben
t3	keine	PIAT	kein
t4	ahnubg	TRUNC	<unknown>
t5	wo	PWAV	wo
t6	der	ART	d
t7	hinfähr-.-	TRUNC	<unknown>

WebLicht-  
Toolchain:

SFS: To TCF Converter  
Document Type  
Language: German

IMS: Tokenizer  
Sentences  
Tokens

IMS: TreeTagger  
Part of Speech: STTS Ta  
Lemmas

# The problem

What's up,  
Deutschland?

A: Rufst an wenn du  
köpenick bist!

B: Ja

B: Wir sehn uns ja gleich

A: Jo

B: Ersatzverkejr

B: Ich hab keine ahnubg  
wo der hinfähr-.-

'I have no idea, where it  
[that one] is going to-.-'

WebLicht

<https://weblicht.sfs.uni-tuebingen.de/weblicht/>

TreeTagger

Annotation Layers: 

language = de

Simple view

- text
- sentences
- Table view
  - tokens
  - POStags
  - lemmas

token ID	tokens	POStags	lemmas
t1	Ich	PPER	ich
t2	hab	VAIMP	haben
t3	keine	PIAT	kein
t4	ahnubg	TRUNC	<unknown>
t5	wo	PWAV	wo
t6	der	ART	d
t7	hinfähr	ADJA	<unknown>
t8	-.-	TRUNC	<unknown>

# The problem

## Problems on several levels of the processing process:

- **Tokenization problems:** The tokens created in the tokenization step do not represent relevant units of the linguistic structure (e.g., due to speedwriting phenomena)
- **Categorization problems:** There's an adequate tag in the tag set but the tagger can't assign it (e.g., in the case of norm-deviating colloquial & dialect spellings)
- **Category problems:** The tagger can't assign an adequate tag because there's no adequate tag in the tag set (e.g., for emoticons, action words, addressings, hashtags, clitics which are typical of dialogical language in informal registers...)

# Ways to solve the problem

## Variant A:

***Normalization ⇒ PoS tagging with standard tools***

### **Open issues:**

- 1) categories for elements that are missing in PoS tagsets for edited text
- 2) adapt tools for automatic normalization

## Variant B:

***No normalization; PoS tagging of the original data***

### **Open issues:**

- 1) categories for elements that are missing in PoS tagsets for edited text
- 2) improve tokenizers & taggers



# Designing a basic PoS tag set for German CMC

- **Initiative in CLARIN-D** (2012-13) for “**updating**” the **canonical STTS** through adapting it for genres which its original creators didn't have in focus (Zinsmeister et al. 2014) – e.g.:
    - historical corpora
    - spoken language corpora
    - learner corpora
    - CMC
  - **Discussions in the DFG network *empirikom*** (2010-2014, <http://www.empirikom.net>) on how to make NLP tools fit for automatically processing & annotating CMC corpora
- ⇒ Idea: Let's set up a **community shared task** on NLP for CMC in order to encourage the developers of NLP tools to adapt their tools & tagging models for CMC
- ⇒ <https://sites.google.com/site/empirist2015/home>  
(supported by **GSCL**)



# What requirements should a basic PoS tag set for CMC meet?

- It should be **compatible with established PoS tag sets**  
⇒ interoperability with other (types of) corpora
- For categories which occur in CMC but which are not CMC-*specific*: try to be **compatible with PoS categories in other (non-CMC) genres** (⇒ interoperability of corpora; interesting research questions)
- For categories which are specific to CMC: **Keep it simple** so that the use of the categories can easily be learned
- As long that there's no consensus in the linguistic community about how to integrate CMC elements into part-of-speech typologies: Don't try to install one (and force people to use it ... because *they won't*) – instead, **design your categories as theory-free as possible.**

# “STTS 2.0”: A basic PoS tag set for German CMC

- **Basis:** The “Stuttgart Tübingen Tagset” (STTS):  
de-facto standard for German (focused on PoS tags for  
the language occurring in edited text / newspaper texts)  
(Schiller et al. 1999)

- |                          |                          |
|--------------------------|--------------------------|
| 1. Nomina (N)            | 7. Adverbien (ADV)       |
| 2. Verben (V)            | 8. Konjunktionen (KO)    |
| 3. Artikel (ART)         | 9. Adpositionen (AP)     |
| 4. Adjektive (ADJ)       | 10. Interjektionen (ITJ) |
| 5. Pronomina (P)         | 11. Partikeln (PTK)      |
| 6. Kardinalzahlen (CARD) |                          |

# “STTS 2.0”: A basic PoS tag set for German CMC

- **Basis:** The “Stuttgart Tübingen Tagset” (STTS): de-facto standard for German

PTKZU	“zu” vor Infinitiv	<i>zu [gehen]</i>
PTKNEG	Negationspartikel	<i>nicht</i>
PTKVZ	abgetrennter Verbzusatz	<i>[er kommt] an, [er fährt] rad</i>
PTKANT	Antwortpartikel	<i>ja, nein, danke, bitte</i>
PTKA	Partikel bei Adjektiv oder Adverb	<i>am [schönsten], zu [schnell]</i>
TRUNC	Kompositions-Erstglied	<i>An- [und Abreise]</i>
VVFIN	finites Verb, voll	<i>[du] gehst, [wir] kommen [an]</i>
VVIMP	Imperativ, voll	<i>komm [!]</i>
VVINF	Infinitiv, voll	<i>gehen, ankommen</i>
VVIZU	Infinitiv mit “zu”, voll	<i>anzukommen, loszulassen</i>
VVPP	Partizip Perfekt, voll	<i>gegangen, angekommen</i>
VAFIN	finites Verb, aux	<i>[du] bist, [wir] werden</i>
VAIMP	Imperativ, aux	<i>sei [ruhig !]</i>
VAINF	Infinitiv, aux	<i>werden, sein</i>
VAPP	Partizip Perfekt, aux	<i>gewesen</i>
VMFIN	finites Verb, modal	<i>dürfen</i>
VMINF	Infinitiv, modal	<i>wollen</i>
VMPP	Partizip Perfekt, modal	<i>[er hat] gekonnt</i>

Structure of STTS  
tags: main category  
> subcategory

# “STTS 2.0”: A basic PoS tag set for German CMC

- **Basis:** The “Stuttgart Tübingen Tagset” (STTS): de-facto standard for German (focused on PoS tags for the language occurring in edited text / newspaper texts) (Schiller et al. 1999)
- **“STTS 2.0”:** canonical STTS extended with new categories, but still downward-compatible with STTS (1999)
- **Compatible with the extended STTS for spoken language** which is used for PoS tagging the FOLK corpus of spoken German at IDS Mannheim (for phenomena which are *not* in the canonical STTS and which also occur in spoken language)

# “STTS 2.0”: A basic PoS tag set for German CMC

Tag	Beschreibung	Beispiele
ADJA	attributives Adjektiv	<i>[das] große [Haus]</i>
ADJD	adverbiales oder prädikatives Adjektiv	<i>[er fährt] schnell</i> <i>[er ist] schnell</i>
ADV	Adverb	<i>schon bald heute, jetzt</i>
APPR	Präposition, Zirkumposition links	<i>in [der Stadt], ohne [mich]</i>
APPRART	Präposition mit Artikel	<i>im [Haus], zur [Sache], vom, überm, fürm</i>
APPO	Postposition	<i>[ihm] zufolge, [der Sache] wegen</i>
APZR	Zirkumposition rechts	<i>[von jetzt] an</i>
ART	bestimmter oder unbestimmter Artikel	<i>der, die, das, ein, eine</i>
CARD	Kardinalzahl	<i>zwei [Männer], [im Jahre] 1994</i>
FM	Fremdsprachliches Material	<i>[Er hat das mit] A big fish [Übersetzt]</i>
ITJ	Interjektion	<i>mhm, ach, ja</i>
ONO	Onomatopoeikon	<i>being, miau, abich</i>
DM	Diskursmarker	prototypisch: <i>well, obwohl, nur also</i> als Einheiten mit projektivem Potential im Vorvorfeld von V2-Sätzen
KOUI	unterordnende Konjunktion mit „zu“ und Infinitiv	<i>um [zu leben] entsteht [zu fragen]</i>
KOUS	unterordnende Konjunktion mit Satz (VL-Stellung)	<i>weil, dass, damit wenn, ob</i>
KON	nebenordnende Konjunktion	<i>und, oder, aber</i>
KOKOM	Vergleichspartikel ohne Satz	<i>als, wie</i>
NN	Appellativa	<i>Tisch, Herr, [das] Reisen</i>
NE	Eigennamen	<i>Hans, Hamburg, HSV</i>
PD\$	substituierendes Demonstrativpronomen	<i>dieser, jener</i>
PDAT	attribuierendes Demonstrativpronomen	<i>jener [Mensch]</i>
PI\$	substituierendes Indefinitpronomen	<i>keiner, viele, man, niemand</i>
PIAT	attribuierendes Indefinitpronomen ohne Determiner	<i>kein [Mensch] irgendn [Bis]</i>
PIDAT	attribuierendes Indefinitpronomen mit Determiner	<i>[ein] wenig [Wasser] [die] beiden [Brüder]</i>
PPER	inflexives Personalpronomen	<i>ich, er, ihm, mich, dir</i>
PPOS\$	substituierendes Possesivpronomen	<i>meins, deines</i>
PPOSAT	attribuierendes Possesivpronomen	<i>mein [Buch], deine [Mutter]</i>
PREL\$	substituierendes Relativpronomen	<i>[der Hund,] der</i>
PRELAT	attribuierendes Relativpronomen	<i>[der Mann,] dessen [Hund]</i>
PRF	reflexives Personalpronomen	<i>sich, einander, dich, mir</i>
PW\$	substituierendes Interrogativpronomen	<i>wer, was</i>
PWAT	attribuierendes Interrogativpronomen	<i>welche [Farbe]</i>
PWAV	adverbiales Interrogativ- oder Relativpronomen	<i>warum, wo, wann worüber, wobei!</i>
PAV	Pronominaladverb	<i>dafür, dabei, deswegen, trotzdem</i>
PTKZU	„zu“ vor Infinitiv	<i>zu [gehen]</i>
PTKNEG	Negationspartikel	<i>nicht</i>

Tag	Beschreibung	Beispiele
PTKVZ	abgetrennter Verbzusatz	<i>[er kommt] an, [er fährt] Rad</i>
PTKANT	Antwortpartikel	<i>ja, nein, danke, bitte</i>
PTKA	Partikel bei Adjektiv oder Adverb	<i>am [schönsten], zu [schnell]</i>
PTKIFS	Intensitäts-, Fokus- oder Gradpartikel	<i>sehr [schön], höchst [eigenartig], nur [sie], voll [geil]</i>
PTKMA	Modal- oder Adörnungspartikel	<i>[Das ist] ja / vielleicht [blos] [ist das] denn [wichtig so?] [Das weiß halt] [technisch] einfach</i>
PTKMWL	Partikel als Teil eines Mehrwort-Lexems	<i>keine mehr, noch mal, schon wieder</i>
TRUNC	Kompositions-Estglied	<i>An- [und Abreise]</i>
VFIN	finites Verb, voll	<i>[du] gehst, [wir] kommen [an]</i>
VVIMP	Imperativ, voll	<i>komm [!]</i>
VVINF	Infinitiv, voll	<i>gehen, ankommen</i>
VVIZU	Infinitiv mit „zu“, voll	<i>anzukommen, loszulassen</i>
VVPP	Partizip Perfekt, voll	<i>gegangen, angekommen</i>
VAFIN	finites Verb, aux.	<i>[du] bist, [wir] werden</i>
VAIMP	Imperativ, aux.	<i>sei [ruhig]</i>
VAINF	Infinitiv, aux.	<i>werden, sein</i>
VAPP	Partizip Perfekt, aux.	<i>gewesen</i>
VMFIN	finites Verb, modal	<i>dürfen</i>
VMINF	Infinitiv, modal	<i>wollen</i>
VMPP	Partizip Perfekt, modal	<i>[er hat] gekommt</i>
VVPPER	Kontraktion: Vollverb + Inflexives Personalpronomen	<i>schwebste, machste</i>
VMPPER	Kontraktion: Modalverb + Inflexives Personalpronomen	<i>wilste, darfst, musste</i>
VAPPER	Kontraktion: Auxiliaverb + Inflexives Personalpronomen	<i>hast, bist, isst</i>
KOUSPPER	Kontraktion: unterordnende Konjunktion mit Satz (VL-Stellung) + Inflexives Personalpronomen	<i>wenns, weils, obse</i>
PPERPPER	Kontraktion: Inflexives Personalpronomen + Inflexives Personalpronomen	<i>ichs, dus, ers</i>
ADVART	Kontraktion: Adverb + Artikel	<i>son, some</i>
EMO\$C	Emoticon, als Zeichenfolge dargestellt (Typ „ASCII“)	<i>:-) :-( ^_&amp;O</i>
EMO\$M	Emoticon, als Grafik-Kond dargestellt (Typ „Image“)	kodiert (Beispiel aus WhatsApp): <i>emoj1QsmilingFaceWithSmilingEyes emoj1QkissingCatFaceWithClosedEyes</i>
AKW	Aktionswort	<i>Yach! feu, großer lol!</i>
H\$T	Hashtag	<i>[Kets war super] #aufab</i>
ADR	Adressierung	<i>@mother [ Wie is set so?]</i>
URL	Uniform Resource Locator	<i>http://www.tu-dortmund.de</i>
EML	E-Mail-Adresse	<i>petew@in@web.de</i>
XY	Nichtwort, Sonderzeichen enthaltend	<i>D&amp;XW3</i>
,	Komma	<i>,</i>
!	Satzbeendende Interpunktion	<i>! ? ! ; ;</i>
()	sonstige Satzzeichen; Satzintem	<i>- [ [ [</i>



# “STTS 2.0”: A basic PoS tag set for German CMC

PoS tag	Category	Examples
<i>I. Tags for phenomena which are specific for CMC / social media discourse:</i>		
<b>EMO ASC</b>	ASCII emoticon	:-) :-( ^^ O.O
<b>EMO IMG</b>	Graphic emoticon	😄 🍌 😜
<b>AKW</b>	Interaction word	*lach*, freu, grübel, *lol*
<b>HST</b>	Hash tag	Kreta war super! #urlaub
<b>ADR</b>	Addressing term	@lothar: Wie isset so?
<b>URL</b>	Uniform resource locator	http://www.tu-dortmund.de
<b>EML</b>	E-mail address	peterklein@web.de
<i>II. Tags for phenomena which are typical for spontaneous spoken language in colloquial registers:</i>		
<b>VV PPER</b>	Tags for types of colloquial contractions which are frequent in CMC (APPRART is already existing in STTS 1999)	schreibste, machste
<b>APPR ART</b>		vorm, überm, fürn
<b>VM PPER</b>		willste, darfst, musst
<b>VA PPER</b>		haste, biste, isses
<b>KOUS PPER</b>		wenns, weils, obse
<b>PPER PPER</b>		ichs, dus, ers
<b>ADV ART</b>		son, sone
<b>PTK IFG</b>	‘Intensitätspartikeln’, ‘Fokuspartikeln’, ‘Gradpartikeln’	<u>sehr</u> schön, <u>höchst</u> eigenartig, <u>nur</u> sie, <u>voll</u> geil
<b>PTK MA</b>	Modal particles	Das ist <u>ja</u> / <u>vielleicht</u> doof. Ist das <u>denn</u> richtig so? Das war <u>halt</u> echt nicht einfach.
<b>PTK MWL</b>	Particle as part of a multi-word lexeme	keine <u>mehr</u> , <u>noch</u> mal, <u>schon</u> wieder
<b>DM</b>	Discourse markers	<u>weil</u> , <u>obwohl</u> , <u>nur</u> , <u>also</u> , ... <i>with V2 clauses</i>
<b>ONO</b>	Onomatopoeia	boing, miau, zisch

# “STTS 2.0”: A basic PoS tag set for German CMC

PoS tag	Category	Examples
---------	----------	----------

## *I. Tags for phenomena which are specific for CMC / social media discourse:*

<b>EMO ASC</b>	ASCII emoticon	:-) :-( ^^ O.O
<b>EMO IMG</b>	Graphic emoticon	  
<b>AKW</b>	Interaction word	*lach*, freu, grübel, *lol*
<b>HST</b>	Hash tag	Kreta war super! <u>#urlaub</u>
<b>ADR</b>	Addressing term	<u>@lothar</u> : Wie isset so?
<b>URL</b>	Uniform resource locator	http://www.tu-dortmund.de
<b>EML</b>	E-mail address	peterklein@web.de

## *II. Tags for phenomena which are typical for spontaneous spoken language in colloquial registers:*

<b>VV PPER</b>	Tags for types of colloquial contractions which are frequent in CMC (APPRART is already existing in STTS 1999)	schreibste, machste
<b>APPR ART</b>		vorm, überm, fürn
<b>VM PPER</b>		willste, darfste, musste
<b>VA PPER</b>		haste, biste, isses
<b>KOUS PPER</b>		wenns, weils, obse



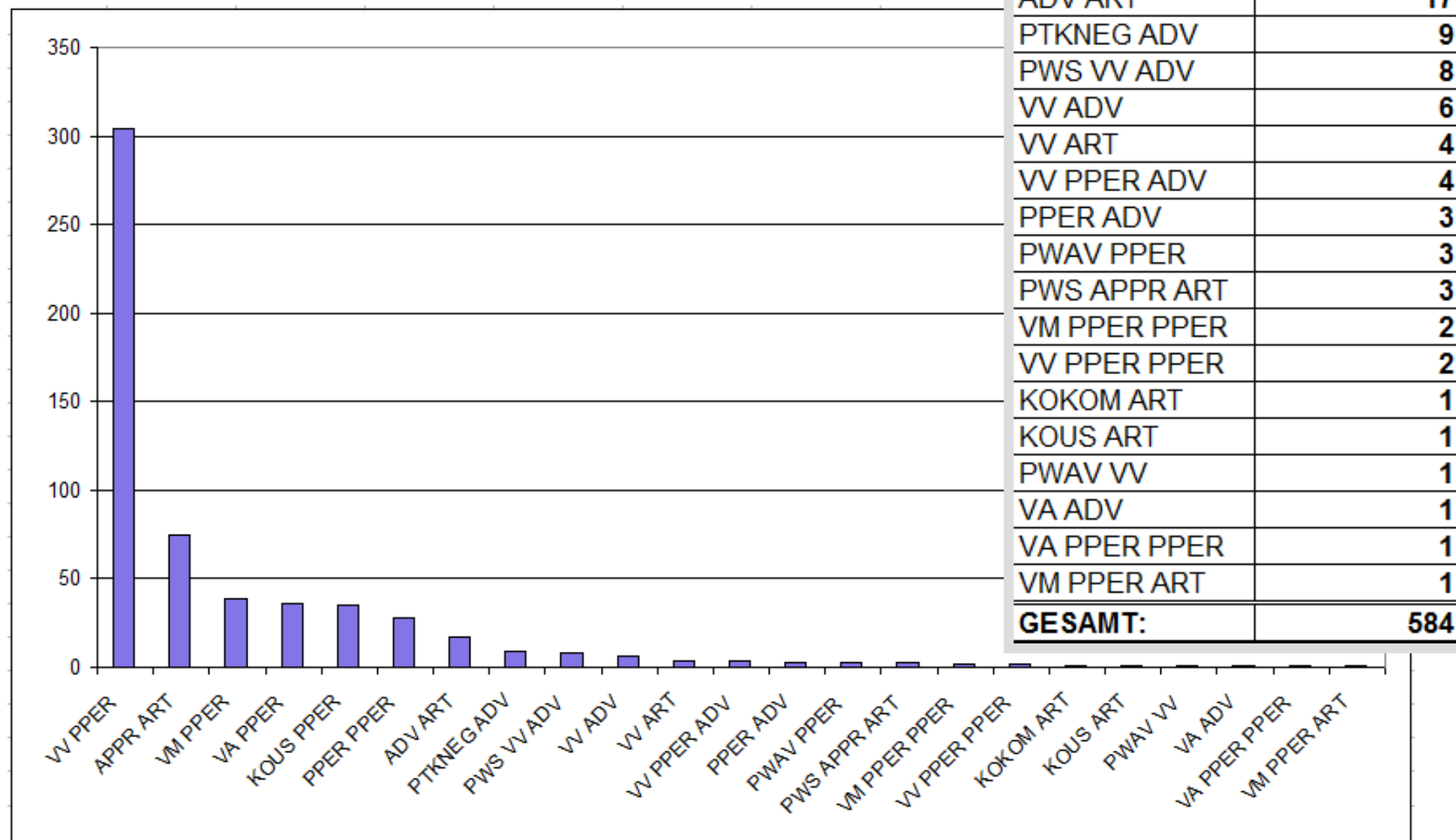
<b>AKW</b>	Interaction word	*lach*, freu, grübel, *lol*
<b>HST</b>	Hash tag	Kreta war super! <u>#urlaub</u>
<b>ADR</b>	Addressing term	<u>@lothar</u> : Wie isset so?
<b>URL</b>	Uniform resource locator	http://www.tu-dortmund.de
<b>EML</b>	E-mail address	peterklein@web.de

**II. Tags for phenomena which are typical for spontaneous spoken language in colloquial registers:**

<b>VV PPER</b>	Tags for types of colloquial contractions which are frequent in CMC (APPRART is already existing in STTS 1999)	schreibste, machste
<b>APPR ART</b>		vorm, überm, fürn
<b>VM PPER</b>		willste, darfst, musste
<b>VA PPER</b>		haste, biste, isses
<b>KOUS PPER</b>		wenns, weils, obse
<b>PPER PPER</b>		ichs, dus, ers
<b>ADV ART</b>		son, sone
<b>PTK IFG</b>	'Intensitätspartikeln', 'Fokuspartikeln', 'Gradpartikeln'	<u>sehr</u> schön, <u>höchst</u> eigenartig, <u>nur</u> sie, <u>voll</u> geil
<b>PTK MA</b>	Modal particles	Das ist <u>ja</u> / <u>vielleicht</u> doof. Ist das <u>denn</u> richtig so? Das war <u>halt</u> echt nicht einfach.
<b>PTK MWL</b>	Particle as part of a multi-word lexeme	keine <u>mehr</u> , <u>noch</u> mal, <u>schon</u> wieder
<b>DM</b>	Discourse markers	<u>weil</u> , <u>obwohl</u> , <u>nur</u> , <u>also</u> , ... <i>with V2 clauses</i>
<b>ONO</b>	Onomatopoeia	boing, miau, zisch

# Contractions in chats

‘social chat’ subcorpus of the Dortmund chat corpus:  
21 logfiles / 104.094 tokens, including 584  
occurrences of colloquial contractions



# Tag set and annotation guidelines @EmpiriST2015

## EmpiriST 2015

[Diese Site durchsuchen](#)

[Navigation](#)  
[Subtasks & deadlines](#)  
[Data sets](#)  
**[Annotation Guidelines](#)**  
[Task Force](#)  
[Sitemap](#)

[GSCL Shared Task: Automatic Linguistic Annotation of Computer-Mediated Communication / Social Media >](#)

### Annotation Guidelines

The training data that will be provided as a gold standard have been manually tokenized and tagged according to the following guidelines:

- Beißwenger, Michael; Bartz, Thomas; Storrer, Angelika; Westpfahl, Swantje (2015): **Tagset und Richtlinie für das Part-of-Speech-Tagging von Sprachdaten aus Genres internetbasierter Kommunikation.** Guideline document from the Empirikom shared task on automatic linguistic annotation of internet-based communication (*EmpiriST 2015*). (21 pages).  
PDF: [EmpiriST\\_Guideline-PoS.pdf](#)
- Beißwenger, Michael; Bartsch, Sabine; Evert, Stefan; Würzner, Kay-Michael (2015): **Richtlinie für die manuelle Tokenisierung von Sprachdaten aus Genres internetbasierter Kommunikation.** Guideline document from the Empirikom shared task on automatic linguistic annotation of internet-based communication (*EmpiriST 2015*). (29 pages).  
PDF: [EmpiriST\\_Guideline-Tokenisierung.pdf](#)

When citing these documents, please use the bibliographic information given above and refer to the URL <http://sites.google.com/site/empirist2015/>.

**Overview: The part of speech tagset used for annotations:**

Extensions to STTS (1999) are highlighted with blue background colour:

Tag	Description (German)	Examples
ADJA	attributives Adjektiv	<i>[das] große [Haus]</i>
ADJD	adverbiales oder prädikatives Adjektiv	<i>[er fährt] schnell [er ist] schnell</i>
ADV	Adverb	<i>schon, bald, heute, jetzt</i>
APPR	Präposition, Zirkumposition links	<i>in [der Stadt], ohne [mich]</i>
APPRART	Präposition mit Artikel	<i>im [Haus], zur [Sache], vorm, überm, fürn</i>
APPO	Postposition	<i>[ihm] zufolge, [der Sache] wegen</i>

PoS tagset + annotation guidelines available on the website of the GSCL/ Empirikom shared task on automatic linguistic annotation of CMC ([EmpiriST2015](#)).

<https://sites.google.com/site/empirist2015/home/>

# ChatCorpus2CLARIN: Project background

Curation project of the CLARIN-D F-AG 1 “German Philology”



**Duration:** May 2015 – February 2016

CLARIN-D

**Project team:** Michael Beißwenger (U Dortmund), Angelika Storrer, Eric Ehrhardt (U Mannheim), Harald Lungen (IDS), Axel Herold (BBAW) + other colleagues at IDS and BBAW

**The task:** Re-modeling of the Dortmund Chat Corpus and samples of other CMC resources compliant with existing standards for the representation of corpora in the Digital Humanities. Integration into the CLARIN-D infrastructures at BBAW and IDS.

## **Main goal:**

- Pave the way for the inclusion of linguistically annotated CMC resources into the CLARIN-D corpus infrastructures and create the prerequisites for investigating linguistic peculiarities of CMC with state-of-the art corpus technology.

# ChatCorpus2CLARIN: Project back

Curation project of the CLARIN-D F-AG 1 “German F



[http://www.clarin-d.de/  
de/kurationsprojekt-1-3-germanistik](http://www.clarin-d.de/de/kurationsprojekt-1-3-germanistik)



Home Accessing ▾ Analysing ▾ Preparation ▾ Disciplines ▾ About ▾ Help ▾

## ChatCorpus2CLARIN: Integration of the Dortmund Chat Corpus into CLARIN-D

### Project content

In the third curation project of the [CLARIN-D working group 1 “German Philology” \(F-AG 1\)](#) an existing corpus of computer-mediated communication (CMC), the Dortmund Chat Corpus, and samples of other CMC resources will be restructured to conform to current standards for the representation of corpora in the Digital Humanities context. The main goal of this work is to pave the way for the inclusion of linguistically annotated CMC resources into CLARIN-D corpus infrastructures and to create the prerequisites for investigating linguistic peculiarities of CMC with state-of-the art corpus technology. To this end, the project will (1) transform the metadata and the annotations of the chat corpus into a TEI-compliant format, (2) enrich the data by further linguistic annotations, and (3) integrate the resulting resource into the CLARIN-D Corpus Infrastructures at the [Institute for the German Language \(IDS\)](#) and the [Berlin-Brandenburg Academy of Sciences \(BBAW\)](#).

The integration in CLARIN-D will allow for a systematic corpus-based analysis of CMC discourse as compared to the language of edited text (as represented in the text corpora at BBAW and IDS) and of spoken conversations (as represented in the spoken language corpora at IDS).

### The Dortmund Chat Corpus

The data for the Dortmund Chat Corpus (Beißwenger & Storrer 2008; Beißwenger 2013) was built at TU Dortmund University. The goal of the corpus project was to create a resource for researching the peculiarities and linguistic variation in written computer-mediated communication. The corpus comprises 478 logfile documents with about 140,000 postings and about 1 million tokens of German chats from different application contexts (social chats, advisory chats, chats in the context of learning and teaching, moderated chats in media context). The corpus has been annotated using an XML format (“ChatXML”) that represents (1) the basic structure and properties of chat logfiles and postings, (2) selected “netspeak” phenomena such as emoticons, interaction words, addressing terms, nicknames and acronyms, (3) selected metadata about the chat users. Since 2005, the corpus has been made available at <http://www.chatcorpus.tu-dortmund.de> as an XML version for download and offline querying and as an HTML version for online browsing. It has been widely used as a resource for studying and teaching the peculiarities of German CMC discourse. In the CLARIN-D context, it has been used as one of the resources of the curation project: „[Linguistic Annotation of Non-standard Varieties – Guidelines and Best Practices](#)” of [working group 7](#).

# The corpus


## Dortmund Chat Corpus

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

478 logfile documents with 140,240 user postings or 1M words of German chat discourse.

Resource for the analysis of **linguistic variation in chats** including chats from different social/institutional contexts (social chats, advisory chats, learning and teaching, moderated chats in the media context).

Annotated in a home-grown XML format ('**ChatXML**'): (1) basic structure of chat logfiles and postings, (2) selected CMC phenomena, (3) selected metadata.

 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.

# Other corpora / data sets in the project focus

- German [WhatsApp Corpus](#) („What's up, Deutschland?“)
- German [Wikipedia corpus](#) in DeReKo
- German [News Corpus](#) in DeReKo
- DWDS [Blog Corpus](#)
- DWDS [German Reference Corpus of CMC](#) (DeRiK)



Jetzt deine **WhatsApp**-Nachrichten für die Wissenschaft spenden!  
17x Gutscheine gewinnen!

**What's up, Deutschland?**

[www.whatsup-deutschland.de](http://www.whatsup-deutschland.de)

Sprachwissenschaftler\_innen an den Universitäten Leipzig, Dortmund, Dresden, Duisburg-Essen, Hannover, Koblenz-Landau und Mannheim wollen untersuchen, wie man in Deutschland WhatsApp-Nachrichten schreibt.

Auf [www.whatsup-deutschland.de](http://www.whatsup-deutschland.de) findest du:

- Anleitungen zur Nachrichtenspende
- unsere Datenschutzrichtlinien
- mehr Informationen zum Gewinnspiel

Foto „Dating in the Rain“ bearbeitet von Garry Knight (by 202).  
Foto „Teen Texting“ bearbeitet von Social Poster (by mrcos 20).

The advertisement features a man in a plaid shirt holding a smartphone in the top right corner. Below this, a QR code is displayed. The main image shows two people sitting on a bench outdoors, one holding an umbrella and both looking at a smartphone. The text is in German, promoting a project where WhatsApp messages are donated to science for a chance to win 17 vouchers. It lists participating universities and provides a website for more information and instructions.

# Work packages in the project

- TEI representation (⇒ “CLARIN-D schema”)
- CLARINification, legal issues + licensing
- enrich the data with additional linguistic annotations (PoS, normalised spellings, ...)



# The vision

After its integration into the CLARIN-D infrastructure the resource will be characterized by the following added values:

- Advanced accessibility and retrieval options;
- interoperability with other corpus resources that are represented in TEI and with annotation and analysis tools that support the TEI format;
- advanced querying options (PoS tags, normalized spellings);
- interoperability with other corpus resources that have been tagged with STTS;
- advanced options for corpus-based analyses on the peculiarities of CMC discourse as compared to the language of edited text and of spoken language, using the text and speech corpora which are already available in the corpus infrastructures of BBAW and IDS.

# PoS annotation of the corpus: workflow

## 1. Automatic tokenisation, PoS annotation & lemmatisation

of the chat corpus with tools + tagging models from the BMBF project „Schreibgebrauch“ at U Saarbrücken (Horbach et al. 2014, Horbach et al. 2015) <http://www.schreibgebrauch.de>

PoS tag set: previous version of “STTS 2.0” (Bartz et al. 2014)

Representation of the tagging results as additions to the ChatXML format.

### *Standard PoS taggers:*

Accuracy on Chat Corpus: ~71% (vs. 97% accuracy on Newspaper)

### *Tagging models from the “Schreibgebrauch” project:*

Average accuracy on Chat Corpus: **83.5%**

## 2. Manual post-processing of the tagging results using OrthoNormal in FOLKER (preview version 1.2) with an import/export filter for PoS tagged ChatXML (defined by Thomas Schmidt/IDS)

# Manual post-processing of PoS tagging results with *OrthoNormal*

OrthoNormal 0.9 [Neue Transkription]

File Edit View Help

Spr... Transkriptionstext

45 Lant... :)))

46 Lant... Na , zori ? :))

47 marc... den ?

48 quaki \* g [G] \*

49 quaki dange [danke] lantonie

50 Lant... Ich habe heute einen SMS von Tigaaaaelse bekommen .

51 Lant... \* erzähl \*

52 Pharao lanto redet wie ne [eine] bewährungshelferin [Bewährungshelferin]

53 zora zora [Zora] freut sich über ihr zeugniss :)))

54 quaki \* aufpluster \*

55 syste... Thor... *betritt den Raum.*

56 marc... ich mal wieder nich [nicht] ...

57 quaki was hast denn zori ??

58 quaki erzähl

59 syste... *stoeps kommt aus dem Raum ((Number\_of\_the\_beast))herein.*

60 Lant... Das hast du [Du] dir [Dir] verdient , zori ?

61 Tomc... oh man wat [was] fürn krawall [Krawall] hier draußen ... \* guck \*

Wort Normal Lemma POS p(POS)

Wort	Normal	Lemma	POS	p(POS)
lanto			NE	
redet		reden	VVFIN	
wie		wie	KOKOM	
ne	eine	eine	ART	
bewährungsh...	Bewährungsh...	Bewährungsh...	NN	
zora	Zora	Zora	NE	
freut		freuen	VVFIN	
sich		sich	PRF	
über		über	APPR	
ihr		ihr	PPOSAT	
zeugniss			NN	
:)))			EMOASC	
*		*	AWIND	
aufpluster			ADJD	
*		*	AWIND	
ich		ich	PPER	
mal		mal	ADV	
wieder		wieder	ADV	
nich	nicht		PTKNEG	
...		...	\$(	
was		was	PWS	
hast		haben	VAFIN	
denn		denn	ADV	
zori			NE	
??			\$.	
erzähl		erzählen	VVIMP	
Das		die	PDS	
hast		haben	VAFIN	
du	Du	du	PPER	
dir	Dir	du	PPER	

Modus:  Normalisieren  Tagging  XML

Automatisches Weiterrücken

(Overview of the FOLK tools: Schmidt 2012)

# Using <w> for the representation of PoS information in our TEI schema

```
<post type="standard" who="#A04" auto="false" rend="color:green">
  <p>
    <w type="VVFIN">dachte</w>
    <w type="PPER">ich</w>
    <w type="ADV">auch</w>
    <w type="ADV">immer</w>
    <w type="$((">,</w>
    <name type="nickname" corresp="#A09">
      <w type="NE">monk</w>
    </name>
    <w type="$. ">..</w>
    <w type=",$("&>*</w>
    <w type="AKW">heul</w>
    <w type=",$("&>*</w>
  </p>
</post>
```

CLARIN-D TEI schema (documentation):

[http://wiki.tei-c.org/index.php/SIG:CMC/  
CLARIN-D schema draft for  
representing CMC in TEI \(2015\)](http://wiki.tei-c.org/index.php/SIG:CMC/CLARIN-D_schema_draft_for_representing_CMC_in_TEI_(2015))

ineli26: dachte ich auch immer, monk .. \*heul\*  
*I was always thinking the same, monk .. \*crying\**

# References

- Bartz, Thomas; Beißwenger, Michael; Storrer, Angelika (2014): Optimierung des Stuttgart-Tübingen-Tagset für die linguistische Annotation von Korpora zur internetbasierten Kommunikation: Phänomene, Herausforderungen, Erweiterungsvorschläge. In: Journal for Language Technology and Computational Linguistics 28 (1), 157-198.  
[http://www.jlcl.org/2013\\_Heft1/7Bartz.pdf](http://www.jlcl.org/2013_Heft1/7Bartz.pdf)
- Beißwenger, Michael (2013): Das Dortmunder Chat-Korpus. In: Zeitschrift für germanistische Linguistik 41 (1), 161-164.  
Extended version: [http://www.linse.uni-due.de/tl\\_files/PDFs/Publicationen-Rezensionen/Chatkorpus\\_Beisswenger\\_2013.pdf](http://www.linse.uni-due.de/tl_files/PDFs/Publicationen-Rezensionen/Chatkorpus_Beisswenger_2013.pdf)
- 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) 3.  
<http://jtei.revues.org/476> (DOI: 10.4000/jtei.476).
- Beißwenger, Michael; Bartz, Thomas; Storrer, Angelika; Westpfahl, Swantje (2015): Tagset und Richtlinie für das PoS-Tagging von Sprachdaten aus Genres internetbasierter Kommunikation. Guideline Document, Dortmund 2015.  
<https://sites.google.com/site/empirist2015/home/annotation-guidelines>
- Horbach, Andrea; Steffen, Diana; Thater, Stefan; Pinkal, Manfred (2014): Improving the Performance of Standard Part-of-Speech Taggers for Computer-Mediated Communication. Proceedings of KONVENS 2014, 171-177.
- Horbach, Andrea; Thater, Stefan; Steffen, Diana; Fischer, Peter M.; Witt, Andreas; Pinkal, Manfred (2015): Internet Corpora: A Challenge for Linguistic Processing. In: Datenbank-Spektrum 15 (1), 41-47.
- Schiller, Anne; Teufel, Simone; Stöckert, Christine (1999): Guidelines für das Tagging deutscher Textcorpora mit STTS (Kleines und großes Tagset). University of Stuttgart: Institut für maschinelle Sprachverarbeitung.
- Schmidt, Thomas (2012): EXMARaLDA and the FOLK tools – two toolsets for transcribing and annotating spoken language. In: Proceedings of the Eighth conference on International Language Resources and Evaluation (LREC'12), Istanbul, Turkey: European Language Resources Association (ELRA). [http://www.lrec-conf.org/proceedings/lrec2012/pdf/529\\_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2012/pdf/529_Paper.pdf).
- Zinsmeister, Heike; Heid, Ulrich; Beck, Kathrin Beck (Eds., 2014): Das STTS-Tagset für Wortartentagging - Stand und Perspektiven. Special issue of the Journal for Language Technology and Computational Linguistics. <http://www.jlcl.org>