



# Wortartentagging der Tübinger Ressourcen nach STTS

Erfahrungen mit verschiedenen Textgenres

**STTS-Workshop Stuttgart**

**24.09.2012, Kathrin Beck, Erhard Hinrichs, Heike Telljohann, Yannick Versley**



---

## STTS-getaggte Tübinger Ressourcen (1)

- **POS-Tags manuell bearbeitet**

- TüBa-D/S (Tübinger Baumbank des Deutschen / Spontansprache)
  - ca. 38 000 Sätze, Quelle: Verbmobil-Dialoge
- TüBa-D/Z (Tübinger Baumbank des Deutschen / Zeitungskorpus)
  - ca. 65 500 Sätze, Quelle: die tageszeitung (taz)
- Sample aus TüBa-D/DC (Diachrones Korpus)
  - ca. 3 800 Sätze aus insg. 6 Texten unterschiedlicher Epochen
  - Quelle: Ausgewähltes Material aus dem Projekt Gutenberg-DE
  - Nur intern verwendet
    - (zur Evaluation der automatischen Annotation in TüBa-D/DC)



---

## STTS-getaggte Tübinger Ressourcen (2)

- **POS-Tags automatisch annotiert**

- TüPP-D/Z (Tüb. Partiiell Gearstes Korpus des Dt. / Zeitungskorpus)
  - > 200 Mio. Wörter, Quelle: Textsammlung aus die tageszeitung (taz)
  - Kombination aus tnt-Modellen für verschiedene Genres
- web-news
  - 1,7 Mia. Wörter, Quelle: Nachrichten- und Blogsites im WWW
  - RFTagger + Konversion nach STTS+morph
- TüBa-D/DC (Tübinger Baumbank des Deutschen/Diachrones Korpus)
  - > 250 Mio. Wörter, Quelle: Projekt Gutenberg-DE
  - TreeTagger



---

## Abweichungen vom originalen STTS-Tagset in den manuell bearbeiteten Korpora

- **TüBa-D/S**

- zusätzlich: BS (Buchstabe)
- umbenannt: PAV (Pronominaladverb,  
z. B. *davon, deswegen, hierfür*)  
→ PROP (pronominale Form einer PP)

- **TüBa-D/Z**

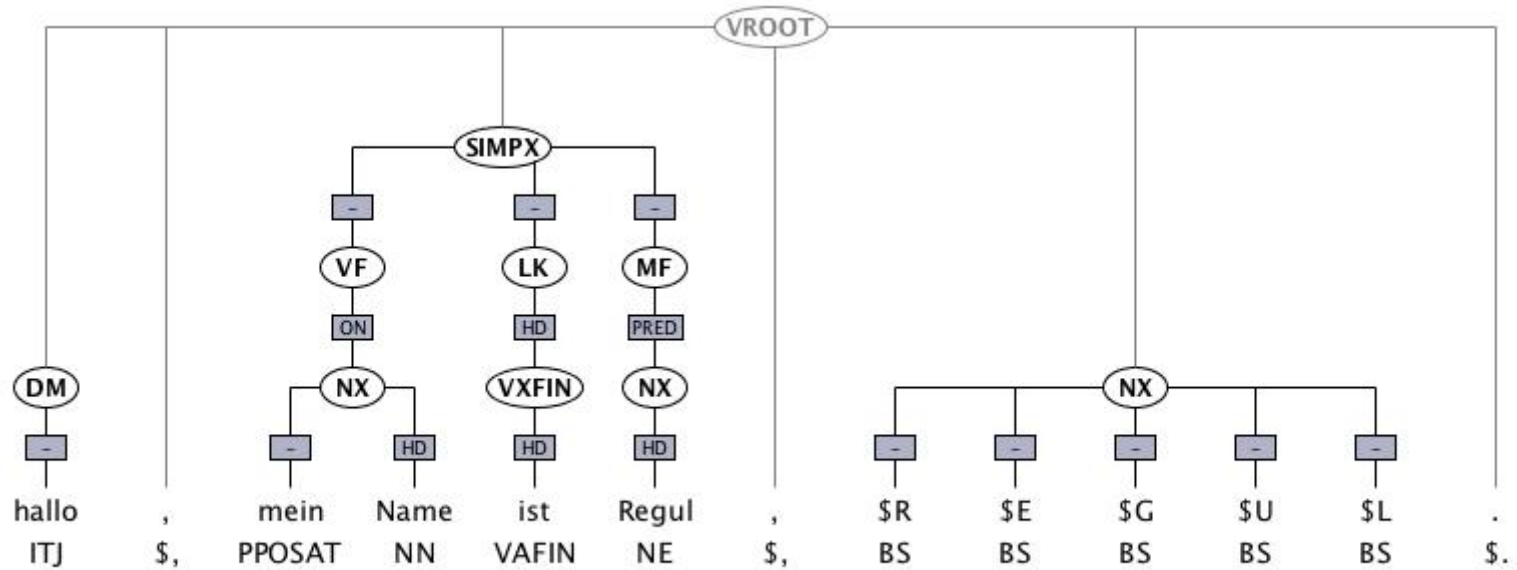
- umbenannt: PAV → PROP (pronominale Form einer PP)

- **TIGER Corpus**

- umbenannt: PAV → PROAV (Pronominaladverb)
- nicht verwendet: PIDAT (attr. Indefinitpronomen m. Determiner)



## Baumbeispiel aus der TüBa-D/S





---

## Anwendungsunterschiede von POS-Tags (Beispiele)

- **PIAT vs. PIDAT**

- TIGER Corpus      Bsp.: **viele (PIAT) Menschen**  
**mehr (PIAT) Fragen**

- TüBa-D/Z      Bsp: **viele (PIDAT) Menschen**  
**mehr (PIAT) Fragen**

- **als APPR vs. KOKOM**

- TIGER Corpus      Bsp.: *Powell galt nach Umfragen*  
**als (APPR) aussichtsreichster**  
*Gegner Clintons.*

- TüBa-D/Z      Bsp.: *Cuomo gilt **als (KOKOM)***  
*exzellenter Redner.*



---

## Anwendungsunterschiede von POS-Tags (Beispiele)

- **FM vs. NE am Beispiel von fremdsprachlichen Filmtiteln**

- TIGER Corpus

Bsp.: *Tous (NE) les (NE) matins (NE) du (NE) monde (NE).*

- TüBa-D/Z

Bsp.: *Knockin' (FM) on (FM) Heaven's (FM) Door (FM).*



## Beispiel 1 für POS-Tag-Korrektur in TüBa-D/DC:

- Johann Wolfgang von Goethe, Die Leiden des jungen Werther (1774)

TEXT	TOKENS	POSTAGS		POSTAGS
t_504	Das	PDS		PDS
t_505	bewog	WFIN		WFIN
t_506	den	ART		ART
t_507	verstorbenen	ADJA		ADJA
t_508	Grafen	NN		NN
t_509	von	APPR		APPR
t_510	M.	NE		NE
t_511	,	\$,		\$,
t_512	einen	ART		ART
t_513	Garten	NN		NN
t_514	auf	APPR		APPR
t_515	einem	PIS		PIS
t_516	der	ART		ART
t_517	Hügel	NN		NN
t_518	anzulegen	WIZU		WIZU
t_519	,	\$,		\$,
t_520	die	PRELS		PRELS
t_521	mit	APPR		APPR
t_522	der	ART		ART
t_523	schönsten	ADJA		ADJA
t_524	Mannigfaltigkeit	NN		NN
t_525	sich	PRF		PRF
t_526	kreuzen	WINF	→	YYY-WFIN
t_527	und	KON		KON
t_528	die	ART		ART
t_529	lieblichsten	ADJA		ADJA
t_530	Täler	NN		NN
t_531	bilden	WINF	→	YYY-WFIN
t_532	.	\$.		\$.



## Beispiel 2 für POS-Tag-Korrektur in TüBa-D/DC:

- Philipp Melanchthon, Die Augsburgerische Konfession (1530)

TEXT	TOKENS	POSTAGS		POSTAGS
t_1321	Hie	NE	→	YYY-ADV
t_1322	werden	VAFIN		VAFIN
t_1323	verworfen	WPP		WPP
t_1324	die	PRELS	→	YYY-PDS
t_1325	,	\$.		\$.
t_1326	so	ADV	→	YYY-PRELS
t_1327	lehren	WVINF	→	YYY-WVFIN
t_1328	,	\$.		\$.
t_1329	daß	KOUS		KOUS
t_1330	diejenigen	PDS		PDS
t_1331	,	\$.		\$.
t_1332	so	ADV	→	YYY-PRELS
t_1333	einst	ADV		ADV
t_1334	seind	ADV	→	YYY-VAFIN
t_1335	fromm	ADJD		ADJD
t_1336	worden	VAPP		VAPP
t_1337	,	\$.		\$.
t_1338	nicht	PTKNEG		PTKNEG
t_1339	wieder	ADV		ADV
t_1340	fallen	WVINF		WVINF
t_1341	mugen	WVINF	→	YYY-VMFIN
t_1342	.	\$.		\$.



## Beispiel 3 für POS-Tag-Korrektur in TüBa-D/DC:

- Gottfried von Straßburg, Tristan (1210)

TEXT	TOKENS	POSTAGS	POSTAGS
t_280	Tiure	NE	YYY-ADJD
t_281	unde	VFIN	YYY-KON
t_282	wert	ADJD	ADJD
t_283	ist	VAFIN	VAFIN
t_284	mir	PPER	PPER
t_285	der	ART	ART
t_286	man	PIS	YYY-NN
t_287	,	,\$	,\$
t_288	der	PRELS	PRELS
t_289	guot	ADJD	YYY-NN
t_290	und	KON	KON
t_291	übel	ADJD	YYY-NN
t_292	betrahten	ADJA	YYY-WINF
t_293	kan	NE	YYY-VMFIN
t_294	,	,\$	,\$
t_295	der	PRELS	PRELS
t_296	mich	PPER	PPER
t_297	und	KON	KON
t_298	iegeſichen	VFIN	YYY-PIDAT
t_299	man	PIS	YYY-NN
t_300	nâch	NE	YYY-APPR
t_301	sînem	ADJA	YYY-PPOSAT
t_302	werde	VAFIN	YYY-NN
t_303	erkennen	WINF	WINF
t_304	kan	XY	YYY-VMFIN
t_305	.	,\$.	,\$.



## Erweiterte STTS-Tagsets im PCFG-Parsing

-PCFG-Parsing braucht feinere Unterscheidung von Wortarten und Knoten  
(Schiehlen, 2004; Dubey, 2005; Versley 2005; Versley&Rehbein 2009)

-Unterschiedliche Verfeinerungen des STTS-Tagsets

STTS	Sch04	Dub05	Ver05	VR09
DET	DET	DET_SB	DET_nsm	DET_der
ADJA	ADJA	ADJA_SB	ADJA_e	
NN	NN_SB		NN_sm	
VAFIN	VAFIN_haben			
VVFIN	VVFIN_OA_OC	VVFIN	VVFIN_a	
APPR	APPR	APPR_OA	APPR_a	



---

**Danke  
fürs  
Zuhören.**