

DH 1: Sprache und Text

# Annotation: Tools und Pipelines

**Andreas Blombach, Stephanie Evert**

Lehrstuhl für Korpus- und  
Computerlinguistik

<https://www.linguistik.phil.fau.de>



Friedrich-Alexander-Universität  
Philosophische Fakultät und  
Fachbereich Theologie

# Wiederholung

## Korpusannotation: Wortebene

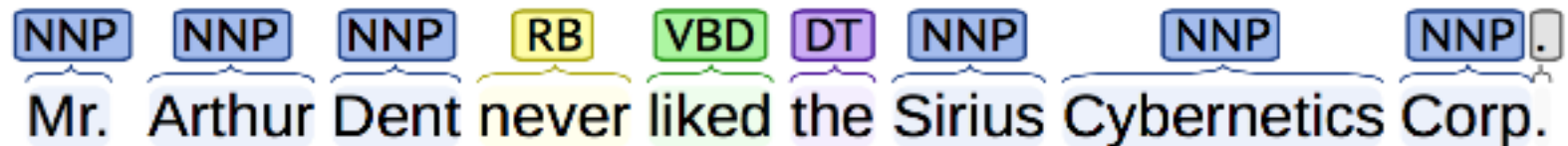
- Jedem (laufenden) Wort wird eine Kategorie zugeordnet  
→ sog. **Tagging** (= Etikettierung)
  - Voraussetzung: Text muss in Wörter zerlegt sein
- **Tokenisierung**
  - **Token** = Wort, Zahl, Symbol (😎), Satzzeichen, ...
  - im Gegensatz zu **Typen** = verschiedene Wörter
- Kann schwieriger sein, als man vermuten würde ...

@Mia1234 #semibk [1] Das schließt direkt an die vorige Frage von @DieMaJa22 an. In jedem Fall gibt es (wie auch in der Sitzung ... @Mia1234 #semibk [2]am BspChats gezeigt) starkeHinweise darauf, dass(wie auch imRealLife) diverseFaktoren die sprVariation beeinflussen: <http://tinyurl.com/3umxkuh>

<https://sites.google.com/site/empirist2015/> (Beißwenger et al. 2016)

## Annotation auf Wortebene

- Zentral: Wortartenannotierung = **POS-Tagging**
  - Substantiv (**noun**), Adjektiv, Verb, Adverb, Pronomen, Präposition, Konjunktion, Zahl, Satzzeichen, ...
  - engl. POS = **part of speech**
- Tagset = Kategorienschema
  - meist feinere Unterschiede: Sg/Pl, inf./fin./imp., ...



- auch: Lemmatisierung (hier kein Tagset!), semantische Kategorien, emotionale Valenz, Schwierigkeitsgrad (CEFR), ...

# Deutsch: STTS-Tagset

<b>ADJA</b>	attributives Adjektiv
<b>ADJD</b>	adverbiales / prädikatives Adjektiv
<b>ADV</b>	Adverb <i>schon, bald, doch</i>
<b>APPR</b>	Präposition / Zirkumposition links
<b>APPRART</b>	Präposition mit Artikel fusioniert <i>zum</i>
<b>APPO</b>	Postposition <i>zufolge, wegen</i>
<b>APZR</b>	Zirkumposition rechts <i>von ... an</i>
<b>ART</b>	bestimmter oder unbestimmter Artikel
<b>CARD</b>	Kardinalzahlen (Ordinalzahl = ADJA)
<b>FM</b>	Fremdsprachliches Material
<b>ITJ</b>	Interjektion <i>ahm, ach, tja</i>
<b>KOUI</b>	unterordnende Konj. mit <i>zu</i> + Inf
<b>KOUS</b>	unterordnende Konjunktion mit Satz
<b>KON</b>	nebenordnende Konjunktion <i>und, oder</i>
<b>KOKOM</b>	Vergleichskonjunktion <i>als, wie</i>
<b>NN</b>	normales Nomen
<b>NE</b>	Eigenname
<b>PDS</b>	substituierendes Demonstrativpron.
<b>PDAT</b>	attribuierendes Demonstrativpron.
<b>PIS</b>	substituierendes Indefinitpron.
<b>PIAT</b>	attrib. Indefinitpron. ohne Determiner
<b>PIDAT</b>	attrib. Indefinitpron. mit Determiner
<b>PPER</b>	Personalpronomen (nicht reflexiv)
<b>PPOSS</b>	substituierendes Possessivpronomen
<b>PPOSAT</b>	attribuierendes Possessivpronomen
<b>PRELS</b>	substituierendes Relativpronomen
<b>PRELAT</b>	attribuierendes Relativpronomen

<b>PRF</b>	reflexives Personalpronomen
<b>PWS</b>	substituierendes Interrogativpron.
<b>PWAT</b>	attribuierendes Interrogativpronomen
<b>PWAV</b>	adverbiales Interrogativ-/Relativpron.
<b>PAV</b>	Pronominaladverb <i>dafür, deswegen</i>
<b>PTKZU</b>	<i>zu</i> vor Infinitiv
<b>PTKNEG</b>	Negationspartikel <i>nicht</i>
<b>PTKVZ</b>	abgetrennter Verbzusatz <i>kommt ... an</i>
<b>PTKANT</b>	Antwortpartikel <i>ja, nein, danke</i>
<b>PTKA</b>	Partikel bei Adjektiv/Adverb <i>am, zu</i>
<b>TRUNC</b>	Kompositions-Erstglied <i>Unter- und ...</i>
<b>VVFIN</b>	finites Verb, voll (= lexikalisch)
<b>VVIMP</b>	Imperativ, voll
<b>VVINFINF</b>	Infinitiv, voll
<b>VVIZU</b>	Infinitiv mit <i>zu</i> , voll
<b>VVPP</b>	Partizip Perfekt, voll
<b>VAFIN</b>	finites Hilfsverb
<b>VAIMP</b>	Imperativ, Hilfsverb
<b>VAINFINF</b>	Infinitiv, Hilfsverb
<b>VAPP</b>	Partizip Perfekt, Hilfsverb
<b>VMFIN</b>	Finites Modalverb
<b>VMINFINF</b>	Infinitiv, Modalverb
<b>VMPP</b>	Partizip Perfekt, Modalverb
<b>XY</b>	Nichtwort mit Sonderzeichen <i>3:7, H2O</i>
<b>\$,</b>	Komma <i>,</i>
<b>\$.</b>	Satzbeendende Interpunktion <i>.?!;:</i>
<b>\$(</b>	sonstige Satzzeichen (intern) <i>- [ ] ( )</i>

<b>CC</b>	Coordinating conjunction
<b>CD</b>	Cardinal number
<b>DT</b>	Determiner
<b>EX</b>	Existential <i>there</i>
<b>FW</b>	Foreign word
<b>IN</b>	Preposition / subordinating conjunction
<b>IN/that</b>	Subordinating conjunction <i>that</i>
<b>JJ</b>	Adjective (positive)
<b>JJR</b>	Adjective (comparative)
<b>JJS</b>	Adjective (superlative)
<b>LS</b>	List item marker
<b>MD</b>	Modal verb
<b>NN</b>	Noun, singular or mass
<b>NNS</b>	Noun, plural
<b>NP</b>	Proper noun, singular
<b>NPS</b>	Proper noun, plural
<b>PDT</b>	Predeterminer
<b>POS</b>	Possessive ending ('s)
<b>PP</b>	Personal pronoun
<b>PP\$</b>	Possessive pronoun
<b>RB</b>	Adverb
<b>RP</b>	Particle
<b>SYM</b>	Symbol (mathematical/scientific)
<b>TO</b>	<i>to</i> (any usage) <i>fly to Paris, ready to go, ...</i>
<b>UH</b>	Interjection
<b>#</b>	Pound sign £
<b>\$</b>	Dollar sign \$

<b>VB</b>	Verb <i>be</i> , base form
<b>VBD</b>	Verb <i>be</i> , past tense
<b>VBG</b>	Verb <i>be</i> , gerund/progressive
<b>VBN</b>	Verb <i>be</i> , past participle
<b>VBP</b>	Verb <i>be</i> , non-3rd pers. sg. present
<b>VBZ</b>	Verb <i>be</i> , 3rd pers. sg. present tense
<b>VH</b>	Verb <i>have</i> , base form
<b>VHD</b>	Verb <i>have</i> , past tense
<b>VHG</b>	Verb <i>have</i> , gerund/progressive
<b>VHN</b>	Verb <i>have</i> , past participle
<b>VHP</b>	Verb <i>have</i> , non-3rd pers. sg. present
<b>VHZ</b>	Verb <i>have</i> , 3rd pers. sg. present tense
<b>VV</b>	Lexical verb, base form
<b>VVD</b>	Lexical verb, past tense
<b>VVG</b>	Lexical verb, gerund/progressive
<b>VVN</b>	Lexical verb, past participle
<b>VVP</b>	Lexical verb, non-3rd pers. sg. present
<b>VVZ</b>	Lexical verb, 3rd pers. sg. present tense
<b>WDT</b>	Wh-determiner
<b>WP</b>	Wh-pronoun
<b>WP\$</b>	Possessive wh-pronoun
<b>WRB</b>	Wh-adverb
<b>SENT</b>	Sentence-final punctuation . ! ?
<b>,</b>	Comma ,
<b>:</b>	Colon, semi-colon : ;
<b>( )</b>	Comma ( [ ]
<b>` ' "</b>	Comma " ... ' ' " "

## Universal-Dependencies-Tags (sprachübergreifend)

Open class words	Closed class words	Other
<u>ADJ</u>	<u>ADP</u>	<u>PUNCT</u>
<u>ADV</u>	<u>AUX</u>	<u>SYM</u>
<u>INTJ</u>	<u>CCONJ</u>	<u>X</u>
<u>NOUN</u>	<u>DET</u>	
<u>PROPN</u>	<u>NUM</u>	
<u>VERB</u>	<u>PART</u>	
	<u>PRON</u>	
	<u>SCONJ</u>	

- ADJ: adjective
- ADP: adposition
- ADV: adverb
- AUX: auxiliary
- CCONJ: coordinating conjunction
- DET: determiner
- INTJ: interjection
- NOUN: noun
- NUM: numeral
- PART: particle
- PRON: pronoun
- PROPN: proper noun
- PUNCT: punctuation
- SCONJ: subordinating conjunction
- SYM: symbol
- VERB: verb
- X: other

## Segmente und Strukturen

- Erkennung von speziellen Wortfolgen (Segmenten bzw. *spans*) und ihre Kategorisierung
- z.B. Eigennamen (**NER** = *named entity recognition*)

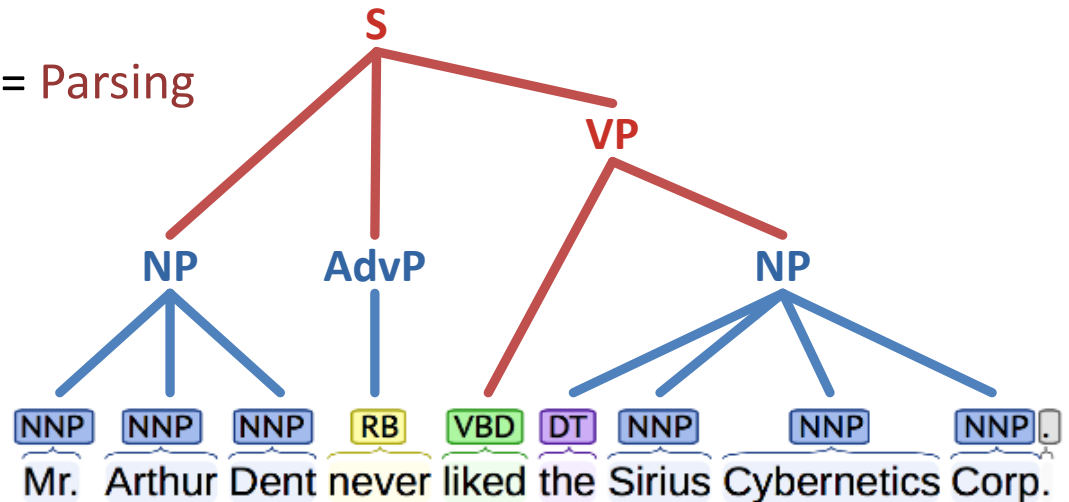
Person Organization  
Mr. Arthur Dent never liked the Sirius Cybernetics Corp.

- Fallen Ihnen noch weitere Beispiele für interessante Segmente ein?
- Kann auch als Tagging operationalisiert werden
  - **B-PERS** **I-PERS** ○ ○ ○ **B-ORG** **I-ORG** **I-ORG**  
Mr. Arthur Dent never liked the Sirius Cybern. Corp.

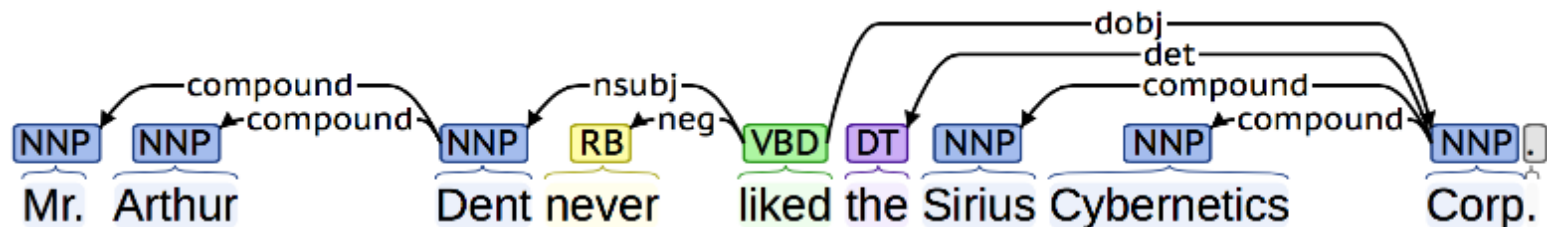


## Strukturen: Parsing

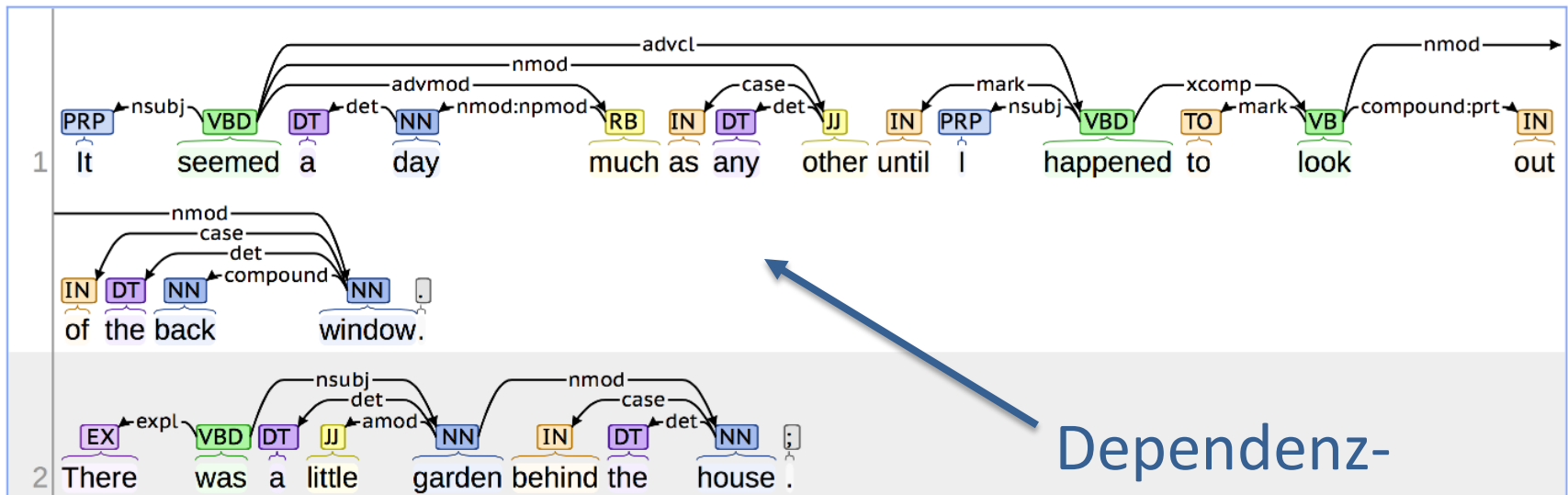
- Erkennung der Satzstruktur = **Parsing**
- Phrasenstruktur als baumförmige Hierarchie



- alternativ: „minimale“ Phrasen als flache Segmente → **Chunk-Parsing**
- **Dependenz-Parsing** findet direkte Abhängigkeiten zwischen Wörtern



## Beispiel: Syntaktische Analyse



Dependenz-  
Graph

Zum Ausprobieren:

- <http://corenlp.run/>
- <https://explosion.ai/demos/displacy>

# Manuelle Annotation

- Kleine Korpora werden oft manuell annotiert
  - z.B. digitale Editionen, Reden eines Präsidenten, ...
- **Annotationsschema** und -kategorien (**Tagset**)
- **Richtlinien** (**Guidelines**)
  - detaillierte Beschreibung und Abgrenzung der Zielkategorien (z.B. für [STTS](#))
  - zusätzlich: Beispielsammlung für schwierige Einzelfälle
- Annotationswerkzeuge (meist Web-basiert)
  - z.B. INCEpTION (<https://inception-project.github.io>), Prodigy (<https://prodi.gy>)
- **Inter-Annotator Agreement** (IAA)
  - wichtig! – überprüft Reliabilität und Validität der Annotation
  - Flüchtigkeitsfehler vs. systematische Differenzen
  - Adjudikation für Endfassung der Annotation

## Automatische Annotation

- Für größere Korpora ist eine manuelle Annotation zu teuer und zeitaufwendig
- Auch in den Digital Humanities ...
  - Romane von Charles Dickens ca. 4 Mio. Wörter
  - Deutsches Gutenberg-Archiv > 100 Mio. Wörter
  - Early English Books (EEBO) > 500 Mio. Wörter
  - Times Online 1780–1900 ca. 4.000 Mio. Wörter

# Automatische Annotation

- Erfolgreichster Ansatz: **maschinelle Lernverfahren**
  - ab ca. 1990 Einsatz von statistischen Modellen („**statistical revolution**“)
  - aktuell große Fortschritte mit **Deep Learning**
- **Trainingskorpus** (manuell annotiert)
  - wichtig: Konsistenz der Annotationen (→ IAA)
  - Flüchtigkeitsfehler scheinen weniger problematisch
- Evaluation auf separatem Testkorpus
  - Gefahr der Überanpassung an das Trainingskorpus
  - zusätzliches **development set** für Optimierung der Lernverfahren (**tuning**)
  - Kreuzvalidierung (**cross-validation**) nutzt alle Daten für Training & Evaluation
- Weiterführend: <https://web.stanford.edu/~jurafsky/slp3/>

# Repräsentationsformat: XML

```
<?xml version="1.0" encoding="UTF-8"?>
<corpus>
  <story title="The Garden">
    <p>
      <s>
        <token pos="PP" lemma="it">It</token>
        <token pos="VBD" lemma="seem">seemed</token>
        <token pos="DT" lemma="a">a</token>
        <token pos="NN" lemma="day">day</token>
        <token pos="RB" lemma="much">much</token>
        <token pos="IN" lemma="as">as</token>
        <token pos="DT" lemma="any">any</token>
        <token pos="JJ" lemma="other">other</token>
        <token pos="IN" lemma="until">until</token>
        <token pos="PP" lemma="I">I</token>
        ...
      </s>
    </p>
  </story>
</corpus>
```

## Repräsentationsformat: Vertical text (.vrt)

```
<corpus>
<text title="The Garden" author="Stefan Evert" author_sex="male"
      date="1991-08-05">
<p num="1">
<s>
It      PP    it
seemed VBD    seem
a       DT    a
day     NN    day
much    RB    much
as      IN    as
any     DT    any
other   JJ    other
until   IN    until
I       PP    I
...
</s>
</p>
</text>
</corpus>
```

## Repräsentationsformat: CoNLL-Format(e)

# story: "The Garden"

# paragraph #1

1	It	PP	it
2	seemed	VBD	seem
3	a	DT	a
4	fine	JJ	fine
5	day	NN	day
6	.	SENT	.

1	There	EX	there
2	was	VBD	be
3	an	DT	a
4	elephant	NN	elephant
5	.	SENT	.

# this is the end of the file

aktuell: CoNLL-U (<https://universaldependencies.org/format.html>)



# Übersicht: Tools

## Manuelle Annotation: Tools

- WebAnno / **INCEpTION** (linguistischer Fokus):
  - <https://webanno.github.io/webanno/documentation/>
  - <https://www.youtube.com/user/webanno>
  - <https://inception-project.github.io>
  - <https://youtube.com/playlist?list=PL5Hz5pttaj96SlXHGRZf8KzIYvpVHIoL->
- **prodigy** (linguistischer Fokus):
  - <https://prodi.gy>
- **CATMA** (literaturwissenschaftlicher Fokus)
  - z.B. Annotation von wörtlicher und indirekter Rede
  - <https://fortext.net/routinen/lerneinheiten/manuelle-annotation-mit-catma>

## Automatische Annotation: komplette Pipelines (1)

- Stanford **CoreNLP** (<https://stanfordnlp.github.io/CoreNLP/>)
  - langlaufendes Projekt, Java
  - Tokenisierung, POS-Tagging, Lemmatisierung, NER, Parsing, Koreferenzauflösung, Sentiment-Analyse, ...
- **Stanza** (<https://stanfordnlp.github.io/stanza/>)
  - Python, Deep Learning, Interface zu CoreNLP (z.B. für Koreferenzauflösung relevant)
  - Tokenisierung, POS-Tagging, Lemmatisierung, NER, Dependenzparsing, Sentiment-Analyse
- **spaCy** – „fastest in the world“ (<https://spacy.io>)
  - Python, Deep Learning
  - Tokenisierung, POS-Tagging, Lemmatisierung, NER, Dependenzparsing

## Automatische Annotation: komplette Pipelines (2)

- Apache **OpenNLP** (<https://opennlp.apache.org/>)
  - Java
  - Tokenisierung, POS-Tagging, Lemmatisierung, NER, Dependenzparsing, Koreferenzauflösung
- **UDPipe** (<http://ufal.mff.cuni.cz/udpipe>)
  - C++/Python, als Bibliothek für diverse C++, C#, Python, Perl und Java verfügbar
  - *UD* steht für *Universal Dependencies*
  - Tokenisierung, POS-Tagging, Lemmatisierung, Dependenzparsing

# Automatische Annotation: Tokenisierung und Tagging

- reine **Tokenisierer**
  - Python: [SoMaJo](#) (DE, EN)
  - generischer Tokenisierer: [Unitok](#)
  - Tokenisierer von NLTK ist bestenfalls mittelmäßig
  - wichtig: Tokenisierung und weitere Verarbeitung müssen kompatibel sein!
- Part-of-speech-**Tagger** (oft mit eigenem Tokenisierer)
  - [TreeTagger](#) (schnell, einfach, viele Sprachen, inkl. Lemmatisierung)
  - [RNNTagger](#) (Deep-Learning-Nachfolger des TreeTaggers, Python, inkl. Lemm.)
  - Python: [SoMeWeTa](#) (DE, EN, FR)
  - Twitter data (EN): [TweetNLP](#)
  - und viele weitere spezialisierte Tokenisierer und Tagger für diverse Sprachen
- Eigene Pipeline im Webservice erstellen:  
[https://weblicht.sfs.uni-tuebingen.de/weblichtwiki/index.php/Main\\_Page](https://weblicht.sfs.uni-tuebingen.de/weblichtwiki/index.php/Main_Page)