# Morphology in CLARIN-D

Daniël de Kok

# Introduction

A whirlwind introduction:

- CLARIN-D tools: WebLicht, TüNDRA
- Resources: corpora with morphology
- Mostly oriented towards inflectional morphology

# WebLicht

WebLicht is a web application for creating and running NLP pipelines

# Services

- Centers provide RESTful annotation services
  - Input: Text Corpus Format (TCF)
  - Output: TCF with the added layers
- Centers create metadata for their annotations services and put them in their repository

# WebLicht architecture

# Morphology services

- Currently available (morphological tagging):
  - German: **Stuttgart Morphology (RFTagger)**, SMOR
  - Dutch: **Alpino**
  - English: MorphAdorner
- Adding new services for morphology:
  - Since WebLicht is decentralized, any CLARIN center could add additional morphology services.
  - If some interesting tool is missing, let us know!

# Stuttgart morphology (German)

- HMM tagger specialized for large, feature-rich tag sets.
- Trained on the Tiger treebank.
- Uses a supplementary lexicon.
- Outputs morphological tags in the TIGER morphology scheme:
  - Part-of-speech
  - Gender
  - Case
  - Number
  - Degree
  - Person
  - Tense
  - Mood
  - Finiteness

# Alpino (Dutch)

- Wide-coverage dependency parser for Dutch.
- But also has:
    - An extensive lexicon with subcategorization frames.
    - A guesser for unknown words.
- Eventual frames are decided by:
    - Filtering by n-best tagging.
    - The parse selected by the disambiguation model.

# Resources for German

- Semi-automatically annotated
  - Tiger treebank
  - TüBa-D/Z
- Automatically annotated
  - TüBa-D/W

# Tiger treebank

- ~50,000 sentences
- Newspaper text (Frankfurter Rundschau)
- Semi-automatically annotated
- Annotations:
    - STTS part-of-speech tags
    - Lemmas
    - Inflectional morphology
    - Constituency structure
    - Dependency conversion (subset hand-annotated)

# TüBa-D/Z

- ~95,500 sentences
- Newspaper text (*taz*)
- Semi-automatically annotated
- Annotations:
  - STTS part-of-speech tags
  - Lemmas
  - Inflectional morphology
  - Constituency structure
  - Dependency conversion
  - Anaphora and coreference relations
  - Subset with GermaNet word senses
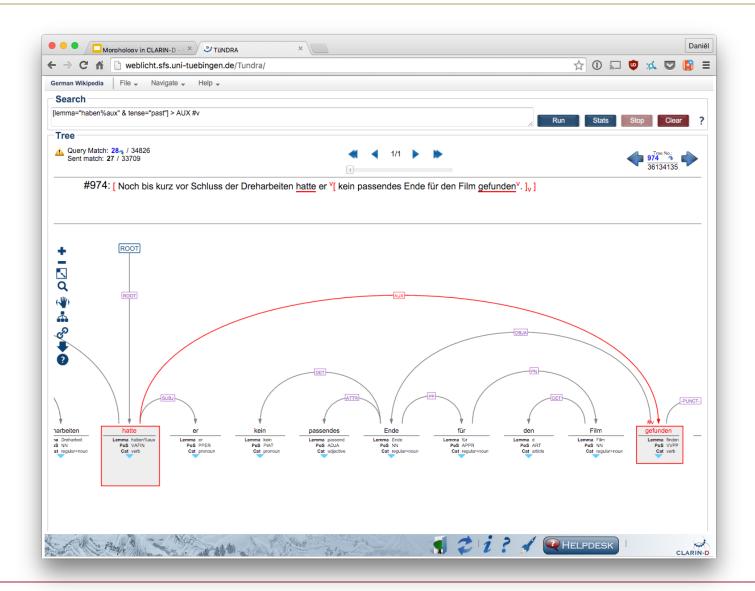  - Named entity class

# TüBa-D/W

- 36.1 million sentences
- German Wikipedia
- Automatically annotated
- Annotations:
  - STTS part-of-speech tags
  - Lemmas
  - Inflectional morphology
  - Dependency structure
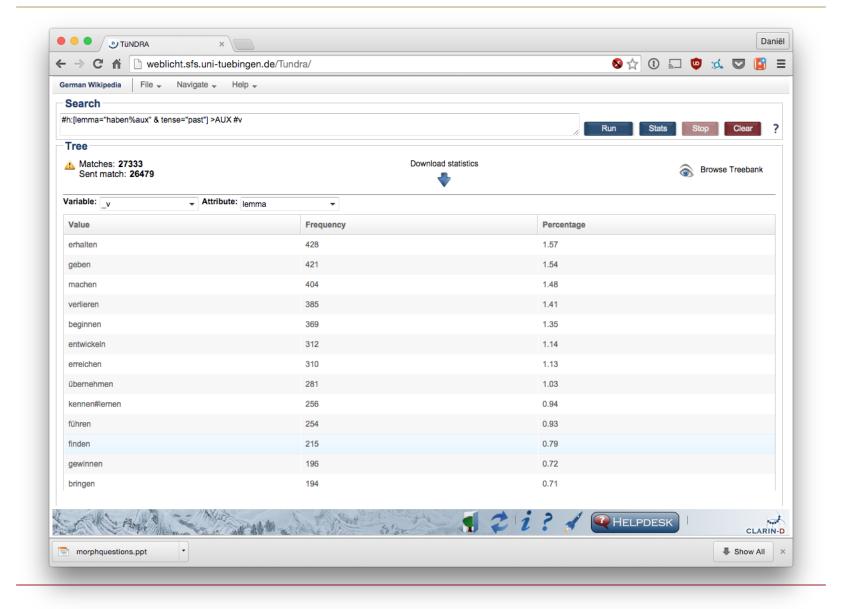- Processed using WebLicht :-)

# TüBa-D/W

TüBa-D/W is fully searchable using the TüNDRA treebank viewer

# Links

WebLicht:

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

TüNDRA:

https://weblicht.sfs.uni-tuebingen.de/Tundra/