# Interactive analysis of multi-layer linguistic corpora with ANNIS

Florian Zipser*, Thomas Krause*, Arne Neumann+
*Humboldt-Universität zu Berlin, +Universität Potsdam

http://annis-tools.org

## Goal

- Higher usability through simplified query syntax (AQL)
- Improvements of AQL for new operators: Equality and Inequality for values, OR Operator and explicit variable naming
- Simple statisitcal evaluations direclty in ANNIS with frequency analysis

## Background

- linguistic phenomena can be spread over different annotation layers, only the combination of annotations will find them
- several multi-layer corpora have been developed for this purpose: TüBa-D/Z (Telljohann et al. 2009), PCC (Stede & Neumann 2014) or Falko (Reznicek et al. 2012)
- annotation layers differ in content and structure: token annotations, constituency trees, dependency trees, rhetorical structure annotation, spoken data, etc.



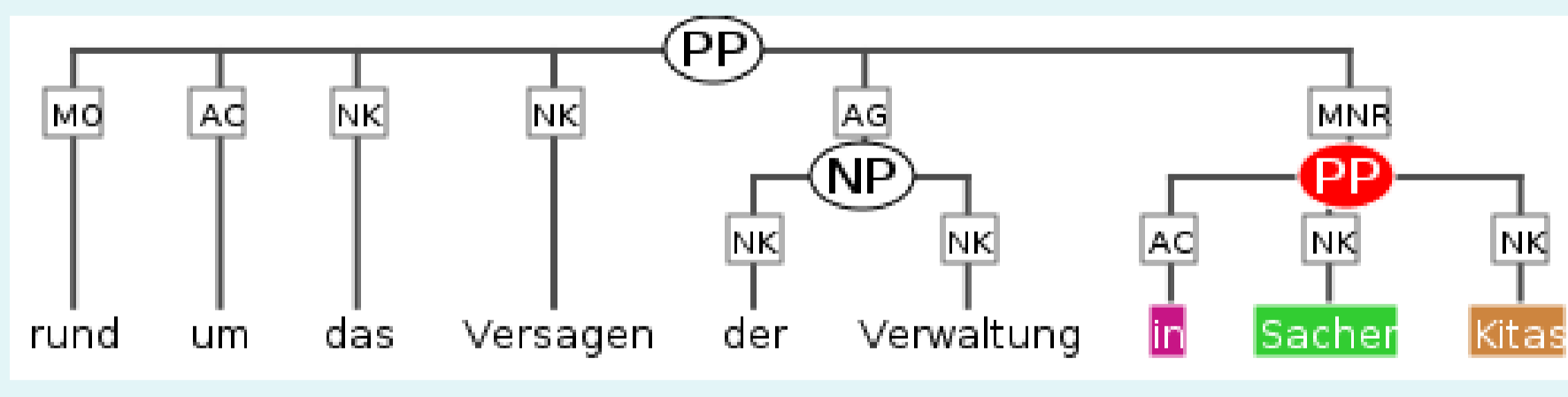→ **We need a unified search and visualization system**

## ANNIS

- not limited to a specific type of annotation
- not limited to a single corpus
- same query language for different corpora (AQL)
- specialized and configurable visualizations for different annotation layers
- support for existing corpora from many formats via SaltNPepper (Zipser & Romary 2010)
- export for further statistical evaluation: CSV, plain text or ARFF (WEKA (Hall et al. 2009))

## Query language improvements

- **simplified query syntax (better readability)**
  *Prepositional phrase with a complex preposition (preposition + noun + noun)*

- so far*:
```
cat="PP" & pos & pos="NN" & pos="NN"
& #1 >[func="AC"] #2
& #2 . #3
& #3 . #4
& #1 >[func="NK"] #3
& #1 >[func="NK"] #4
```



- simplified*:
```
cat="PP" >[func="AC"] pos . pos="NN" . pos="NN"
& #1 >[func="NK"] #3
& #1 >[func="NK"] #4
```

- **search for equal/different value**
  *two same/different part-of-speech tags, one directly following the other*
```
pos . pos & #1 == #2
```
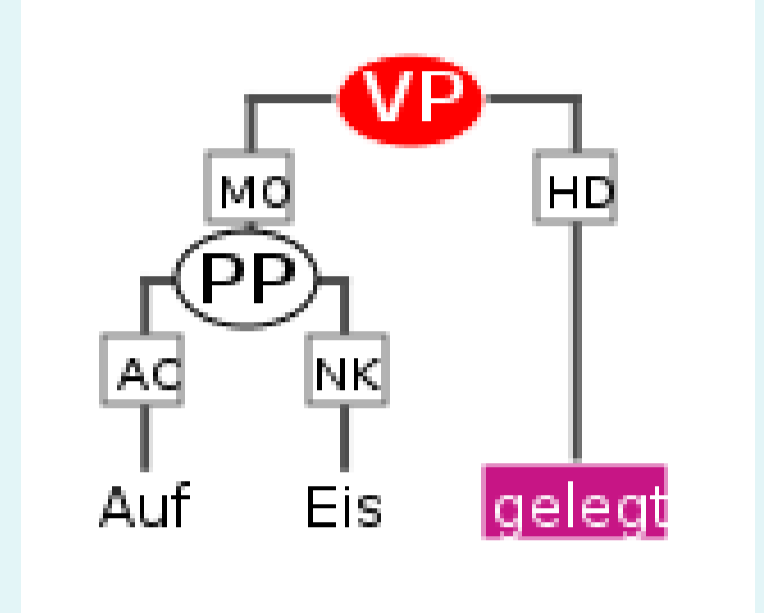```
pos . pos & #1 != #2
```



* queries use the TIGER annotation scheme

## Query language improvements

- **OR operator and variable naming**

  *A verbal phrase, which dominates either an adjective or the token „gelegt" which is often misclassified by automatic POS-taggers*

```
a#cat="VP"
&(b#pos="ADJA" | b#"gelegt")
& #a > #b
```



## Frequency analysis

- **frequency analysis directly in ANNIS**



Example 1:

```
pos . pos="NN"
```

*What are the most frequent part-of-speech tags followed by a noun?*



Example 2:

```
cat="S" > cat="PP" > pos
```

*What are the most frequent part-of-speech tags in a prepositional phrase which is in a sentence?*

## Outlook

- meta data in frequency analysis
- syntax highlighting for AQL
- implement a new and much faster/scalable graph database backend as a replacement for the relational database PostgreSQL

## References

- **M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, I. H. Witten. 2009.** The WEKA Data Mining Software: An Update; SIGKDD Explorations, Volume 11, Issue 1.
- **M. Reznicek, A. Lüdeling, C. Krummes, F. Schwantuschke, M. Walter, K.Schmidt,H. Hirschmann, T. Andreas. 2012.** Das Falko-Handbuch. Korpusaufbau und Annotationen Version 2.01. https://www.linguistik.hu-berlin.de/institut/professuren/korpuslinguistik/forschung/falko/Falko-Handbuch_Korpusaufbau%20und%20Annotationen_v2.01
- **M. Stede, A. Neumann 2014.** Potsdam Commentary Corpus 2.0: Annotation for Discourse Research. Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14). Reykjavik, Island.
- **H. Telljohann, E. W. Hinrichs, S. Kübler, H. Zinsmeister, K. Beck. 2009.** Stylebook for the Tübingen Treebank of Written German (TüBa-D/Z). Universität Tübingen Seminar für Sprachwissenschaft.
- **F. Zipser, L. Romary. 2010.** A Model Oriented Approach to the Mapping of Annotation Formats using Standards. Workshop Language Resource & Language Technology Standards, LREC 2010. Malta, 7-18. http://u.hu-berlin.de/saltnpepper

Web: http://annis-tools.org
Email: annis-admin@ling.uni-potsdam.de