# Extending and Exploiting the Entity Graph for Analysis, Classification and Visualization of German Texts

Julia Suter and Michael Strube, Heidelberg Institute for Theoretical Studies

## In a Nutshell

- Enhancing the entity graph
- Visualizing and analyzing German texts using the Entity graph
- Text classification with entity graph metrics

## Entity Graph

One mode projection $(P_p$ and $P_{pu})$ captures the relations between sentences

## Improvements with Adjacent Syntactic Categories

Evaluation by sentence reordering task on TüBa-D/Z news corpus

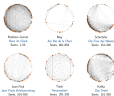| Original Entity Graph | Adjusted Entity Graph |
|---|---|

## Text Analysis by Graph Metrics

- Visualization of entity graph reveals differences in structure, text coherence among different authors and text types
- Graph metrics describe text properties

Text sample of 30 sentences extracted from Project Gutenberg.
Syntactic parsing using ParZu, coreference resolution using CorZu.

## Graph Metrics for German-Literary Texts

## Author and Genre Classification

- Support Vector Classifier
- 24 features extracted from entity graph : centrality measures, edge weights, diameter, maximum flow, edge betweenness, clustering coefficient, etc.
- Comparison to 31 syntactic features

Class adding entity graph features improves the compare / feature space?

| Author Classification | | Genre Classification | |
|---|---|---|---|

- Entity graph features capture information not encoded in syntactic features

## Conclusion

- Preposition modifiers category and indexed weights improve entity graph
- Entity graph visualization and graph metrics are suitable for analyzing texts
- Contribution of new information to author and genre classification