# Speech Rhythm and Syntax in Poetry and Prose

# Thomas Haider[1,2], Debby Trzeciak[1,3], Gerrit Kentner[1,4]

[1] MPI für Empirische Aesthetik, [2] Universität Stuttgart, [3] TU Darmstadt, [4] Goethe Universität Frankfurt

thomas.haider@ae.mpg.de, debby.trzeciak@ae.mpg.de, gerrit.kentner@ae.mpg.de

**CLP**

Prosody undisputedly affects the choice of syntactic constructions and the order of constituents within a sentence (Anttila 2016). The influence of prosody on syntax is most obviously attested in metered poetry where strict metrical rules and poetic license influence word order and grammaticality (Donat 2010). However, for German, it is largely unclear which prosodic factors affect sentence construction, and how strong their influence is on grammatical encoding.

To analyze the interaction of syntax with meter and rhythm, we train Conditional Random Fields (CRF) with sklearn-crfsuite to annotate large scale for *part-of-speech* (POS), *binary meter* (BM), and *ternary rhythm* (TR). For POS tagging we rely on gold annotation from DTA (deutsches-textarchiv) and TIGER, according to the STTS tagset. For BM and TR we manually annotate 3600 lines of school canon poetry (158 poems) with Cohen κ > .90 on syllable level.

BM includes binary syllable prominence (+/-) and foot boundaries (|). TR segments the verse into *rhythmic groups* at caesuras (|) and in these segments allows for main accent (2), side accent (1), and no accent (0). See example (1) for an annotated line. We use the CRF models to annotate the German Poetry Corpus DLK (Haider & Eger, 2019) (74k poems, 1.6m lines, 11m token) on all three levels. We achieve up to 95% accuracy for POS and BM, while TR still has substantial 85% accuracy. The TR model and humans similarly confuse main and side accents.

(1) <l met="-+|-+|-+|-+|-+|" rhythm="01020|20102|">Geduckte Hütten, Pfade wirr verstreut,</l>

In our experiments we investigate:

1. the likelihood of a certain POS class to be stressed or unstressed, which allows us to establish a stress hierarchy like Antilla et al. (2018), but obtaining differing results: NOUN > VERB$_{modal}$ > VERB$_{full}$ > ADJ > ADV > FUNC. We agree however, that function words (FUNC, e.g. KONJ, ART) are seldom stressed, while nouns are usually stressed.

2. how words change their stress based on their context. We find that the established hierarchy reiterates for contextual dependence. If a word is preceded by a conjunction, then likelihood of stress is higher. Nouns rarely lose their prominence, and adverbs, which are quite balanced, also show a balanced context dependence.

3. the interaction of enjambement (line end without punctuation) with part-of-speech and line measures. We find no obvious differences (beyond tendencies) in POS transitions between lines, but shorter lines and hexameter prefer enjambement, while the alexandrine disprefers it.

4. the distinction of prose vs. poetry with a regularized linear discriminant analysis to interpret feature weights. Classifying single sentences reaches barely over random baseline while whole documents can be classified on POS n-grams (.93 Acc.) and rhythmic groups (.84 Acc.), suggesting infrequent recurring types of syntactic inversion (Gopidi & Alam 2019).

**References:** Anttila, A. (2016). Phonological effects on syntactic variation. Anttila, A. et al. (2018). Sentence stress in presidential speeches. Donat, S. (2010). Deskriptive Metrik. Gopidi, A. & Alam, A. (2019). Computational Analysis of the Historical Changes in Poetry and Prose. Haider, T. & Eger, S. (2019). Semantic Change and Emerging Tropes in a Large Corpus of New High German Poetry.