

1 General information

German Named Entity Recognizer

Gate group¹, Department of Computer Science, University of Sheffield, UK

Status: stable

Languages covered: German

2 Implemented NLP services

The German IE module consists of the following main language processing tools: tokeniser, sentence splitter, named entity recogniser and orthomatcher.

The named entity recogniser identifies and categorizes entity names (such as persons, professions, organizations, and location names), temporal expressions (dates and times), and numerical expressions.

For this purpose, it uses three types of processing resources: a gazetteer, a part of speech tagger and a rule grammar module. The gazetteer consists of lists such as cities, organizations, days of the week, etc. It not only consists of entities, but also of names of useful indicators, such as typical company designators (e.g. 'Verein.'), titles, etc. The gazetteer lists are compiled into finite state machines, which can match text tokens. The grammar component allows the encoding of rules that operate on the output of the gazetteer in order to annotate text spans with the relevant named entity types. There is also a grammar module for the analysis of nominal compounds. The orthomatcher establishes co-reference relations between textual objects.

The text spans and annotations are exported into an RDF ontology, in which the named entity types such as Organization and Person constitute classes, and the text spans instances of these classes.

3 Language resource standards

No language resource standards are required for input.

4 Linguistic data encoding

The service does not assume any data categories as input.

It will run over a text corpus containing documents from a large variety of formats:

plain text

¹ <http://www.gate.ac.uk>

HTML

SGML

XML

RTF

PDF (not all)

Microsoft Word (not all)