# A combined presentation and demonstration of the INESS treebanking infrastructure

Who? Petter Haugereid
Department of Linguistic, Literary and Aesthetic Studies
University of Bergen

When? DELPH-IN meeting
St. Wendel, July 29 – August 2, 2013

# INESS – *Infrastructure for the Exploration of Syntax and Semantics*

- Carried out by
- The University of Bergen (Norway)
  - Helge Dyvik, Victoria Rosén, Koenraad De Smedt, Gunn Inger Lyse , Gyri Smørdal Losnegaard , Martha Thunes, Petter Haugereid
- Uni Computing (a division of Uni Research, Bergen)
  - Paul Meurer
- Funded by
- The Research Council of Norway
- The University of Bergen
- Duration: 2010–2016

# INESS – Main activities

Main activities:

- The development of a large, deep LFG parsebank for Norwegian
- Implementation and operation of a comprehensive open treebanking environment for building, hosting and exploring treebanks

# Hosting treebanks

- INESS currently provides the most comprehensive web-based treebanking services available
- A normal browser sufficient for:
  - accessing, searching and downloading treebanks
  - annotation of LFG-based parsebanks:
    - computer-aided manual disambiguation
    - text cleanup
    - handling of unknown words

# Treebanks hosted by INESS

- The Icelandic Parsed Historical Corpus (IcePaHC)
- 73,014 sentences
- The German Tiger treebank
- 50,472 sentences with dependency annotation
- 9,221 with LFG annotation
- The dependency part of the Bulgarian BulTreeBank
- 11,900 sentences
- HPSG treebanks
- WeScience: 48,000 sentences
- DeepBank: 37,000 sentences
- and other smaller corpora

# INESS – Parallel corpora

The INESS infrastructure offers tools for the development and exploitation of parallel corpora

- Excerpts of the novel *Sofies verden* (*Sophie's World*)
  - 26 aligned language pairs
- A document from the *Acquis communautaire* (EU legislation)
  - 21 aligned language pairs

# The INESS Norwegian treebank – Development

- Obtained by parsing automatically with the LFG grammar NorGram on the XLE parsing platform
- `http://iness.uib.no/redmine/projects/inesspublic/wiki/NorGram_documentation`
- Consists of a number of different text types in fiction and non-fiction
- Part of the treebank is manually disambiguated with the LFG Parsebanker (Rosén et al., 2009)

# The INESS Norwegian treebank – Data

- Manually (at least partially) disambiguated:
- 8,405 sentences (81,934 words)
- Fully disambiguated:
- 6,966 sentences (65,599 words)

# Stochastic disambiguator

- A stochastic disambiguator is implemented based on the manually disambiguated sentences
- Used to rank the analyses of a given sentence
- Corpus searches are limited to the top ranked analyses in the parsebank
- The manually selected discriminants are saved and reused in case the corpus is reparsed

# Treebank selection page

# DEMO

http://iness.uib.no/iness/main-page