# A Norwegian Language Grid

## (A 'Vision' from the Semi-Outside)

**Stephan Oepen**

Universitetet i Oslo, Institutt for Informatikk

`oe@ifi.uio.no`

(CLARIN Norge — June 18, 2010)

# D-SPIN: Language Resources & Technology On-Line

**Pros and Cons of a Web-Based SOA:**

$\vdots$

– *Not applicable for huge amounts of data.*

$\vdots$

# The IFI Language Technology Group

## Table of Contents

| | | |
|---|---|---|
| Gordana Ilić Holen | Doctoral Fellow | Coreference Resolution |
| Elisabeth Lien | Doctoral Fellow | Textual Inference |
| Jan Tore Lønning | Professor | Computational Semantics |
| Stephan Oepen | Professor | Grammar-Based Processing |
| Woodley Packard | Doctoral Fellow | Joint Disambiguation |
| Erik Velldal | Post-Doctoral Fellow | Classification |
| Gisle Ytrestøl | Doctoral Fellow | Incremental Parsing |
| Aleksander Øhrn | Adjunct Professor | Information Retrieval |
| Lilja Øvrelid | Post-Doctoral Fellow | Data-Driven NLP |
| NN | Associate Professor | Empirical Methods |
| NN | Doctoral Fellow | High-Quality Research |

# The IFI Language Technology Group

## Table of Contents

| Gordana Ilić Holen | Doctoral Fellow | Cor *Spring 2009* ution |
| Elisabeth Lien | Doctoral Fellow | Te *Fall 2009* ce |
| Jan Tore Lønning | Professor | Computational Semantics |
| Stephan Oepen | Professor | Grammar-Based Processing |
| Woodley Packard | Doctoral Fellow | Joi *Spring 2010* on |
| Erik Velldal | Post-Doctoral Fellow | ( *Fall 2009* |
| Gisle Ytrestøl | Doctoral Fellow | Incremental Parsing |
| Aleksander Øhrn | Adjunct Professor | Inf *Spring 2010* val |
| Lilja Øvrelid | Post-Doctoral Fellow | Da *Fall 2010* LP |
| NN | Associate Professor | En *Fall 2010* ods |
| NN | Doctoral Fellow | High *Fall 2010* earch |

# The IFI Language Technology Group

# An Example: Syntacto-Semantic Analysis of Wikipedia

## General Idea

- Enabling technology: Wikipedia as a corpus and a knowledge source;

- e.g. research in linguistics, lexical acquisition, ontology learning, etc.
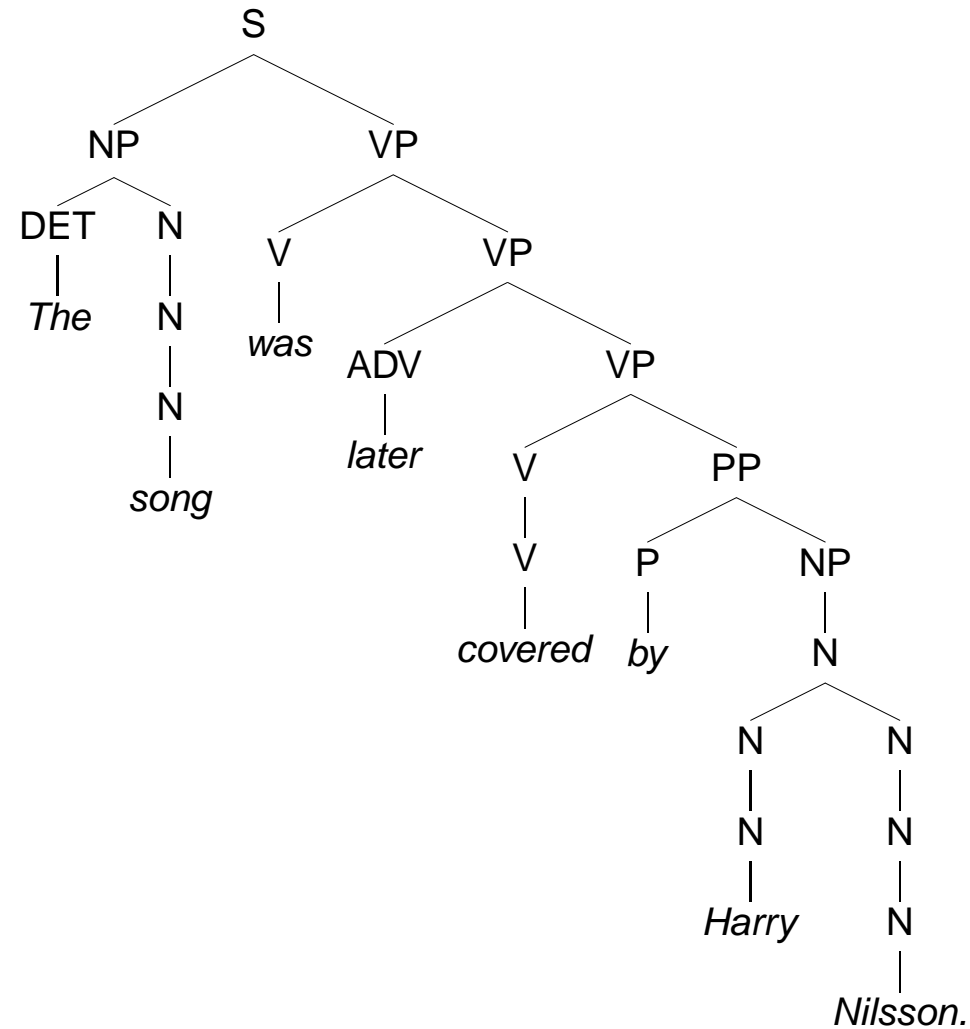
## Approach & Technology

- Semi-automated 'deep' linguistic annotation, from pre-existing parser;

- gold-standard annotation of domain-specific subset: ~250,000 words.
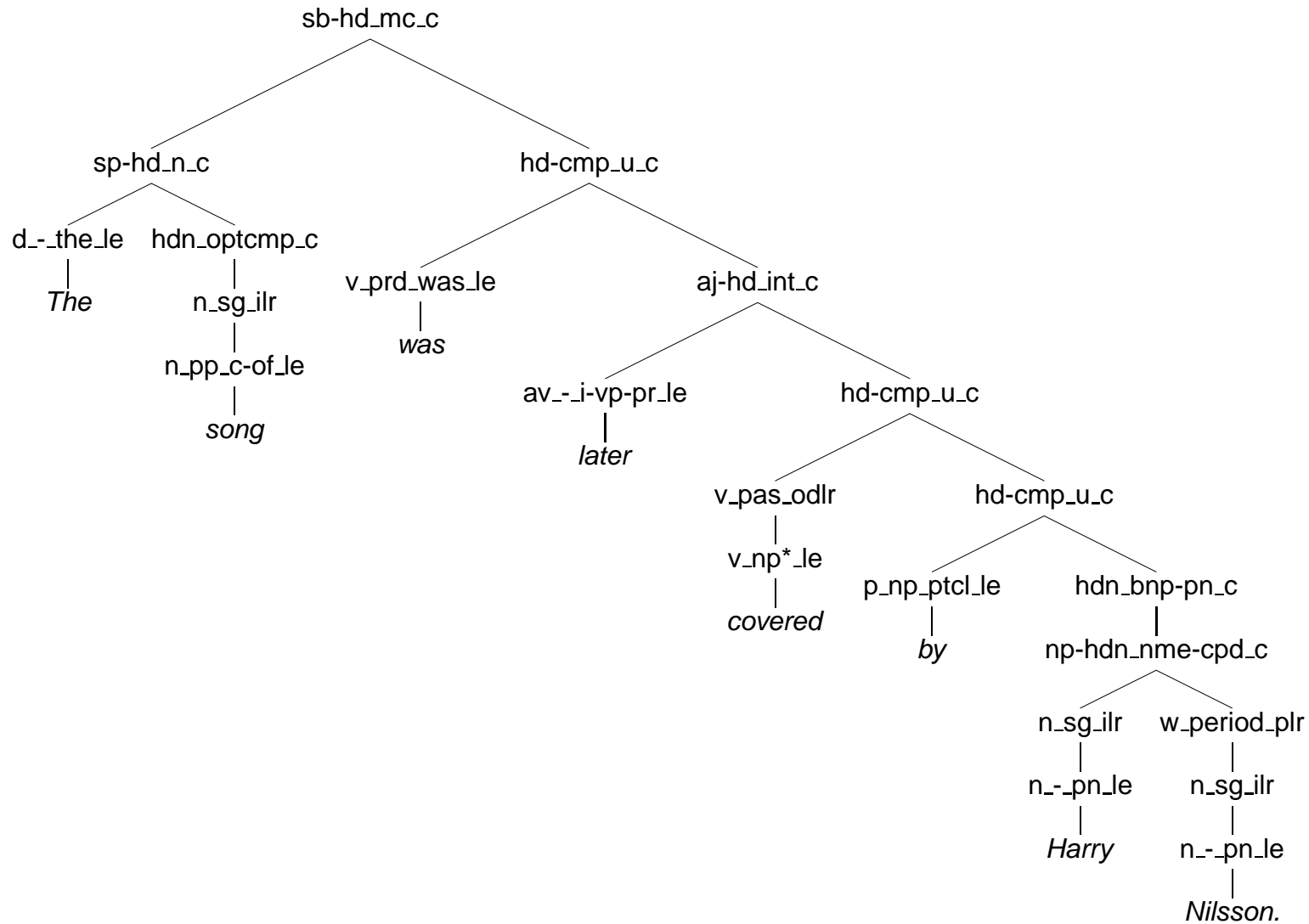
## More Information (Download Site)

`http://www.delph-in.net/wikiwoods`

# Syntactic Annotation: 'Classic' Constituent Tree

# Syntactic Annotation: HPSG Derivation

sb-hd_mc_c
- sp-hd_n_c
  - d_-_the_le — *The*
  - hdn_optcmp_c
    - n_sg_ilr
      - n_pp_c-of_le — *song*
- hd-cmp_u_c
  - v_prd_was_le — *was*
  - aj-hd_int_c
    - av_-_i-vp-pr_le — *later*
    - hd-cmp_u_c
      - v_pas_odlr
        - v_np*_le — *covered*
      - hd-cmp_u_c
        - p_np_ptcl_le — *by*
        - hdn_bnp-pn_c
          - np-hdn_nme-cpd_c
            - n_sg_ilr
              - n_-_pn_le — *Harry*
            - w_period_plr
              - n_sg_ilr
                - n_-_pn_le — *Nilsson.*

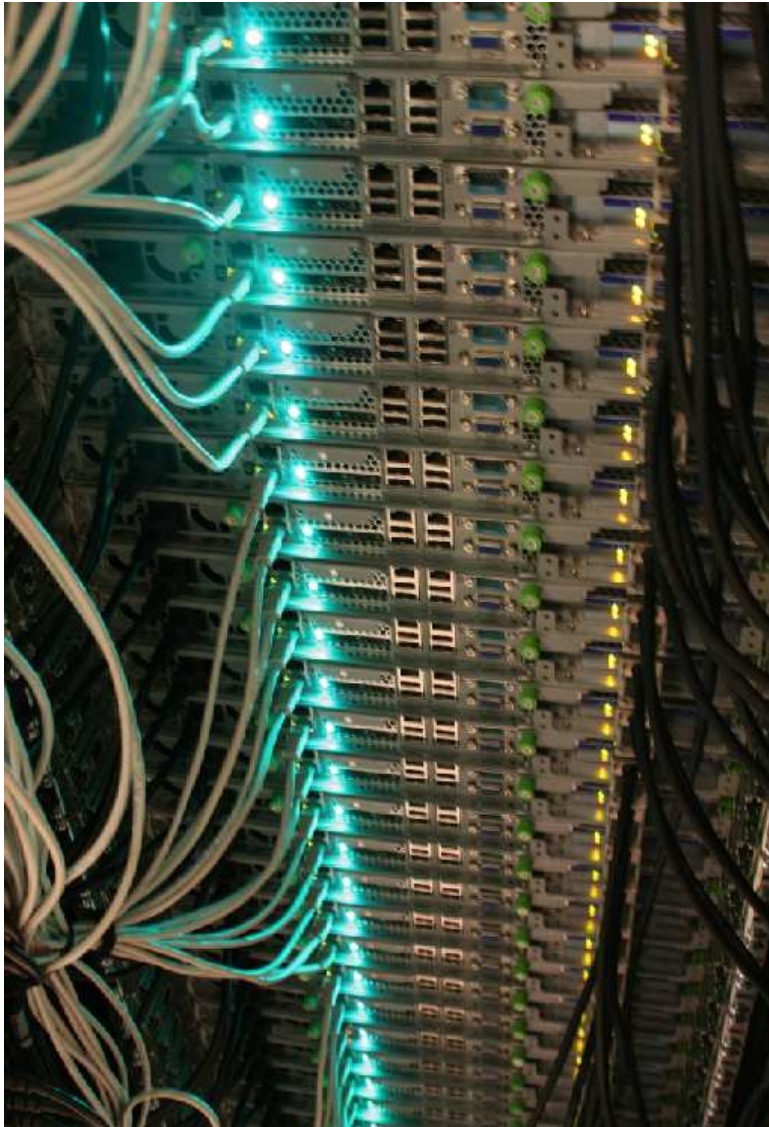# Semantic Annotation: Predicate−Argument Structure

*The song was later covered by Harry Nilsson.*

$\langle\, h_1,$

$h_3{:}\_\text{the\_q}(x_5, h_6, h_4), h_7{:}\_\text{song\_n\_of}(x_5\{\text{PERS 3}, \text{NUM } sg\}, \_\_),$

$h_9{:}\_\text{cover\_v\_1}(e_2\{\text{SF } prop, \text{TENSE } past, \text{MOOD } ind\}, x_{11}, x_5),$

$h_9{:}\_\text{later\_a\_1}(\_\_, e_2),$

$h_{16}{:}\text{compound\_name}(\_\_, x_{11}, x_{17}),$

$h_{19}{:}\text{proper\_q}(x_{17}, h_{20}, h_{21}), h_{22}{:}\text{named}(x_{17}\{\text{PERS } 3, \text{NUM } sg\}, Harry),$

$h_{13}{:}\text{proper\_q}(x_{11}, h_{14}, h_{15}), h_{16}{:}\text{named}(x_{11}\{\text{PERS } 3, \text{NUM } sg\}, Nilsson)$

$\{\, h_{20} =_q h_{22},\, h_{14} =_q h_{16},\, h_6 =_q h_7 \,\}\, \rangle$

# Semantic Annotation: Predicate−Argument Structure

*The song was later covered by Harry Nilsson.*

$\langle h_1,$

$h_3{:}\_the\_q(x_5, h_6, h_4), h_7{:}\_song\_n\_of(x_5\{\text{PERS } 3, \text{NUM } sg\}, \_\_),$

$h_9{:}\_cover\_v\_1(e_2\{\text{SF } prop, \text{TENSE } past, \text{MOOD } ind\}, x_{11}, x_5),$

$h_9{:}\_later\_a\_1(\_\_, e_2),$

$h_{16}{:}compound\_name(\_\_, x_{11}, x_{17}),$

$h_{19}{:}proper\_q(x_{17}, h_{20}, h_{21}), h_{22}{:}named(x_{17}\{\text{PERS } 3, \text{NUM } sg\}, Harry),$

$h_{13}{:}proper\_q(x_{11}, h_{14}, h_{15}), h_{16}{:}named(x_{11}\{\text{PERS } 3, \text{NUM } sg\}, Nilsson)$

$\{ h_{20} =_q h_{22}, h_{14} =_q h_{16}, h_6 =_q h_7 \} \rangle$

→ 1.3 million content articles, 55 million utterances, ~900 million tokens;

→ ~85 % parsing coverage, ~83 % of analyses totally or nearly correct.

# Semantic Annotation: Predicate–Argument Structure

*ter covered by Harry Nilsson.*

$ng\_n\_of(x_5\{\text{PERS } 3, \text{NUM } sg\}, \_\_),$

$\text{TENSE } past, \text{MOOD } ind\}, x_{11}, x_5),$

$\_{11}, x_{17}),$

$h_{22}\text{:named}(x_{17}\{\text{PERS } 3, \text{NUM } sg\}, Harry),$

*~120,000 cpu hours (six days);*
*~130 gigabytes compressed data;*
*→ subject extraction present
in one of 15 utterances;*
*→ ~90 % in relative clauses.*

# A Candidate Role Model: BioPortal at UiO

A Norwegian Language Grid (In Fifteen Minutes) (8)

# Imagine: Language Resources & Technology Portal

> ## Motivation
>
> • Reduce technology barriers: on-line demonstrators *and* processing;
>
> • unified, Web-based point of entry; balance ease of use and flexibility.

> ## Core Components
>
> • **Data**  *Språkbanken*, ELRA, LDC, and others; user-contributed data;
>
> • **Tools**  text extraction (PDF, {HT|W|X}ML, et al.), segmentation, morphology, tagging, chunking, parsing, search, concordancing, etc.

> ## Scalability
>
> • Built on top of national HPC infrastructure: NoTur, NorStore, NorGrid.

# Imagine: Language Resources & Technology Portal

## Motivation

- Reduce technology barriers: on-line demonstrators *and* processing;

- unified, Web-based point of entry; balance ease of use and flexibility.

## Core Components

- **Data**    *Språkbanken*, ELRA, LDC, and others; user-contributed data;

- **Tools**    text extraction (PDF, {HT|W|X}ML, et al.), segmentation,

*Preferably mostly through bottom-up, grass-roots process:
plurarilty of approaches: different frameworks and methods;
some convergence needed: exchange formats and interfaces;
starting points: UIMA, Language Grid, D-SPIN, and others.*

# More Concretely: Short-Term Initiatives

**High-Performance LRT User Group**

- UniNett Sigma ($\Sigma$) looking to establish discipline-specific user groups;

- group-internal functions: exchange experience, coordinate activities;

- interface function to $\Sigma$: give feedback on user experience and needs;

- at least one annual meeting $\rightarrow$ contact `oe@ifi.uio.no` if interested.

**'Deep' Parsing Portal at UiO (`http://www.delph-in.net`)**

- Existing international network on multi-lingual 'deep' parsing (HPSG);

- Fall 2010, seek NoTur support and cpu allocations to establish portal.

# Credits

NoTur and NorStore (via UniNett Sigma);

The UiO Scientific Computation Group;

The Norwegian Taxpayers.