# LAP: The CLARINO Language Analysis Portal

Emanuele Lapponi, Stephan Oepen, Arne Skjærholt, and Erik Velldal

University of Oslo,
Department of Informatics

October 16, 2015

- Introduction & high level goals
- Design & Implementation: Galaxy
- LAF as a model for tool interchange in LAP
- Tool integration & versioning with the LAP tree
- Reaching out to other research communities
- Current state of development and future plans

LAP:

- A portal providing easy access to NLP tools
- Unlike other processing environments:
    - LAP is not web-service based
    - Tools run on a high-capacity compute cluster
    - Annotation representation and interchange format

LAP:

- A portal providing easy access to NLP tools
- Unlike other processing environments:
    - LAP is not web-service based
    - Tools run on a high-capacity compute cluster
    - Annotation representation and interchange format
- Part of an ongoing PhD project that investigates how NLP can benefit SSH research

A web-app platform for accessing and configuring tools, organizing datasets and annotations, and share results.

**Galaxy**

# Galaxy

A web-app platform for accessing and configuring tools, organizing datasets and annotations, and share results.

**Tools**

search tools

**Import**

**Segmentation**

Tokenizer: Rule-Based Sentence Segmenter determine 'sentences' (top-level utterances)

PEPP: Regular Expression–Based Tokenizer determine word-like units

NLTK Punkt Sentence Segmenter determine 'sentences' (top-level utterances)

NLTK Tokenizer determine word-like units

**Tagging**

HunPOS: Part of Speech Tagger determine word classes

**Parsing**

MaltParser: Linear-Time Dependency Parsing determine bi-lexical syntactic dependencies

Bohnet & Nivre (2012) Joint Part of Speech Tagger and Parser determine word classes and bi-lexical syntactic dependencies

**Giellatekno**

**Oslo-Bergen Tagger**

**Export**

Export as Tab-Separated Values CoNLL-style annotations

Export to VISL CG-3 Format Constraint Grammar annotations
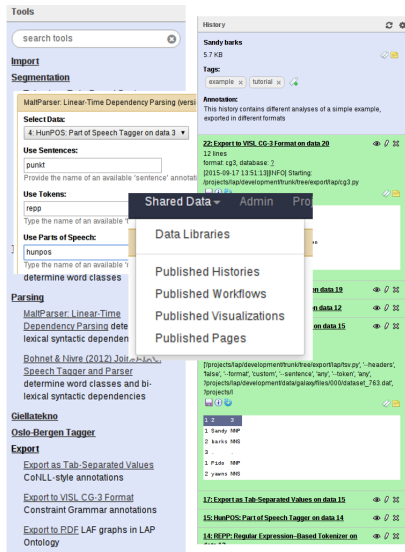
Export to PDF LAF graphs in LAP Ontology

4

A web-app platform for accessing and configuring tools, organizing datasets and annotations, and share results.

A web-app platform for accessing and configuring tools, organizing datasets and annotations, and share results.

A web-app platform for accessing and configuring tools, organizing datasets and annotations, and share results.

A GUI to build (potentially complex) workflows

Requirements:

- Stand-off

- Scalable in terms of coverage of linguistic information and data volume

- Granular: on-demand access of relevant annotations

# Tool interchange format

Data model: **L**inguistic **A**nnotation **F**ramework [2]

- Stand-off: text regions linked to a graph that describes them

- Content agnostic

- Flexible structure

Implementation: MongoDB

- Records describing the structural LAF elements:
  Regions, Nodes and Edges
- Flexible data-access

# Tool interchange format

*Not* in competition with richer end-user formats!

A LAP tool is made of:

- binaries for the actual annotator (i.e. the B&N parser)
- a *wrapper* that communicates with MongoDB

# Tool integration and versioning: the LAP tree

A LAP tool is made of:

- binaries for the actual annotator (i.e. the B&N parser)
- a *wrapper* that communicates with MongoDB

Which means:

- Different programming languages
- Different virtual machines and interpreters
- Different versions

The LAP Tree

- A version controlled repository of the core LAP parts
  (i.e. those that transcend Galaxy and the OS)
- Easily relocatable
- Enables reproducibility of experiments performed with
  historical versions of tools

If we build it, will they come? [3]

# Reaching out to other research communities

Our position:

- Start out with actual research questions
- Work jointly with SSH researchers
- Investigate how and to what degree this work can be generalized into workflows

## Reaching out to other research communities

Our position:

- Start out with actual research questions
- Work jointly with SSH researchers
- Investigate how and to what degree this work can be generalized into workflows

So far:

- Joint work with Political Scientists
- Data-driven analysis of plenary debate speeches in the European Parliament [1]

Talk of Europe

▸ A project that aims at curating EP datasets to linked data

Our contribution:

▸ State-of-the-art syntacto-semantic annotations in rdf triples
  (and possible ontological means to connect them to the ToE graph)

# Current state of development and future plans

LAP, currently:

- A feide- and eduGAIN-accessible development instance
- HPC-ready tools for English, Sami and Norwegian
- Tabulated, cg3 and rdf export
- Basic user documentation

## Current state of development and future plans

Short- to mid-term goals:

- ► Broaden the range of processing types (e.g. deep semantic parsing)
- ► Preprocessing tools for e.g. xml-datasets
- ► Export interfaces with other CLARINO platforms such as Corpuscle and Glossa

# Thank you!

https://lap.hpc.uio.no/

📄 B. Høyland, J.-F. Godbout, E. Lapponi, and E. Velldal.
Predicting party affiliations from European Parliament debates.
In *Proceedings of the 52nd Meeting of the Association for Computational Linguistics: Workshop on Language Technologies and Computational Social Science*, page 56 – 60, Baltimore, MD, USA, 2014.

📄 N. Ide and K. Suderman.
The Linguistic Annotation Framework: A standard for annotation interchange and merging.
*Language Resources and Evaluation*, (forthcoming), 2013.

📄 J. v. Zundert.
If you build it, will we come? Large scale digital infrastructures as a dead end for digital humanities.
*Historical Social Research*, 37(3), 2012.