

# The Language Analysis Portal Large-Scale, State-of-the-Art Language Analysis for Everyone

Emanuele Lapponi, Milen Kouylekov, Stephan Oepen

Department of Informatics, University of Oslo

lap-developers@ifi.uio.no

## A Digital Sea Change: 'Text as Data'

We constructed a corpus of digitized texts containing about 4% of all books ever printed. Analysis of this corpus enables us to investigate cultural trends quantitatively. [...] This approach can provide insights about fields as diverse as lexicography, the evolution of grammar, collective memory, the adoption of technology, the pursuit of fame, censorship, and historical epidemiology. [...] If you tried to read only the entries from the year 2000 alone, [...] it would take eighty years. The sequence of letters is a thousand times longer than the human genome.

[Jean-Baptiste Michel et al., 2010; Science]



## A Digital Sea Change: 'Text as Data'

We constructed a corpus of digitized texts containing about 4% of all books ever printed. Analysis of this corpus enables us to investigate cultural trends quantitatively. [...] This approach can provide insights about fields as diverse as lexicography, the evolution of grammar, collective memory, the adoption of technology, the pursuit of fame, censorship, and historical epidemiology. [...] If you tried to read only the entries from the year 2000 alone, [...] it would take eighty years. The sequence of letters is a thousand times longer than the human genome.

[Jean-Baptiste Michel et al., 2010; Science]

Google Books or Europeana, for Example: Digitizing the World Libraries





- USIT — 5-APR-17 (oe@ifi.uio.no) -

#### Raw Text

were series by the series of the set of the



LAP: The UiO Language Analysis Portal (3)

#### Raw Text

never for filled have to estimate the set of the set o



Science



USIT — 5-APR-17 (oe@ifi.uio.no) -

#### Raw Text

as an end black the set of the s



Utility



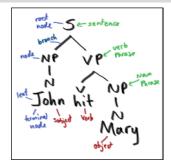
USIT — 5-APR-17 (oe@ifi.uio.no) -

#### Raw Text

there, be filled a lower is a lower of fille  $p_{ij}$  as a first result from of the set  $p_{ij}$  and  $p_{ij}$  is a set of the set o



#### Refinement



#### Science



USIT — 5-APR-17 (oe@ifi.uio.no) -



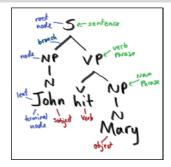


#### Raw Text

anow we find the set of the set



#### Refinement



#### Science



USIT — 5-APR-17 (oe@ifi.uio.no) -





#### Raw Text

mere for Boch hand is sell un plate on a self-the self-and then the transmission of the self-and the self-and the self-and then it of model sever. In Counter to the of the self-and self-and then it of model sever. In Counter to the self-and self-and the several is in the self-and the self-and self-and self-and self-and the several is public to the several several several several self-and the several several several several several several several to the several sever



Refinement







Science



USIT — 5-APR-17 (oe@ifi.uio.no) -





#### Raw Text

anow for Book have a still we fight  $\rho_{0}$  as a still we strend that an effect stages of the start of the start of the start of the start is equivalent to the start of the start is proved by the start of the start start of the start







Refinement







Science



- USIT — 5-APR-17 (oe@ifi.uio.no) -

#### **Textual Building Blocks (All Too Often Taken for Granted)**

- Sentence Splitting headlines, bullet lists, interactions with mark-up, etc.
- Tokenization punctuation marks, abbreviations, dashes and slashes, etc.



#### **Textual Building Blocks (All Too Often Taken for Granted)**

- Sentence Splitting headlines, bullet lists, interactions with mark-up, etc.
- Tokenization punctuation marks, abbreviations, dashes and slashes, etc.

#### The Olympic Committee already regrets choosing Russia for this Winter Olympics .



#### **Textual Building Blocks (All Too Often Taken for Granted)**

- Sentence Splitting headlines, bullet lists, interactions with mark-up, etc.
- Tokenization punctuation marks, abbreviations, dashes and slashes, etc.

The Olympic Committee already regrets choosing Russia for this Winter Olympics . the Olympic Committee already regret choose Russia for this Winter Olympics .



#### **Textual Building Blocks (All Too Often Taken for Granted)**

- Sentence Splitting headlines, bullet lists, interactions with mark-up, etc.
- Tokenization punctuation marks, abbreviations, dashes and slashes, etc.

The Olympic Committee already regrets choosing Russia for this Winter Olympics the Olympic Committee already regret choose Russia for this Winter Olympics . DT NNP NNP RB VBZ VBG NNP IN DT NN NN



USIT — 5-APR-17 (oe@ifi.uio.no)

#### **Textual Building Blocks (All Too Often Taken for Granted)**

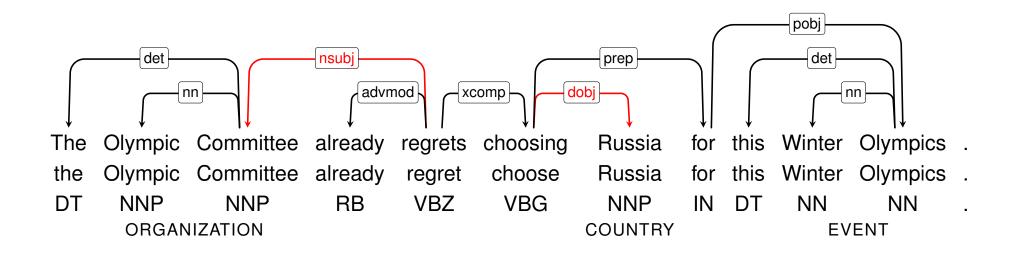
- Sentence Splitting headlines, bullet lists, interactions with mark-up, etc.
- Tokenization punctuation marks, abbreviations, dashes and slashes, etc.

#### The Olympic Committee already regrets choosing Russia for this Winter Olympics the Olympic Committee already regret choose Russia for this Winter Olympics . DT NNP NNP RB VBZ VBG NNP IN DT NN NN ORGANIZATION COUNTRY **EVENT**



#### **Textual Building Blocks (All Too Often Taken for Granted)**

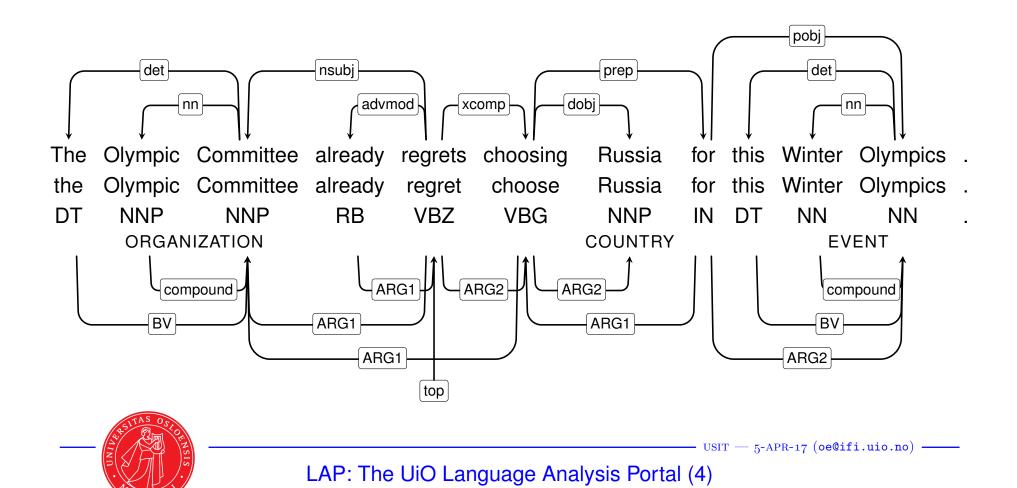
- Sentence Splitting headlines, bullet lists, interactions with mark-up, etc.
- Tokenization punctuation marks, abbreviations, dashes and slashes, etc.





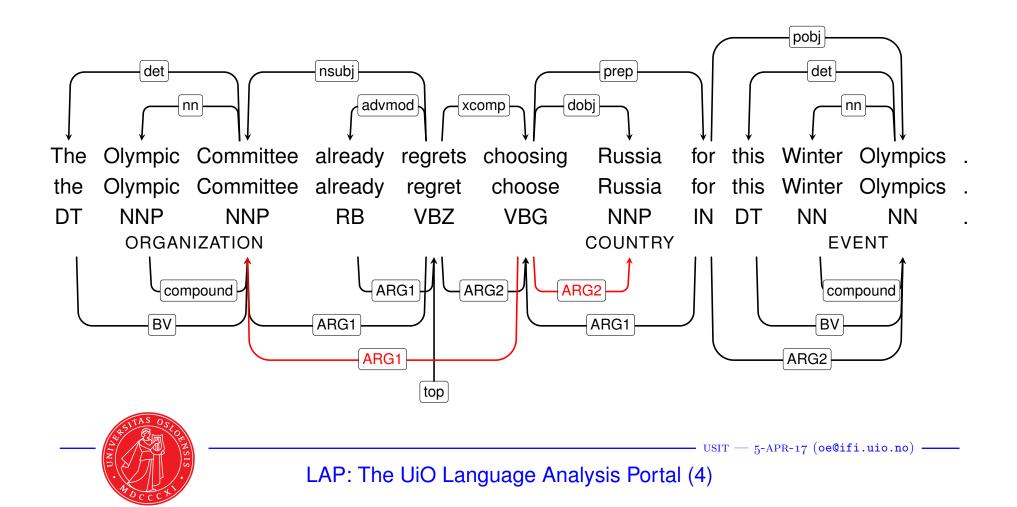
#### **Textual Building Blocks (All Too Often Taken for Granted)**

- Sentence Splitting headlines, bullet lists, interactions with mark-up, etc.
- Tokenization punctuation marks, abbreviations, dashes and slashes, etc.



#### **Textual Building Blocks (All Too Often Taken for Granted)**

- Sentence Splitting headlines, bullet lists, interactions with mark-up, etc.
- Tokenization punctuation marks, abbreviations, dashes and slashes, etc.



## In Support of Digital Humanities and Social Sciences

#### **Predicting Party Affiliations from European Parliament Debates**

Bjørn Høyland Department of Political Science University of Oslo bjorn.hoyland@stv.uio.no Jean-François Godbout Department of Political Science University of Montreal godboutj@umontreal.ca

Emanuele Lapponi Department of Informatics University of Oslo emanuel@ifi.uio.no Erik Velldal Department of Informatics University of Oslo erikve@ifi.uio.no

#### Abstract

This paper documents an ongoing effort to assess whether party group affiliation of participants in European Parliament debates can be automatically predicted on the basis of the content of their speeches, premise for the success of such an approach is that differences in ideology and belief systems are reflected in differences in choice of words in plenary debates. Another premise is that a shared belief system translates to the same choice of party group. As discussed below, systematic differences in prediction performance in the data can be used

(2014 Workshop on Language Technologies and Computational Social Science)



## **A Trivial Example: Benefits of Linguistic Analysis**





USIT — 5-APR-17 (oe@ifi.uio.no) -

## **A Trivial Example: Benefits of Linguistic Analysis**



- 250,000 speech transcripts from the Norwegian Parliament: 1998–2016;
- rich and structured metadata available; e.g. all speeches from one party.



## Seeking to 'Amplify' the Linguistic Signal







LAP: The UiO Language Analysis Portal (7)

## Seeking to 'Amplify' the Linguistic Signal



- Morphological analysis through the Oslo–Bergen Tagger;
- ignore everything but nouns; use lemmata ('base forms').



USIT — 5-APR-17 (oe@ifi.uio.no) -



#### Motivation

- Lower barrier to entry for non-technical users: state-of-the-art tools;
- unified, Web-based point of entry; balance ease of use and flexibility.

#### **Core Elements**

- **Tools** Pre-installed and -configured, multi-lingual, common tasks;
- Interfaces Portal access; data injection and export; parameterization;
- **Plug & Play** 'mix and match' tools: unified LAP interchage format.

#### Scalability

• Submit to Norwegian eInfrastructure: 11,000<sup>+</sup> cores at your fingertips.



#### Motivation

Incremental, demand-driven approch to tool integration;

 $\rightarrow$  so far, languages and tools represented in Norway.

#### **Core Elements**

- **Tools** Pre-installed and -configured, multi-lingual, common tasks;
- Interfaces Portal access; data injection and export; parameterization;
- **Plug & Play** 'mix and match' tools: unified LAP interchage format.

#### Scalability

• Submit to Norwegian eInfrastructure: 11,000<sup>+</sup> cores at your fingertips.



#### Motivation

Incremental, demand-driven approch to tool integration;  $\rightarrow$  so far, languages and tools represented in Norway.

#### Core Flements

In parallel, reach out to candidate consumers (SSH researchers); jointly implement 'use case' studies, e.g. Høyland, et al. (2014).

• Plug & Play 'mix and match' tools: unified LAP interchage format.

#### Scalability

• Submit to Norwegian eInfrastructure: 11,000<sup>+</sup> cores at your fingertips.



#### Motivation

Incremental, demand-driven approch to tool integration;  $\rightarrow$  so far, languages and tools represented in Norway.

#### Core Flements

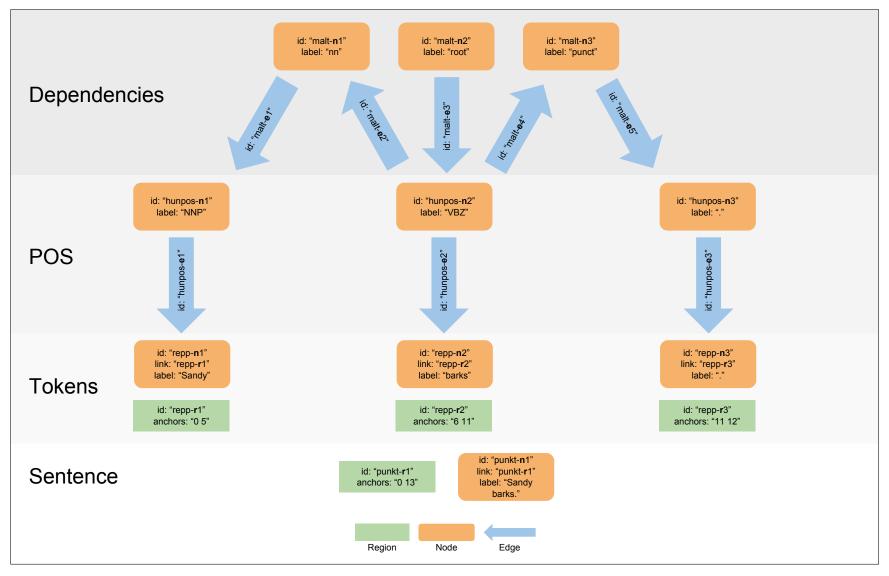
In parallel, reach out to candidate consumers (SSH researchers); jointly implement 'use case' studies, e.g. Høyland, et al. (2014).

• Plug & Play 'mix and match' tools: unified LAP interchage format.

CLARINO: Funded by Norwegian Research Council (2013–2018).

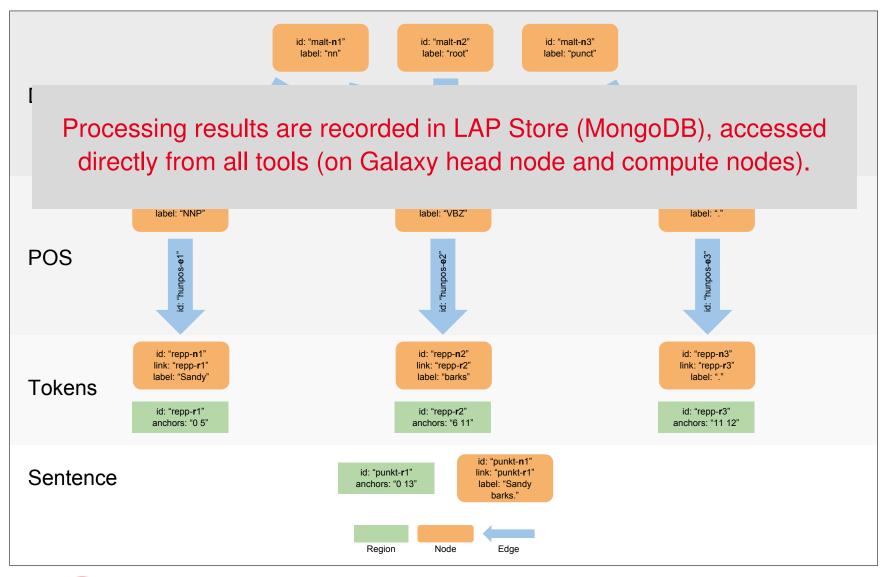
• Submit to Norwegian eInfrastructure: 11,000<sup>+</sup> cores at your fingertips.





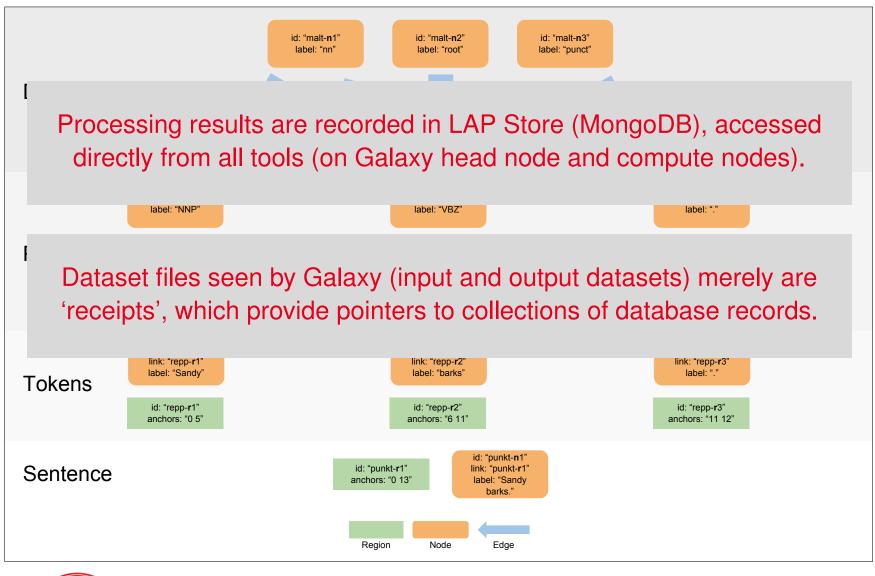


LAP: The UiO Language Analysis Portal (9)



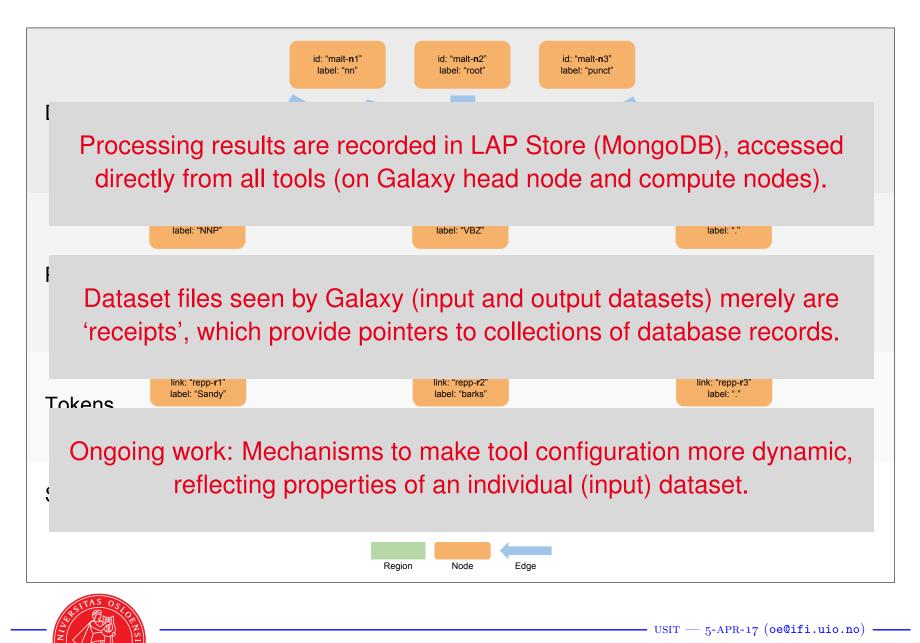


LAP: The UiO Language Analysis Portal (9)





LAP: The UiO Language Analysis Portal (9)





# http://lap.clarino.uio.no/



USIT — 5-APR-17 (oe@ifi.uio.no) -

# http://lap.clarino.uio.no/

If we build it ... Will you come?



USIT — 5-APR-17 (oe@ifi.uio.no)

# http://lap.clarino.uio.no/



If we build it ... Will you come?



USIT — 5-APR-17 (oe@ifi.uio.no)

# http://lap.clarino.uio.no/



If we build it ... Will you come?





USIT — 5-APR-17 (oe@ifi.uio.no)