

The Language Analysis Portal



Design, implementation and use

Emanuele Lapponi (he/him)

Language Technology Group

Institute for Informatics

University of Oslo

@emanlapponi



The digital Social Sciences and Humanities (SSH)

- Interest in computational methods has grown substantially in the last decade
- One strain of digital SSH: research built on "more data than you can read"
- NLP tools are becoming an integral part of research outside of CS and Linguistics

CLARIN

- Develop an infrastructure to facilitate NLP ❤ SSH research
- Part of the CLARINO mandate: build a portal for NLP applications focusing on users with no NLP background



CLARIN - European Research Infrastructure for Language Resources and Technology

CLARIN makes digital language resources available to scholars, researchers, students and citizen-scientists from all disciplines, especially in the humanities and social sciences, through single sign-on access. CLARIN offers long-term solutions and technology services for deploying, connecting, analyzing and sustaining digital language data and tools. CLARIN supports scholars who want to engage in cutting edge data-driven research, contributing to a truly multilingual European Research Area. [Read more...](#)



Common Language Resources and Technology Infrastructure Norway

Centers

CLARINO Bergen Center
Språkbanken
Text Laboratory
LAP
Termportalen
Trolling

CLARIN 2019 annual conference in Leipzig

Five representatives from Norway attended the CLARIN 2019 Annual Conference in Leipzig from September 30 to October 2, 2019. Scott Rettberg was the first invited speaker at the event. Koenraad De Smedt attended the National Coordinators' Forum.



Undertaking such a project raises questions:

- How should a system the facilitates their use be designed?
- Can off-the-shelf tools be integrated?
- What's a good/scalable/coherent way to do it?
- And what's the point?

Part 1

Political Science & Text-as-data 



Classifier evaluation scores as a quantity of interest

- In a parliamentary setting, the contents of a speech should reflect the ideology of the speaker.
- Use party classification scores for parliamentary analysis!

For example: Measure polarization

- Higher accuracy → more polarized*
- Hirst et al** do it for the Canadian Parliament
 - and find that position is a much stronger signal than policy!

...Or ask more involved questions

- Are newcomers to the EU parliament joining groups based on ideology or convenience?***

* Classification Accuracy as a Substantive Quantity of Interest: Measuring Polarization in Westminster Systems

** Party Status as a Confound in the Automatic Classification of Political Speech by Ideology.

*** Lost in Translation? Predicting Party Group Affiliation from European Parliament Debates

Text-as-data pre-processing cookbook

- 🏅 Grimmer and Stewart (2013)* survey the text-as-data field, warning of the dangers of "one-size-fits-all" experiments
- 😐 They also introduce a recipe for pre-processing, with emphasis on stemming and normalization
- 😞 Dismissal of expensive techniques that "do little to enhance performance", e.g. context, disambiguation, lemmatization and so on

* Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts

Text-as-data pre-processing cookbook

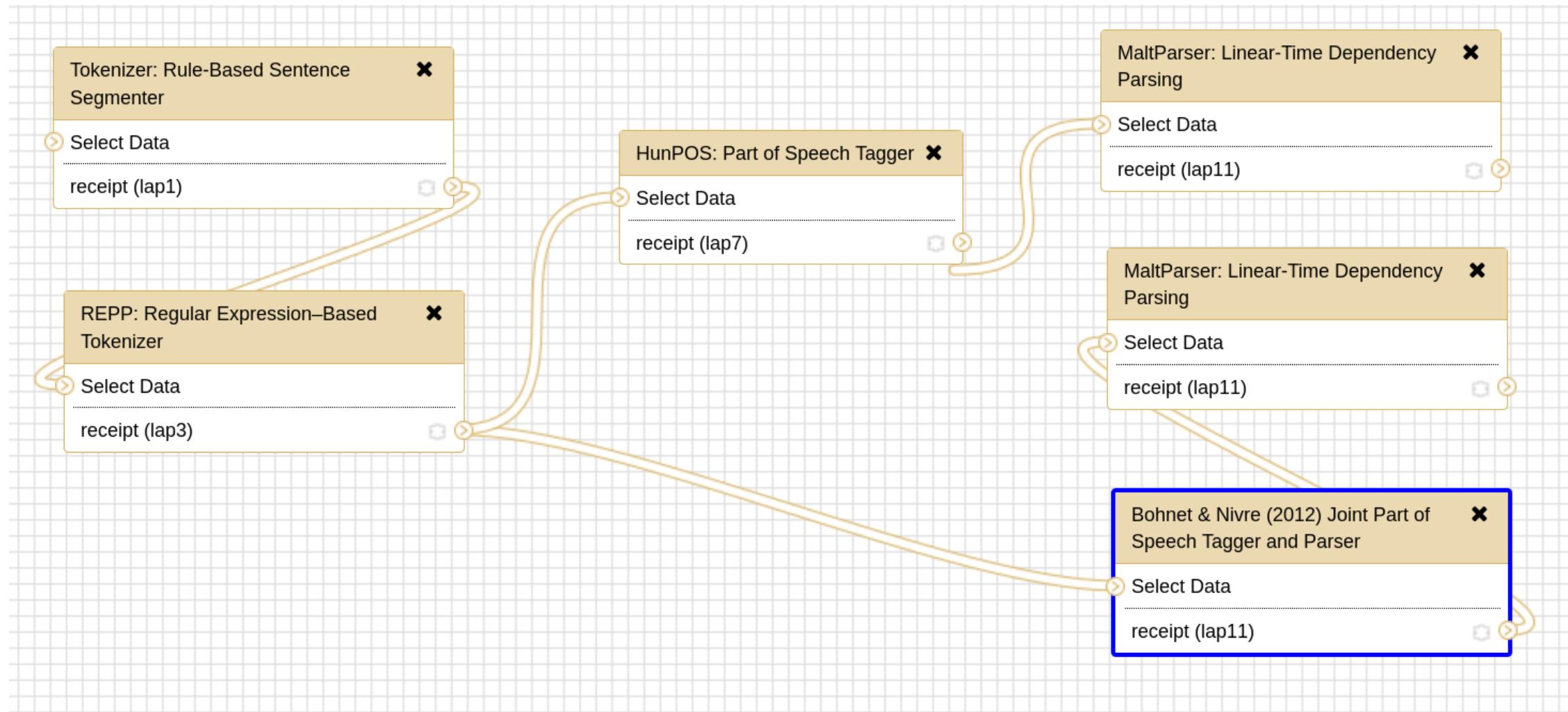
- The recipe becomes gospel
- Even attempts to address the effects of pre-processing on political analysis provide alternate recipes using the same ingredients*
- Experiments miss out on potentially useful (and contentfully significant) tools and techniques!

* Text Preprocessing For Unsupervised Learning: Why It Matters, When It Misleads, And What To Do About It

Part 2

LAP & LXF 

The Language Analysis Portal



Lap eXchange Format

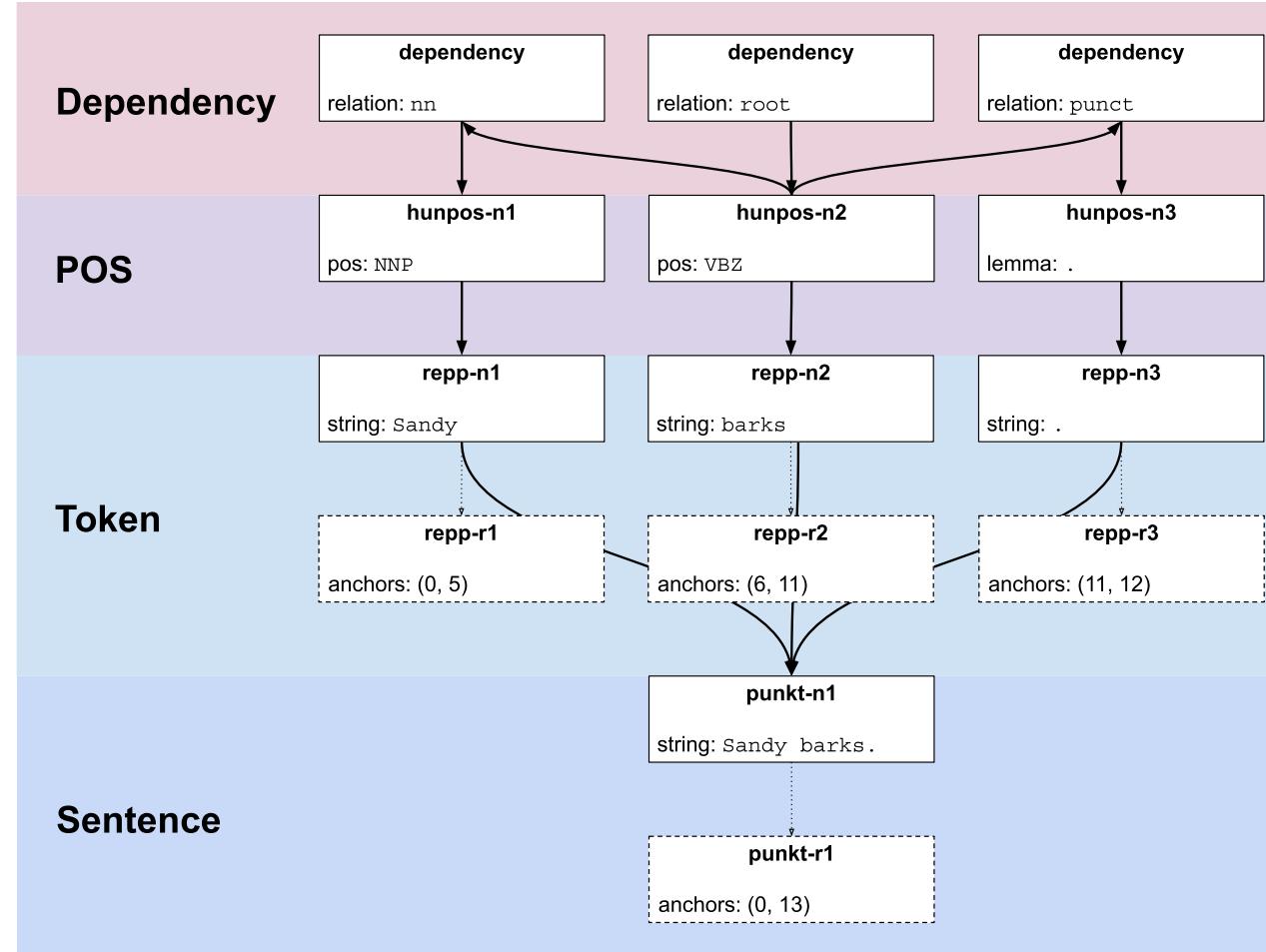
- Scalable to linguistic annotations beneath and above the token level
- Applicable to a wide range of linguistic phenomena and text annotations
- Scalable to large data volumes

The Language Annotation Framework*

An abstract data model: represent arbitrary text **annotations** in a graph. **Nodes** and **edges** organize the structure and the relations of the annotations, and refer to **regions** in the original data.

- Purpose: full interoperability across annotation formats.
- Machine readable instantiations of the data model exist, but are focused on *representation* rather than *interchange*.

- LXF builds on and expands LAF, focusing on intermediate graph representations and building a graph incrementally
- Achieves scalability, both representational and computational, while staying faithful to the original design



Part 4

The Talk of Norway 💡

<https://github.com/ltgoslo/talk-of-norway>

- 18 years of speeches from the Norwegian Parliament, in a novel combination of 80+ meta data variables and language, sentence, token, lemma, PoS and morphology annotations
- An open, plug-and-research resource for text-as-data experimentation

F1 as a quantity of interest, in Norway

- Marked improvements with better normalization, disambiguation and context
- Mixing linguistic and meta-data features helps, too
- Pre-processing does matter for the substance of the analysis

	system	SV	AP	SP	KRF	V	H	FRP	MACRO	ACC
majority class	—	—	—	—	—	—	—	—	0.05	0.24
meta only	0.19	0.38	0.36	0.31	0.24	0.16	0.34	0.28	0.28	0.30
unigram	stem	0.57	0.58	0.58	0.51	0.51	0.58	0.62	0.57	—
	token	0.65	0.65	0.66	0.61	0.62	0.65	0.69	0.65	—
	lemma	0.64	0.65	0.65	0.62	0.62	0.64	0.68	0.64	—
	lemma/pos	0.66	0.66	0.67	0.64	0.64	0.66	0.70	0.66	—
	+meta	0.69	0.69	0.72	0.68	0.69	0.68	0.73	0.70	—
n-gram	stem	0.63	0.66	0.65	0.60	0.61	0.65	0.68	0.65	—
	token	0.67	0.69	0.69	0.66	0.66	0.67	0.71	0.68	—
	lemma	0.68	0.69	0.69	0.66	0.67	0.68	0.72	0.69	—
	lemma/pos	0.69	0.70	0.71	0.67	0.69	0.69	0.73	0.70	—
	+meta	0.71	0.72	0.73	0.71	0.72	0.70	0.75	0.72	—

	SV	A	Sp	KrF	V	H	FrP
SV	71.4	15.4	1.8	2.2	1.5	4.4	3.3
A	3.9	80.3	2.6	2.8	1.4	5.5	3.4
Sp	2.9	14.3	72.1	3.0	1.3	3.3	3.1
KrF	3.1	12.2	3.6	70.1	1.5	6.2	3.4
>	3.4	15.5	3.2	4.2	64.7	6.4	2.6
H	3.1	15.1	1.6	3.3	1.3	69.6	6.0
FrP	3.4	12.4	3.6	1.4	1.6	13.6	64.0

	SV	A	Sp	KrF	V	H	FrP
SV	71.2	10.1	2.0	2.7	2.0	5.5	6.4
A	6.1	66.8	3.4	4.3	2.0	7.9	9.5
Sp	4.1	12.4	66.1	4.1	1.4	4.5	7.5
KrF	4.0	9.3	2.2	65.7	1.2	7.9	9.7
>	4.7	6.6	3.0	2.4	69.5	6.0	7.8
H	3.8	8.9	1.8	2.8	1.2	69.0	12.5
FrP	2.1	4.5	1.0	2.0	0.6	6.6	83.1

Conclusions ✨

- Heterogenous workflows of off-the-shelf NLP tools can be made easily accessible without sacrificing scalability and flexibility
- Easy access to NLP pre-processing tools benefits political research and analysis

  LAP & LXF

- A modular system for HPC-ready, user friendly language annotations
- Scalable to large volumes of data and heterogenous linguistic (and non-linguistic!) annotation



Talk of Norway

- An open resource for parliamentary analysis, informed by and developed with Political Scientists
- The first ToN experiments show substantive differences in the analyses built on party classification performance

2013 ●
Towards large-scale language
analysis in the cloud
E Lapponi, E Velldal, NA Vazov, S Oepen

2014 ●
Predicting party affiliations from
European Parliament debates
B Høyland, JF Godbout, E Lapponi, E Velldal

2014 ●
Off-Road LAF: Encoding and
Processing Annotations in NLP
Workflows.
E Lapponi, E Velldal, S Oepen, RL Knudsen

Representation and Interchange of
Linguistic Annotation. An In-Depth,
Side-by-Side Comparison of Three
Designs
RE de Castilho, N Ide, E Lapponi, S
Oepen, K Suderman, E Velldal, Marc
Verhagen

2017 ●
The Talk of Norway: a richly
annotated corpus of the Norwegian
parliament, 1998–2016
E Lapponi, MG Søyland, E Velldal, S Oepen

2018 ●
The Language Analysis Portal:
Design, implementation and use
E Lapponi

Defense! 

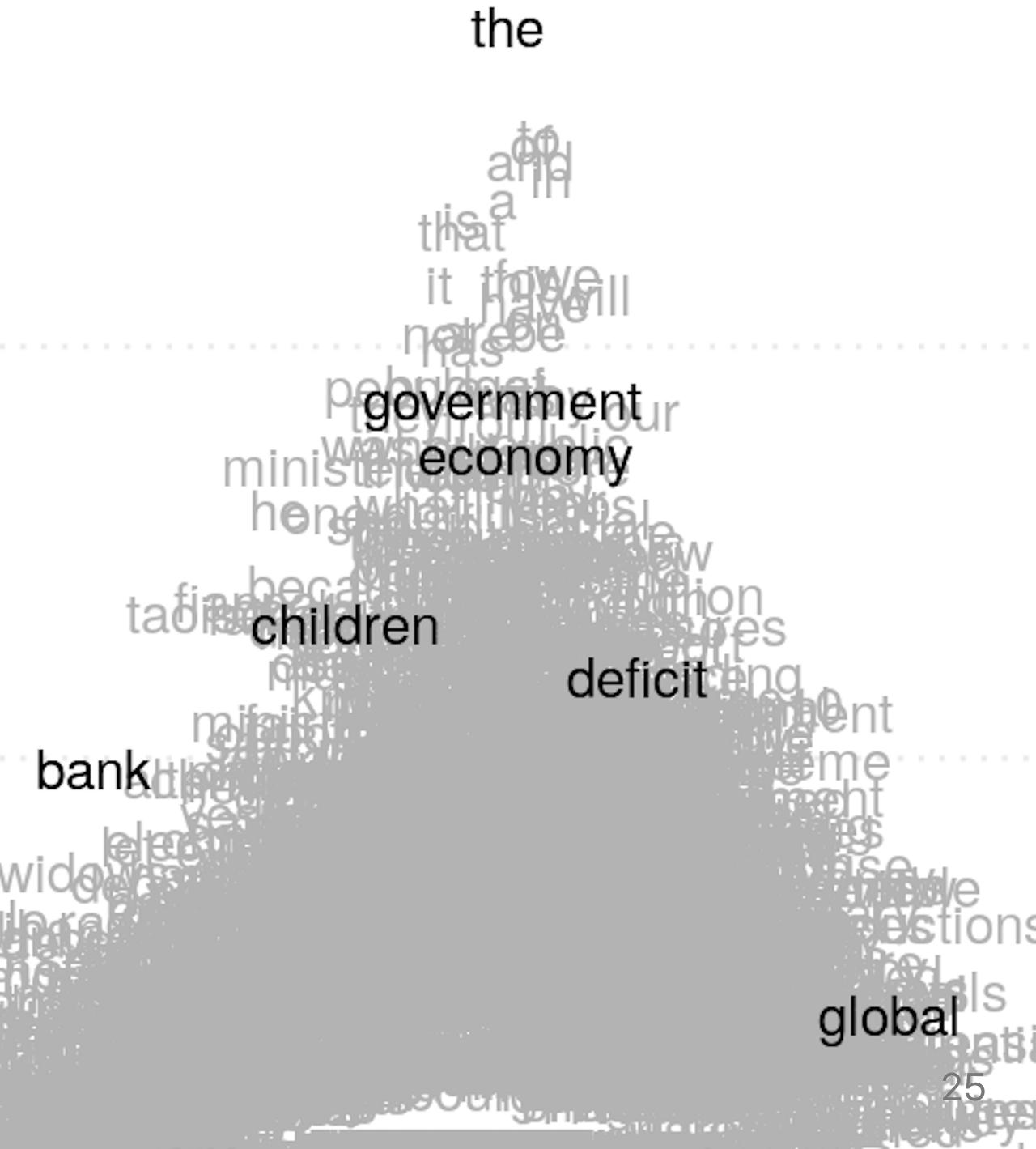
extras

Director's cut 

Using NLP techniques for political analysis

- For example: place words on a left-right political axis*

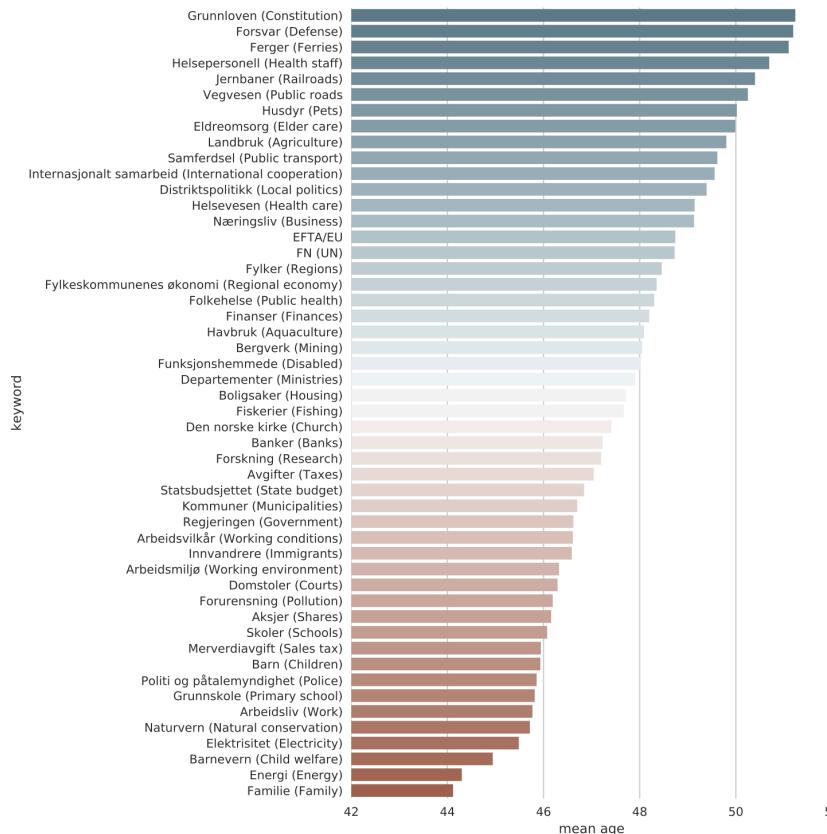
* A Scaling Model for Estimating Time-Series Party Positions from Texts



Getting our feet wet

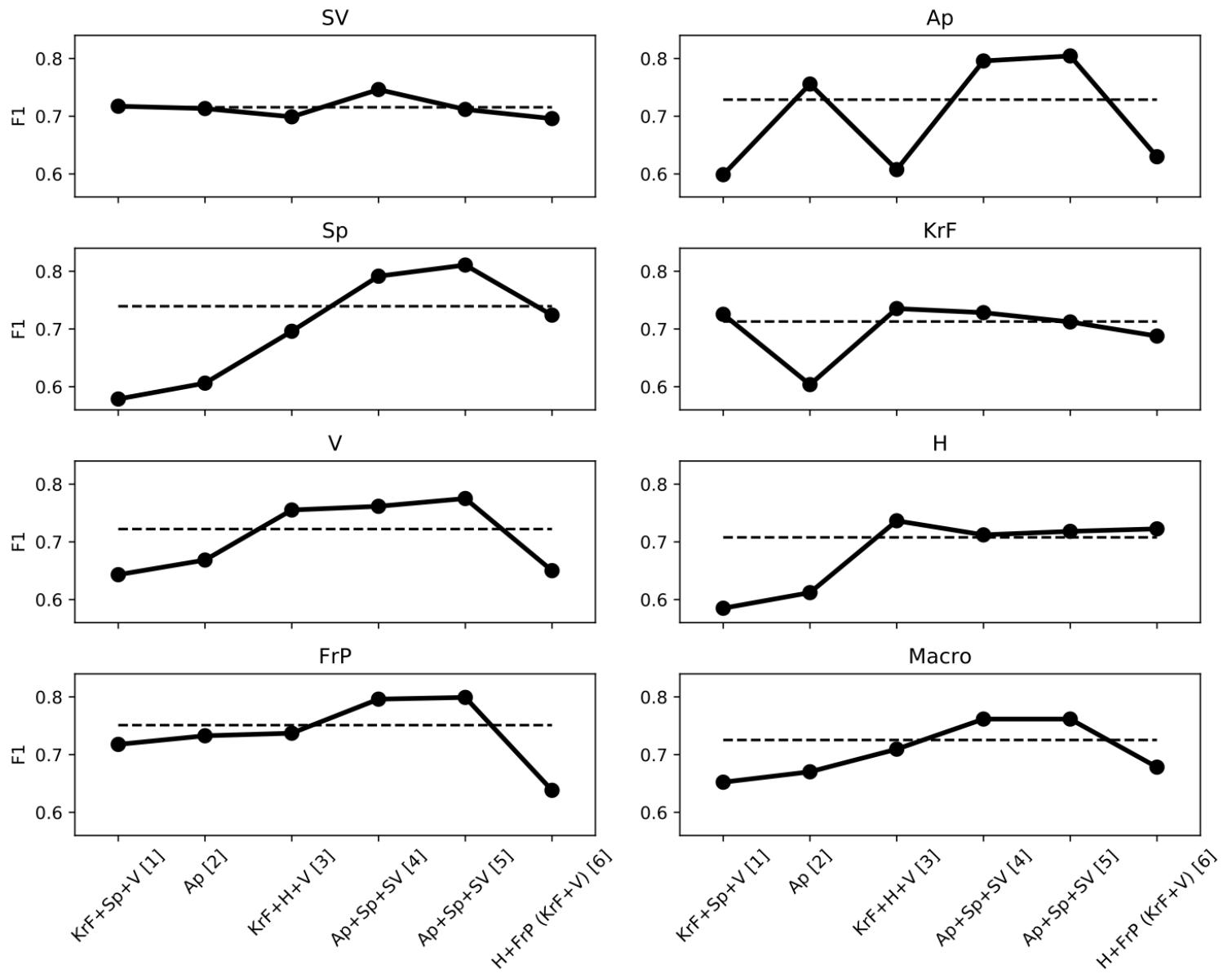
- Preliminary experiments, going off the recipe
- Signs that linguistic features might contribute positively to these methods

	Baseline	stem	dep/stem
Acc	0.394	0.476	0.492
P	0.065	0.439	0.458
R	0.166	0.399	0.393
F_1	0.094	0.418	0.423

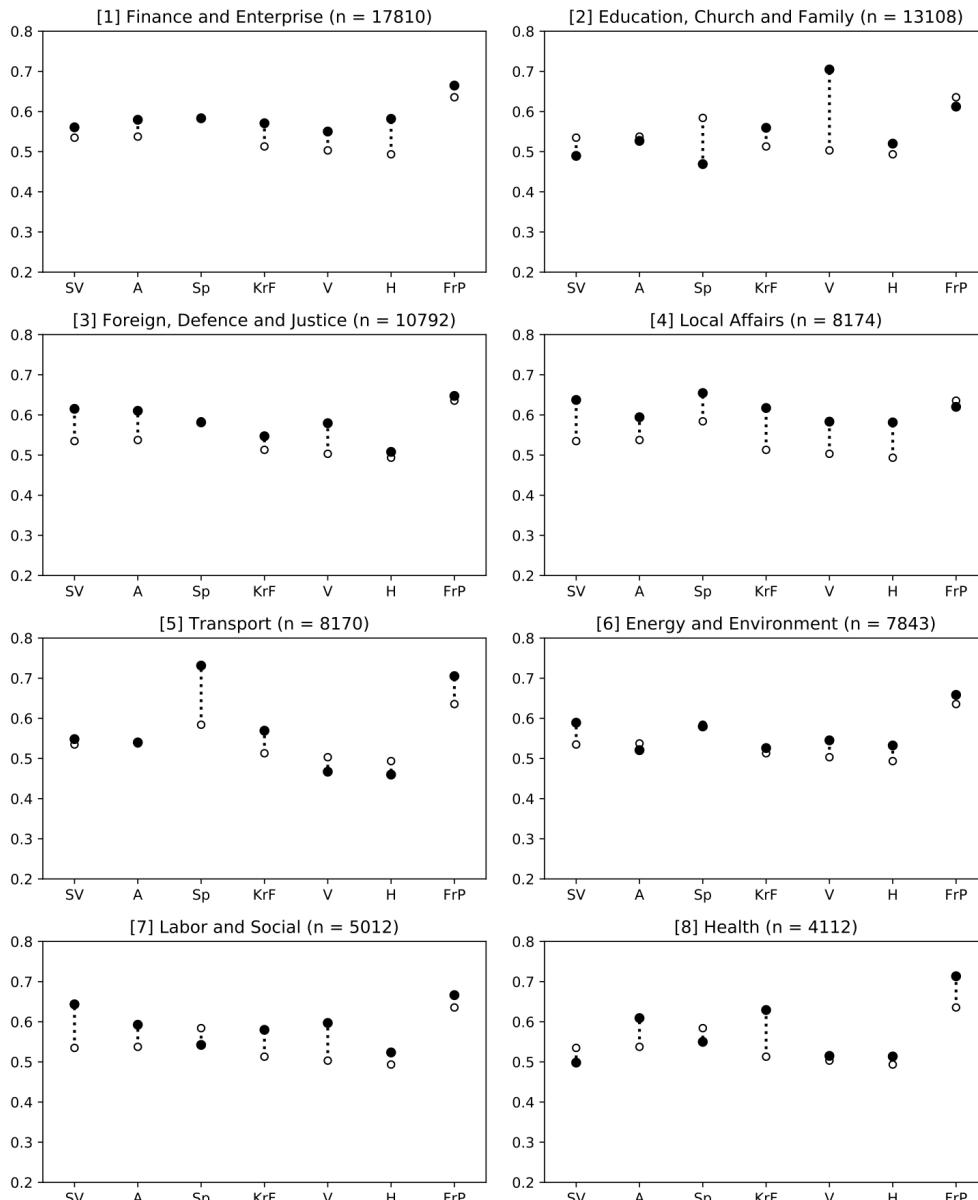


	SV	FrP
1	ramme videregående skole <i>affect high schools</i>	bevilge mer midler <i>allocate more funds</i>
2	tøyle grådige markedskrefter <i>control greedy market forces</i>	bygge ny vei <i>build new roads</i>
3	være liten plass <i>be little room</i>	putte mer penger <i>put more money</i>
4	vurdere ulike virkemidler <i>consider different means</i>	bygge firefelts vei <i>build four-lane highways</i>
5	medføre doblet slaktevekt <i>cause doubling of the carcass weight</i>	effektivisere norsk økonomi <i>streamline the norwegian economy</i>
6	nå ambisiøse mål <i>reach ambitious goals</i>	innføre lokal kontantstøtte <i>introduce local child care allowance</i>
7	være betydelige kostnader <i>be substantial costs</i>	bruke nytt verktøy <i>use new tools</i>
8	redusere forurensende utslipp <i>reduce pollutant emissions</i>	gi betydelige samfunnsgevinster <i>give significant social gains</i>
9	gjøre større utslippskutt <i>do bigger emission cuts</i>	utføre kriminelle handlinger <i>commit criminal acts</i>
10	få mindre høyrepolitikk <i>get less right-wing politics</i>	sikre økt pasientbehandling <i>ensure increased patient treatment</i>

On their own, the meta data variables are useful for statistics (what do younger vs. older MEP talk about?), and can be used in combination with the linguistic annotations to quickly go from raw data to policies



Text as data



What is an NLP tool?

- Tokenizers, PoS taggers, lemmatizers, syntactic parsers +++
- Come in many flavors, adhere to different traditions, have different pros and cons (depending on the task)
- It's a good idea to test different tools before settling for one for a given task!

Tool interoperability and interchange formats

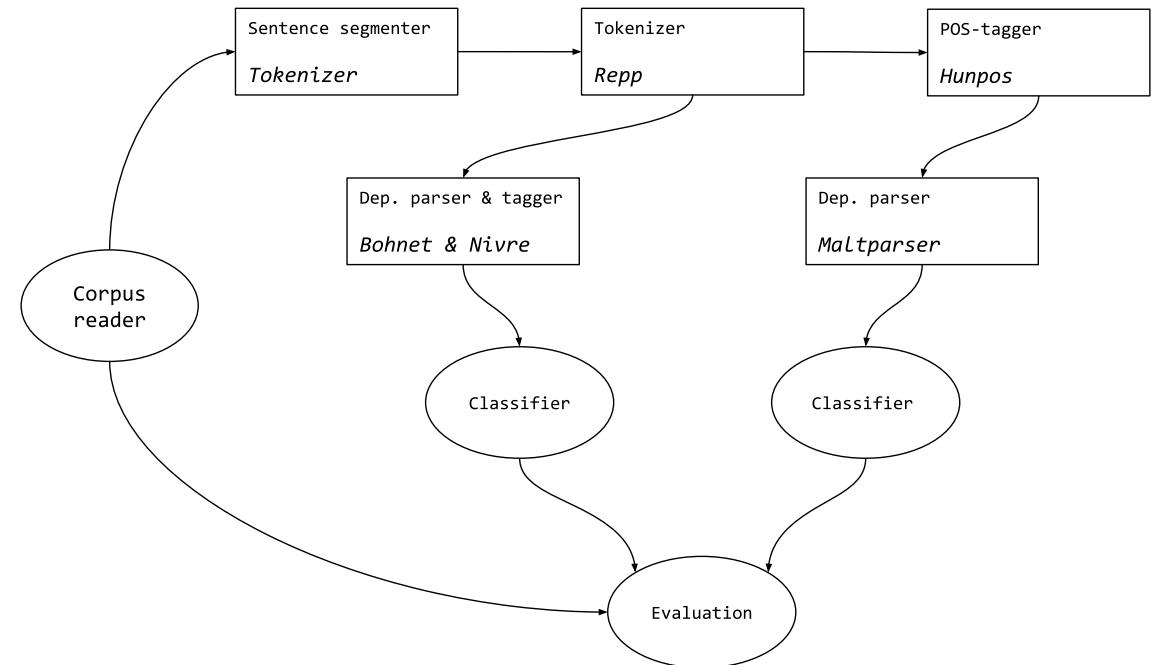
```
Sandy NNP  
barks NNS
```

```
1 Sandy _ _ NNP _ _ _ _ _  
2 barks _ _ NNS _ _ _ _ _  
3 . _ _ . _ _ _ _ _
```

```
1 Sandy _ _ NNP _ 2 nn  
2 barks _ _ NNS _ 0 null _ _  
3 . _ _ . _ 2 punct _ _
```

What is a workflow?

- Combine tools to get annotations for some downstream task (linguistic or otherwise)
- Tools better be compatible



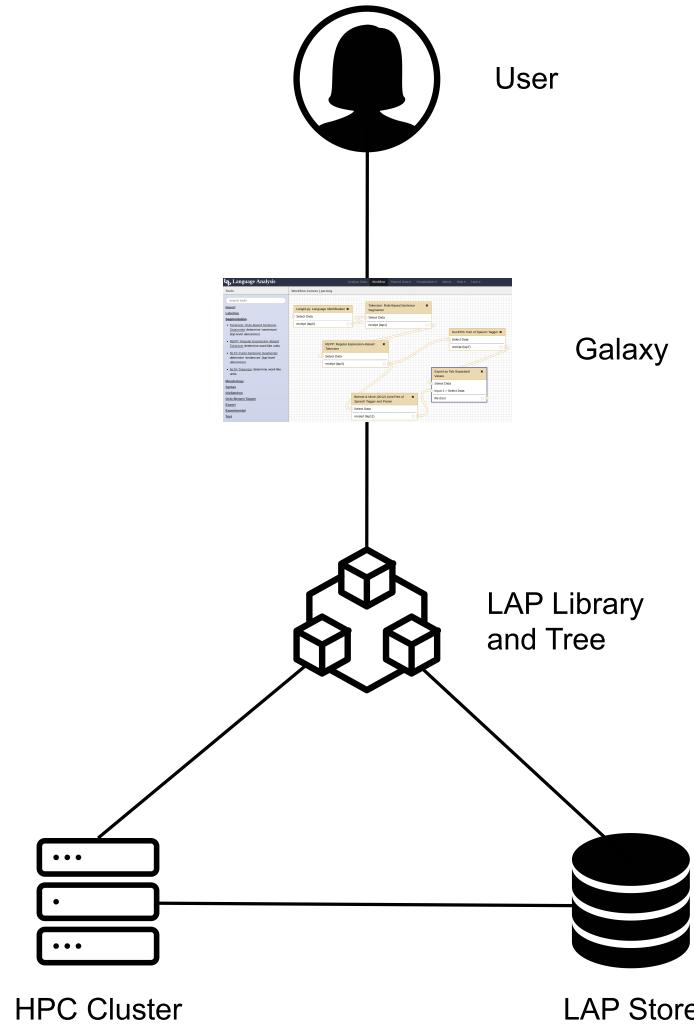
What is a portal?

- No local installation is required
- Data upload and results download
- Graphical user interface to configuring and running tools
- Not developer tools like **NLTK** and **SpaCy**!

What is big data?

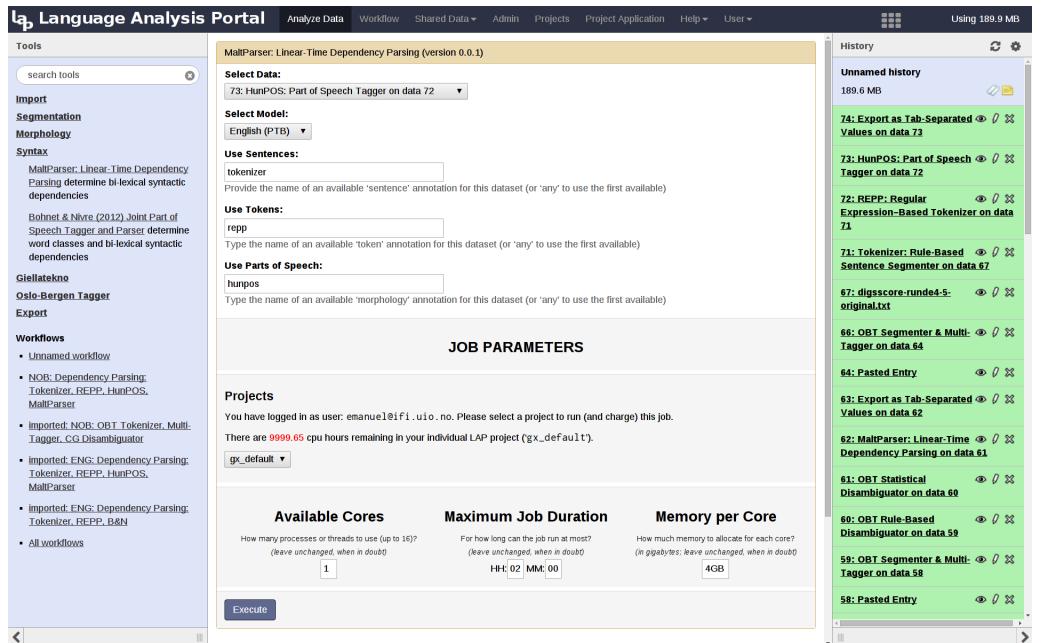
- Big enough to make processing it on a laptop impractical

LAP Architecture

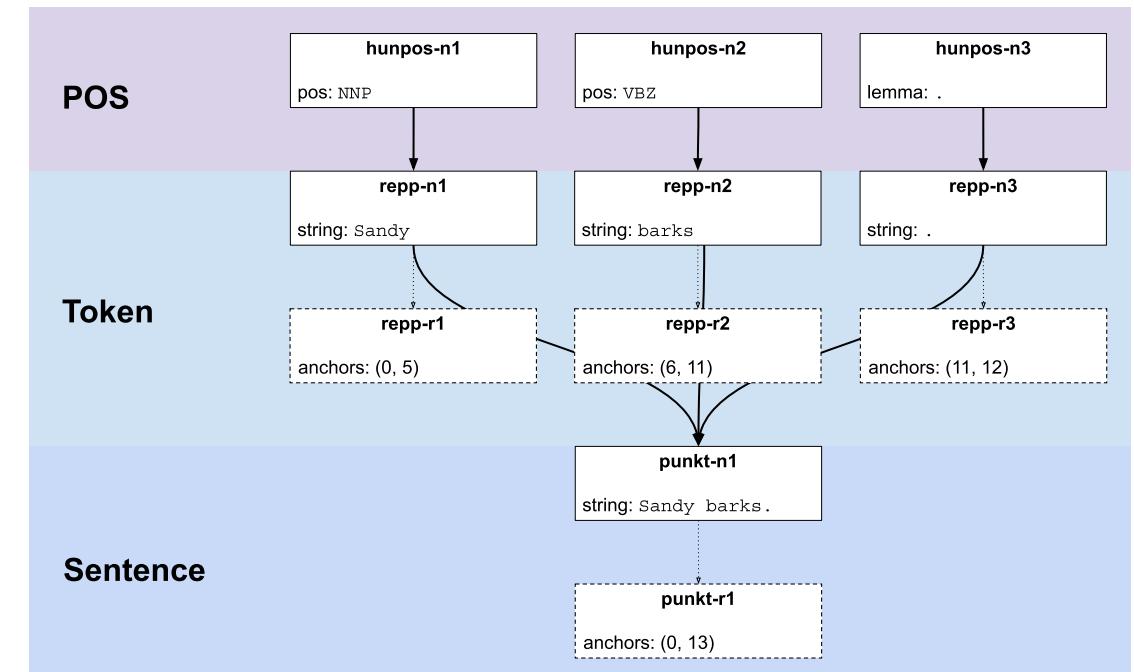


Graphical user interface

- LAP is GUI agnostic
- Blazing the trail for using Galaxy for NLP



```
{
  "annotations": {
    "morphology": [
      "hunpos"
    ],
    "sentence": [
      "nltk_punkt"
    ],
    "token": [
      "repp"
    ]
  },
  "annotators": {
    "hunpos": "26a43d0c-f3db-11e8-b17a-00259075dac6",
    "nltk_punkt": "1d112232-f3db-11e8-a3ee-00259075db92",
    "repp": "20aaddc0-f3db-11e8-b01c-b083fed3d77f"
  },
  "media": {
    "text": "0d3b5012-f3db-11e8-91a5-b083fed3d77f"
  },
  "receipt_origin": "hunpos"
}
```



Our perspective

- Thoroughly navigating the NLP sea before doing political analysis means no political analysis
- Building NLP-for-SSH infrastructures might just be a good idea
- Doing it without navigating the SSH sea likely is not: if we build it, will they come?