

# The ParlaMint Tutorial

– demo of a CLARIN teaching material –

Darja Fišer\* and Kristina Pahor de Maiti\*\*

\*University of Ljubljana; Jožef Stefan Institute; Institute of Contemporary History, Slovenia

\*\*CY Cergy Paris University, France

# Today's plan

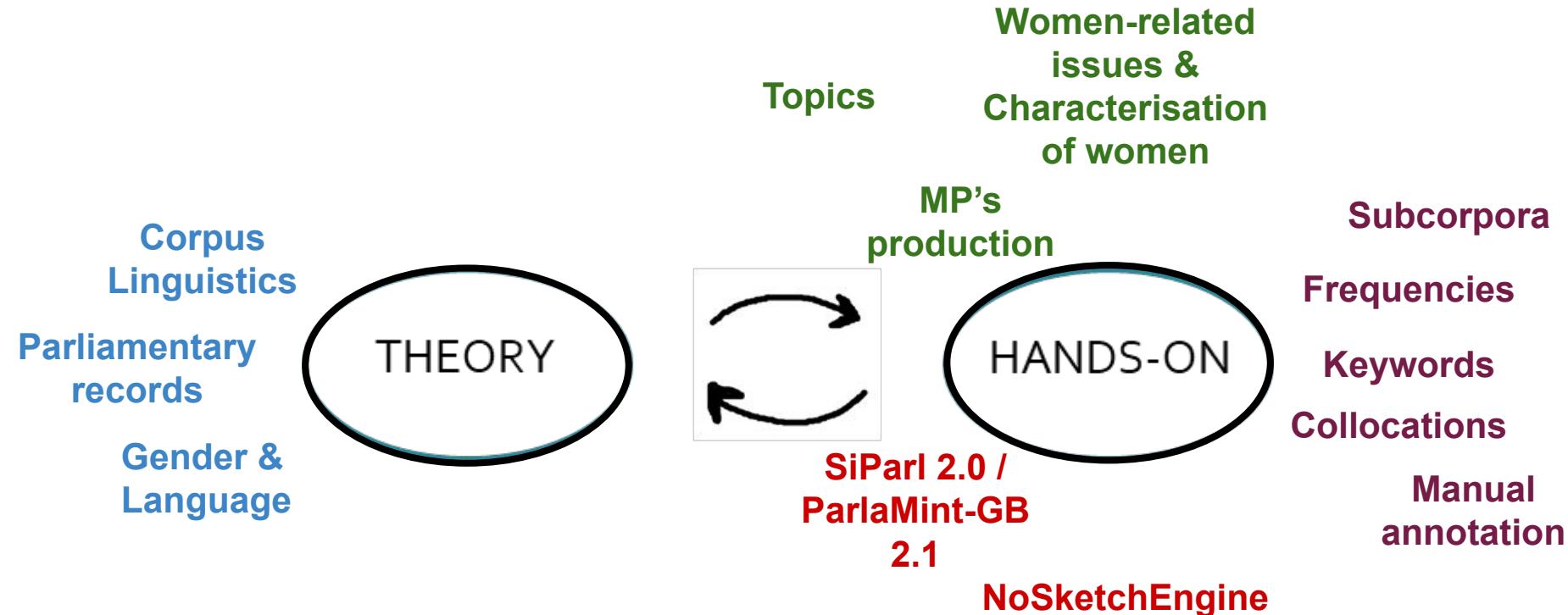
- Demo of the tutorial *Voices of the parliament*
  - tutorial overview
  - showcasing analyses adapted to the ParlaMint-GB data
    - step-by-step instructions + screen casts in the tutorial
    - special parameter tweaks for this analysis provided on the slides
- Some reflections on the tutorial creation process
  - aim & target audience
  - structure
  - sustainability & adaptability

# Tutorial demo



Image by Joakim Honkasalo

# OVERVIEW



Voices of the Parliament – A Corpus Approach to Parliamentary Discourse Research

SiParl – Slovene parliamentary debates 1990–2018

# THEORETICAL PART

CORPORA and CONCORDANCERS	PARLIAMENTARY DISCOURSE, RECORDS and CORPORA	GENDER, LANGUAGE and POLITICS
<ul style="list-style-type: none"><li>• What are language corpora?</li><li>• How is data encoded?</li><li>• What is a token?</li><li>• What are different levels of annotation?</li><li>• Why use linguistic annotations?</li><li>• What is lemmatization?</li><li>• What is PoS tagging?</li><li>• What is metadata?</li><li>• What is a concordancer?</li><li>• Which are different concordancers available?</li><li>• What are the most popular corpus analysis techniques?</li></ul>	<ul style="list-style-type: none"><li>• What are parliamentary corpora?</li><li>• What are the most often included metadata?</li><li>• What are the specifics of parliamentary discourse?</li><li>• How faithful are the records to the actual spoken word?</li><li>• Why is it important to know your data?</li></ul>	<ul style="list-style-type: none"><li>• Is there a women-specific debating style?</li><li>• Which topics do women and men debate in the parliaments?</li><li>• What are the dangers of performing analyses when knowing a parameter, such as gender, in advance?</li></ul>

# Data encoding

Attributes **ana** and **msd** both encode the morphosyntactic characteristics of the word in focus (in blue). The **ana** attribute contain Slovene-specific tags according to MULTTEXT-East Specifications, while the **msd** attribute contains the Universal Dependency tagset which greatly simplifies cross-lingual comparison.

The attribute values Q and PART both stand for *particle*. The difference between the two tags is in the tagset used.

The attribute **lemma** indicates the basic word form of the token, that is of the word in focus (in blue).

The **word form** also known as a token from the running text in the corpus.

```
<s>
<w ana="mte:Q" msd="UpoTag=PART" lemma="zlasti">zlasti</w>
<w ana="mte:Sg" msd="UpoTag=ADP|Case=Gen" lemma="glede">glede</w>
<w ana="mte:Ncmsg" msd="UpoTag=NOUN|Case=Gen|Gender=Masc|Number=Sing"
    lemma="nadzor">nadzora</w>
<w ana="mte:Va-r3s-n"
    msd="UpoTag=AUX|Mood=Ind|Number=Sing|Person=3|Polarity=Pos|Tense=Pres|VerbForm=Fin"
    lemma="biti">je</w>
<w ana="mte:Pd-fsn" msd="UpoTag=DET|Case=Nom|Gender=Fem|Number=Sing|PronType=Dem"
    lemma="ta">ta</w>
<w ana="mte:Ncfsn" msd="UpoTag=NOUN|Case=Nom|Gender=Fem|Number=Sing"
    lemma="stvar">stvar</w>
<w ana="mte:Rgp" msd="UpoTag=ADV|Degree=Pos" lemma="zelo">zelo</w>
<w ana="mte:Agpsn" msd="UpoTag=ADJ|Case=Nom|Degree=Pos|Gender=Fem|Number=Sing"
    lemma="kočljiv">kočljiva</w>
<pc ana="mte:Z" msd="UpoTag=PUNCT">.</pc>
</s>
```

The corpus also encodes syntactic parses and named entities but since they are not used in this tutorial they were omitted from this illustrative example.

One pair of the opening and closing structural tags which, in our case, indicate punctuation (**pc**), word (**w**) or sentence (**s**).

# Faithfulness of the records

## *Official records*

*"Madam President, I think we all agreed that equal pay is, first of all, a matter of elementary fairness, and the fact that we still have this average pay gap of 16-17% also shows, among other things, the lack of respect for women, which is so deeply rooted in traditions and stereotypes. /.../. "*

*(Věra Jourová, Member of the Commission)*

## *Verbatim transcription*

***"Ladies and Gentlemen, yes gentlemen also are here still. Am... I think we all agreed that equal pay is, first of all, a matter of elementary fairness, and that the fact that still we have this pay gap am... average 16-17% this also shows, among other things, the lack of respect to women, which is so deeply rooted in traditions and stereotypes. /.../. "***

*(Věra Jourová, Member of the Commission)*

# ParlaMint

- The ParlaMint corpora
  - comparable and uniformly annotated corpora of parliamentary sessions
    - more about the [ParlaMint project](#)
    - half a billion words or 5 million speeches produced by around 11 thousand speakers
  - records of 17 national parliaments
    - Belgium, Bulgaria, Croatia, the Czech Republic, Denmark, France, Great Britain, Hungary, Iceland, Italy, Latvia, Lithuania, the Netherlands, Poland, Romania, Slovenia, and Turkey
  - (at least) 2015–mid 2020
    - all corpora divided into *Reference* (until 2019-10) and *COVID* (from 2019-11) subcorpora
  - richly annotated
    - linguistic annotations (lemma, PoS, morphological and syntactic features, NE)
    - metadata (SPEECH: date, session, title ...; SPEAKER: name, party, role, type, gender, birth)

# ParlaMint-GB\*

ID	Lang	Houses	Terms	From	To	Yrs	Mill. words/Yr	Mill. Words
GB	en	Lower+Upper	4	2015–01	2021–03	6.3y	17.25	109.30
ID	Parties	C/O	Orgs	Prsns	Gender	MP	URL	IMG
GB	31	5	2	1901	1901	1865	1901	1029
ID	Spchs	By KnwnSpks	By NChairs	ByMPs	Heads	Notes	COVID (2019-11)	REF
GB	552,103	549,71	537,928	547,305	31,389	165,65	109940	442,163

[Erjavec et al. 2022](#)

\*Now available at the [new noSketchEngine concordancer](#) as well as in the [SketchEngine concordancer](#)

# PRACTICAL PART

- Pre-task: creation of subcorpora (CQL)
- Task 1: production of speeches by M and F MPs (frequencies)
- Task 2: topics debated by F and M MPs (keywords)
- Task 3: women-related issues by F and M MPs and characterisation of women and men (collocations, UD tags)

# Pre-task: subcorpora creation

- <speech speaker\_gender="F|M" & house="Lower house" & subcorpus="COVID"/>
- <speech speaker\_gender="F|M" & house="Lower house" & from>="2018-07-03" & from<="2019-10-31"/>

	<b>REF (tokens)</b>	<b>COVID</b>
<b>lower_F</b>	4,950,325	4,023,907
<b>lower_M</b>	10,033,622	9,243,264
<b>lower_FM</b>	14,983,947	13,267,171

# TASK 1: Production of speeches

SPKRS	M	%	F	%	TOTAL
COV	453	66.4	229	33.6	682
REF	435	67.7	208	32.3	643
SPCHS					
COV	50,966	72.3	19,517	27.7	70,483
REF	66,788	69.3	29,638	30.7	96,426
TKNS					
COV	9,243,264	69.7	4,023,907	30.3	13,267,171
REF	10,033,622	67.0	4,950,325	33.0	14,983,947

- What is the proportion of MPs per gender?
- How often the MPs take the floor?
- How much MPs' talk?

# TASK 1: frequencies

Word list options 

Corpus: ParlaMint-GB 2.1 (British parliament) 

Subcorpus: COVID\_lower\_F  [info](#) [create new](#) 

Search attribute: speech.speaker\_name 

use n-grams. Value of n: from 2  to 2  

hide/nest sub-n-grams

**Filter options:**

Filter word list by: Regular expression:  

Minimum frequency:  

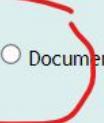
Maximum frequency: 0 (0 = no maximum frequency)

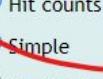
Whitelist: [Izberite datoteko](#) Nobena dat...ka ni izbrana [Clear](#)

Blacklist: [Izberite datoteko](#) Nobena dat...ka ni izbrana [Clear](#) [format](#)

Include non-words

**Output options:**

Frequency figures:  Hit counts  Document counts  ARF 

Output type:  Simple 

Keywords

Reference (sub)corpus: ParlaMint-GB 2.1 (British parliament)

Prefer: rare words  common words

Change output attribute(s)

[http://www.clarin.si/noske/parlamint21.cgi/wordlist\\_form?corpname=parlamint21\\_gb](http://www.clarin.si/noske/parlamint21.cgi/wordlist_form?corpname=parlamint21_gb)

# TASK 1: Production of speeches

AVG SPCH LENGTH	M	F	TOTAL
COV	181.36	206.17	188.23
REF	150.23	167.02	155.39
AVG No. of SPCH/DAY			
COV	278.50	106.65	385.15
REF	335.62	148.93	484.55

- How long are the MPs' speeches?
- How active are the MPs in both periods on a daily basis?

# TASK 2: Topic preference

parlamint21_gb COVID_lower_F	parlamint21_gb COVID_lower_M	parlamint21_gb REF_last15mths_lower_F	parlamint21_gb REF_last15mths_lower_M
<b>hysteroscopy</b>	interconnector	nda	journalistic
moat	highland	pml	seahorse
pesticide	courtroom	lantern	ofs
<b>domicile</b>	orthopaedic	defibrillator	a40
<b>pregnant</b>	alas	abnormality	proposition
<b>maternity</b>	forbearance	<b>pornography</b>	beetle
<b>sibling</b>	le	<b>asc</b>	highland
pedicab	proposition	lara	borrower
priory	distinguished	<b>baby</b>	motorsport
explicable	flatten	bladed	unauthorised
loo	characteristically	<b>post-natal</b>	m26
<b>bame</b>	evidential	<b>incest</b>	manifestly
vascular	chairman	<b>pornographic</b>	gtr
marble	magistrate	moorland	alacrity
<b>anti-racist</b>	cyber	<b>menopause</b>	fertiliser
slp	diocese	<b>underpaid</b>	fella
<b>midwifery</b>	e3	<b>upskirting</b>	sharia
<b>unequal</b>	agility	<b>autistic</b>	mellifluous
<b>pregnancy</b>	coalfield	<b>cervical</b>	nationalisation
<b>transgender</b>	occasionally	crossbow	encampment

- Topic analysis of top 50 keywords
  - [a–z.\*]: to skip NE
- **Reproduction**
- **Equality**

# TASK 2: keywords

[http://www.clarin.si/noske/parlamint21.cgi\\_wordlist\\_form?corpname=parlamint21\\_gb](http://www.clarin.si/noske/parlamint21.cgi_wordlist_form?corpname=parlamint21_gb)

The screenshot shows the ParlaMint-GB 2.1 (British parliament) wordlist interface. Key elements include:

- Corpus:** ParlaMint-GB 2.1 (British parliament)
- Subcorpus:** COVID\_lower\_F (highlighted with a red circle)
- Search attribute:** lemma
- Filter options:**
  - Filter word list by: Regular expression: [a-z].\*
  - Minimum frequency: 5 (highlighted with a red circle)
  - Maximum frequency: 0 (no maximum frequency)
  - Whitelist: Izberite datoteko Nobena dat...ka ni izbrana (Clear)
  - Blacklist: Izberite datoteko Nobena dat...ka ni izbrana (Clear, format)
- Output options:**
  - Frequency figures: Hit counts (selected), Document counts, ARF
  - Output type: Simple, Keywords (selected)
  - Reference (sub)corpus: ParlaMint-GB 2.1 (British parliament)
  - Prefer: rare words (slid to the left), common words (1), Change output attribute(s)

COVID Key

[https://www.clarin.si/noske/run.cgi/wordlist?corpname=parlamint21\\_gb&viewmode=kwic&attrs=word%2Cdep&ctxattrs=word&structures=speech&refs=%3Dspeech.from%2C%3Dspeech.speaker\\_name&pagesize=50&gdexconf=&attr\\_tooltip=nott&wlmaxitems=100&wlsort=f&subcnorm=freq&corpname=parlamint21\\_gb&reload=&usesubcorp=Covid\\_lower\\_F&wlattr=lemma&usengrams=0&ngrams\\_n=2&ngrams\\_max\\_n=2&nest\\_ngrams=0&wlpat=%5Ba-z%5D\\*&wlminfreq=5&wlmaxfreq=0&wlfile=&wlblacklist=&wlnums=frq&wltype=keywords&ref\\_corpname=parlamint21\\_gb&ref\\_usesubcorp=Covid\\_lower\\_M&simple\\_n=1](https://www.clarin.si/noske/run.cgi/wordlist?corpname=parlamint21_gb&viewmode=kwic&attrs=word%2Cdep&ctxattrs=word&structures=speech&refs=%3Dspeech.from%2C%3Dspeech.speaker_name&pagesize=50&gdexconf=&attr_tooltip=nott&wlmaxitems=100&wlsort=f&subcnorm=freq&corpname=parlamint21_gb&reload=&usesubcorp=Covid_lower_F&wlattr=lemma&usengrams=0&ngrams_n=2&ngrams_max_n=2&nest_ngrams=0&wlpat=%5Ba-z%5D*&wlminfreq=5&wlmaxfreq=0&wlfile=&wlblacklist=&wlnums=frq&wltype=keywords&ref_corpname=parlamint21_gb&ref_usesubcorp=Covid_lower_M&simple_n=1)

REF Key

[https://www.clarin.si/noske/run.cgi/wordlist?corpname=parlamint21\\_gb&viewmode=kwic&attrs=word%2Cdep&ctxattrs=word&structures=speech&refs=%3Dspeech.from%2C%3Dspeech.speaker\\_name&pagesize=50&gdexconf=&attr\\_tooltip=nott&wlmaxitems=100&wlsort=f&subcnorm=freq&corpname=parlamint21\\_gb&reload=&usesubcorp=REF\\_last15mths\\_lower\\_F&wlattr=lemma&usengrams=0&ngrams\\_n=2&ngrams\\_max\\_n=2&nest\\_ngrams=0&wlpat=%5Ba-z%5D\\*&wlminfreq=5&wlmaxfreq=0&wlfile=&wlblacklist=&wlnums=frq&wltype=keywords&ref\\_corpname=parlamint21\\_gb&ref\\_usesubcorp=REF\\_last15mths\\_lower\\_M&simple\\_n=1](https://www.clarin.si/noske/run.cgi/wordlist?corpname=parlamint21_gb&viewmode=kwic&attrs=word%2Cdep&ctxattrs=word&structures=speech&refs=%3Dspeech.from%2C%3Dspeech.speaker_name&pagesize=50&gdexconf=&attr_tooltip=nott&wlmaxitems=100&wlsort=f&subcnorm=freq&corpname=parlamint21_gb&reload=&usesubcorp=REF_last15mths_lower_F&wlattr=lemma&usengrams=0&ngrams_n=2&ngrams_max_n=2&nest_ngrams=0&wlpat=%5Ba-z%5D*&wlminfreq=5&wlmaxfreq=0&wlfile=&wlblacklist=&wlnums=frq&wltype=keywords&ref_corpname=parlamint21_gb&ref_usesubcorp=REF_last15mths_lower_M&simple_n=1)

# TASK 3.1: Women-related issues

parlamint21_gb	parlamint21_gb
COVID_lower_F	COVID_lower_M
771.14/mill mentions	169.85/mill mentions
<b>pregnant</b>	<b>man</b>
<b>man</b>	<b>girl</b>
<b>girl</b>	<b>servicemen</b>
<b>young</b>	<b>pregnant</b>
<b>against</b>	<b>sport</b>
<b>equalities</b>	<b>bisexual</b>
<b>who</b>	<b>brave</b>
's	<b>armed</b>
<b>violence</b>	<b>lbt</b>
<b>migrant</b>	<b>young</b>
<b>black</b>	<b>trans</b>
<b>committee</b>	<b>equalities</b>
<b>many</b>	<b>violence</b>
<b>rights</b>	<b>waspi</b>
<b>more</b>	<b>serviceman</b>
<b>face</b>	<b>lesbian</b>
<b>affect</b>	1950
<b>woman</b>	<b>against</b>
<b>particularly</b>	<b>child</b>
<b>and</b>	<b>black</b>

	F (%)	M (%)
<b>collective reference</b>	5.88	<b>17.95</b>
<b>equality &amp; representation</b>	<b>55.88</b>	<b>53.85</b>
<b>reproduction</b>	5.88	<b>10.26</b>
<b>violence &amp; other problems</b>	<b>32.35</b>	17.95

F	M
face, experience, encourage, work, affect, protect	suffer, encourage

F	M
continue training and take up career, run for office, other	participate in sport

- Topic analysis of top 50 collocations for [lemma = "woman"], win.3
- verbs → *encourage*
- content analysis of concordances for *encourage*

# TASK 3.1: collocations

	<b>CONCORDANCES</b> lemma <i>woman</i>	<b>COLLOCATIONS</b> lemma (lc); -3/+3
<b>COVID F</b>	<a href="https://www.clarin.si/noske/run.cgi/first?corpname=parlamint21_gb&amp;reload=&amp;query=&amp;queryselector=lemmarow&amp;lemma=woman&amp;phrase=&amp;word=&amp;char=&amp;cql=&amp;default_atr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsize=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=COVID_lower_F&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=">https://www.clarin.si/noske/run.cgi/first?corpname=parlamint21_gb&amp;reload=&amp;query=&amp;queryselector=lemmarow&amp;lemma=woman&amp;phrase=&amp;word=&amp;char=&amp;cql=&amp;default_atr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsize=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=COVID_lower_F&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=</a>	<a href="https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=COVID_lower_F&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d">https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=COVID_lower_F&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d</a>
<b>COVID M</b>	<a href="https://www.clarin.si/noske/run.cgi/first?corpname=parlamint21_gb&amp;reload=&amp;query=&amp;queryselector=lemmarow&amp;lemma=woman&amp;phrase=&amp;word=&amp;char=&amp;cql=&amp;default_atr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsize=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=COVID_lower_M&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=">https://www.clarin.si/noske/run.cgi/first?corpname=parlamint21_gb&amp;reload=&amp;query=&amp;queryselector=lemmarow&amp;lemma=woman&amp;phrase=&amp;word=&amp;char=&amp;cql=&amp;default_atr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsize=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=COVID_lower_M&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=</a>	<a href="https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=COVID_lower_M&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d">https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=COVID_lower_M&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d</a>
<b>REF F</b>	<a href="https://www.clarin.si/noske/run.cgi/first?corpname=parlamint21_gb&amp;reload=&amp;query=&amp;queryselector=lemmarow&amp;lemma=woman&amp;phrase=&amp;word=&amp;char=&amp;cql=&amp;default_atr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsize=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=REF_last15mths_lower_F&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=">https://www.clarin.si/noske/run.cgi/first?corpname=parlamint21_gb&amp;reload=&amp;query=&amp;queryselector=lemmarow&amp;lemma=woman&amp;phrase=&amp;word=&amp;char=&amp;cql=&amp;default_atr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsize=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=REF_last15mths_lower_F&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=</a>	<a href="https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=REF_last15mths_lower_F&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d">https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=REF_last15mths_lower_F&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d</a>
<b>REF M</b>	<a href="https://www.clarin.si/noske/run.cgi/first?corpname=parlamint21_gb&amp;reload=&amp;query=&amp;queryselector=lemmarow&amp;lemma=woman&amp;phrase=&amp;word=&amp;char=&amp;cql=&amp;default_atr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsize=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=REF_last15mths_lower_M&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=">https://www.clarin.si/noske/run.cgi/first?corpname=parlamint21_gb&amp;reload=&amp;query=&amp;queryselector=lemmarow&amp;lemma=woman&amp;phrase=&amp;word=&amp;char=&amp;cql=&amp;default_atr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsize=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=REF_last15mths_lower_M&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=</a>	<a href="https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=REF_last15mths_lower_M&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d">https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=REF_last15mths_lower_M&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d</a>

	<b>CONCORDANCE</b> of <i>encourage</i> from collocation candidates
<b>F</b>	<a href="https://www.clarin.si/noske/run.cgi/sortx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;o=P_3+3+1%5Blemma_lc%3D%22enco urage%22%5D&amp;atrs=sword%2Ft+1%3C0%7E-3%3C0;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=COV_ID_lower_F&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d">https://www.clarin.si/noske/run.cgi/sortx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;o=P_3+3+1%5Blemma_lc%3D%22enco urage%22%5D&amp;atrs=sword%2Ft+1%3C0%7E-3%3C0;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=COV_ID_lower_F&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d</a>
<b>M</b>	<a href="https://www.clarin.si/noske/run.cgi/sortx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;o=P_3+3+1%5Blemma_lc%3D%22enco urage%22%5D&amp;atrs=sword%2Ft+1%3C0%7E-3%3C0;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=COV_ID_lower_M&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d">https://www.clarin.si/noske/run.cgi/sortx?q=aword%2C%5Blemma%3D%22woman%22%5D&amp;o=P_3+3+1%5Blemma_lc%3D%22enco urage%22%5D&amp;atrs=sword%2Ft+1%3C0%7E-3%3C0;corpname=parlamint21_gb&amp;viewmode=kwic&amp;atrs=word%2Clemma%2Cdep&amp;ctxatrs=word&amp;structs=spech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker_name&amp;lemma=woman&amp;pagesize=50&amp;usesubcorp=COV_ID_lower_M&amp;qdexconf=5&amp;cmfreq=3&amp;cmaxitems=50&amp;cbgrfn=t&amp;cbgrfn=m&amp;cbgrfn=d&amp;csorfn=d</a>

# TASK 3.2: Characterization of women and men

parlamint21_gb	parlamint21_gb
1297 hits for WOMAN etc.	623 hits for MAN etc.
<b>underrepresented</b>	<b>breadwinner</b>
pregnant	island
likely	<b>labourer</b>
twice	<b>vicar</b>
<b>frightened</b>	<b>imam</b>
<b>unaware</b>	likely
<b>reluctant</b>	<b>merchant</b>
polish	twice
<b>concentrated</b>	physically
<b>unhappy</b>	dead
<b>desperate</b>	mother
<b>unable</b>	continent
dead	entire
<b>capable</b>	<b>soldier</b>
<b>fearful</b>	innocent
less	<b>perpetrator</b>
<b>absent</b>	<b>graduate</b>
<b>active</b>	victim
<b>angry</b>	<b>stronger</b>
<b>disadvantaged</b>	father

- Analysis of top 50 collocations for
  - woman|girl|mother|female .\*
  - man|boy|father|male .\*
- in the subject position followed by a copula verb
  - CQL: [lemma="woman|girl|mother|female .\*" & dep="nsubj"] [dep="cop"]
  - (the entire ParlaMint-GB corpus used)

	WOMAN (%)	MAN (%)
<b>positive</b>	25.93	20.83
<b>negative</b>	<b>66.67</b>	25.00
<b>profession</b>	7.41	<b>50.00</b>

# TASK 3.2: UD search + collocations

Entire ParlaMint-GB corpus used!

CQL queries:

- [lemma="woman|girl|mother|female .\*" & dep="nsubj"] [dep="cop"]
- [lemma="man|boy|father|male .\*" & dep="nsubj"] [dep="cop"]

Collocation candidates: lemma (lc);  
0/+3

	CONCORDANCES	COLLOCATIONS
<b>WOMAN</b>	<a href="https://www.clarin.si/noske/run.cgi/first?corpusname=parlament21_gb&amp;reload=&amp;iquery=&amp;queryselector=cqrow&amp;lemma=&amp;phrase=&amp;word=&amp;char=&amp;cql=%5Blemma%3D%22woman%7Cair%7Cmother%7Cfemale+.%22+%26+dep%3D%22nsubj%22%5D+%5Bdep%3D%22cop%22%5D&amp;default_attr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsiz=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=">https://www.clarin.si/noske/run.cgi/first?corpusname=parlament21_gb&amp;reload=&amp;iquery=&amp;queryselector=cqrow&amp;lemma=&amp;phrase=&amp;word=&amp;char=&amp;cql=%5Blemma%3D%22woman%7Cair%7Cmother%7Cfemale+.%22+%26+dep%3D%22nsubj%22%5D+%5Bdep%3D%22cop%22%5D&amp;default_attr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsiz=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name="</a>	<a href="https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%7Cgirl%7Cmother%7Cfemale+.%22+%26+dep%3D%22nsubj%22%5D+%5Bdep%3D%22cop%22%5D&amp;corpusname=parlament21_gb&amp;viewmode=kwic&amp;attr=word%2Clemma%2Cdep%3D&amp;ctxatrs=word&amp;strucs=speech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker.name&amp;pagesize=50&amp;qdexconf=&amp;attr_tooltip=nott&amp;catr=lemma_lc&amp;cfromw=0&amp;ctow=3&amp;cminfreq=5&amp;cminbgr=3&amp;cmaxitems=50&amp;cbqrfn=&amp;cbqrfn=m&amp;cbqrfn=d&amp;csortfn=d">https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22woman%7Cgirl%7Cmother%7Cfemale+.%22+%26+dep%3D%22nsubj%22%5D+%5Bdep%3D%22cop%22%5D&amp;corpusname=parlament21_gb&amp;viewmode=kwic&amp;attr=word%2Clemma%2Cdep%3D&amp;ctxatrs=word&amp;strucs=speech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker.name&amp;pagesize=50&amp;qdexconf=&amp;attr_tooltip=nott&amp;catr=lemma_lc&amp;cfromw=0&amp;ctow=3&amp;cminfreq=5&amp;cminbgr=3&amp;cmaxitems=50&amp;cbqrfn=&amp;cbqrfn=m&amp;cbqrfn=d&amp;csortfn=d</a>
<b>MAN</b>	<a href="https://www.clarin.si/noske/run.cgi/first?corpusname=parlament21_gb&amp;reload=&amp;iquery=&amp;queryselector=cqrow&amp;lemma=&amp;phrase=&amp;word=&amp;char=&amp;cql=%5Blemma%3D%22man%7Cboy%7Cfather%7Cmale+.%22+%26+dep%3D%22nsubj%22%5D+%5Bdep%3D%22cop%22%5D&amp;default_attr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsiz=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name=">https://www.clarin.si/noske/run.cgi/first?corpusname=parlament21_gb&amp;reload=&amp;iquery=&amp;queryselector=cqrow&amp;lemma=&amp;phrase=&amp;word=&amp;char=&amp;cql=%5Blemma%3D%22man%7Cboy%7Cfather%7Cmale+.%22+%26+dep%3D%22nsubj%22%5D+%5Bdep%3D%22cop%22%5D&amp;default_attr=word&amp;fc_lemword_window_type=both&amp;fc_lemword_wsiz=5&amp;fc_lemword=&amp;fc_lemword_type=all&amp;usesubcorp=&amp;fsca_speech.from=&amp;fsca_speech.to=&amp;fsca_speech.speaker_name="</a>	<a href="https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22man%7Cboy%7Cfather%7Cmale+.%22+%26+dep%3D%22nsubj%22%5D+%5Bdep%3D%22cop%22%5D&amp;corpusname=parlament21_gb&amp;viewmode=kwic&amp;attr=word%2Clemma%2Cdep%3D&amp;ctxatrs=word&amp;strucs=speech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker.name&amp;pagesize=50&amp;qdexconf=&amp;attr_tooltip=nott&amp;catr=lemma_lc&amp;cfromw=0&amp;ctow=3&amp;cminfreq=5&amp;cminbgr=3&amp;cmaxitems=50&amp;cbqrfn=&amp;cbqrfn=m&amp;cbqrfn=d&amp;csortfn=d">https://www.clarin.si/noske/run.cgi/collx?q=aword%2C%5Blemma%3D%22man%7Cboy%7Cfather%7Cmale+.%22+%26+dep%3D%22nsubj%22%5D+%5Bdep%3D%22cop%22%5D&amp;corpusname=parlament21_gb&amp;viewmode=kwic&amp;attr=word%2Clemma%2Cdep%3D&amp;ctxatrs=word&amp;strucs=speech&amp;refs=%3Dspeech.from%2C%3Dspeech.speaker.name&amp;pagesize=50&amp;qdexconf=&amp;attr_tooltip=nott&amp;catr=lemma_lc&amp;cfromw=0&amp;ctow=3&amp;cminfreq=5&amp;cminbgr=3&amp;cmaxitems=50&amp;cbqrfn=&amp;cbqrfn=m&amp;cbqrfn=d&amp;csortfn=d</a>

# Creating a tutorial



Image by Alexas Fotos

# Identifying the aim

- Transferring CorpLing techniques to **other SSH disciplines** (other than linguistics)
- Making it relevant for **trans-national research** (comparable, multilingual corpora)
- Showcasing the potential of **special data type** (parliamentary records)
- Making it available to **non-coders** (corpora containing large amounts of data)
- Making it **accessible** (by using non-commercial resources and tools)
- Making it a **self-study manual** divided into several modules (theory supporting the tasks, mentioning potential errors in interpretation)

# Choosing resources and tools

- Parliamentary corpora in CLARIN
  - available for several languages, comparable (ParlaMint)
  - richly annotated (with metadata)
  - well-documented
  - free and open-access
  - available through free online concordancers
- Several paths tested:
  - Scenario 1: Use several parliamentary corpora for comparative analysis (but suboptimal before ParlaMint; uneven time spans, metadata, annotations, different concordancers)
  - Scenario 2: Use English corpora for internationally understandable examples, but the concordancer did not offer the functionalities needed for the selected research problem
  - Scenario 3: siParl corpus (comprehensive in size and annotations) in NoSketchEngine (offering all major CorpLing analytical techniques)

# Creating the tutorial

- **Adapting RQ** to the target audience (SSH fields)
- Embedding the analysis into **their theoretical and methodological framework** (manual analysis of topics based on the ministries, in a later study on Comparative Agendas Project schema)
- Making it possible to **skip and return** to individual parts in theory and hands-on (by analytical design and formatting)
- Leveraging the **affordances of the online format** (hyperlinks, cross-references, screencasts)
- Exposing it to reviews and tests

# Pitfalls to avoid

- Using commercial tools and resources
- Not checking the maturity of the resource or scheduled updates to the tools
- Making it too specific to hinder adaptability to other data resources
- Not adapting the RQ to the interests of the target audience
- Making it too long or not modular

# Why use/adapt the Voices Of The Parliament tutorial?

- **Teach powerful corpus linguistic techniques**
  - No programming skills needed
  - No previous experience with concordancers needed
- **Show elaborated hands-on examples supported by theory**
  - (Screencasts available)
  - Rich in cross-references to dig deeper
- **Help your students assimilate the methods to the point where they can think of possible uses on other data**
  - Similar parliamentary corpora
  - Other discourse types
- **Wide target audience**
  - Researchers outside of linguistics
  - Students in linguistics
  - Anyone interested in multidisciplinary research

Try it out, use it in the classroom, create spin-offs  
and let us know what you think.

Thank you for your  
attention!