

Use Case 1

Gender in Parliamentary Discourse

Background

While more women than ever are being elected to parliaments around the world, equality is still a long way off, and current progress is far too slow. Most parliaments are still heavily male-dominated, and some have no women members of parliament at all. Even where women are present in greater numbers, glass ceilings often remain firmly in place. (Source: [Women in Parliament](#) by the Inter-Parliamentary Union)

CLARIN Resources and services used

- [CLARIN Resource Families – parliamentary corpora](#)
- [the CLARIN.SI repository](#)
- [the noSketch Engine concordancer](#)

Citizen scientist

Do female speakers in the Slovenian and Croatian parliaments speak more or less than their male counterparts?



Student

Is the language of female parliamentary speakers similar in Slovenia and Croatia?



Step-by-step guide

- 1 Search the [Parliamentary CLARIN Resource Family](#) for relevant Slovenian and Croatian corpora. In this walkthrough, we'll use the Croatian and Slovenian ParlaMETER corpora, since they are roughly comparable in terms of time span, linguistic annotation and speaker metadata, but you can also use any of the other parliamentary corpora.

Croatian parliamentary corpus ParlaMeter-hr 1.0

Size: 14.1 million tokens
Annotation: tokenised, MSD-tagged, lemmatised, named entities
Licence: CC-BY

Croatian

The corpus contains minutes of the National Assembly of the Republic of Croatia and currently covers its VIth mandate from 15 November 2016 to 21 November 2018. The corpus contains speaker metadata (gender, age, education, party affiliation).

The corpus is available for download from the CLARIN.SI repository and through the concordancers [KonText](#) and [noSketchEngine](#), as well as through a [dedicated webpage](#).

Concordancer

Download

Slovenian parliamentary corpus siParl 1.0

Size: 227.8 million tokens
Annotation: tokenised, PoS-tagged, lemmatised
Licence: CC BY

Slovenian

The corpus contains Slovenian parliamentary debates from 1990 to 2018. It differs from the SlovParl 2.0 corpus (listed below) in that it contains only basic meta-data about the speakers, a typology of sessions and structural and editorian annotations.

The corpus is available for download from the CLARIN.SI repository and through the concordancers [KonText](#) and [noSketchEngine](#).

Concordancer

Download



- For both corpora, check their descriptions to see that:
 - In terms of linguistic annotation, both corpora are **annotated for syntactic and morphological features** (“MSD-tagged”), **lemmatized** and **marked for named entities**.
 - In terms of extra-linguistic annotation, both corpora are marked for **speaker metadata (gender, age, education, party affiliation)**.
 - The CC-BY licence shows that the corpus is publicly available, either for **download** or **on-line querying**.

- Let's start by analysing the Slovenian corpus. First click on [Slovenian parliamentary corpus ParlaMeter-sl 1.0](#) in the CLARIN Resource Families. This takes you to the record for this corpus in the CLARIN.SI repository:

- The CLARIN.SI repository shows how the corpus has to be cited to ensure proper authorship attribution, and offers a **persistent identifier** for the resource – <http://hdl.handle.net/11356/1208>.

- 5 The corpus can be queried via two concordancers – **KonText** and **noSketch Engine**. Both offer very versatile search environments in which complex queries can be narrowed down on the basis of the **speaker metadata (age, party affiliation, etc)**.

- 6 Let's query the corpus by using the **noSketch Engine**. In the repository, click on the downward arrow next to “noSketch” and then select search.

The screenshot shows the noSketch Engine interface for the ParlaMeter-si (parliament) corpus. The left sidebar contains navigation links: Home, Search, Word list, Corpus info, My jobs, and User guide. The main content area displays the following information:

ParlaMeter-si (parliament)
 Korpus parlamentarnih razprav Republike Slovenije: Mandat VII (2014-08 - 2018-06) // Corpus of parliamentary debates of the Republic of Slovenia: Mandate VII (2014-08 - 2018-06)

Counts	General info	Lexicon sizes	Tags legend	Lempos suffixes
Tokens: 40,987,516	Corpus description: Document	word: 263,007	samostalnik: S.*	samostalnik: -s
Words: 34,882,499	Language: Slovenian	lempos: 109,066	glagol: G.*	glagol: -g
Sentences: 1,833,147	Encoding: UTF-8	tag_en: 1,080	pridevnik: P.*	pridevnik: -p
Paragraphs: 133,287	Compiled: 12/31/2018 22:16:01	tag: 1,080	pristov: R.*	pristov: -r
Documents: 1,338	Tagset: Description	lc: 228,682	zaimek: Z.*	zaimek: -z
		norm: 228,682	predlog: D.*	predlog: -d
		lemma: 104,247	veznik: V.*	veznik: -v
		lemma_lc: 100,467	členek: L.*	členek: -l
			medmet: M.*	medmet: -m

- 7 Let's recall our task: we're interested how parliamentary speakers are represented in the corpus in terms of gender. In other words, how many words of the total 34,882,499 are spoken by female parliament speakers and how many by male speakers?

- 8 We can figure this out by creating a **word list** and narrowing it down to the “Female” subcorpus.

The screenshot shows the noSketch Engine interface for creating a word list. The left sidebar contains navigation links: Home, Search, Word list, Corpus info, My jobs, and User guide. The main content area displays the following information:

Word list options

Corpus: ParlaMeter-si (parliament)

Subcorpus: **Female** (circled in red) [info](#) [create new](#)

Search attribute: word

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: 5

Maximum frequency: 0 (0 = no maximum frequency)

Whitelist: no file selected

Blacklist: no file selected [format](#)

Include non-words

- 9 After clicking on **Make word list**, we get the [result](#) for female speakers. We repeat the procedure for the “Male” subcorpus and [see](#) that the male speakers say 2.5 times more words than their female counterparts.

Repeat the procedure for the Croatian ParlaMeter corpus.

- What is the gender division in terms of words between male and female speakers in this corpus?
- Is the difference greater or smaller than that in the Slovenian corpus?

Additional Task

We can also construct word lists for individual word classes, such as nouns, verbs, adjectives, etc.

Which are the most frequent nouns used by the speakers in the Slovenian corpus?

- a. Under search attribute, change from “word” to “tag_en”. This specifies that you’re searching for parts of speech rather than individual words.
- b. In the filter option “Regular expression”, write N.*. This specifies that you’re searching for all nouns.
- c. Under output options, select “lemma” under “Change output attributes”. This ensures that all inflectional variants are all accounted for under a single base form of the word.
- d. Click [here](#) to see the result for such a query.

Repeat the procedure for the Croatian ParlaMeter corpus. Are the results similar to the Slovenian corpus?

Research bite

In the Slovenian ParlaMeter corpus, the most frequent topics among the female speakers are *health* and *labour*, *family* and *social affairs*, which are followed by *public administration and education*, *science and sport*.

Most of the 100 top-ranking keywords uttered by female speakers, on the other hand, could not be classified into a single topic because they were used either to achieve a *stylistic effect*, were general words that were used in *multiple topics*, such as descriptive adjectives or legal terms, or *ideological expressions*, all of which indicate a more discursive, debating *style* of the male speakers, but could also stem from the fact that the leading roles in that term were predominantly held by male members of parliament (Source: [Parlamenteer – a Corpus of Contemporary Slovene Parliamentary Proceedings](#) by Darja Fišer, Nikola Ljubešić and Tomaž Erjavec).

Use Case 2

Creating a linguistically annotated corpus of 19th century English novels

Background

The digital humanities provide a new conception of the world of literature. Not only is this world larger – the sheer volume of the material we can access is unprecedented – but it is open to levels of analysis that could never be achieved by human brainpower alone. Hierarchies and themes fade into the background as patterns and networks emerge. These methods simultaneously divide texts into new categories and connect them to each other to form new wholes. (Source: [When computers read: Literary analysis and digital technology](#) by Sarah Jones)

CLARIN Resources and services used

- [Virtual Language Observatory \(VLO\)](#)
- [Language Resource Switchboard](#)
- [WebLicht](#)

Teacher

How can my students create an annotated corpus of 19th century English novels from scratch in an easy-to-use online environment?



Researcher

Can you help me find resources and tools to research the stylistic differences between 19th century female and male novelists?

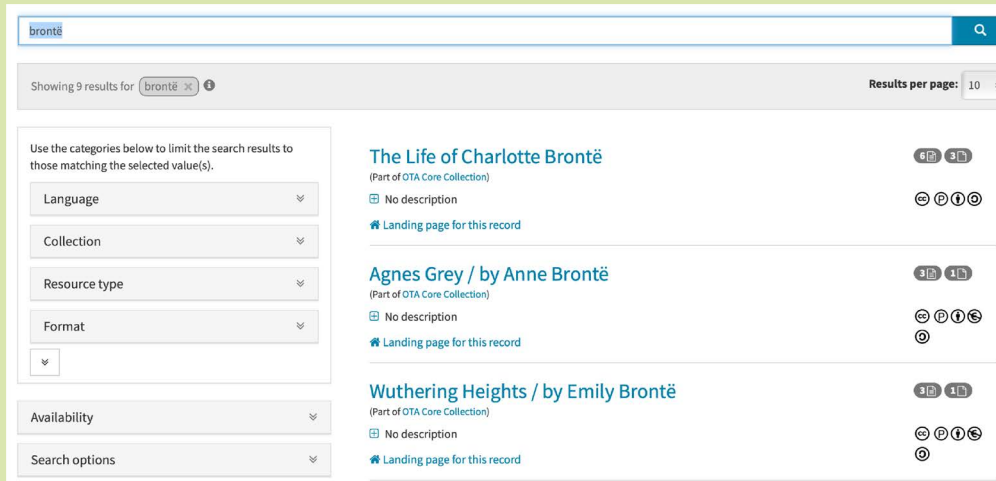


Step-by-step guide

The **Language Resource Switchboard (LRS)** aims at bridging the gap between resources (as identified in the [VLO](#), [Federal Content Search](#), and the [CLARIN Virtual Collection](#)) and tools that can process these resources in one way or another. For a given resource in question, it identifies all tools that can process the resource. It then sorts the tools in terms of the tasks they perform, and presents a task-oriented list to the user. Users can then select and invoke the tool of their choosing. (Source: Adapted from [The Language Resource Switchboard](#) by Claus Zinn)

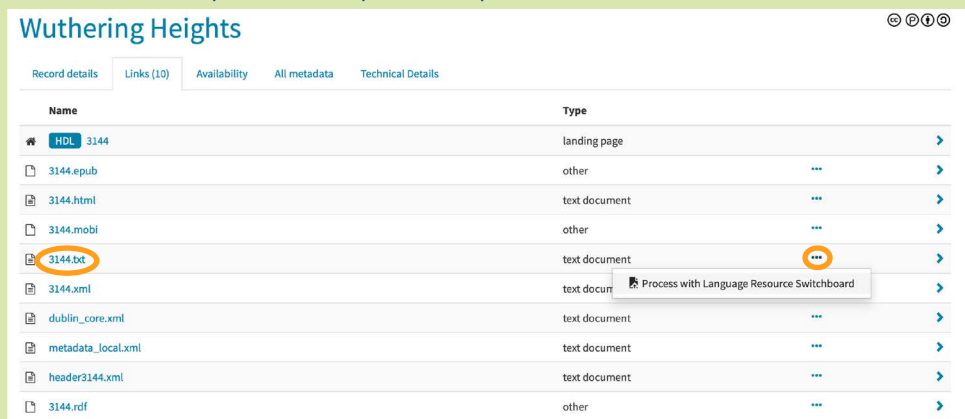


1 Search the VLO with the simple query Brontë.

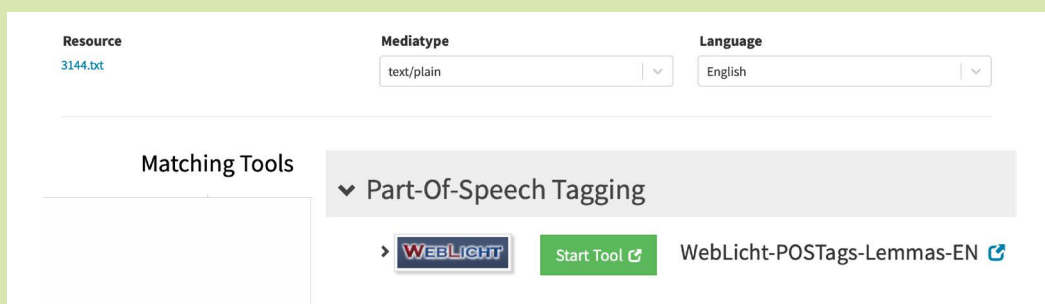


This query gives you VLO records for 19th century English novels by the Brontë sisters, such as Wuthering Heights by Emily Brontë, The Tenant of Wildfell Hall by Anne Brontë, and Jane Eyre by Charlotte Brontë.

2 In each VLO record under the “Links” tab, we can find the complete novels in the form of .txt files. Each file can be processed through the Language Resource Switchboard by clicking on “...” next to the .txt file (in our case, 3144.txt).



3 The first step in linguistic annotation typically involves **part-of-speech tagging**, with which each word in a corpus is assigned a part of speech, like *noun*, *verb*, and *adjective*. In the LRS, we see that part-of-speech tagging is performed by WebLicht.



4 In the WebLich application, we select *PoS tags/lemmas* under “Available Annotations for English Plain text”

The screenshot shows the WebLich application interface. On the left, under 'Available Annotations for: English Plain Text', the 'Pos Tags/Lemmas' option is selected. The 'Visualization Area' contains the text: 'Your results will be visualized here once you have: 1. Chosen one of the annotations on the left 2. Run the tools by clicking Run Tools below. Tip: hovering over an annotation type will give a detailed overview of the output.' Below this, the 'Input and Chain Selection' section shows a table of tool configurations:

Title (Plain Text)	SSS: To TCF Converter	SSS: Stanford Tokenizer	SSS: JSLR POS Tagger	SSS: MorphAdorner Lemmatizer
Wuthering Heights	Language: English	Sentence	Part of Speech: Penn Treebank 1	Lemmas
by Emily Brontë (author)	Document type: TCF	tokens		
Chapter II	TCF Version: 0.4			
1824 -- have just returned from a visit to my land-lord -- the solitary neighbour that I shall be troubled with	Text			

5 After clicking on *Run Tools*, the entire *Wuthering Heights* novel becomes tagged for parts of speech.

The screenshot shows the WebLich application interface after running tools. The 'Query' field contains the text: 'Enter either a TIGERSearch query, or simply a word in quotation marks.' Below the query field, the 'Visualization' area displays the text: '-- I have just returned from a visit to my land-lord -- the solitary neighbour that I shall be troubled with .' The text is tagged with parts of speech (PoS) and lemmas. For example, 'I' is tagged as 'PRP', 'have' as 'VBP', 'just' as 'RB', 'returned' as 'VBN', 'from' as 'IN', 'a' as 'DT', 'visit' as 'NN', 'to' as 'TO', 'my' as 'PRP', 'land-lord' as 'NN', '--' as 'DT', 'the' as 'DT', 'solitary' as 'JJ', 'neighbour' as 'NN', 'that' as 'IN', 'I' as 'PRP', 'shall' as 'MD', 'be' as 'VB', 'troubled' as 'VBN', 'with' as 'IN', and '.' as 'P'.

The annotated novel can now be queried like a regular corpus either for simple words or by using the [TIGERSearch corpus query language](#). To find all the adjectives in the newly tagged corpus, type `[pos = /J.*/]` in the Query field. Make sure to enclose the query in square brackets. Try visualizing the results. Which are the most and least common adjectives in the novel? Hint: In the Statistics visualisation under “Add/remove columns”, try adding the values PoS and lemma.

WebLich also allows you to create additional annotation chains (“New Chain”). Try to repeat the task above by also tagging [The Tenant of Wildfell Hall](#) for **parts of speech** as well as for **named entities**. By creating several annotation chains in this way, you are able to create a full-fledged linguistically annotated corpus, consisting of several novels which were originally in simple plain text.



Research bite

By analysing a corpus of novels by Charles Dickens, Malmberg et al. (2019) have studied how fictional dialogue is used by the author to create a sense of realism and authenticity. The authors have shown that Dickens consistently writes dialogue characterised by linguistic features that fictional and real people share (e.g., question fragments, set expressions conveying politeness and vagueness), which contributes to a sense of naturalness to speech in fiction. By contrast, the range of frequent word combinations in fictional dialogue is more limited than that in spoken fiction, so it is possible that literature adds an iconic or heightened meaningful effect onto these forms. Particularly, shorter lexical combinations (e.g. *I mean, you know* and *I don't know*) are less frequent in fiction. (Source: [Speech-bundles in the 19th-century English novel](#) by Michaela Malmberg, Viola Wiegand, Peter Stockwell, and Anthony Hennessey)