# 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?
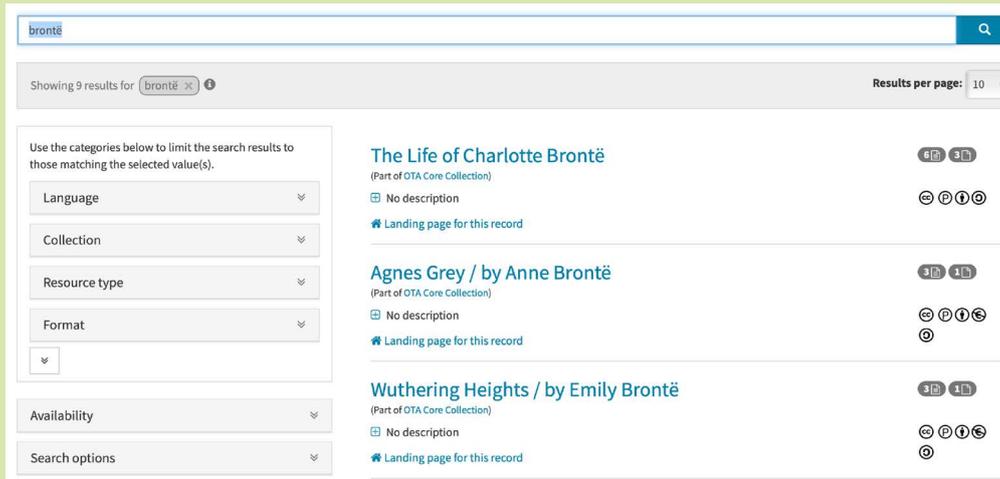
### Step-by-step guide

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)
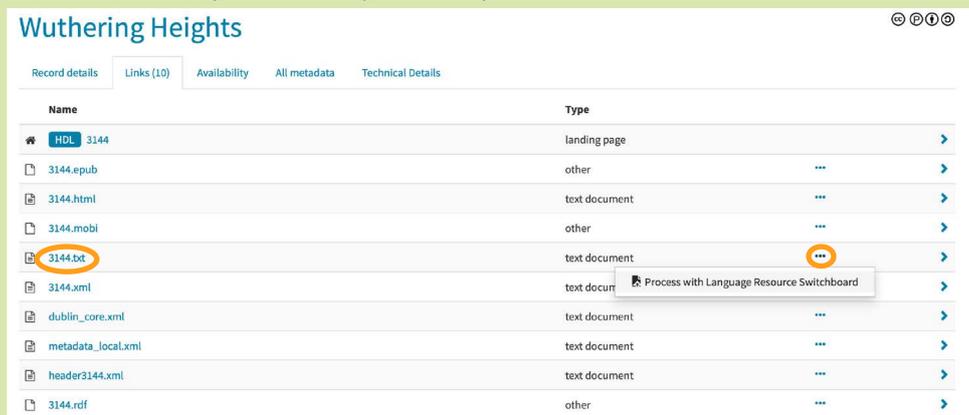
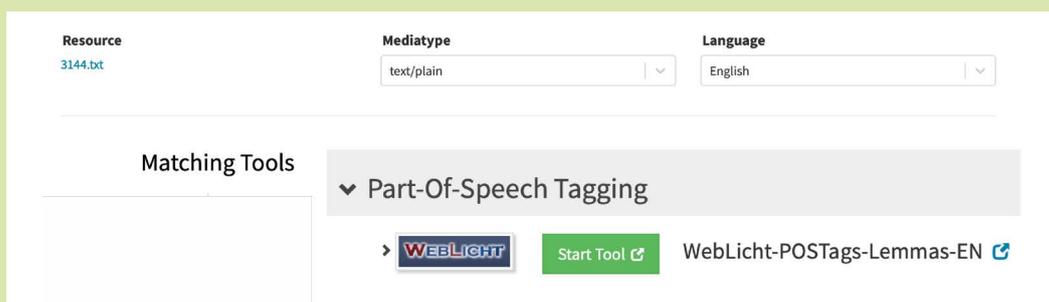**1** Search the VLO with the simple query <u>Brontë</u>.



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

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



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



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.

WebLicht 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, Mahlberg 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)