

# U-Compare Syntactic Chunking Service

## 1. BASIC INFORMATION

### *Service name*

U-Compare Syntactic Chunking Service

### *Overview and purpose of the tool*

This is a web service that identifies and categorises syntactic chunks in Galician text. Also identifies sentence and tokens, and assigns parts-of-speech and lemmas to tokens.

### *A short description of the algorithm*

This web service is based on a UIMA-based workflow, created using the U-Compare text mining system<sup>1</sup>. The workflow was exported from U-Compare as a web service using the built in functionality (Kontonastios et al., In Press). The workflow was created as part of the work to increase the number of interoperable tools operating on different European languages (Ananiadou et al, 2011).

The workflow consists of a single UIMA compliant tool:

1) Freeling<sup>2</sup> parsed web service (performs syntactic chunking)<sup>3</sup>, configured to operate on Galician (service provided by the PANACEA project<sup>4</sup>)

## 2. TECHNICAL INFORMATION

### *Software dependencies and system requirements*

This is a web service that can be run from a browser or accessed programmatically. The only basic requirement is an internet connection.

### *Installation*

There is no installation. The web service can be accessed at the following URL:

[http://nactem001.mib.man.ac.uk:8080/UCompareWebServices/SyntacticChunking\\_FreelingParser](http://nactem001.mib.man.ac.uk:8080/UCompareWebServices/SyntacticChunking_FreelingParser)

---

<sup>1</sup> <http://nactem.ac.uk/ucompare/>

<sup>2</sup> <http://nlp.lsi.upc.edu/freeling/>

<sup>3</sup> <http://registry.elda.org/services/205>

<sup>4</sup> <http://www.panacea-lr.eu/>

The web form available at this URL is shown in Figure 1, with some Galician text entered.

The screenshot shows a web browser window with the address bar containing the URL: `nactem001.mib.man.ac.uk:8080/UCCompareWebServices/SyntacticChunking_FreelingParser`. The page header features the logo of 'The National Centre for Text Mining'. Below the logo, there are three radio buttons: 'XML document', 'inline XML', and 'stand-off annotation', with 'stand-off annotation' selected. A large text area contains the following Galician text: 'O presidente da Xunta, Alberto Núñez Feijóo, confirmou que a Comunidade activou formalmente o protocolo contra a seca, que integra os diferentes plans deseñados polo seu Goberno para facer fronte a este problema, posto que as actuacións se iniciaron "hai meses. Ao tempo, chamou os cidadáns ao consumo responsable da auga. "Pídelles aos galegos esta sensibilidade porque a auga é agora mesmo en Galicia un ben escaso", sentenciou na rolda de prensa posterior ao consello de la Xunta, na que reivindicou a necesidade de aplicar "as maiores doses de responsabilidade que se deron nunca" na Comunidade para facer fronte a esta situación.' Below the text area is a blue 'Run' button. Below the form, there are four sections: 'Service Description', 'Usage', 'References', and 'Application programming interface'. 'Service Description' contains the text: 'Ideas in worklow: freeling snairow parser web service (service provided by the PANACEA project) Supported language(s): Galician (gl)'. 'Usage' contains the text: 'POST request should be sent to use the service' and '1) text -- the value of this parameter is the text to analyze. Expected encoding is UTF-8'. 'References' is empty. 'Application programming interface' contains the code snippet: 'String text = "Hello Mr. John Smith I!"; String parameters = "text=" + URLEncoder.encode(text, "UTF-8") + "&mode=inline"; URL url = new URL(url of the service);'.

Figure 1: Web form for the service

### Execution instructions

The web service can be executed by typing or pasting text into the online form and clicking on the “Run” button.

Alternatively, the web service can be executed from within program code, as explained in the “Usage” and “Application programming interface” boxes of the web form.

A POST request should be used to call the service. The following parameters may be used in the request:

- **text** - the value of this parameter is the text to analyze. Expected encoding is UTF-8. This parameter is obligatory.
- **lang** - This parameter sets the language of the text. If this parameter is not provided, then the value "en" will be used

- **mode** - This parameter sets the format of the annotated information returned by the service. If this parameter is not set, XML output will be produced. The two possible types of output are as follows:
  - **inline** – annotations are encoded as inline XML.
  - **xml** – results are output as an XML document containing the annotations added

The following code example shows how the web service can be called from Java code:

```
//Set the input text
String text = "<Text_to_be_analysed>";
//Set the parameter string
String parameters = "text=" + URLEncoder.encode(text,
"UTF-8") + "&mode=inline";
//Create the URL connection
URL url = new
URL("http://nactem001.mib.man.ac.uk:8080/UCompareWebServi
ces/SyntacticChunking_FreelingParser");
URLConnection connection = url.openConnection();
connection.setDoOutput(true);
//Create Output stream
OutputStreamWriter writer = new
OutputStreamWriter(connection.getOutputStream());
//write parameters to output stream
writer.write(parameters);
writer.flush();

//Read the results returned by the service
BufferedReader reader = new BufferedReader(new
InputStreamReader(connection.getInputStream(), "UTF-8"));
String line;
while ((line = reader.readLine()) != null) {
    System.out.println(line);
}
```

## Input/Output data formats

### Input data formats

The input is plain text, UTF-encoded.

### Output data format

If the service is run from the web interface, then the output is visualized in the interface using colored highlights in the text to show the individual annotations, and one or more tables of information below, each corresponding to a particular type of annotation.

If the service is run programmatically, then the output is provided in XML format. See section 3 for an example.

### Integration with external tools

The API allows the functionality of the web service to be embedded in any application.

## 3. CONTENT INFORMATION

Using the web interface, the output of the service is visualised as shown in Figure 2.

Chunk	labelString
para	SPS00
facер	VMN0000
fronte a	CA
esta	DDOFS0
situación	NCFS000
.	Fp
O presidente da Xunta, Alberto Núñez Feijóo, confirmou que a Comunidade activou formalmente o protocolo contra a seca, que integra os diferentes plans deseñados polo seu Goberno para facer fronte a este problema, posto que as actuacións se iniciaron "hai meses. Ao tempo, chamou os cidadáns ao consumo responsable da auga. "Pídelles aos galegos esta sensibilidade porque a auga é agora mesmo en Galicia un ben escaso", sentenciou na rolda de prensa posterior ao consello de la Xunta, na que reivindicou a necesidade de aplicar "as maiores doses de responsabilidade que se deron nunca" na Comunidade para facer fronte a esta situación.	grup-verb
O presidente da Xunta, Alberto Núñez Feijóo,	sn
O	j-ms
presidente	n-ms
da Xunta	sp-de
a Xunta	sn
a	j-fs
Xunta	w-fs
,	Fc
Alberto Núñez Feijóo	w-ms
,	Fc
confirmou	verb
que	conj-subord
a Comunidade activou formalmente o protocolo contra a seca, que integra os diferentes plans deseñados polo seu Goberno para facer	grup-verb
a Comunidade	sn
a	j-fs
Comunidade	w-fs
activou	verb
formalmente	sadv

Figure 2: Visualisation of web service output

In Figure 2, the different types of annotations are displayed in different tables. At the top of the screen is the bottom of the table containing information about tokens, their parts-of-speech and lemmas. At the bottom of the screen is the “Chunk” table,

showing all the (possibly overlapping) chunks that have been identified by the service. For each chunk, the span of text is displayed, together with the category assigned to the chunk, in the “labelString” column.

An example of the XML output format, which is more suited to programmatic use, is shown in Figure 3. In the XML, different XML element types encode the different types of annotations. In the example shown, the RichToken elements store information about tokens, while the Chunk annotations store information about syntactic chunks. In each XML element, the start and end offsets of the corresponding annotation in the text are encoded in the “begin” and “end” attributes. Other attributes are possible, depending on the element type. The Chunk type stores the chunk category in the “labelString” attribute, while the RichToken type stores the part-of-speech tag and base form assigned to each token, in the “posString” and “base” attributes, respectively.

```
- <Chunk begin="23" end="188" labelString="sn">
  Alberto Núñez Feijóo, confirmou que a Comunidade activou formalmente o protocolo contra a seca, que integra os diferentes plans de
</Chunk>
<Chunk begin="23" end="43" labelString="w-ms">Alberto Núñez Feijóo</Chunk>
<Chunk begin="23" end="43" labelString="grup-nom-ms">Alberto Núñez Feijóo</Chunk>
<RichToken base="alberto_núñez_feijóo" begin="23" end="43" posString="NP00000">Alberto Núñez Feijóo</RichToken>
<Chunk begin="43" end="44" labelString="Fc">,</Chunk>
<RichToken base="," begin="43" end="44" posString="Fc">,</RichToken>
<Chunk begin="45" end="54" labelString="verb">confirmou</Chunk>
<RichToken base="confirmar" begin="45" end="54" posString="VMIS3S0">confirmou</RichToken>
<Chunk begin="55" end="58" labelString="conj-subord">que</Chunk>
<RichToken base="que" begin="55" end="58" posString="CS">que</RichToken>
- <Chunk begin="59" end="188" labelString="grup-verb">
  a Comunidade activou formalmente o protocolo contra a seca, que integra os diferentes plans deseñados polo seu Goberno para facer
</Chunk>
```

Figure 3: XML output example

### 3. LICENCE

a) The web service only is licenced NaCTeM Web Service Licence Agreement (standard non-commercial use) – see “U-Compare-Syntactic-Chunking-Service-Licence.pdf” in the “licences” directory. Please contact us using the details below if you require a commercial licence.

b) The tools used in the workflow on which the web service is based may have their own licences. The NaCTeM Web Service Licence Agreement does NOT apply to these tools.

### 4. ADMINISTRATIVE INFORMATION

#### Contact

For further information, please contact Sophia Ananiadou:

[sophia.ananiadou@manchester.ac.uk](mailto:sophia.ananiadou@manchester.ac.uk)

## 5. REFERENCES

Ananiadou, S., Thompson, P., Kano, Y., McNaught, J., Attwood, T. K., Day, P. J. R., Keane, J., Jackson, D. and Pettifer, S.. (2011). Towards Interoperability of European Language Resources. *Ariadne*, 67.

Kontonatsios, G., Korkontzelos, I., Kolluru, B., Thompson, P. and Ananiadou, S. (In Press). Deploying and Sharing U-Compare Workflows as Web Services. *Journal of Biomedical Semantics*.