

U-Compare Co-Reference Identification Service

1. BASIC INFORMATION

Service name

U-Compare Co-Reference Identification Service

Overview and purpose of the tool

This is a web service that identifies co-reference chains in Romanian text. Also identifies sentences, tokens with parts-of-speech and lemmas, and NP chunks.

A short description of the algorithm

This web service is based on a UIMA-based workflow, created using the U-Compare text mining system¹. 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 the following UIMA-compliant tools

- 1) Language Identifier (RACAI, Romania)
- 2) TTL-Tokenizer (RACAI, Romania)
- 3) TTL-Tagger (RACAI, Romania)
- 4) TTL-Lemmatizer (RACAI, Romania)
- 5) UAIC-NPChunker (UAIC, Romania)
- 6) UAIC-RARE (UAIC, Romania)

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:

¹ <http://nactem.ac.uk/ucompare/>

http://nactem001.mib.man.ac.uk:8080/UCompareWebServices/Coreference_TTL_UAIC

The web form available at this URL is shown in Figure 1, with some Romanian text entered into the text box.

The screenshot shows a web browser window with the URL http://nactem001.mib.man.ac.uk:8080/UCompareWebServices/Coreference_TTL_UAIC. The page header features the logo of 'The National Centre for Text Mining'. Below the header, there are radio buttons for 'XML document', 'inline XML', and 'stand-off annotation', with 'stand-off annotation' selected. A large text area contains a paragraph of Romanian text about William Gibson. Below the text area is a blue 'Run' button. The form is divided into several sections: 'Service Description' (describing the UIMA-based workflow), 'Usage' (specifying a POST request with a text parameter), 'References' (an empty box), and 'Application programming interface' (providing a code snippet for a REST client).

Figure 1: Web form for the web 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/UCompareWebServices/Coreference_TTL_UAIC");
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.

The screenshot shows a web interface for selecting annotation types. At the top, there is a blue header bar. Below it, the text "Select type of annotation" is displayed. Underneath, there are four checkboxes, all of which are checked: "CoreferenceChain", "NpChunkWithHead", "Sentence", and "RichToken". Below the checkboxes is another blue header bar. At the bottom, there is a text box containing a paragraph of text about William Gibson. The text is highlighted with colored boxes: yellow for names and dates, green for locations, blue for titles, and red for specific terms or concepts. The text in the box is: "William Gibson (n. 17 martie 1948) este un scriitor canadian de origine americană. Romanele lui de ficțiune termenul de "cyberspațiu" în povestirea sa " Burning Chrome " (1982) și a popularizat conceptul mai târziu în iconografie a erei informaționale înainte de apariția Internetului în anii '90. Tot lui Gibson îi se datorează termenii virtuale precum jocurile video și World Wide Web. Pentru că familia lui s-a mutat foarte des în timpul cercetării științifico-fantastică. După ce și-a petrecut adolescența într-un internat privat din Arizona, Gibson a evitat să fie preocupat de contracultură și, după ce s-a stabilit în Vancouver, a devenit scriitor profesionist. Primele sale opere și rețelelor de calculatoare asupra oamenilor — o " combinație între calitatea slabă a vieții și calitatea înaltă a tehnologiei, decorurile și personajele dezvoltate în aceste povestiri culminează în primul său roman, Neuromancer, literară cyberpunk. Deși mare parte din reputația lui Gibson a rămas legată de Neuromancer, opera sa a fost înlocuită astfel distopica trilogie Sprawl, Gibson a devenit un autor important al altui subgen literar SF — steampunk Sterling. În anii '90, el a scris trilogia trilogie Bridge, care analizează societatea urbană a viitorului a steampunk Country (2007) și Zero History (2010) — descriu lumea contemporană și l-au plasat pe autor în rândul scriitorilor de SF contemporani".

Figure 2: Visualisation of web service output

In Figure 2, the top of the screen has check boxes corresponding to each type of annotation produced by the workflow – in this case “CoreferenceChain”, “NPChunkWithHead” (corresponding to the NP chunks), “Sentence” and “RichToken” annotations (the latter of which allow both part-of-speech tags and lemma information to be associated with tokens). Checking one or more of the boxes will cause the annotations to become highlighted in the view of the text below. In figure 2, all different boxes are selected. Each type of annotation is highlighted using a different colour, hence the several highlighting colours used in Figure 2.

Below the text, the different types of annotations added by the workflow are shown in tabular format, with each type of annotation in a separate table.

Using the web interface, it is also possible to view the output in XML format, by selecting one of the other radio buttons above the text input box shown in Figure 1. Each different annotation type is encoded using a different XML element type, and the attributes of the annotation (e.g., part-of-speech tags) are encoded as attributes of the XML element.

3. LICENCE

a) The web service only is licenced NaCTeM Web Service Licence Agreement (standard non-commercial use) – see “U-Compare-Co-reference-Identification--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

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*.