

# U-Compare Lemmatisation Service

## 1. BASIC INFORMATION

### *Service name*

U-Compare Lemmatisation Service

### *Overview and purpose of the tool*

This is a web service that identifies sentences and tokens in plain text. Parts of speech and lemmas are assigned to tokens. Language is automatically identified amongst the supported languages (French, Romanian and English), and language-specific processing is carried out.

### *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 the following UIMA-compliant tools. All are provided by the Research Institute for Artificial Intelligence (RACAI), Romania.<sup>2</sup>

- 1) Language Identifier
- 2) TTL-Tokenizer
- 3) TTL-Tagger
- 4) TTL-Lemmatizer

## 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.

---

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

<sup>2</sup> <http://www.racai.ro/webservices/>

## Installation

There is no installation. The web service can be accessed at the following URL:  
[http://nactem001.mib.man.ac.uk:8080/UCompareWebServices/POS\\_Tagging\\_MLRS](http://nactem001.mib.man.ac.uk:8080/UCompareWebServices/POS_Tagging_MLRS)

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

The screenshot shows a web browser window with the address bar containing the URL `nactem001.mib.man.ac.uk:8080/UCompareWebServices/Lemmatisation_TTL`. The page title is "wikipedia french". The main content area features a logo for "The National Centre for Text Mining" and a form with the following elements:

- Radio buttons for output format:  XML document,  inline XML,  stand-off annotation.
- A text area containing the following French text: "Plusieurs athlètes suisses ayant des chances de médailles manquent cependant leurs objectifs. Par exemple, le cycliste Fabian Cancellara, champion olympique en titre, est septième du contre-la-montre sur route après une chute lors de la course en ligne. La délégation suisse ne remplit pas l'objectif principal de Swiss Olympic, qui est d'être parmi les 25 meilleures nations. Elle gagne moins de médailles que lors des quatre éditions précédentes des Jeux d'été. Les athlètes suisses obtiennent 6 diplômes, contre 13 en 2008. Gian Gilli juge ces résultats insuffisants."
- A blue "Run" button.
- Four informational boxes: "Service Description", "Usage", "References", and "Application programming interface".

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/UCompareWebServic
es/Lemmatisation_TTL");
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.

### Select type of annotation

Sentence  RichToken

Plusieurs athlètes suisses ayant des chances de médailles manquent cependant leurs objectifs ; Par exemple , le cycliste Fabian Cancellara, champion olympique en titre, est septième du contre-la-montre sur route après une chute lors de la course en ligne ; La délégation suisse ne remplit pas l'objectif principal de Swiss Olympic ; Elle gagne moins de médailles que lors des quatre éditions précédentes des Jeux d'été ; Les athlètes suisses obtiennent 6 diplômes , contre 13 en 2008 ; Gian Gilli juge ces résultats insuffisants.

#### Sentence

Plusieurs athlètes suisses ayant des chances de médailles manquent cependant leurs objectifs.  
 Par exemple, le cycliste Fabian Cancellara, champion olympique en titre, est septième du contre-la-montre sur route après u  
 La délégation suisse ne remplit pas l'objectif principal de Swiss Olympic, qui est d'être parmi les 25 meilleures nations.  
 Elle gagne moins de médailles que lors des quatre éditions précédentes des Jeux d'été.  
 Les athlètes suisses obtiennent 6 diplômes, contre 13 en 2008.  
 Gian Gilli juge ces résultats insuffisants.

RichToken	posString	base
Plusieurs	Ai-mp	plusieurs
athlètes	Ncfp	athlète
suisses	Ncfp	suisse
ayant	Vapp	avoir
des	Dg-fp	de_le
chances	Ncfp	chance
de	Spd	de
médailles	Ncfp	médaille
manquent	Vmsp3p	manquer
cependant	R	cependant
leurs	Ds3mp	leur
objectifs	Ncmp	objectif
.	PERIOD	.
Par	Sp	par
exemple	Ncms	exemple

**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 “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, only “RichToken” annotations are highlighted.

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. In Figure 2, it can be seen that there are two tables, one for “Sentence” annotations and one for “RichToken” annotations. For each annotation type, the information associated with each annotation is shown in a row of the table. For sentences, the information comprises only the text covered by the sentence annotation. For “RichToken” annotations, the information additionally includes the part-of-speech tag assigned to the token (in the “posString” column) and the lemma (in the “base” column).

An example of the XML output format, which is more suited to programmatic use, is shown in Figure 3. In the XML, the start and end offsets of each annotation in the text are encoded in the “begin” and “end” attributes. For the “RichToken” type, the “posString” and “base” attributes of the annotations encode the part-of-speech and lemma of each token, respectively.

```

- <result>
- <Sentence begin="0" end="94">
  Plusieurs athlètes suisses ayant des chances de médailles manquent cependant leurs objectifs.
  </Sentence>
  <RichToken base="plusieurs" begin="0" end="9" posString="Ai-mp">Plusieurs</RichToken>
  <RichToken base="athlète" begin="10" end="18" posString="Ncfp">athlètes</RichToken>
  <RichToken base="suisse" begin="19" end="26" posString="Ncfp">suisses</RichToken>
  <RichToken base="avoir" begin="27" end="32" posString="Vapp">ayant</RichToken>
  <RichToken base="de_le" begin="33" end="36" posString="Dg-fp">des</RichToken>
  <RichToken base="chance" begin="37" end="44" posString="Ncfp">chances</RichToken>
  <RichToken base="de" begin="45" end="47" posString="Spd">de</RichToken>
  <RichToken base="médaille" begin="48" end="57" posString="Ncfp">médailles</RichToken>
  <RichToken base="manquer" begin="58" end="66" posString="Vmsp3p">manquent</RichToken>
  <RichToken base="cependant" begin="67" end="76" posString="R">cependant</RichToken>
  <RichToken base="leur" begin="77" end="82" posString="Ds3mp">leurs</RichToken>
  <RichToken base="objectif" begin="83" end="92" posString="Ncmp">objectifs</RichToken>
  <RichToken base="." begin="92" end="93" posString="PERIOD">.</RichToken>
- <Sentence begin="94" end="254">
  Par exemple, le cycliste Fabian Cancellara, champion olympique en titre, est septième du contre-la-montre st
  </Sentence>
  <RichToken base="par" begin="94" end="97" posString="Sp">Par</RichToken>
  <RichToken base="exemple" begin="98" end="105" posString="Ncms">exemple</RichToken>
  <RichToken base="," begin="105" end="106" posString="COMMA">,</RichToken>
  <RichToken base="le" begin="107" end="109" posString="Da-ms">le</RichToken>
  <RichToken base="cycliste" begin="110" end="118" posString="Af-fs">cycliste</RichToken>
  <RichToken base="Fabian" begin="119" end="125" posString="Np">Fabian</RichToken>

```

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-Lemmatization-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*.