

U-Compare Paragraph-Breaking Service

1. BASIC INFORMATION

Service name

U-Compare Paragraph-Breaking Service

Overview and purpose of the tool

This is a web service that identifies paragraphs in plain text. The tool will work for most European languages.

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 (Kontonastsios 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, i.e.:

- 1) MLRS Paragraph Breaker² (University of Malta)

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/Paragraph_breaking_MLRS

¹ <http://nactem.ac.uk/ucompare/>

² <http://mlrs.research.um.edu.mt/index.php?page=1>

The web form available at this URL is shown in Figure 1. Some text in the Basque language has been pasted in.

The screenshot shows a web browser window with the following details:

- URL:** nactem001.mib.man.ac.uk:8080/UCompareWebServices/Paragraph_breaking_MLRS
- Page Content:**
 - Logo:** The National Centre for Text Mining logo, featuring three overlapping documents and a green 'U' shape.
 - Section:** Examples
 - Options:** XML document, inline XML (selected), stand-off annotation
 - Text:** Euskal Autonomia Erkidegoko (EAE) enplegurako zerbitzu publikoen bulegoetan erregistraturiko langabeen kopurua 6.198 lagunetan igo zen joan den urtarrilean, hau da, abenduan baino %3,67 gehiago, eta horrela, guztira 175.281 langabe daude, Espainiako Enplegu eta Gizarte Segurantza Ministerioak jakinarazi duenez.
 - Text Analysis:** Beste alde batetik, 2012ko urtarrilarenk alderatzen bada, aurtengo lehen hilabetean 19.414 langabe gehiago izan dira EAEn, hau da, %12,46ko hazkundea izan da urte batetik bestera.
 - Text Statistics:** Lurraldeka, Araban hazi zen gehien langabezia (%4,56), 1.213 langabe gehiago zenbatuta, eta guztira 27.797 daude; Gipuzkoan 2.035 langabe gehiago izan ziren urtarrilean (%4,27) eta guztira 49.723 daude; Bizkaian %3,11 igo zen langabezia (2.950 langabe gehiago) eta guztira 97.761 daude. Urte artean, Araban %14,78 igo da langabezia (3.580 langabe gehiago), Bizkaian %11,52 (10.098 gehiago) eta Gipuzkoan %13,04 (5.736 gehiago).
- Buttons:** Run

Service Description	Usage
Web service created by exporting UIMA-based workflow from the U-Compare text mining system. Functionality: Identifies paragraphs in plain text Tools in workflow: MI-BC Paragraph Splitter (University of Malta)	POST request should be sent to use the service 1) text -- the value of this parameter is the text to analyze. Expected encoding is UTF-8

References	Application programming interface
	String text = "Hello Mr. John Smith !"; String parameters = "text=" + URLEncoder.encode(text, "UTF-8") + "&mode=inline"; URI url = new URI(url of the service);

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/UCompareWebServi
ces/POS_Tagging_MLRS");
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 user interface for visualizing service output. At the top, there is a blue header bar. Below it, a white form with a blue border. The form has a title "Select type of annotation" and a checkbox labeled "Paragraph" which is checked. The main area contains a text box with highlighted annotations. The text is in Spanish and discusses the number of annotations and their distribution across different regions. Below the text box, there is a table with two rows, each containing a title and a detailed description of the annotations found in the text.

Paragraph	Euskal Autonomia Erkidegoko (EAE) enplegurako zerbitzu publikoan bulegoetan erregistraturiko langabeen kopurua 6.198 lagunetan igo zen joan den urtarrilean, hau da, abenduan baino %3,67 gehiago, eta horrela, guztira 175.281 langabe daude, Espainiako Enplegu eta Gizarte Segurantza Ministerioak jakinarazi duenez.
Beste alde batetik, 2012ko urtarrilarenak alderatzen bada, aurtengo lehen hilabetean 19.414 langabe gehiago izan dira EAEn, hau da, %12,46ko hazkundeia izan da urte batetik bestera.	Lurraldeka, Araban hazi zen gehien langabezia (%4,56), 1.213 langabe gehiago zenbatuta, eta guztira 27.797 daude; Gipuzkoan 2.035 langabe gehiago izan ziren urtarrilean (%4,27) eta guztira 49.723 daude; Bizkaian %3,11 igo zen langabezia (2.950 langabe gehiago) eta guztira 97.761 daude. Urte artean, Araban %14,78 igo da langabezia (3.580 langabe gehiago), Bizkaian %11,52 (10.098 gehiago) eta Gipuzkoan %13,04 (5.736 gehiago)

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 only “Paragraph” annotations are present. Checking the box determines whether the paragraph annotations are highlighted in the text below.

Below the text, the annotations added by the workflow are shown in tabular format. In Figure 2, there is a single table, which indicates the span of text covered by each “Paragraph” annotation (one paragraph per row of the table).

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.

```
- <result>
  - <Paragraph begin="0" end="312">
    Euskal Autonomia Erkidegoko (EAE) enplegurako zerbitzu publikoen bulegoetan erregistra
    horrela, guztira 175.281 langabe daude, Spainiako Enplegu eta Gizarte Segurantza Ministe
  </Paragraph>
  - <Paragraph begin="318" end="498">
    Beste alde batetik, 2012ko urtarrilarekin alderatzen bada, aurtengo lehen hilabetean 19.414 ↗
  </Paragraph>
  - <Paragraph begin="504" end="930">
    Lurraldetik, Araban hazi zen gehien langabezia (%4,56), 1.213 langabe gehiago zenbatuta, e
    Bizkaian %3,11 igo zen langabezia (2.950 langabe gehiago) eta guztira 97.761 daude. Urtean
    %13,04 (5.736 gehiago)
  </Paragraph>
</result>
```

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-Paragraph-Breaking-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*.

Rosner, M., Attard, A., Thompson, P., Gatt, A. and Ananiadou, S. (2011). Extending a Tool Resource Framework with U-Compare. In *Proceedings of the 5th Language & Technology Conference (LTC'2011)*.