



2017-IPR-I-000-8664

**Multilingual Named Entity Recognition
and Classification**

<p>Position for:</p> <p>Trainee</p>	<p><u>Short description of activity:</u></p> <p>As the science and knowledge service of the Commission, the mission of DG Joint Research Centre is to support EU policies with independent evidence throughout the whole policy cycle.</p> <p>The JRC is located in 5 Member States (Belgium, Germany, Italy, the Netherlands and Spain). Further information is available at: http://www.jrc.ec.europa.eu</p> <p><u>Short description of activity:</u></p> <p>The JRC's <i>Europe Media Monitor</i> (EMM) team carries out research and development in the field of highly multilingual text mining (Language Technology; Computational Linguistics) for the purposes of media monitoring. EMM gathers an average of 300,000 online news articles per day in over 70 languages and analyses them to help its large international user community understand and use this enormous amount of media information. The <i>Europe Media Monitor</i> EMM is publicly accessible and widely used. The EMM team has produced over 200 international peer-reviewed publications. The team has also produced and distributes a number of highly multilingual Language Technology resources.</p> <p>The <i>Text and Data Mining Unit</i> (I3) of the European Commission's <i>Joint Research Centre</i> (JRC) in Ispra, Italy, is looking for a trainee to support the JRC's <i>Europe Media Monitor</i> (EMM) team in its effort to improve its Named Entity Recognition and Classification (NERC) tools, especially for multi-word entities such as organisation and event names. EMM gathers and analyses reports from traditional and social media in dozens of languages by clustering related news items; categorising them; extracting information such as entities (persons, organisations, locations), events (who did what to whom, where and when), quotations by and about people; identifying sentiment; as well as linking related news clusters over time and across languages. Methods used are mostly hybrid: machine learning tools are used to gather evidence, learn vocabulary and rules,</p>
--	---

but the results are usually controlled and optimised through human intervention. EMM is used by European Institutions, by national authorities in EU Member States, by international organisations and by the public. The public EMM applications [NewsBrief](#), [NewsExplorer](#) and [MedISys](#) can be accessed freely by the general public. EMM is part of the [JRC's Competence Centre on Text Mining and Analysis](#).

As of now, the EMM team has accumulated several very large independent sets of multi-word entities and their monolingual and multilingual name variants. Some of the entities are classified according to an entity type hierarchy, while others are not. **The successful trainee will help to improve the current tools to recognise multi-word entities, classify entities, merge the various lists of entities and their variants into one single repository, and integrate the NERC tools with the EMM processing chain.** The trainee is also expected to contribute to writing a scientific publication on the work carried out.

Qualifications:

Essential:

- a degree (or an almost completed degree) in computational linguistics, computer science or related areas;(Applications from students currently preparing a thesis for a University degree are eligible. The thesis should match with the subject of the project call).
- Java programming skills;
- good working knowledge of English. (B2 level)

Advantage:

- knowledge of further foreign languages;
- proven advanced programming skills, especially in Java;
- good knowledge of Language Technology-related tools and methods;
- proven ability to work independently and as part of a team.

In your application, please provide clear information on your skill set, by elaborating on the above-mentioned list of requirements and by listing your level of languages and your computer / programming skills.

	<p><u>For general eligibility requirements, please read the rules governing the traineeship scheme of the JRC:</u></p> <p>https://ec.europa.eu/jrc/en/working-with-us/jobs/temporary-positions/jrc-trainees</p>
Institute/Directorate Unit	Directorate <i>Competences</i> I03 – <i>Text and Data Mining Unit</i>
Indicative duration	5 months
Preferred starting date	As soon as possible
JRC Site	Ispra
Country	Italy
<u>JRC contact details</u>	<p>For any technical problems with your application, please contact:</p> <p>JRC-ESRA@ec.europa.eu</p>