# Linking News Content Across Languages

**Ralf Steinberger**

European Commission – Joint Research Centre

21027 Ispra (VA), Italy

http://langtech.jrc.it/  –  http://press.jrc.it/overview.html

`Ralf.Steinberger@jrc.it`

## 1 Introduction

Organisations and individuals that need to monitor what the media say about certain issues face an extreme information overload, especially if they are interested in the news written in more than one language. News aggregators sometimes pre-filter potentially user-relevant articles or automatically group related articles into clusters. However, the enormous amount of available online information calls for further automatic information processing to enable users to sieve through even larger amounts of textual data in less time and to navigate and explore the document collections efficiently.

## 2 NewsExplorer

NewsExplorer is a freely available news analysis system that offers such functionality in 19 languages. NewsExplorer integrates various text analysis applications including clustering, multi-label document classification, named entity recognition, name variant matching across languages and writing systems, topic detection and tracking, and more. The purpose of this presentation is to present this news exploration and analysis system and to especially address the multilinguality issue and the cross-lingual functionality of the application. References to prior art will be made, where appropriate.

## 3 News Data and the EMM family of applications

NewsExplorer is part of the Europe Media Monitor (EMM) family of applications (http://press.jrc.it/overview.html). EMM gathers a daily average of 80,000 news articles from about 2,200 web news sources in 43 languages. *NewsBrief* and the Medical Information System *MedISys* classify the news, cluster related articles and alert users of breaking news when unexpected spikes are detected. *EMM-Labs* gives access to data visualisation tools and to the results of a collection of advanced text processing tools such as relation extraction, event scenario template filling, and various types of social networks. The freely available EMM online applications attract between one and two Million hits per day.

## References

Steinberger Ralf, Bruno Pouliquen & Camelia Ignat (2008). *Using language-independent rules to achieve high multilinguality in text mining*. In: Françoise Fogelman-Soulié, Domenico Perrotta, Jakub Piskorski & Ralf Steinberger (eds): Mining Massive Data Sets for Security. IOS-Press, Amsterdam, Holland.

Pouliquen Bruno & Ralf Steinberger (2009). *Automatic Construction of Multilingual Name Dictionaries*. In: Cyril Goutte, Nicola Cancedda, Marc Dymetman & George Foster (eds.): Learning Machine Translation. MIT Press, NIPS series.

Pouliquen Bruno, Ralf Steinberger & Camelia Ignat (2003). *Automatic Annotation of Multilingual Text Collections with a Conceptual Thesaurus*. In: Proceedings of the EUROLAN Workshop Ontologies and Information Extraction at the Summer School The Semantic Web and Language Technology - Its Potential and Practicalities. Bucharest, Romania.

Pouliquen Bruno, Marco Kimler, Ralf Steinberger, Camelia Ignat, Tamara Oellinger, Ken Blackler, Flavio Fuart, Wajdi Zaghouani, Anna Widiger, Ann-Charlotte Forslund, Clive Best (2006). *Geocoding multilingual texts: Recognition, Disambiguation and Visualisation*. Proceedings of the 5th International Conference on Language Resources and Evaluation LREC, pp. 53-58. Genoa, Italy.

Pouliquen Bruno, Ralf Steinberger & Clive Best (2007). *Automatic Detection of Quotations in Multilingual News*. In: Proceedings of the International Conference Recent Advances in Natural Language Processing RANLP. Borovets, Bulgaria.

Pouliquen Bruno, Olivier Deguernel & Ralf Steinberger (2008). *Story tracking: linking similar news over time and across languages*. In: Proceedings of the CoLing'2008 workshop: Multi-source, multilingual information extraction and summarization, Manchester, August 2008.
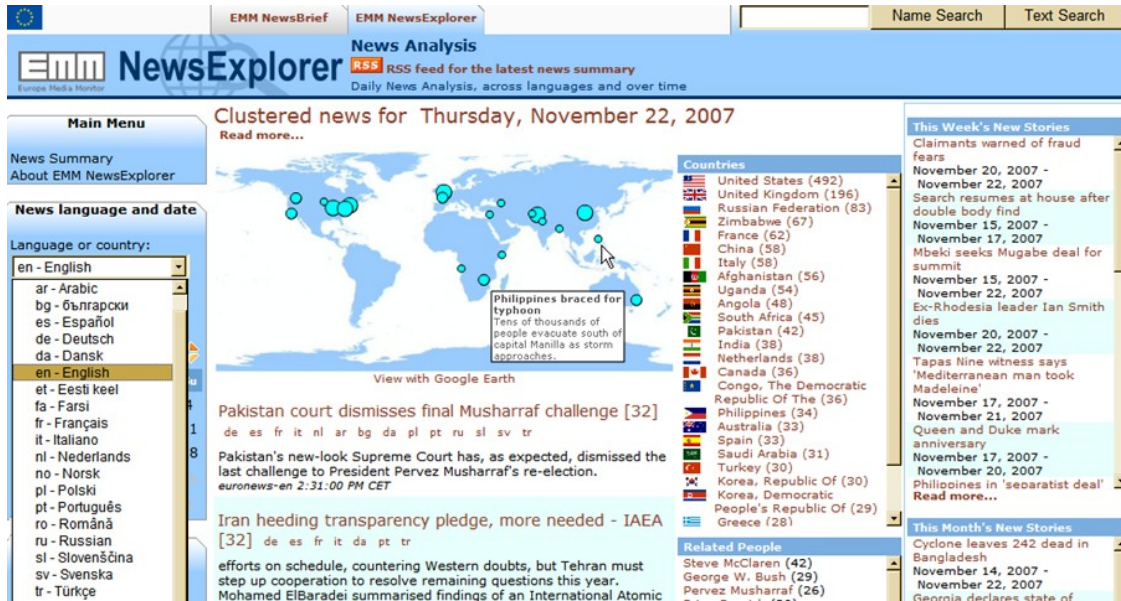
**Figure 1.** Screenshot of NewsExplorer, showing a map with the location of today's news, the largest English language news clusters, links to related news in the other 18 languages, lists of countries, people and organizations mentioned in the news that day, and lists of the biggest 'stories' (daily news clusters linked over time) this week, month and year.
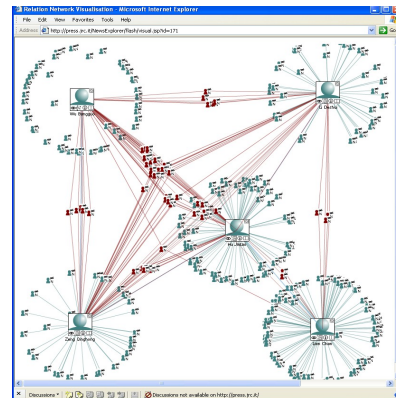


**Figure 2.** Screenshot of (part of) NewsExplorer's page on Pakistani President Pervez Musharraf. That page shows automatically collected name variants (including variants in different scripts such as Arabic, Farsi, and Russian), titles, lists of related persons, lists of related news clusters, quotations by and about Musharraf, multi-day 'stories' in which he is mentioned, and more. Relations between two or more persons can be visualised graphically.