# Challenges and solutions for multilingual text mining

lrec2010
malta
17-23 may

Ralf Steinberger

http://langtech.jrc.ec.europa.eu/
http://press.jrc.it/overview.html
Ralf.Steinberger@jrc.ec.europa.eu

# What this presentation is about

- Mostly Information Extraction, but not only
- Mostly rule-based approaches, but not only

- Show the benefits of – and the need for – multilingual text processing.

- Challenge: it is an enormous effort to develop these tools.

  - For N languages, N times the effort of developing tools for one language?

- Question: Are there ways to minimise this effort?

  - Literature review
  - Our own insights

# Acknowledgements    (alphabetical order)

## JRC colleagues (incl. former):

- Martin Atkinson
- Maud Ehrmann
- Flavio Fuart
- Erik van der Goot
- Camelia Ignat (now ECHA)
- Mijail Kabadjov
- **Bruno Pouliquen** (now WIPO)
- Hristo Tanev
- Vanni Zavarella
- …

## Multilingual system developers:

- Kalina Bontcheva (Sheffield University);
- Khalid Choukri (ELRA/ELDA);
- Gregory Grefenstette (Exalead);
- Frédérique Segond, Caroline Hagège and Claude Roux (Xerox Research Centre Europe);
- Aarne Ranta (Gothenburg University);
- Gregor Thurmair (Linguatec);
- Jacques Vergne (Caen University);
- Eric Wehrli (Geneva University).

# Thank you !

# Agenda

- European Commission – Joint Research Centre:
  - Who we are
  - Our customers and users (motivation)
- Importance of multilinguality
- Examples of multilingual EMM functionality
  - Multilingual media monitoring (publicly accessible at http://press.jrc.it/overview.html )



- **How to minimise the effort of producing multilingual applications?**
- Language resources that would facilitate the development of highly multilingual applications.

# Joint Research Centre - Who we are

**BRUSSELS (BE)**
The Directorate General (**DG**)
The Institutional and Scientific Relations Directorate (**ISR**)
The Programme and Resource Management Directorate (**PRM**)

**GEEL (BE)**
The Institute for Reference Materials and Measurements (**IRMM**)

**KARLSRUHE (DE)**
The Institute for Transuranium Elements (**ITU**)

**ISPRA (IT)** Download the Ispra site Brochure (English - Italian)
The Institute for the Protection and Security of the Citizen (**IPSC**)
The Institute for Environment and Sustainability (**IES**)
The Institute for Health and Consumer Protection (**IHCP**)
The Ispra site Directorate (**IS**)
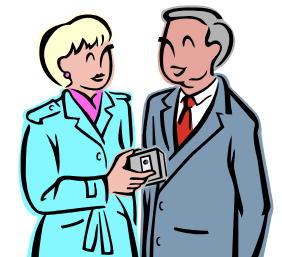
**PETTEN (NL)**
The Institute for Energy (**IE**)

**SEVILLE (E)**
The Institute for Prospective Technological Studies (**IPTS**)

- European Commission
  (scientific-technical arm of public administration)
- Non-commercial
- Relatively small team working on Language Technology

- **European Commission** (most DGs) and other EU Institutions
- **EU Agencies**:
  - e.g. Public Health (ECDC), **Food** Safety (EFSA), **Chemicals** Bureau (ECHA), etc.
- **EU Member State organisations**: e.g.
  - Public **Health**,
  - **law enforcement** authorities,
  - **parliaments**,
  - crisis management/**humanitarian**
- **International and extra-European organisations**: e.g.
  - various UN organisations
  - Centres for **Disease** Prevention and Control in the **US**, **Canada**, **China**, …
- **The public**:
  - Ca. 30,000 anonymous **internet** users of publicly accessible EMM systems.
  - Combined between 1 and 2 Million hits per day

# Agenda

- European Commission – Joint Research Centre:
  - Who we are
  - Our customers and users (motivation)
- **Importance of multilinguality**
- Examples of multilingual EMM functionality
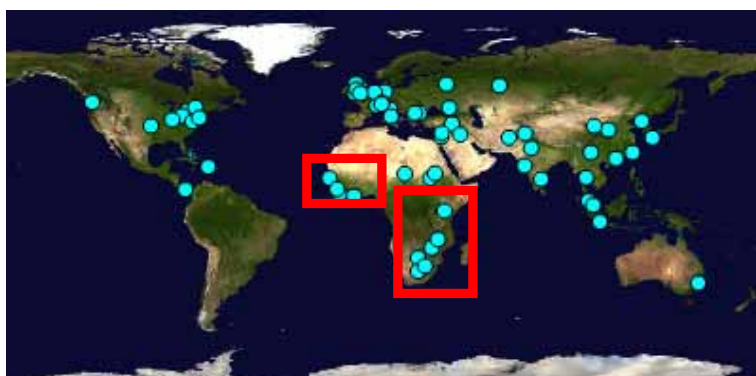  - Multilingual media monitoring (publicly accessible at http://press.jrc.it/overview.html )



- How to minimise the effort of producing multilingual applications?
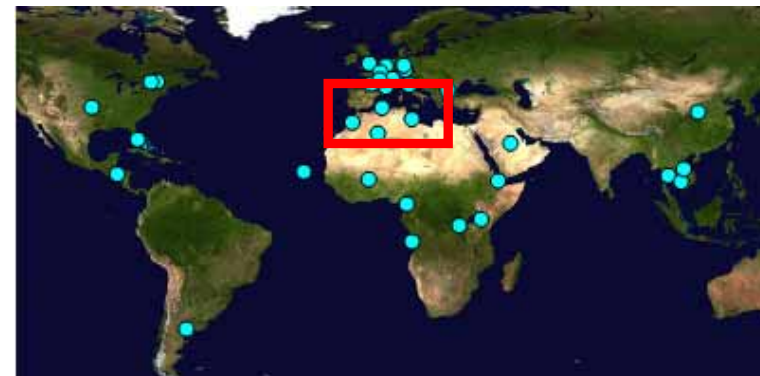- Language resources that would facilitate the development of highly multilingual applications.

# Multilinguality: complementary news coverage in various languages

## Locations mentioned in MedISys medical articles across languages – complementary coverage



Italian  -  German

English  -  French

Spanish -  Portuguese

# Multilinguality: Gathering *more information* about people

## Alexander Litvinenko
### Information about this person

| Names | Key Titles and Phrases | External resources |
|---|---|---|
| Alexander Litvinenko (Eu,nl) | russo (it,pt - 349) | |
| Alexander Litwinenko (de) | agent russe (fr - 134) | |
| Alexandre Litvinenko (fr) | ruso (es - 208) | |
| Aleksandr Litvinenko (fi,no) | agenten (de,sv - 134) | |
| Aleksander Litvinenko (nl,sv) | kritikers (de - 79) | |
| Александра Литвиненко (ru) | agent (en,sv - 130) | |
| Александр Литвиненко (ru) | russa (it,pt - 76) | |
| Alexander Livtinenko (it) | agent secret russe (fr - 39) | |
| Alexander V. Litvinenko (en) | russe (de,fr - 73) | |
| Alexandr Litvinenko (it) | former russian agent (en - 20) ← | |
| Alexander Litvineko (es) | morte di (it - 45) | |
| Alexandre Livinenko (fr) | ryske agenten (sv - 13) | |
| Alexander Litveneneko (en) | kritiker (de - 19) ← | |
| 亞歷山大·利特維年科 (zh) | 43 ans (fr - 17) | |
| Oleksandr Lytvynenko (en) | russi (it - 14) | |
| Olexandre Litvinenko (fr) | russian (en - 15) | |
| Aleksandar Litvinjenko (hr) | omicidio di (it - 11) ← | |
| Alexander Litvinenk (it) | officer (en - 13) ← | |
| アレクサンダー・リトビネンコ (ja) | former (en - 16) | |
| Alexander Walterowitsch Litwinenko (de) | | |

*Image obtained automatically from Wikipedia*

Read Wikipedia entry

Steinberger Ralf & Bruno Pouliquen (2009). **Cross-lingual Named Entity Recognition.** In: Satoshi Sekine & Elisabete Ranchhod (eds.): Named Entities - Recognition, Classification and Use, Benjamins Current Topics, Volume 19, pp. 137-164. John Benjamins Publishing Company.

# Multilinguality: less news bias and more transparency

# Multilinguality: More information about relations between people

Social networks produced on the basis of many languages are **more complete** and **less biased**.



**Associated People**

Salman Bashir (1.9)
Джон Негропонте (1.5)
Nawaz Sharif (1.2)
Раджа Омар Хатаб (1.2)
Мухоммада Али Джинны (1.2)
Зульфикара Али Бхутто (1.2)
Iftikhar Muhammad Chaudhry (1.1)
Christian College (1.1)
Tariq Azeem (1.1)
Benazir Bhutto (1.1)
Malik Mohammad Qayyum (1.0)
Chaudhry Shujaat Hussain (1.0)
Furqan Bahadur (1.0)
Javed Cheema (0.9)
Shaukat Aziz (0.9)
Abdul Rashid Ghazi (0.9)
Amir Mir Lahore (0.9)
Amin Fahim (0.9)
Гордон Джонроу (0.9)
Wajihuddin Ahmed (0.9)
Mohammed Ali Durrani (0.9)
Rashid Qureshi (0.9)
Qazi Hussain Ahmed (0.9)



live

Pouliquen Bruno, Ralf Steinberger, Jenya Belyaeva (2007). **Multilingual multi-document continuously updated social networks**. Proceedings of the Workshop *Multi-source Multilingual Information Extraction and Summarization* (**MMIES'2007**) held at **RANLP'2007**, pp. 25-32. Borovets, Bulgaria, 26 September 2007.(**PDF**)

Hristo Tanev (2007). **Unsupervised Learning of Social Networks from a Multiple-Source News Corpus**. Proceedings of the Workshop *Multi-source Multilingual Information Extraction and Summarization* (**MMIES'2007**) held at **RANLP'2007**, pp. 33-40. Borovets, Bulgaria, 26 September 2007. (**PDF**)

# Language coverage of various media analysis tools (March 2010)

The seventh international conference on Language Resources and Evaluation, LREC 2010, Valletta, Malta, 19-21 May 2010          12

## News aggregators



in early 2008: 34, 17, 43 languages

## News analysis systems



in early 2008: the same

# Agenda

- European Commission – Joint Research Centre:
  - Who we are
  - Our customers and users (motivation)
- Importance of multilinguality
- **Examples of multilingual EMM functionality**
  - Multilingual media monitoring (publicly accessible at  http://press.jrc.it/overview.html )



- How to minimise the effort of producing multilingual applications?
- Language resources that would facilitate the development of highly multilingual applications.

# What we do:  1. News gathering

The seventh international conference on Language Resources and Evaluation, LREC 2010, Valletta, Malta, 19-21 May 2010                    14

- **EMM news gathering engine**
  - Monitors ~ 2,500 news sources
  - Gathers  ~100,000  news articles per day
  - Clusters and categorises news
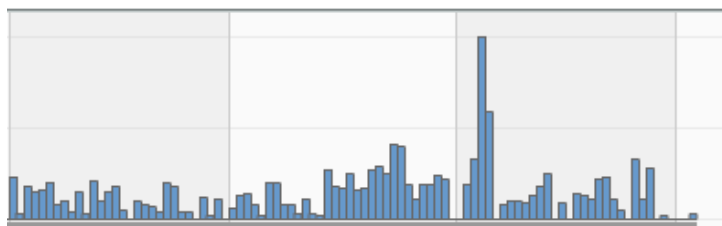  - In 50 languages
  - Feeds news into the public media monitoring applications





Geographical distribution of EMM news sources

Steinberger Ralf, Bruno Pouliquen & Erik van der Goot (2009). **An Introduction to the Europe Media Monitor Family of Applications**. In: Fredric Gey, Noriko Kando & Jussi Karlgren (eds.): Information Access in a Multilingual World - Proceedings of the SIGIR 2009 Workshop (SIGIR-CLIR'2009), pp. 1-8. Boston, USA. 23 July 2009.

# What we do: 2. Deeper news analysis (~20 languages)

- Breaking news detection; alerting; **tracking topics** over time;

- **Named entity recognition** and disambiguation (persons, organisations, locations);

- **Name variant** matching;

- **Quotation** recognition;

- **Social network** generation;

- Multi-label categorisation;

- Linking related clusters across languages;

- **Event scenario** template filling (6 languages);

- …

Rice [-said]: I would many times over liberate Iraq again from Saddam Hussein, CBSnews 19-MAR-10

Abdul Rahman [ said]: Facts that were created after the overthrow of the Saddam Hussein regime will not be easy to maintain because there will be no US umbrella, guardian 04-MAR-10

# NewsExplorer – Multilingual daily news overview

# NewsExplorer – Aggregation of clusters into longer 'stories'

# Detection and visualisation of events (violence/disasters/humanitarian/...)

**Objective**: global crisis monitoring.     **Languages**: En, Fr, Es, It, Pt, Ru + (Ar)



Atkinson Martin, Jakub Piskorski, Bruno Pouliquen, Ralf Steinberger, Hristo Tanev & Vanni Zavarella (2008). **Online-monitoring of security-related events**. In Proceedings of the 22nd International Conference on Computational Linguistics (**CoLing'2008**). Manchester, UK, 18-22 August 2008. (PDF)

# Detection and visualisation of events (violence/disasters/humanitarian/...)

EMM-Labs

**В Турции террористы захватили 15 заложников**
Turkey Lat/Lon: 39.9293,32.8533 Size: 18

Last update 2009-05-13T17:35+0200
Двое неизвестных в масках совершили в среду вооруженный налет на филиал банка в турецком курортном городе Кушадасы на побережье Эгейского...
Event Type:Kidnapping/Hostage Taking. Severity :0 Killed, 0 Injured and 15 .Victims were no one killed and no one injured.

На турецком курорте грабители взяли в заложники посетителей банка

**Освобождены захваченные в турецком банке заложники**
news-open 13.05.2009 17:40:00 *i*
Перевести:[ar] [bg] [zh] [hr] [cs] [da] [nl] [en] [fi] [fr] [de] [el] [hi] [it] [ja] [ko] [no] [pl] [pt] [ro] [es] [sv]
Турецкая полиция обезвредила налетчика, захватившего в заложники более 10 человек в банке в курортном городе Кушадасы на побережье Эгейского моря, передают местные СМИ. Полиция пошла на штурм банка....

**Сдался полиции преступник, захвативший 10 заложников в банке Кушадасы**
news-meta 13.05.2009 17:21:00 *i*
Турецкая полиция обезвредила налетчика, захватившего в заложники более 10 человек в банке в курортном городе Кушадасы на побережье Эгейского моря,......

**Полиция Турции обезвредила налетчика, захватившего заложников в банке**
rian 13.05.2009 16:54:00 *i*
Драма с заложниками продолжалась более шести часов - с 11.00 местного времени (12.00 мск). По сообщению телеканалов NTV и CNN-Turk, полиции удалось убедить налетчика сдаться. Подробности проведенной операции пока неизвестны....

**Google**™   This page was automatically translated from Russian.
View original web page or mouse over text to view original language.

// World news // Wednesday, May 13, 2009          Archive   Search

**In Turkey, an unknown masked seized 13 hostages in the bank**

publication time: 14:29                                    Images print download submit
Last Updated: 19:08

The police stormed the bank in the western Turkish city of Kusadasi, where the hostages were 13 people. Robbers arrested, ITAR-TASS reported with reference to state television.

# Agenda

- European Commission – Joint Research Centre:
  - Who we are
  - Our customers and users (motivation)
- Importance of multilinguality
- Examples of multilingual EMM functionality
  - Multilingual media monitoring (publicly accessible at  http://press.jrc.it/overview.html )



- **How to minimise the effort of producing multilingual applications?**
- Language resources that would facilitate the development of highly multilingual applications.

# Insights collected from various teams

Effort for N languages = N times effort for one language? How to save effort?

- Depends on required level of analysis

- Complex applications may require deeper linguistic analysis

- But: even simple means can take you relatively far

- Recognition easier than generation (e.g. agreement)

- Typical and most common approach for rule-based systems:
  - Develop in one language
  - Reuse resources and adapt to new languages
  - E.g. Gamon et al. (1997); Rayner & Bouillon (1996), Pastra et el. (2002); Carenini et al. (2007); Maynard et al. (2003)

# Insights collected from various teams; Guidelines; Ideas (1)

1. Use **Unicode** (Maynard et al. 2002)

2. Use **virtual keyboards** to enter foreign language data (Maynard et al. 2002)

3. **Modularity** (Pastra et al. 2002; Maynard et al. 2002)

# Insights collected from various teams (2)

The seventh international conference on Language Resources and Evaluation, LREC 2010, Valletta, Malta, 19-21 May 2010      24

4. **Shared token classes**, ideally based on **surface features** (Bering et al. 2003)
   - E.g. case information, includes-number, includes hyphen; string length

5. **Uniform input and output** structures (Carenini et al. 2007; Bering et al. 2003)
   - Data format
   - When possible: same part-of-speech classes, same grammatical categories

```
The_DD  N-terminal_JJ region_NN had_VHD high_JJ
P12571010A13
Several_JJ sequences_NNS were_VBD identified_V
P12576309A07
Our_PNG findings_NNS indicate_VVB that_CST CRC
P12582233A05
The_DD corresponding_VVGJ mRNA_NN of_II 3.5_MC
P12586375A07
A_DD few_JJ examples_NNS of_II heterologous_JJ
```

| The | DT | the |
| TreeTagger | NP | TreeTagger |
| is | VBZ | be |
| easy | JJ | easy |
| to | TO | to |
| use | VB | use |
| . | SENT | . |

# Insights collected from various teams (3)

The seventh international conference on Language Resources and Evaluation, LREC 2010, Valletta, Malta, 19-21 May 2010    25

6. **Simplicity of rules and the lexicon** (Carenini et al. 2007; Vergne 2002)

- E.g. Identification of subject-verb pairs for 5 languages (Vergne 2002)
  - <200 dictionary elements per language
  - Case information
  - Regular expressions matching certain combinations of word endings

- E.g. Chunker for 23 languages (Vergne 2009), using only
  - String length information
  - Word frequency

- E.g. under-specification (agreement; order of modifiers) in recognition of quotation and speakers (Pouliquen et al. 2007)
  - "…" said the <u>former</u>  <u>56-year-old</u>  <u>British</u>  <u>Prime Minister</u>  *Tony Blair*

# Insights collected from various teams (4)

7.  **Share resources between languages** (lexica, gazetteers, grammar rules)
    (Bering et al. 2003)

    - Language-independent rule set + language-specific rules, e.g. for date recognition
        - 20.10.2010   (generic)
        - 20[th] of October of the year 2010   (language-specific)
          (Ignat et al. 2003 also cover this and more cases with generic rules)

8.  **Use theory-neutral date types** (Pastra et al. 2002; Maynard et al. 2002)
    - For the Language Engineering architecture GATE

# Insights collected from various teams (5)

9. **Adhere to a grammar theory**, e.g.

- Bender & Flickinger (2005) use **HPSG** for general-purpose grammars;
- Gamon et al. (1997) adhere to **Universal Grammar** for generic Microsoft NLP grammar;
- Wehrli (2007) uses **Chomsky's generative grammar** to build parsers;
- Ranta (2009, e.g. pp. 47ff; LREC tutorial) works within the **Grammatical Framework** GF.

## Benefits mentioned:

- The mere existence of an abstract syntax implies grammar sharing;
- Creating a generic grammar and parameterise it to handle many languages;
- Shared grammars by language group;
- Treating some linguistic phenomena in a systematic way (e.g. clitics; morphological agreement; phrase ordering; …);
- Generating starter grammars based on linguistic features (Bender & Flickinger 2005).

# The promising contribution of Machine Learning (ML)

- Idea: self-learning software learns rules and vocabulary, e.g. from examples

- Success story: **Statistical Machine Translation** (SMT);
  - Google: 57 languages (incl. 6 alpha versions); 1596 language pairs (May 2010).

- ML for **Named Entity Recognition** (NER) (Nadeau & Sekine 2009)
  - Supervised ML: train on previously annotated corpora
    - Challenge: Corpora annotation is labour-intensive and expensive
  - Semi-supervised, weakly-supervised ML:
    - Use set of human-provided seeds to start learning process
    - Use boot-strapping to increase number of patterns and resources
  - Unsupervised ML:
    - E.g. words appearing synchronously in news articles (Shinyama & Sekine 2004).

# Machine Learning versus manual intervention

- Issue: how to combine ML methods with manual intervention, e.g. to correct
  - E.g. Output of SVM or HMM is difficult to modify manually.

→ EMM approach:
  - Use hand-crafted rules;
  - Use external knowledge sources (dictionaries) when available;
  - Use bootstrapping and ML methods to enhance these dictionaries;
  - Empirical testing.

- Benefit:
  - Keep control over the recognition performance
  - Spend less time per language than when using ML: 1 week – 3 months per language to
    - add news sources,
    - translate Boolean category definitions
    - Add linguistic IE resources (to recognise persons, organisations, locations, quotations, dates)
    - Evaluation

# EMM guidelines

- Keep applications simple!
  - Under-specify (constraints are time-consuming to produce and may hinder you in other languages)
  - Use bags of words without specifying agreement and order, if possible
  - Don't disambiguate if you can avoid it
  - No grammar theory

- Use language-independent rules, if possible
- Use as little language-specific resources as possible
  (POS taggers, parsers, dictionaries, …)
- Modularity: keep language-specific resources outside the rules
  → plug in any new language
- Do not use language pair-specific resources
  - NewsExplorer covers 20 languages, 190 language pairs

Times Square evacuated as police defuse
car bomb [50]  de es fr it nl ar bg da et
fa no pl pt ro ru sl sv tr

Steinberger Ralf, Bruno Pouliquen & Camelia Ignat (2008). **Using language-independent rules to achieve high multilinguality in Text Mining**. In: Fogelman-Soulié Françoise, Domenico Perrotta, Jakub Piskorski & Ralf Steinberger (eds.): Mining Massive Data Sets for Security. pp. 217-240. IOS Press,

# Concretely, …

- **<u>What we (JRC) cannot use:</u>**
  - Syntactic parsers
  - Part-of-speech taggers
  - Full dictionary for any of the languages

- <u>What we do use:</u>
  - Targeted world lists
    - Name titles, etc.
    - Gazetteers of place names and locations
    - Sentiment words
    - Reporting verbs
    - Stop word lists!!!
  - Light-weight suffix stripping rules and generation of morphological variants
  - Boolean combinations of category-defining words (for document classification)
  - The output of our own NER tools
  - Machine Learning and boot-strapping

# Agenda

- European Commission – Joint Research Centre:
  - Who we are
  - Our customers and users (motivation)
- Importance of multilinguality
- Examples of multilingual EMM functionality
  - Multilingual media monitoring (publicly accessible at http://press.jrc.it/overview.html )



- How to minimise the effort of producing multilingual applications?
- Language resources that would facilitate the development of highly multilingual applications.

# EMM insight: Need for multilingual resources

→ Need for **uniform highly multilingual resources** for information extraction, etc.
  Similar to Multext and Multext-East, EuroWordNet, …

- *Parallel* lexical resources, e.g.
  - Multilingual **gazetteers** for geo-coding (e.g. GeoNames)
  - Multilingual **dictionaries** using the same features and format (MulText, Euro-WordNet)
- Multilingual **taggers and parsers** producing the same output type and output format.
- Annotated **multilingual parallel text corpora** (both for training and for testing) (e.g. JRC-Acquis)
- **Single access point for licensing** issues (e.g. ELRA / LDC)
- Ideally **freely available** (see also: G. Grefenstette's recommendation at recent FLaReNet workshop: http://www.flarenet.eu/?q=node/347)

| Surface form | Lemma | POS | translation |
|---|---|---|---|
| étrangères | étranger | ADJ | foreign, stranger |
| libération | libération | N | liberation |
| mauvaise | mauvais | ADJ | bad |
| avantage | avantage | N | advantage |
| représentent | répresenter | V | represent |

# Summary

# Summary

- Benefits of multilinguality: capture complementary coverage across languages
  - Contents
  - Opinions

- Major challenge: time needed to develop resources and applications

- Proposals how to keep the effort down:

  - Unicode
  - Virtual keyboards
  - Modularity
  - Shared token classes (surface features)
  - Uniform input and output structures
  - Simplicity of rules and the lexicon
  - Share resources between languages

  - Theory-neutral / adhere to grammar theory
  - Use Machine Learning
  - Under-specification
  - Minimise use of language-specific resources
  - Avoid language pair-specific resources

# Summary (2)

- Useful language resources to achieve high multilinguality
    - *Uniform and parallel* dictionaries, corpora and tools.

    - Surely, there are more unpublished insights → please share them with me.

        Ralf.Steinberger@jrc.ec.europa.eu

- Different applications require different means and there are many ways of doing things.

# Challenges and solutions for multilingual text mining



Ralf Steinberger

http://langtech.jrc.ec.europa.eu/
http://press.jrc.it/overview.html
Ralf.Steinberger@jrc.ec.europa.eu

# Summary of our approach

- MMDSS paper (2008) for details of the ideas, description of 9 text mining applications, limitations and language-specific exceptions:

> Steinberger Ralf, Bruno Pouliquen & Camelia Ignat (2008). **Using language-independent rules to achieve high multilinguality in Text Mining**. In: Fogelman-Soulié Françoise, Domenico Perrotta, Jakub Piskorski & Ralf Steinberger (eds.): Mining Massive Data Sets for Security. pp. 217-240. **IOS Press**, Amsterdam, The Netherlands.

- LRE journal article (submitted): example of the effort required to add a new language (African Bantu language Swahili):

> Steinberger Ralf, Sylvia Ombuya, Mijail Kabadjov, Bruno Pouliquen, Leonida Della Rocca, Jenya Belyaeva, Monica De Paola & Erik van der Goot (submitted). **Expanding a multilingual media monitoring and information extraction tool to a new language: Swahili.**

# Agenda

- European Commission – Joint Research Centre:
  - Who we are
  - Our customers and users (motivation)
- Importance of multilinguality
- Examples of multilingual EMM functionality
  - Multilingual media monitoring (publicly accessible at  http://press.jrc.it/overview.html )



- How to minimise the effort of producing multilingual applications?
- What do these guidelines mean for real-life applications? Sample solutions.
- Language resources that would facilitate the development of highly multilingual applications.

# Example 1: Quotation recognition

- ## <u>Motivation:</u> detect and display quotations by and about entities
  - Incl. usage for quotation network, sentiment analysis, etc.

- ## Solution should be applicable for many languages (currently 13)
  - Not considering syntax and part-of-speech

## <u>Proposed solution:</u>
- Write simple patterns, that use
  - Previously extracted information on named entities and their titles
  - Simple vocabulary lists  (reporting verbs, modifiers, determiners)
  - Lists of quotation marks

## <u>Sample pattern:</u>

- "QUOTE" [,] *reporting-verb* [*modifier*] [*determiner*] [*title*]  *name*
  e.g. **"blah blah"**, *said again the journalist* **John Smith.**



**NewsExplorer**

**Quotes from - English**

[bg] [fr] [pt] [es] [it] [de] [ru] [sv] [nl]

*[said:]:* In addition, we will seek to ratify the Comprehensive Test Ban Treaty and negotiate a treaty to end the production of fissile material for use in nuclear weapons. (Reuters) Published:

**Quotes about - English**

[tr] [ru] [es] [pt] [it] [fr] [nl] [bg] [ro] [de]

*Hillary Clinton [ said]:* Both President Obama and I have made clear, both last year and again this year, that we do not believe any action by the Congress is appropriate, and we oppose it, *voanews 05-MAR-10*

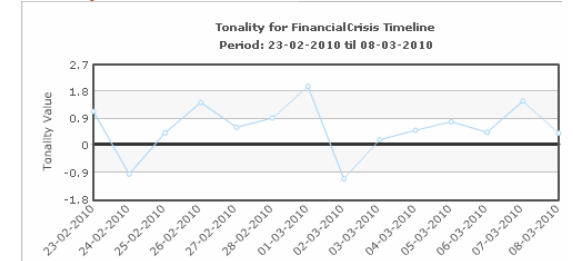Pouliquen Bruno, Ralf Steinberger & Clive Best (2007). **Automatic Detection of Quotations in Multilingual News**. In: Proceedings of the International Conference Recent Advances in Natural Language Processing (**RANLP'2007**), pp. 487-492. Borovets, Bulgaria, 27-29.09.2007.

# Example 2: Sentiment analysis (opinion mining)

- **Motivation**: detect and display opinion on entities / themes
  - Incl. changes over time
  - Incl. differences across languages / sources
- Solution should be applicable to ~20 languages:
  - Not considering syntax and part-of-speech
- **Proposed solution**:
  - Use simple sentiment word lists
  - Occurring within a word window around the entity, but
  - Separating good/bad news from pos./neg. sentiment
    - By not considering category-defining words, or
    - By setting average sentiment for news category as baseline
  - Challenge: create multilingual sentiment dictionaries







Balahur-Dobrescu Alexandra & Ralf Steinberger (2009). **Rethinking sentiment analysis in the news: from theory to practice and back**. 'Workshop on Opinion Mining and Sentiment Analysis' (**WOMSA**), held at the 2009 **CAEPIA-TTIA** 13th Conference of the Spanish Association for Artificial Intelligence, pp. 1-12. Sevilla, Spain, 13.11.2009.

# Example 3: Multi-document summarisation

- <u>Motivation:</u> summarise redundant news data (clusters)
- Incl. update summaries (daily and live – every 10 minutes)
- Solution should be applicable to ~20 languages:
  - Extractive summarisation
  - Not considering syntax

- <u>Proposed solution:</u>
  - Use Lexical Semantic Analysis (LSA) for sentence selection (similar to Gong & Liu, 2002)
  - Using as LSA input: word n-grams + entity mentions + non-disambiguated terms from the multilingual MeSH thesaurus



**Wednesday, January 13, 2010**

Thousands feared dead after Haiti quake [107]   de  es  fr  it  nl  bg  da  et  fa  no  pl  pt
ro  ru  sl  sv  tr

There are growing fears that thousands of people were killed in a massive earthquake which struck the Caribbean state of Haiti.
*RTERadio 10:42:00 AM CET*

Timeline
Story Tracking: Quake-hit Haiti cries for help

Steinberger Josef, Mijail Kabadjov, Bruno Pouliquen, Ralf Steinberger & Massimo Poesio (2009). **WB-JRC-UT's Participation in TAC 2009: Update Summarization and AESOP Tasks.** In: Proceedings of the Text Analysis Conference 2009 (TAC'2009). National Institute of Standards and Technology, Gaithersburg, Maryland USA, 16-17 November 2009.

# Example 4: Dealing with inflection in lookup



**Names**

Tony Blair (Eu,yo)
توني بلير (ar)
Тони Блэр (ru)
Тони Блеър (bg)
Tonyja Blaira (sl)
توني بلر (fa)
Tonyjem Blairom (sl)
Tony Blairi (et)
طوني بلير (ar)
Toni Blair (de,sl)
Tonyju Blairu (sl)
Anthony Charles Lynton Blair (da,sl)
Tony Blaira (pl,sl)
Tonny Blair (en,pt)
Тони Блер (mk,sr)
Tony Blairile (et)
Tony BLAIR (eo)
□□□□ □□□□□□□ (ml)
Toy Blair (de)
Tonijs Blērs (lv)

- **Lookup of known names** from database
  - Currently over 1,200,000 names and known variants

  - Pre-generate morphological variants (Slovene example):

Tony(a|o|u|om|em|m|ju|jem|ja)?\s+Blair(a|o|u|om|em|m|ju|jem|ja)

Tony Blair / Tonyju Blairu / Tony Blairt / Tonijs Blērs / Тони Блэра / توني بلير / Tony Blairi / طوني بلير

Steinberger Ralf & Bruno Pouliquen (2009). **Cross-lingual Named Entity Recognition.** In: Satoshi Sekine & Elisabete Ranchhod (eds.): Named Entities - Recognition, Classification and Use, Benjamins Current Topics, Volume 19, pp. 137-164. John Benjamins Publishing Company.