

DELOS
Network of Excellence in Digital Libraries

CLEF

GeoCLEF 2007: the CLEF 2007 Cross-Language Geographic Information Retrieval Track Overview

Thomas Mandl¹, Fredric Gey¹, Giorgio Di Nunzio⁶, Nicola Ferro⁶,
Ray Larson², Mark Sanderson³, Diana Santos⁴, Christa Womser-Hacker³, Xing Xie⁵

¹University of California, Berkeley, CA, USA gey@berkeley.edu, ray@ims.berkeley.edu
²Department of Information Studies, University of Sheffield, Sheffield, UK, m.sanderson@sheffield.ac.uk
³Information Science, University of Hildesheim, GERMANY, mandl@womser@uni-hildesheim.de
⁴Linguatca, SINTEF ICT, NORWAY, Diana.Santos@sintef.no
⁵Microsoft Research Asia, Beijing, China, Xingx@microsoft.com
⁶Department of Information Engineering, University of Padua, Italy, dinunzio@ferro@dei.unipd.it

CLEF 2007, Budapest, September 19-21, 2007

GeoCLEF Administration

Joint effort of

- Fredric Gey (U. California at Berkeley)
- Diana Santos (Linguatca, SINTEF ICT, Norway)
- Mark Sanderson (U. Sheffield)
- Nicola Ferro, Giorgio Di Nunzio (U. Padua)
- Xing Xie (Microsoft Research, Asia)
- Thomas Mandl, Christa Womser-Hacker (U. Hildesheim)

○ and many others ...

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 2

Overview

More on the geographic search task ...

- Thomas Mandl: Overview
 - Query Classification Subtask
- Mark Sanderson: *Topic Creation and Relevance Assessment*
- Giorgio Di Nunzio: *Results*
- Diana Santos: *Approaches and Interpretation*

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 3

Geographic Information Systems

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 4

Initial Aim of GeoCLEF

- Aim: to evaluate retrieval of multilingual documents with an emphasis on geographic search (GIR)
 - "find me news stories about riots near Dublin"

(Fred Gey @ CLEF Workshop 2005)

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 5

Participation

CLEF Year	2005	2006	2007
Nr. of Participants	11	17	13
Nr. of submitted Experiments	117	149	108

New: Query Classification Task
6 participants

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 6

Search Task 2007

- Three languages
- 600,000 + docs
- 25 topics (75 in three years now)

- Intention behind topics
 - geographically challenging

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 7

Reliability?

- 25 topics are sufficient under most circumstances to reliably order systems
(Sanderson & Zobel 2005)

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 8

Partial Swap Rate Analysis

Number of Topics	Average Correlation
1	0.25
2	0.40
5	0.70
10	0.85
15	0.90
20	0.95
25	1.00

Average Correlation between system rankings of full and partial topic set for German Mono-lingual task 2006

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 9

Search Task

- How much and which geo knowledge and reasoning is necessary?
- Each year, keyword based systems do well on the task

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 10

Query Classification Task

- Goal: find geo queries in a log of real queries
- New in 2007

- Organized by Xing Xie (Microsoft Research Asia, Beijing, China)

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 11

Data

- Query log from the MSN search engine
 - in English
 - 800.000 queries (collected August 2006)
 - 500 queries were labelled and used for evaluation
 - 100 queries for training
 - 400 for testing

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 12

Task

- Find queries with a geographic scope
 - Extract where component
 - Extract geo-relation-type
 - Extract what component
 - Classify what type {information, yellow page, map}

Example:

Lottery in Florida

<local>YES

<what>lottery

<what-type>information

<where>Florida, US

<geo-relation>in

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 13

Geo-relation-types

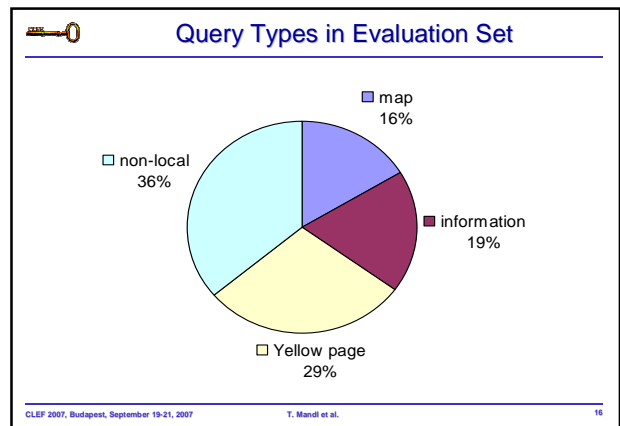
- 27 classes
- Examples:
 - In
 - On
 - Near
 - Along
 - Distance
 - North_of
 - North_west_of
 - North_to
 - ...

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 14

Evaluation Set

- Choose 800 queries randomly from the query set.
- Remove the typos and the ambiguous queries from the 800 ones manually.
- Select the queries with special geo-relations from the remainder queries in the query set manually and add them to the evaluation set.
- Select 500 queries for the final evaluation set.

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 15



Evaluation Metrics


- Three assessors
 - individually assessed all system answers
 - reached an agreement
- Fully Correct classified query instances
- Recall, precision and combined F1-Score

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 17

Results

Team	Precision	Recall	F1
Ask	0.625	0.258	0.365
Csum	0.201	0.197	0.199
Linguit	0.112	0.038	0.057
Miracle	0.428	0.566	0.488
Talp	0.222	0.249	0.235
Xldb	0.096	0.08	0.088


CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 18

 **Approaches**

- ❑ Gazeteers for location identification
 - Large base of geo names
- ❑ Pre-defined Rules


- ❑ Issues
 - Low Performance
 - Few training classes for many geo-types

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 19

 **geoCLEF topic creation**


Mark Sanderson

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 20

 **Topics**


- ❑ 25 adhoc topics
- ❑ Developed
 - One third in English
 - One third in German
 - One third in Portuguese

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 21

 **Classic CLEF topic creation**


- ❑ Developed topics in local language
- ❑ Other geoCLEF partners translated topics and checked for relevance in their collection.

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 22

 **Motivation behind topic design**


- ❑ Text only often wins
- ❑ How to tackle that
 - Imprecise regions
 - Regions surrounding, but not including a point
 - Regions where local knowledge is important

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 23

 **Imprecise regions**


- ❑ "Documents describing the damage caused by acid rain in the countries of northern Europe"

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 24

 Surrounding...


- "Find information about social problems afflicting places in greater Lisbon."
- "Find documents mentioning airplane crashes close to Russian cities"

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 25

 Local knowledge


- "To be relevant, a document must describe a whisky made, or a whisky distillery located, on a Scottish island."

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 26

 Problems


- Always hard to find topics that work well across languages

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 27

 Assessment - DIRECT


- Mostly volunteers
 - Hildesheim
 - SINTEF
 - Fred in Berkeley
- Some funding
 - Sheffield - Tripod project

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 28

 Participation (1/2)

Participant	Monolingual Tasks			Bilingual Tasks			TOTAL
	DE	EN	PT	X2DE	X2EN	X2PT	
catalunya		5					5
cheshire	1	1	1	3	3	3	12
csusm	6	6	5		4	4	25
depok*					6		6
groningen		5					5
hagen	5			5			10
hildesheim	4	4					8
icl		4					4
linguit*		4					4
moscow*		2					2
msasia		5					5
valencia		12					12
xldb		5	5				10
TOTAL	16	53	11	8	13	7	108

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 29

 Participation (2/2)

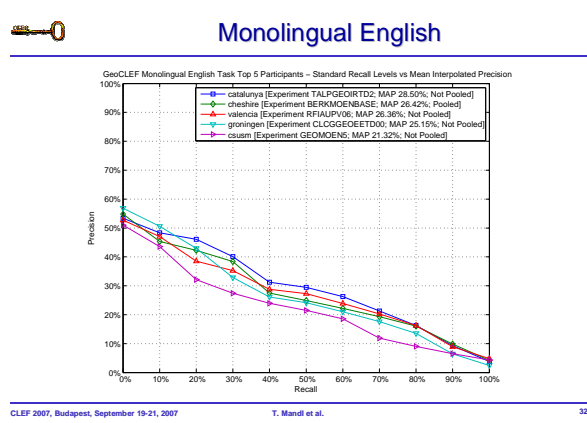
Track	Source Language					TOTAL
	DE	EN	ES	ID	PT	
Bilingual X2DE		6	1		1	8
Bilingual X2EN	1		5	6	1	13
Bilingual X2PT	1	1	5			7
TOTAL	2	7	11	6	2	28

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 30

Monolingual Tasks

- English
- German
- Portuguese

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 31



Monolingual English

Track	Rank	Part.	Experiment DOI	MAP
Monolingual English	1 st	catalunya	10.2415/GC-MONO-EN-CLEF2007.CATALUNYA.TALPGEORITD2	28.50%
	2 nd	cheshire	10.2415/GC-MONO-EN-CLEF2007.CHESHIRE.BERKMOENBASE	26.42%
	3 rd	valencia	10.2415/GC-MONO-EN-CLEF2007.VALENCIA.RPIAUPV06	26.36%
	4 th	groningen	10.2415/GC-MONO-EN-CLEF2007.GRONINGEN.CLOGGEOETD00	25.15%
	5 th	caism	10.2415/GC-MONO-EN-CLEF2007.CSUSM.GEOMDENS	21.32%
	Diff.			33.68%

POS, lemmas, named entities

graphical thesaurus

World Gazetteer

WordNet

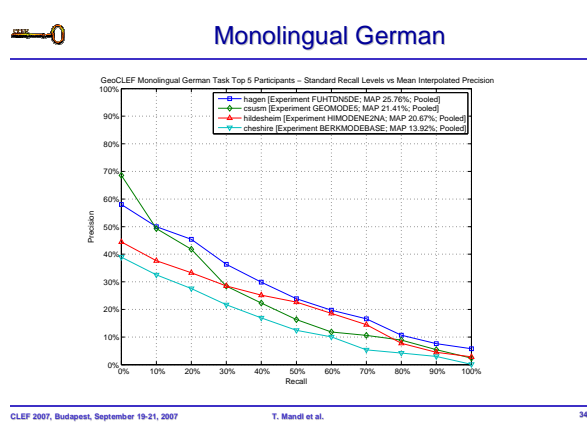
LingPipe

World Gazetteer, GeoNET, Wikipedia, WordNet

Ling Pipe

mixed txidf (Lucene)

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 33



Monolingual German

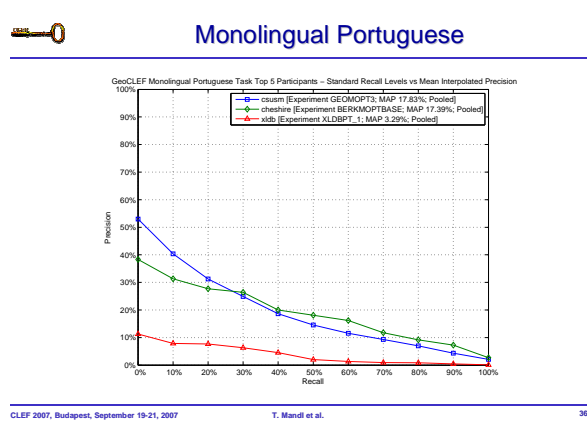
Track	Rank	Part.	Experiment DOI	MAP
Monolingual German	1 st	hagen	10.2415/GC-MONO-DE-CLEF2007.HAGEN.FUHTINDSE	25.76%
	2 nd	caism	10.2415/GC-MONO-DE-CLEF2007.CSUSM.GEOMD04	21.41%
	3 rd	hildesheim	10.2415/GC-MONO-DE-CLEF2007.HILDESHEIM.HIMODENE2NA	20.67%
	4 th	cheshire	10.2415/GC-MONO-DE-CLEF2007.CHESHIRE.BERKMOENBASE	13.92%
	5 th			

Semantic analysis for GIR txidf (Zebra DBMS)

GeoNet name server BM25, LM, InL2 (Terrier)

LingPipe txidf (Lucene)

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 35



Monolingual Portuguese

Track	Rank	Part.	Experiment DOI	MAP
Monolingual Portuguese	1 st	csusm	10.2415/GC-MONO-PT-CLEF2007_CSUSM_GEOMDPT3	17.83%
	2 nd	cheshire	10.2415/GC-MONO-PT-CLEF2007_CHESHIRE_BERKHOPTBASE	17.39%
	3 rd	xlkb	10.2415/GC-MONO-PT-CLEF2007_XLKB_XLDBPT_1	3.29%
	4 th			
	5 th			
	Diff.			441.95%

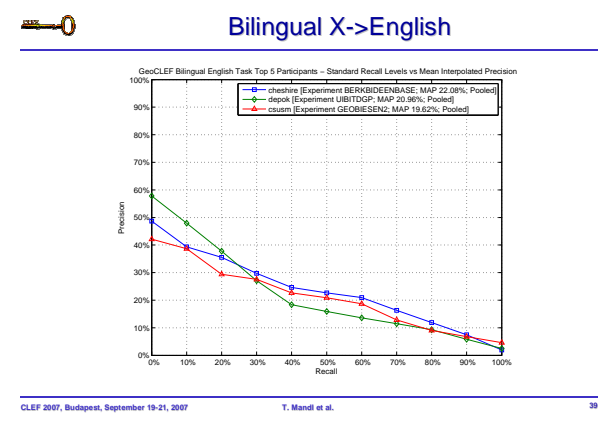
System with different modules
Geographic Ontology + Text Mining
MG4J + mixed Okapi BM25

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 37

Bilingual Tasks

- X -> English
- X -> German
- X -> Portuguese

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 38



Bilingual X->EN

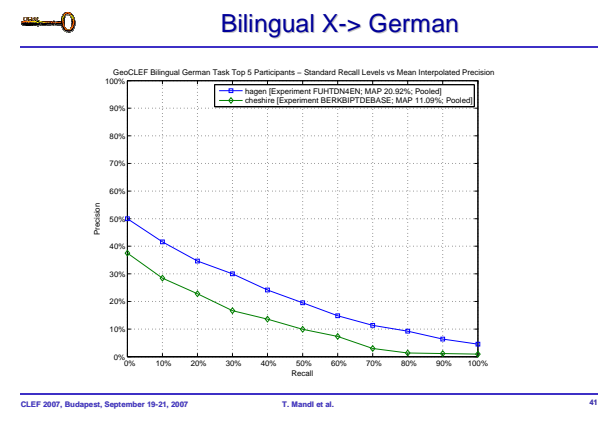
Track	Rank	Part.	Experiment DOI	MAP
Bilingual English	1 st	cheshire	10.2415/GC-BILI-XZEN-CLEF2007_CHESHIRE_BERKHOPTBASE	22.08%
	2 nd	depok*	10.2415/GC-BILI-XZEN-CLEF2007_DEPOK_UIBITDGP	20.96%
	3 rd	csusm	10.2415/GC-BILI-XZEN-CLEF2007_CSUSM_GEOBIEN2	19.62%
	4 th			
	5 th			
	Diff.			12.54%

LEC Power Translator
World Gazteer
LR

Free MT tool
Geonames + Wikipedia
LM (Lemur)

Spanish Toponymy from Euro Parl.
GeoNet name server
BM25, LM, InL2 (Terrier)

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 40



Bilingual X-> German

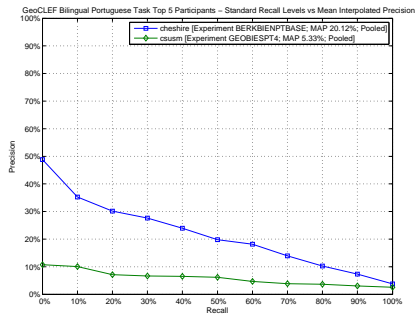
Track	Rank	Part.	Experiment DOI	MAP
Bilingual German	1 st	hagen	10.2415/GC-BILI-XZDE-CLEF2007_HAGN_FUHTDN4EN	20.92%
	2 nd	cheshire	10.2415/GC-BILI-XZDE-CLEF2007_CHESHIRE_BERKBIPTDBASE	11.09%
	3 rd			
	4 th			
	5 th			
	Diff.			88.64%

Prompt Web service MT
Semantic analysis for GIR
tfxidf (Zebra DBMS)

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 42



Bilingual X-> Portuguese



Bilingual X-> Portuguese

Track	Rank	Part.	Experiment DOI	MAP
Bilingual Portuguese	1 st	cheshire	10.2415/GC-BILI-X2PT-CLEF2007_CHESHIRE_BERKBENPTBASE	20.12%
	2 nd	csusm	10.2415/GC-BILI-X2PT-CLEF2007_CSUSM_GEOBIESPT4	5.33%
	3 rd			
	4 th			
	5 th			
	Diff.			277.49%

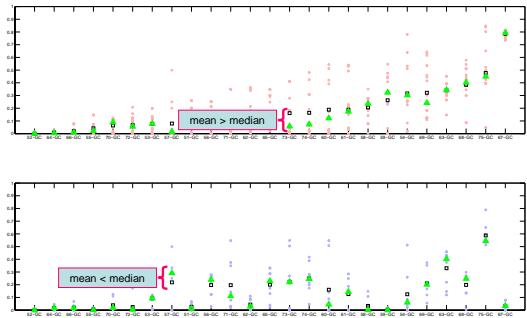


More Analyses

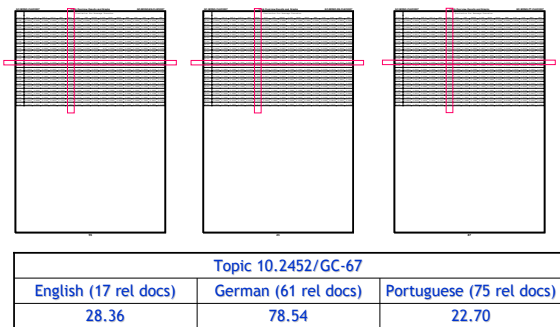
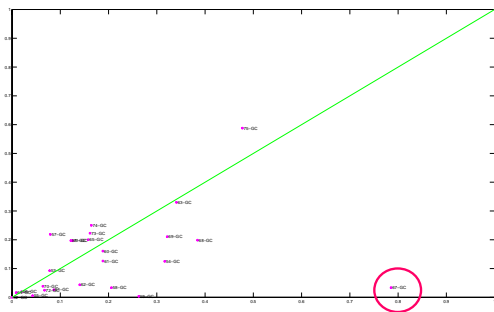
- Monolingual vs Bilingual Analyses
 - Mean of average precision for each topic of a task
 - Median of average precision for each topic of a task




Mean vs Median (mono, bili German)



Monolingual Mean VS Bilingual Mean German



 **GeoCLEF 2007**

```


<stop lang="en">
<name>10.245267</GC/name>
<title>F1 circuits where Ayrton Senna competed in 1994</title>
<desc>Find documents that mention circuits where the Brazilian driver Ayrton Senna participated in 1994. The name and location of the circuit is required</desc>
<name>Documents should indicate that Ayrton Senna participated in a race in a particular stadion, and the location of the race track.</name>
<stop>

<stop lang="de">
<name>10.245267</GC/name>
<title>Formel 1 Rennstrecken, auf denen Ayrton Senna 1994 gefahren ist</title>
<desc>Dokumente, die Rennstrecken erwähnen, auf denen der Brasilianische Fahrer Ayrton Senna 1994 gefahren ist. Name und Ort müssen erwähnt sein.</desc>
<name>Dokumente sollten angeben, dass Ayrton Senna an einem Rennen auf einer bestimmten Strecke teilgenommen hat. Der Ort der Strecke sollte genannt sein.</name>
<stop>

<stop lang="pt">
<name>10.245267</GC/name>
<title>Pistas em que Ayrton Senna correu em 1994</title>
<desc>Em que circuitos de fórmula 1 Ayrton Senna competiu em 1994</desc>
<name>Documentos que indiquem que Ayrton Senna participou em corridas num determinado autódromo, em que a localização deste é também mencionada.</name>
<stop>

```

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 49


 **Overview of GeoCLEF 2007**

- IR techniques
- IE/NLP techniques
- GIR techniques

- Systems
- Resources
- Experiments


- Translation
- General comments

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 50

 **IR techniques and systems**

- automatic and manual query expansion
- blind relevance feedback
- INL2 term weighting model
- divergence from randomness framework
- latent Dirichlet allocation model
- logistic regression
- query decomposition
- vector space model
- stemming
- systems: Lemur (2), Lucene (4), MG4J, Terrier (2), Zebra DBMS


CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 51

 **NLP techniques and systems**

- named entity recognition
- WordNet-based expansion
- semantic analysis
- part-of-speech tagging
- use of a QA system for subqueries
- decomposing (for German)
- WordNet lemmatizer

- systems: LingPipe, Annie (GATE)

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 52

 **GIR/GIE techniques**

- location disambiguation
 - geographic unique strings
- location normalization
- query expansion based on geographical terms
- query expansion based on a geographic ontology
- heuristic geographically informed filtering
 - removing candidates
 - using shape files for close or near geographical relations
- separate geographical indexes
- geographic cooccurrence model (based on Wikipedia)
- geographic relation finder

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 53

 **GIR/GIE resources**

- geonames gazetteer
- World Gazetteer
- Getty TGN
- GNIS
- GeoWorldMap World gazetteer
- ADL Feature Type thesaurus
- GKB 2.0
- Spanish toponymy list

- plus Wikipedia and WordNet

CLEF 2007, Budapest, September 19-21, 2007 T. Mandl et al. 54



Translation (only queries)

- Systran (PT-EN)
- Prompt (ES-EN)
- Promt (EN-DE)
- LEC Power Translator (all possible pairs)
- transfer MT system (ES-EN, ES-PT)
- Toggletext (ID-EN)



Experiments

- In addition to the particular experiments of each group concerning their original approaches
- Separate indexes for geographic and non-geographic information
- Apply different techniques to different parts of the query
- Use T, TD or TDN



Comments

- Very few papers mention distinguished treatment for different kinds of topics
 - although some discuss differences between topics
 - only one provides some per-topic analysis
- Most participants have best results with text only!
- Most participants use NER,
 - but no NER evaluation in itself has been reported: does it really help? or does it also diminish performance?
 - they don't use the output of NER the same way: wide difference in behaviour after NER has been invoked