

Thomas Mandl

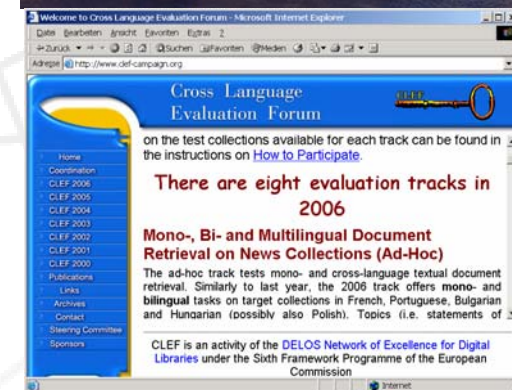

Information Science  
Universität Hildesheim  
mandl@uni-hildesheim.de

# CLEF 2006

-

# Robust Overview

Cross-Language  
Evaluation Forum (CLEF)



# Why Robust?

The user never sees the perspective of an evaluation (=MAP)  
but only the performance on his/her request(s).



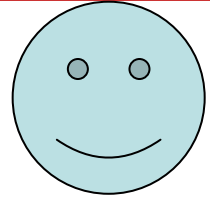
# Why Robust?

«The **unhappy customer**, on average, **will tell 27** other people about their experience. With the use of the internet, whether web pages or e-mail, that number can increase to the thousands ...»

«**Dissatisfied customers tell** an average of **ten** other people about their bad experience. Twelve percent tell up to twenty people.»

→ ***Bad news travels fast.***

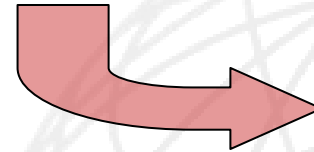
# Why Robust?



On the other hand, satisfied customers will tell an average of *five* people about their positive experience.

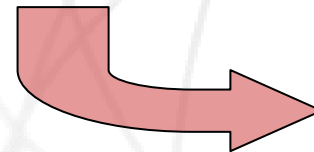
→ ***Good news travels somewhat slower***

Your system should produce less bad news!



improve on  
*hard* topics

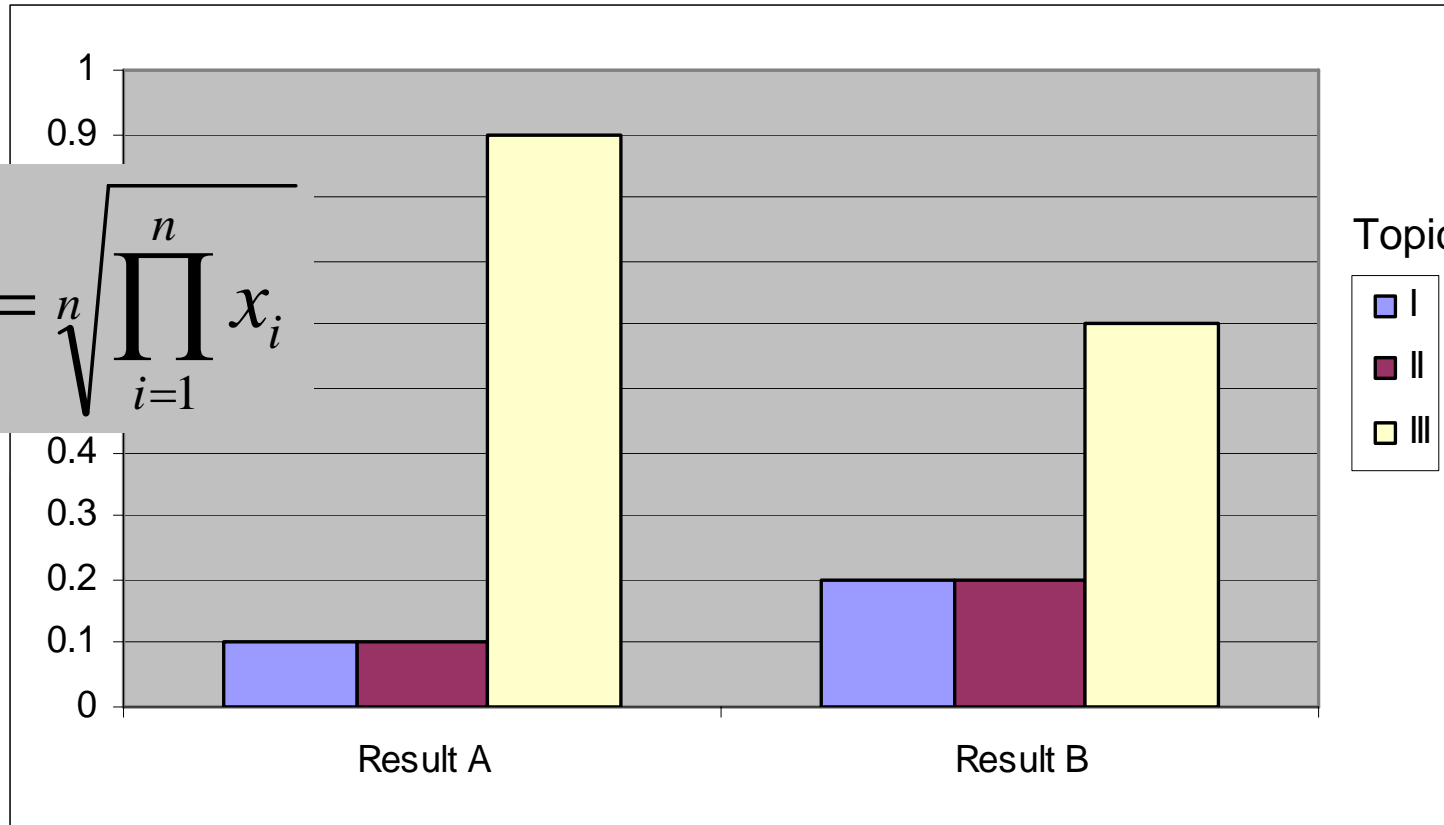
Don't worry too much about the good news



best topics

# Which system is better?

$$\text{geoAve} = \sqrt[n]{\prod_{i=1}^n x_i}$$



GeoAve	A	0.21	GeoAve	B	<b>0.29</b>
MAP	A	<b>0.37</b>	MAP	B	0.33

- Robustness in multilingual retrieval
  - Stable performance over all topics instead of high average performance (like at TREC, just for other languages)
  - Stable performance over all topics for multi-lingual retrieval
  - Stable performance over different languages (so far at CLEF ?)

- „More work needs to be done on customizing methods for each topic“ (Harman 2005)
- Hard, but some work ist done
  - RIA Workshop
  - SIGIR Workshop on Topic Difficulty
  - TREC Robust Track (until 2005)
  - ...

- Six languages
- 1.35 million documents
- 3.6 gigabytes of text

<b>Language</b>	<b>Collection</b>
English	LA Times 94, Glasgow Herald 95
French	ATS (SDA) 94/95, Le Monde 94
Italian	La Stampa 94, AGZ (SDA) 94/95
Dutch	NRC Handelsblad 94/95, Algemeen Dagblad 94/95
German	Frankfurter Rundschau 94/95, Spiegel 94/95, SDA 94
Spanish	EFE 94/95



# Data Collections

CLEF Year	2001	2002	2003
Documents			<i>Missing relevance judgements</i>
Topics	#41-90	#91-140	#141-200
Relevance Judgements			

- Arbitrary split of the 160 topics
  - 60 training topics
  - 100 test topics



- Finding a set of difficult topics  
... was difficult
  - Before the campaign, **no** topics were found which were difficult for more than one language or task at CLEF 2001, CLEF 2002 and CLEF 2003
  - Several definitions of difficulty were tested
- Just as at Robust Track @ TREC
  - Topics are not difficult by themselves
  - but only in interaction with a collection

- Sub-Tasks
  - Mono-lingual
    - Only for the six document languages as topic language
  - Bi-lingual
    - it- > es
    - fr- > nl
    - en- > de
  - Multi-lingual



- Submission of Training data was encouraged
  - Further analysis of topic difficulty
- Overall, systems did better on test topics

- U. Coruna & U. Sunderland (Spain & UK)
- U. Jaen (Spain)
- DAEDALUS & Madrid Univs. (Spain)
- U. Salamanca – REINA (Spain)
- Hummingbird Core Tech. (Canada)
- U. Neuchatel (Switzerland)
- Dublin City U. – Computing (Ireland)
- U. Hildesheim – Inf. Sci. (Germany)

# More than 100 Runs ...

<b>Task</b>	<b>Language</b>	<b>nr test runs</b>	<b>nr training runs</b>	<b>nr groups</b>
<b>mono</b>	<b>en</b>	<b>13</b>	<b>7</b>	<b>6</b>
	<b>fr</b>	<b>18</b>	<b>10</b>	<b>7</b>
	<b>nl</b>	<b>7</b>	<b>3</b>	<b>3</b>
	<b>de</b>	<b>7</b>	<b>3</b>	<b>3</b>
	<b>es</b>	<b>11</b>	<b>5</b>	<b>5</b>
	<b>it</b>	<b>11</b>	<b>5</b>	<b>5</b>
<b>bi</b>	<b>it-&gt;es</b>	<b>8</b>	<b>2</b>	<b>3</b>
	<b>fr-&gt;nl</b>	<b>4</b>	<b>0</b>	<b>1</b>
	<b>en-&gt;de</b>	<b>5</b>	<b>1</b>	<b>2</b>
<b>multi</b>	<b>multi</b>	<b>10</b>	<b>3</b>	<b>4</b>

# Results – Mono-lingual

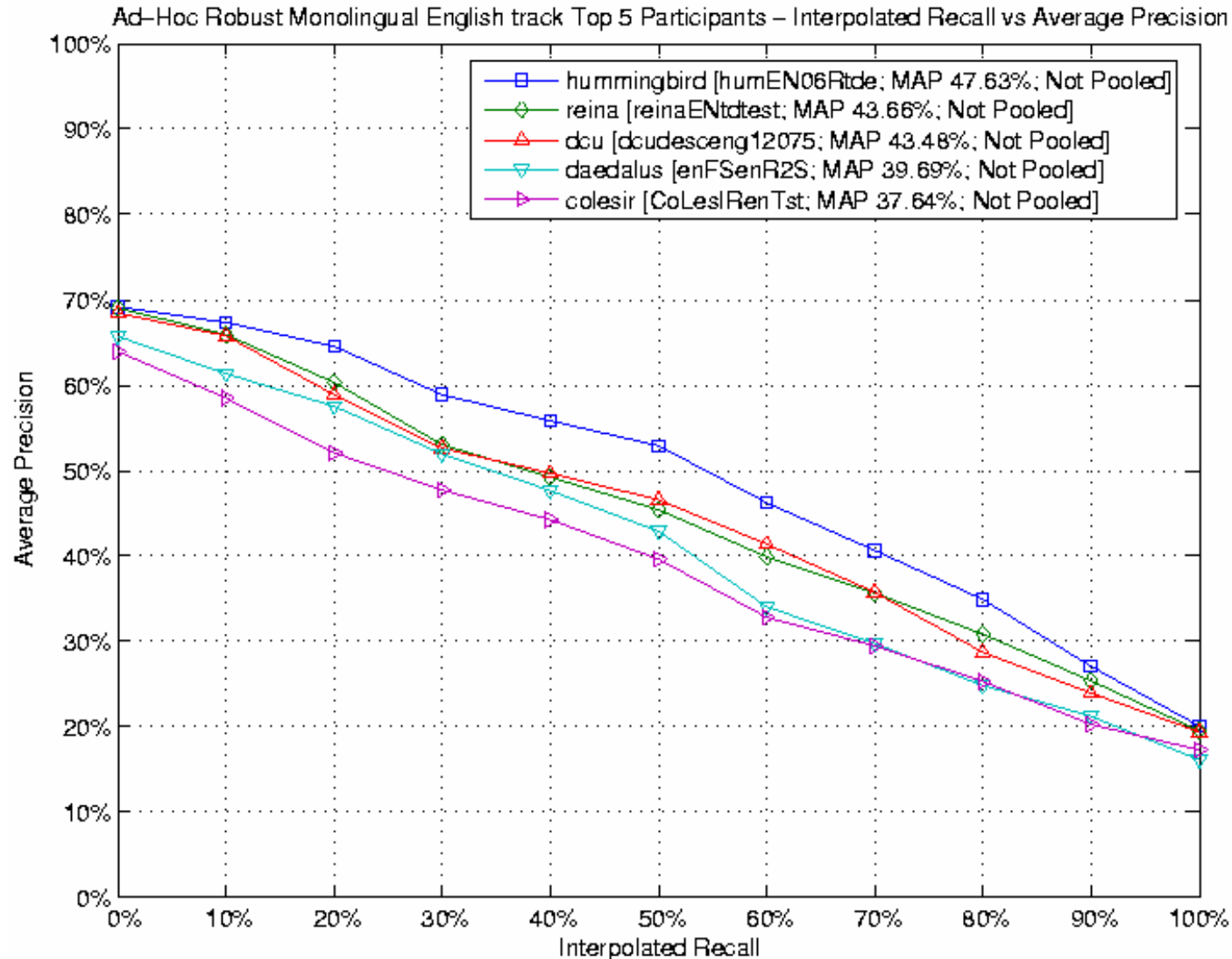
Track	Participant Rank					
	1st	2nd	3rd	4th	5th	Diff.
<b>Dutch</b>	hummingbird	daedalus	colesir			1st vs 3rd
MAP	51.06%	42.39%	41.60%			22.74%
GMAP	25.76%	17.57%	16.40%			57.13%
Run	humNL06Rtd	enFSnlR2S	CoLeslRnlTst			
<b>English</b>	hummingbird	reina	dcu	daedalus	colesir	1st vs 5th
MAP	47.63%	43.66%	43.48%	39.69%	37.64%	26.54%
GMAP	11.69%	10.53%	10.11%	8.93%	8.41%	39.00%
Run	humEN06Rtd	reinaENtdtest	dcudesceng1	enFSenR2S	CoLeslRenTst	
<b>French</b>	unine	hummingbird	reina	dcu	colesir	1st vs 5th
MAP	47.57%	45.43%	44.58%	41.08%	39.51%	20.40%
GMAP	15.02%	14.90%	14.32%	12.00%	11.91%	26.11%
Run	UniNEfr1	humFR06Rtd	reinaFRtdtest	dcudescfr120	CoLeslRfrTst	



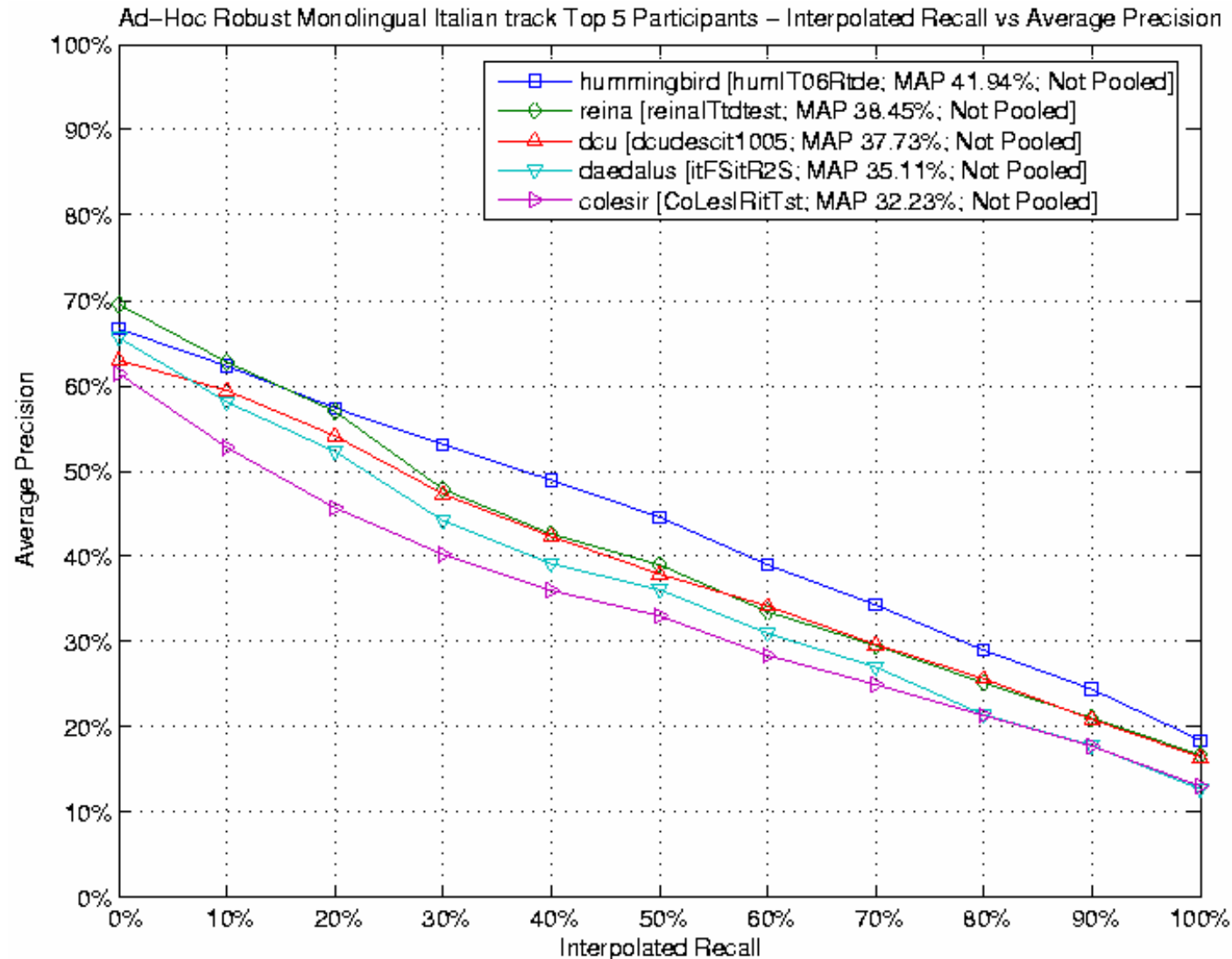
# Results – Mono-lingual

Track	Participant Rank					Diff.
	1st	2nd	3rd	4th	5th	
<b>German</b>	hummingbird	colesir	daedalus			1st vs 3rd
MAP	48.30%	37.21%	34.06%			41.81%
GMAP	22.53%	14.80%	10.61%			112.35%
Run	humDE06Rtd	CoLesIRdeTs	deFSdeR2S			
<b>Italian</b>	hummingbird	reina	dcu	daedalus	colesir	1st vs 5th
MAP	41.94%	38.45%	37.73%	35.11%	32.23%	30.13%
GMAP	11.47%	10.55%	9.19%	10.50%	8.23%	39.37%
Run	humIT06Rtd	reinalTtdtest	dcudes100	itESitR2S	CoLesIRitTst	
<b>Spanish</b>	hummingbird	reina	dcu	daedalus	colesir	1st vs 5th
MAP	45.66%	44.01%	42.14%	40.40%	40.17%	13.67%
GMAP	23.61%	22.65%	21.32%	19.64%	18.84%	25.32%
Run	humES06Rtd	reinaEStdtest	dcudescsp12	esFSesR2S	CoLesIResTst	

# Results – Mono-lingual English

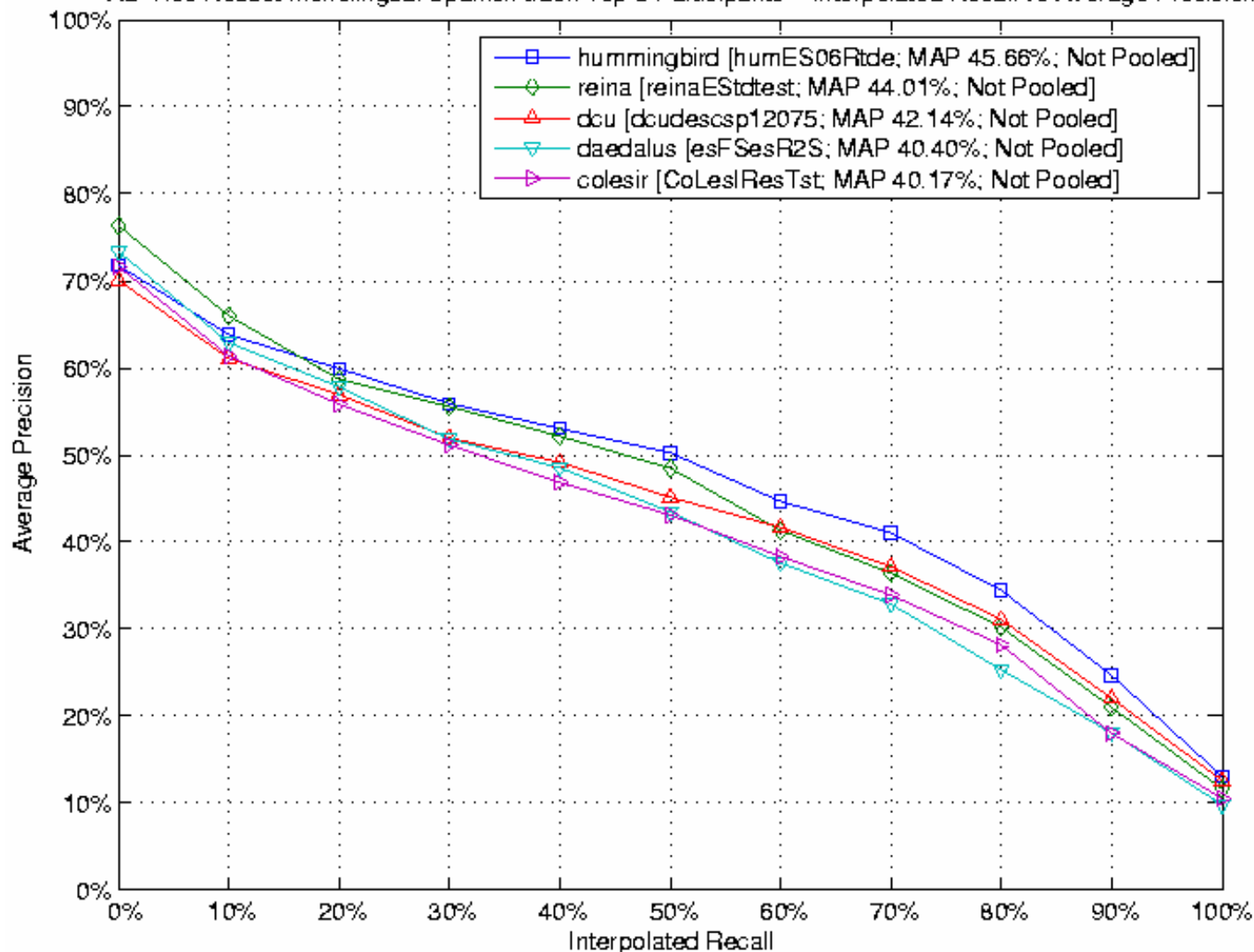


# Results – Mono-lingual Italian



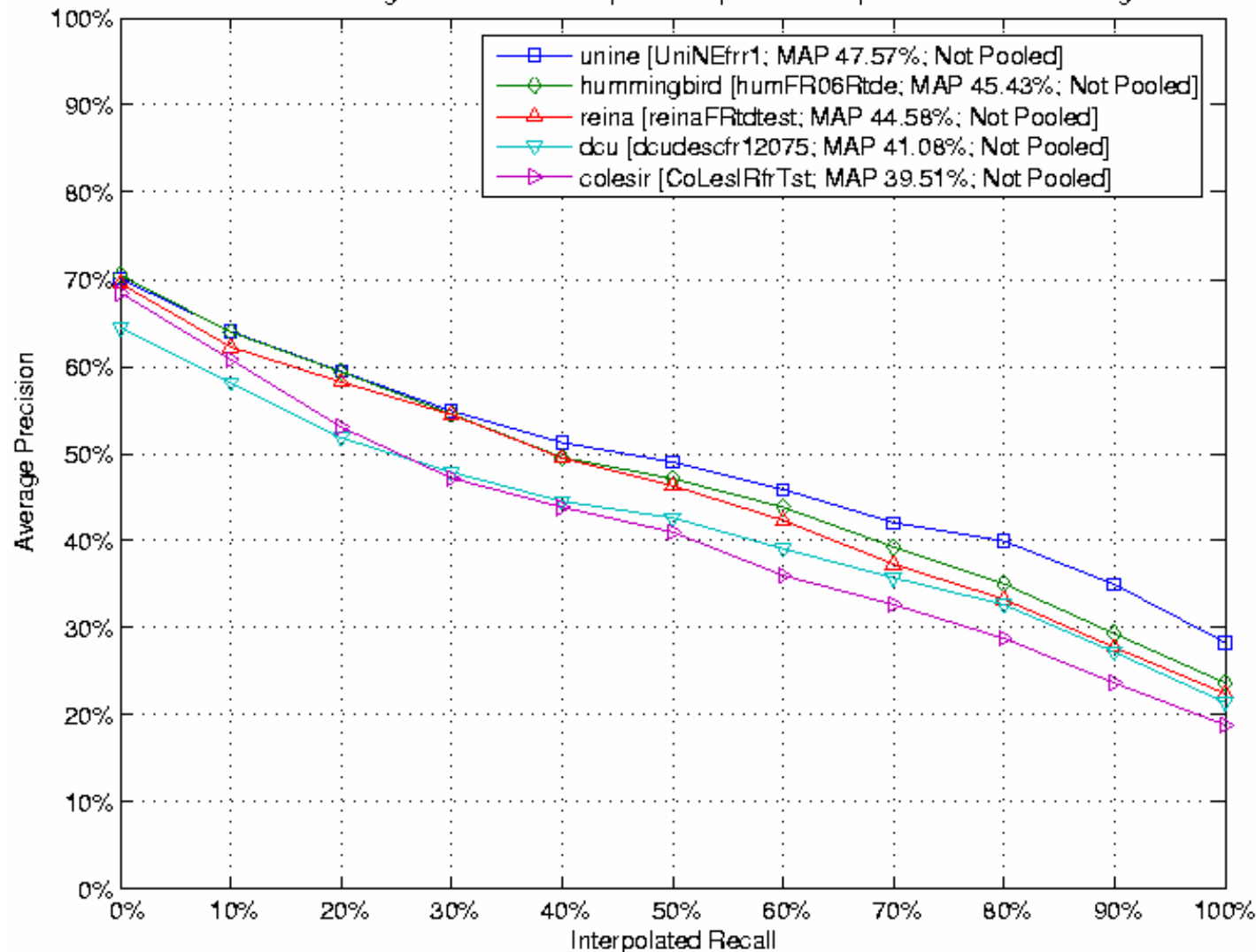
# Results – Mono-lingual Spanish

Ad-Hoc Robust Monolingual Spanish track Top 5 Participants – Interpolated Recall vs Average Precision



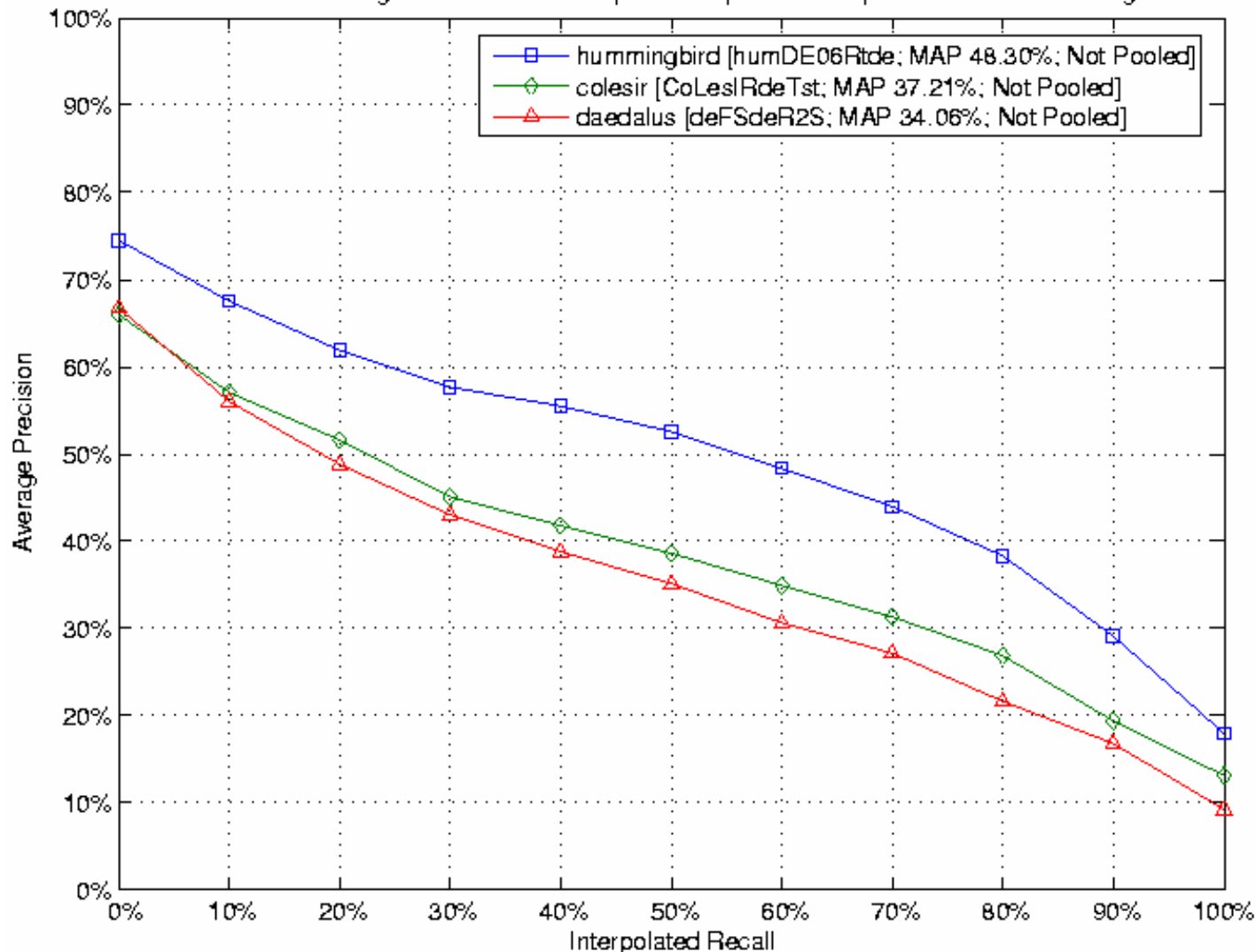
# Results – Mono-lingual French

Ad-Hoc Robust Monolingual French track Top 5 Participants – Interpolated Recall vs Average Precision



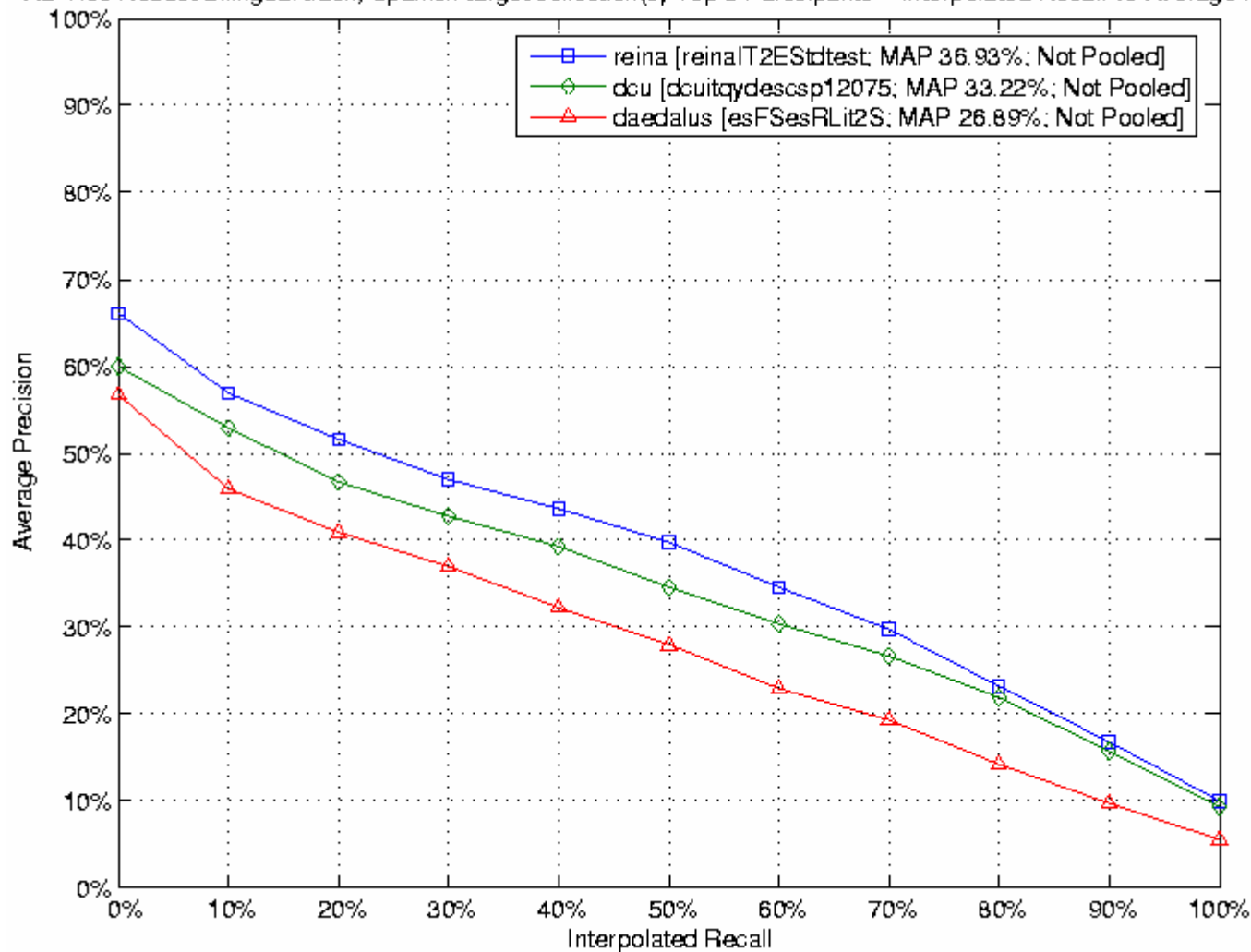
# Results – Mono-lingual German

Ad-Hoc Robust Monolingual German track Top 5 Participants – Interpolated Recall vs Average Precision



# Results – Bi-lingual

Ad-Hoc Robust Bilingual track, Spanish target collection(s) Top 5 Participants – Interpolated Recall vs Average Precision

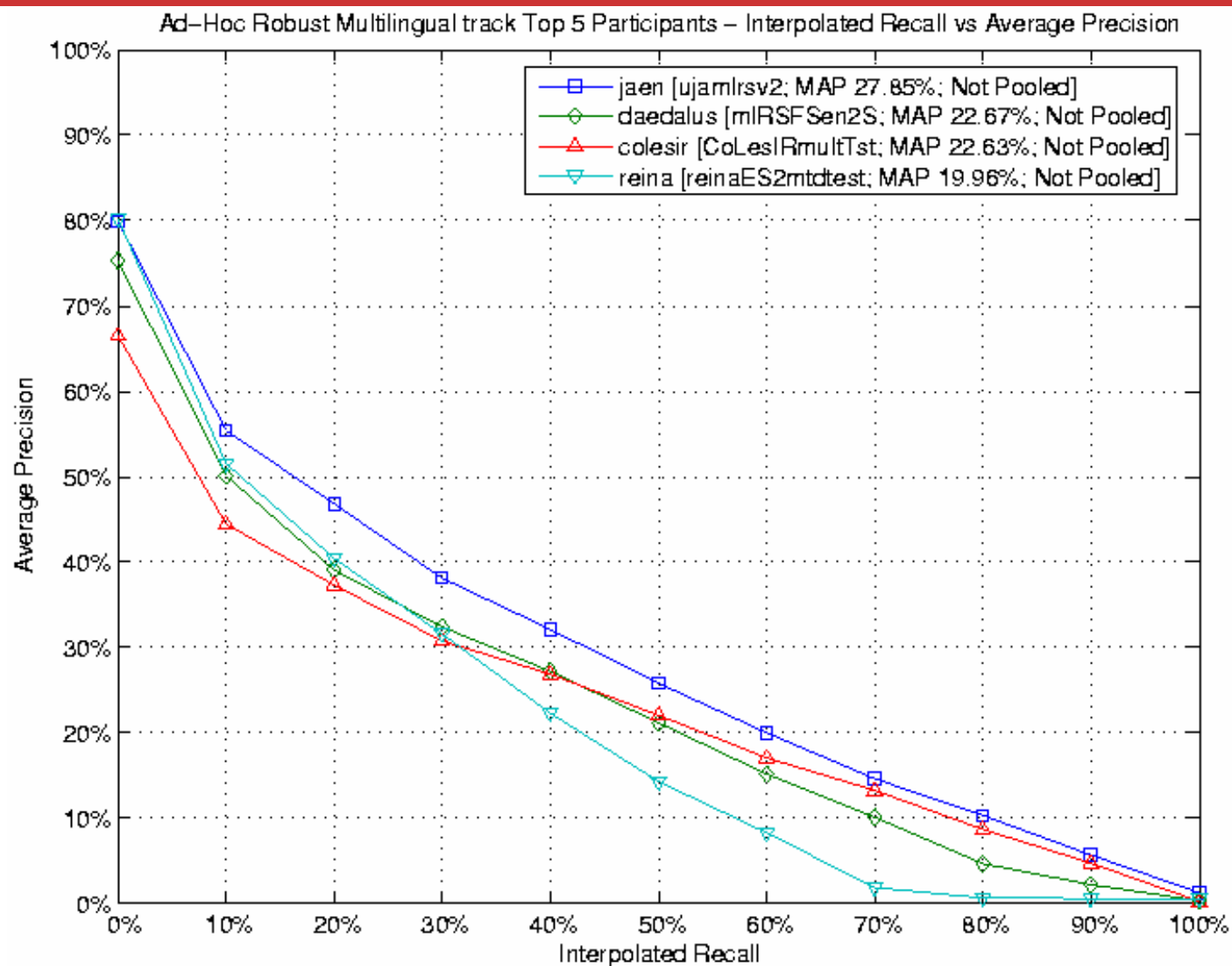


# Results – Multi-lingual

Track	Participant Rank			
	1st	2nd	3rd	4th
<b>Multilingual</b>	jaen	daedalus	colesir	reina
MAP	27.85%	22.67%	22.63%	19.96%
GMAP	15.69%	11.04%	11.24%	13.25%
Run	ujamlrsv2	mIRSFSen2S	COLESIRmultTst	reinaES2mtdtest



# Results – Multi-lingual



- Example: topic 64
  - Easiest topic for mono and bi German
  - Hardest topic for mono Italian
- Example: topic 144
  - Easiest topic for four bi Dutch
  - hardest for mono and bi German
- Example: topic 146
  - Among three hardest topics for four sub-tasks
  - mid-range for all other sub-tasks

- SINAI expanded with terms gathered from a web search engine [Martinez-Santiago et al. 2006]
- REINA used a heuristic for determining hard topics during training. Different expansion techniques were applied [Zazo et al. 2006]
- Hummingbird experimented with other evaluation measures then used in the track [Tomlinson 2006].
- MIRACLE tried to find a fusion scheme which is good for the robust measure [Goni-Menoyo et al. 2006]

- What can we do with the data?
- Have people improved in comparison to CLEF 2001 through 2003?
- Are low MAP scores a good indicator of topic difficulty?

# Acknowledgements

- Giorgio di Nunzio & Nicola Ferro (U Padua)
- Robust Committee
  - Donna Harman (NIST)
  - Carol Peters (ISTI-CNR)
  - Jacques Savoy (U Neuchâtel)
  - Gareth Jones (Dublin City U)
- Ellen Voorhees (NIST)
- Participants

*Thanks for your Attention*

*I am looking forward to the Discussion*

*Please come to the Breakout Session*