

# Overview of the wikipediaMM task at ImageCLEF 2009

Theodora Tsikrika<sup>1</sup> and Jana Kludas<sup>2</sup>

<sup>1</sup>CWI, Amsterdam, The Netherlands

<sup>2</sup>CUI, University of Geneva, Switzerland

Theodora.Tsikrika@cwi.nl, jana.kludas@cui.unige.ch

## Abstract

ImageCLEF's wikipediaMM task provides a testbed for the system-oriented evaluation of multimedia information retrieval from a collection of Wikipedia images. The aim is to investigate retrieval approaches in the context of a large and heterogeneous collection of images (similar to those encountered on the Web) that are searched for by users with diverse information needs. This paper presents an overview of the resources, topics, and assessments of the wikipediaMM task at ImageCLEF 2009, summarises the retrieval approaches employed by the participating groups, and provides a first analysis of the main evaluation results.

## Keywords

ImageCLEF, Wikipedia image collection, image retrieval, evaluation

## 1 Introduction

The wikipediaMM task is an ad-hoc image retrieval task. The evaluation scenario is thereby similar to the classic TREC ad-hoc retrieval task and the ImageCLEF photo retrieval task: simulation of the situation in which a system knows the set of documents to be searched, but cannot anticipate the particular topic that will be investigated (i.e. topics are not known to the system in advance). Given a multimedia query that consists of a title and one or more sample images describing a user's multimedia information need, the aim is to find as many relevant images as possible from the (INEX MM) wikipedia image collection. A multi-modal retrieval approach in that case should be able to combine the relevance of different media types into a single ranking that is presented to the user.

The wikipediaMM task differs from other benchmarks in multimedia information retrieval, like TRECVID, in the sense that the textual modality in the wikipedia image collection contains less noise than the speech transcripts in TRECVID. Maybe that is one of the reasons why, both in last year's task and in INEX Multimedia 2006-2007 (where this image collection was also used), it has proven challenging to outperform the text-only approaches. This year, the aim is to promote the investigation of multi-modal approaches to the forefront of this task by providing a number of resources to support the participants towards this research direction.

The paper is organised as follows. First, we introduce the task's resources: the wikipedia image collection and additional resources, the topics, and the assessments (Sections 2-4). Section 5 presents the approaches employed by the participating groups and Section 6 summarises their main results. Section 7 concludes the paper.

## 2 Task resources

The resources used for the wikipediaMM task are based on Wikipedia data. The collection is the **(INEX MM) wikipedia image collection**, which consists of approximately 150,000 JPEG and PNG Wikipedia images provided by Wikipedia users. Each image is associated with user-generated alphanumeric, unstructured metadata in English. These metadata usually contain a brief caption or description of the image, the Wikipedia user who uploaded the image, and the copyright information. These descriptions are highly heterogeneous and of varying length. Further information about the image collection can be found in [4].



Figure 1: Wikipedia image+metadata example from the wikipedia image collection.

Additional resources were also provided to support the participants in their investigations of multi-modal approaches. These resources are:

**Image similarity matrix:** The similarity matrix for the images in the collection has been constructed by the IMEDIA group at INRIA. For each image in the collection, this matrix contains the list of the top  $K = 1000$  most similar images in the collection together with their similarity scores. The same is given for each image in the topics. The similarity scores are based on the distance between images; therefore, the lower the score, the more similar the images. Further details on the features and distance metric used can be found in [1].

**Image classification scores:** For each image, the classification scores for the 101 MediaMill concepts have been provided by UvA [3]. The UvA classifier is trained on manually annotated TRECVID video data for concepts selected for the broadcast news domain.

**Image features:** For each image, the set of the 120D feature vectors that has been used to derive the above image classification scores [2] has also been made available. Participants can use these feature vectors to custom-build a content-based image retrieval (CBIR) system, without having to pre-process the image collection.

The additional resources are beneficial to researchers who wish to exploit visual evidence without performing image analysis. Of course, participants could also extract their own image features.

## 3 Topics

The topics are descriptions of multimedia information needs that contain textual and visual hints.

### 3.1 Topic Format

These multimedia queries consist of a textual part, the query title, and a visual part, one or several example images.

**<title>** query by keywords

**<image>** query by image content (one or several)

**<narrative>** description of query in which the definitive definition of relevance and irrelevance are given

#### 3.1.1 <title>

The topic <title> simulates a user who does not have (or want to use) example images or other visual constraints. The query expressed in the topic <title> is therefore a text-only query. This profile is likely to fit most users searching digital libraries.

Upon discovering that a text-only query does not produce many relevant hits, a user might decide to add visual hints and formulate a multimedia query.

#### 3.1.2 <image>

The visual hints are example images, which can be taken from outside or inside the wikipedia image collection and can be of any common format. Each topic has at least one example image, but it can have several, e.g., to describe the visual diversity of the topic.

#### 3.1.3 <narrative>

A clear and precise description of the information need is required in order to unambiguously determine whether or not a given document fulfils the given information need. In a test collection this description is known as the narrative. It is the only true and accurate interpretation of a user's needs. Precise recording of the narrative is important for scientific repeatability - there must exist, somewhere, a definitive description of what is and is not relevant to the user. To aid this, the <narrative> should explain not only what information is being sought, but also the context and motivation of the information need, i.e., why the information is being sought and what work-task it might help to solve.

These different types of information sources (textual terms and visual examples) can be used in any combination. It is up to the systems how to use, combine or ignore this information; the relevance of a result does not directly depend on these constraints, but it is decided by manual assessments based on the <narrative>.

### 3.2 Topic Development

The topics in the ImageCLEF 2009 wikipediaMM task have been partly developed by the participants and partly by the organisers. This year the participation in the topic development process was not obligatory, so only 2 of the participating groups submitted a total of 11 candidate topics. The rest of the candidate topics were created by the organisers with the help of the log of an image search engine. After a selection process performed by the organisers, a final list of 45 topics was created.

These final topics range from simple, and thus relatively easy (e.g., “bikes”), to semantic, and hence highly difficult (e.g., “aerial photos of non-artificial landscapes”), with the latter forming the bulk of the topics. Semantic topics typically have a complex set of constraints, need world knowledge, and/or contain ambiguous terms, so they are expected to be challenging for current state-of-the-art retrieval algorithms. We encouraged the participants to use multi-modal approaches since they are more appropriate for dealing with semantic information needs. On average, the 45 topics contain 1.7 images and 2.7 words.

## 4 Assessments

The wikipediaMM task is an image retrieval task, where an image with its metadata is either relevant or not (binary relevance). We adopted TREC-style pooling of the retrieved images with a pool depth of 50, resulting in pools of between 299 and 802 images with a mean and median both around 545. The evaluation was performed by the participants of the task within a period of 4 weeks after the submission of runs. The 7 groups that participated in the evaluation process used the web-based interface that was used last year and which has also been previously employed in the INEX Multimedia and TREC Enterprise tracks.

## 5 Participants

A total of 8 groups submitted 57 runs: CEA (LIC2M-CEA, Centre CEA de Saclay, France), DCU (Dublin City University, School of Computing, Ireland), DEU (Dokuz Eylul University, Department of Computer Engineering, Turkey), IIIT-Hyderabad (Search and Info Extraction Lab, India), LaHC (Laboratoire Hubert Curien, UMR CNRS, France), SZTAKI (Hungarian Academy of Science, Hungary), SINAI (Intelligent Systems, University of Jaen, Spain) and UALICANTE (Software and Computer Systems, University of Alicante, Spain).

Table 1: Types of the 57 submitted runs.

<b>run type</b>	<b># runs</b>
text	26
visual	2
text/visual	29
query expansion	18
relevance feedback	7

Table 1 gives an overview of the types of the submitted runs. This year more multi-modal (text/visual) than text-only runs were submitted. A short description of the participants’ approaches follows.

**CEA (12 runs)** They extended the approach they employed last year by refining the textual query expansion procedure and introducing of a k-NN based visual reranking procedure. Their main aim was to examine whether combining textual and content-based retrieval improves over purely textual search.

**DCU (5 runs)** Their main effort concerned the expansion of the image metadata using the Wikipedia abstracts’ collection DBpedia. Since the metadata is short for retrieval by query text, they expand the query and documents using the Rocchio algorithm. For retrieval, they used the LEMUR toolkit. They also submitted one visual run.

**DEU (6 runs)** Their research interests focussed on 1) the expansion of native documents and queries, term phrase selection based on WordNet, WSD and WordNet similarity functions and 2) a new reranking approach with Boolean retrieval and C3M based clustering.

**IIT-H (1 run)** Their system automatically ranks the most similar images to a given textual query using a combination of the Vector Space Model and the Boolean model. The system preprocesses the data set in order to remove the non-informative terms.

**LaHC (13 runs)** In this second participation, they extended their approach (a multimedia document model defined as a vector of textual and visual terms weighted using tf.idf) by using 1) additional information for the textual part (legend and image bounding text extracted from the original documents), 2) different image detectors and descriptors, and 3) a new text/image combination approach.

**SINAI (4 runs)** Their approach focussed on query and document expansion techniques based on WordNet. They used the LEMUR toolkit as their retrieval system.

**SZTAKI (7 runs)** They used both textual and visual features and employed image segmentation, SIFT keypoints, Okapi BM25 based text retrieval, and query expansion by an online thesaurus. They preprocessed the annotation text to remove author and copyright information and biased retrieval towards images with filenames containing relevant terms.

**UALICANTE (9 runs)** They used IR-n, a retrieval system based on passages and applied two different term selection strategies for query expansion: Probabilistic Relevance Feedback and Local Context Analysis, and their multi-modal versions. They also used the same technique for Camel Case decompounding of image filenames that they used in last year’s participation.

## 6 Results

Table 2 presents the evaluation results for the 15 best performing runs ranked by Mean Average Precision (MAP). DEU’s text-only runs performed best. But as already seen last year, approaches that fuse several modalities can compete with the text-only ones. Furthermore, it is notable that all participants that used both mono-media and multi-modal algorithms achieved their best results with their multi-modal runs. The complete list of results can be found at the ImageCLEF website <http://www.imageclef.org/2009/wikiMM-results>.

Table 2: results for the top 15 runs

	Participant	Run	Modality	FB/QE	MAP	P@10	P@20	R-prec.
1	deuceng	deuwiki2009_205	TXT	QE	0.2397	0.4000	0.3133	0.2683
2	deuceng	deuwiki2009_204	TXT	QE	0.2375	0.4000	0.3111	0.2692
3	deuceng	deuwiki2009_202	TXT	QE	0.2358	0.3933	0.3189	0.2708
4	lach	TXTIMG_100_3_1_1_5_meanstd	TXTIMG	NOFB	0.2178	0.3378	0.2811	0.2538
5	lach	TXTIMG_50_3_1_1_5_meanstd	TXTIMG	NOFB	0.2148	0.3356	0.2867	0.2536
6	cea	cealateblock	TXTIMG	QE	0.2051	0.3622	0.2744	0.2388
7	cea	ceaearyblock	TXTIMG	QE	0.2046	0.3556	0.2833	0.2439
8	cea	ceabofblock	TXTIMG	QE	0.1975	0.3689	0.2789	0.2342
9	cea	ceatlepbblock	TXTIMG	QE	0.1959	0.3467	0.2733	0.2236
10	cea	ceabofblockres	TXTIMG	QE	0.1949	0.3689	0.2789	0.2357
11	cea	ceatlepbblockres	TXTIMG	QE	0.1934	0.3467	0.2733	0.2236
12	lach	TXTIMG_Siftdense_0.084	TXTIMG	NOFB	0.1903	0.3111	0.2700	0.2324
13	lach	TXT_100_3_1_1_5	TXT	NOFB	0.1890	0.2956	0.2544	0.2179
14	lach	TXT_50_3_1_1_5	TXT	NOFB	0.1880	0.3000	0.2489	0.2145
15	ualicante	Alicante-MMLCA	TXTIMG	FB	0.1878	0.2733	0.2478	0.2138

Next, we analyse the evaluation results. In our analysis, we use only the top 90% of the runs to exclude noisy and buggy results. Furthermore, we excluded 3 runs that we considered to be redundant, i.e., they were produced by the same group and achieved the exact same result, so as to reduce the bias of the analysis.

### 6.1 Performance per modality for all topics

Table 3 shows the average performance and standard deviation with respect to modality. On average, the multi-modal runs manage to outperform the mono-media runs with respect to all

examined evaluation metrics (MAP, Precision at 20, and precision after R (= number of relevant documents are retrieved)).

Table 3: Results per modality over all topics.

Modality	MAP		P@20		R-prec.	
	Mean	SD	Mean	SD	Mean	SD
All top 90% runs (46 runs)	0.1751	0.0302	0.2356	0.0624	0.2076	0.0572
TXT in top 90% runs (23 runs)	0.1726	0.0326	0.2278	0.0427	0.2038	0.0328
TXTIMG in top 90% runs (23 runs)	0.1775	0.0281	0.2433	0.0364	0.2115	0.0307

## 6.2 Performance per topic and per modality

To analyse the average difficulty of the topics, we classify the topics based on the average MAP values per topic as follows:

**easy:**  $aMAP > 0.3$

**medium:**  $0.2 < aMAP \leq 0.3$

**hard:**  $0.1 < aMAP \leq 0.2$

**very hard:**  $aMAP < 0.1$ .

Table 4 presents the top 6 topics per class (i.e., easy, medium, hard, and very hard), together with the total number of topics per class. Most of the topics are considered to be hard. This was actually intended during the topic development process where we opted for highly semantic topics that are challenging for current retrieval approaches. Nonetheless, 10 out of 45 topics were of easy and medium difficulty. Only 7 topics were very hard to solve. Therein, topics #97 “woman in pink dress” and #98 “close up of people doing sport” can be considered as unsolvable, since their  $aMAP < 0.05$ .

Table 4: Topics classified based on their difficulty. The top 6 topics are shown per class together with the total number of topics per class.

easy	medium	hard	very hard
(112) hot air balloons	(118) coral reef underwater	(120) yellow flower	(105) snowy street
(88) madonna portrait	(90) satellite image of river	(91) landline telephone	(78) sculpture of an animal
(80) orthodox icons	(110) desert landscape	(99) flowers on trees	(117) earth from space
(108) bird nest	(77) real rainbow	(79) stamp human face	(85) aerial ph. of landscapes
(103) palm trees		(107) red fruit	(89) people laughing
(93) close up antenna		(94) people with dogs	(97) woman in pink dress
6	4	28	7

We also analysed the performance of runs that use only text (TXT) versus text and visual resources (TXTIMG). Figure 2 shows the average performance on each topic for all runs, the text-only and text-visual based ones. The text-based runs outperform the text-visual ones in 22 out of the 45. So, slightly more than half of the topics profit from a multi-modal approach.

## 6.3 Visuality of topics

The “visuality” of topics can be deduced from the performance of text-only and text-visual approaches that we presented in the last section. We consider that if, for a topic, the text-visual approaches improve significantly the MAP over all runs (e.g., by  $diff(MAP) \geq 0.01$ ), then we could consider that to be a visual topic. In the same way, we can define topics as textual, if the text-only approaches improve significantly the MAP over all runs of a topic. Based on this, 15 of the topics can be characterised as textual and 14 as visual. The remaining 16 topics, where no clear improvements are observed, are considered to be neutral.

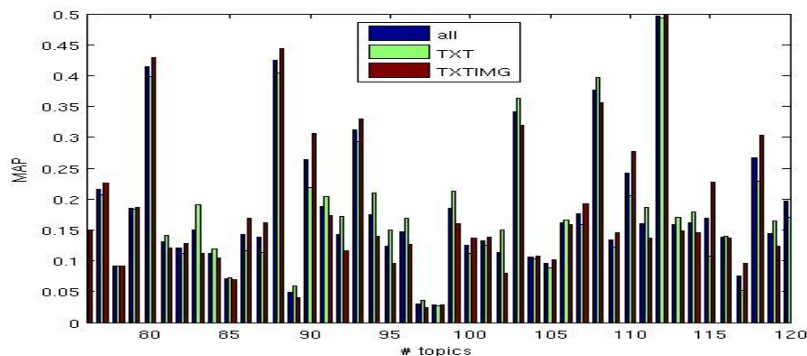


Figure 2: Average topic performance over all, text-only and text/visual runs

Table 5 presents the top 10 topics in each group, as well as some statistics on the topic, their relevant documents, and their distribution over the difficulty. As expected, visual topics have more image examples per topic ( $\#images/topic$ ) than textual ones (1.66 vs. 1.85); however, the neutral topics have an even higher average of 2.06 images per topic. The same tendency is observed in the average number of words in the topics ( $\#words/topic$ ). Short titled topics are better solved with text-only approaches, topics with longer titles tend to be visual or neutral. Therefore, it appears that the latter two groups contain the more complex/semantic topics. The distribution of the textual, visual and neutral topics over the classes expressing their difficulty shows that the visual topics are more likely to fall into the easy/medium class than the textual or neutral ones. The neutral topics seem to contain in general very difficult topics, where neither the text-only approaches nor the text-visual ones could achieve good retrieval results.

Table 5: Top 10 best performing topics for textual and text-visual runs relative to the average over all runs.

	textual	visual	neutral
topics	(83) advertisement for cars (102) building site (94) people with dogs (92) bikes (95) group of dogs (99) flowers on trees (111) pol. campaign poster (103) palm trees (96) cartoon with a cat (119) harbor	(115) notes on music sheet (90) satellite image of river (118) coral reef underwater (110) desert landscape (120) yellow flower (86) situation after katrina (87) airplane crash (117) earth from space (88) madonna portrait (93) close up of antenna	(76) shopping on a market (77) real rainbow (78) sculpture of an animal (79) stamp without human face (81) Greek mythological figures (82) rider on horse (84) advertisement on buses (101) fire (104) street musician (105) snowy street
$\#$ out of 45	15	14	16
$\#images/topic$	1.66	1.85	2.06
$\#words/topic$	2.53	3.00	3.31
$\#reldocs$	35.33	36.28	36.50
$\#words/reldocs$	29.65	44.99	39.24
easy	2	3	1
medium	0	3	1
hard	12	7	9
very hard	1	1	5

## 6.4 Effect of Query Expansion and Relevance Feedback

Finally, we analyse the effect of the application of query expansion (QE) and relevance feedback (FB) techniques. Similarly to the analysis in the previous section, we consider the techniques to be useful for a topic, if they improved significantly the MAP over all runs. Table 6 presents the top 10 best performing topics for these techniques and some statistics. Query expansion is useful for 17 topics and relevance feedback for 11. The statistics show that these techniques can help improve the retrieval results for topics defined without too much detail, e.g., topics having a short

title (#words/topic) and/or a small number of example images (#images/topic).

Table 6: Top 10 best performing topics for textual and text-visual runs relative to the average over all runs.

	QE	FB
topics	(110) desert landscape (118) coral reef underwater (120) yellow flower (109) tennis player on court (92) bikes (82) rider on horse (101) fire (115) notes on music sheet (117) earth from space (119) harbor	(88) madonna portrait (115) notes on music sheet (87) airplane crash (93) close up of antenna (96) cartoon with a cat (79) stamp without human face (116) illustration of engines (118) coral reef underwater (95) group of dogs (104) street musician
# out of 45	17	11
#images/topic	1.94	1.72
#words/topic	2.76	3.18
#reldocs	46.47	40.36
#words/reldocs	37.98	42.74
easy	1	2
medium	2	1
hard	11	8
very hard	3	0

## 7 Conclusions

This year (similarly to 2008), a text-based approach performed best in the wikipediaMM task, even though highly semantic multimedia topics were developed with the aim to encourage and show the potential of multi-modal approaches. It is worth noting though that all of the participants that submitted both mono-media and multi-modal runs achieved their best results with their multi-modal runs. Additionally, we as organisers are really glad to see more than half of the submitted runs being multi-modal.

## 8 Acknowledgements

Theodora Tsikrika was supported by the European Union via the European Commission project VITALAS (contract no. 045389).

## References

- [1] Marin Ferecatu. Image retrieval with active relevance feedback using both visual and keyword-based descriptors. In *Ph.D. Thesis, Universit de Versailles, France*, 2005.
- [2] Jan C. van Gemert, Jan-Mark Geusebroek, Cor J. Veenman, Cees G. M. Snoek, and Arnold W. M. Smeulders. Robust scene categorization by learning image statistics in context. In *Proceedings of the 2006 Conference on Computer Vision and Pattern Recognition Workshop*, page 105, Washington, DC, USA, 2006. IEEE Computer Society.
- [3] Cees G. M. Snoek, Marcel Worring, Jan C. van Gemert, Jan-Mark Geusebroek, and Arnold W. M. Smeulders. The challenge problem for automated detection of 101 semantic concepts in multimedia. In *Proceedings of the 14th annual ACM international conference on Multimedia*, pages 421–430, New York, NY, USA, 2006. ACM Press.
- [4] Thijs Westerveld and Roelof van Zwol. The INEX 2006 multimedia track. In Norbert Fuhr, Mounia Lalmas, and Andrew Trotman, editors, *Advances in XML Information Retrieval: 5th International Workshop of the Initiative for the Evaluation of XML Retrieval, INEX 2006, Revised Selected Papers*, volume 4518, pages 331–344. Springer, 2007.