



## Combining text/image in WikipediaMM task 2009

Christophe Moulin, Cécile Barat, Cédric Lemaître, Mathias Géry, Christophe Ducottet, Christine Largeron

► **To cite this version:**

Christophe Moulin, Cécile Barat, Cédric Lemaître, Mathias Géry, Christophe Ducottet, et al.. Combining text/image in WikipediaMM task 2009. ECDL 2009 - Workshop CLEF, Sep 2009, Corfu, Greece. <ujm-00432319>

**HAL Id: ujm-00432319**

**<https://hal-ujm.archives-ouvertes.fr/ujm-00432319>**

Submitted on 16 Nov 2009

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Combining text/image in WikipediaMM task 2009

Christophe Moulin, Cécile Barat, Cédric Lemaître, Mathias Géry, Christophe Ducottet, Christine Largeron

Université de Lyon, F-69003, Lyon, France  
Université de Saint-Étienne, F-42000, Saint-Étienne, France  
CNRS UMR5516, Laboratoire Hubert Curien  
{christophe.moulin, cecile.barat, cedric.lemaitre, mathias.gery,  
ducottet,largeron}@univ-st-etienne.fr

**Abstract.** This paper reports our multimedia information retrieval experiments carried out for the ImageCLEF track 2009. In 2008, we proposed a multimedia document model defined as a vector of textual and visual terms weighted using a tf.idf approach [5]. For our second participation, our goal was to improve this previous model in the following ways: 1) use of additional information for the textual part (legend and image bounding text extracted from the original documents, 2) use of different image detectors and descriptors, 3) new text / image combination approach. Results allow to evaluate the benefits of these different improvements.

## 1 Introduction

ImageCLEFwiki is a multimedia collection where each document is composed of text and one image. User needs are represented by queries ("topics"), which are also multimedia. Therefore, a multimedia document model is necessary to handle such a collection. In 2008, we proposed a first model that combines text and image information for multimedia retrieval [5]. This year, we improve our model adding textual information and using different detectors and descriptors for the visual information. Moreover we use a linear combination to merge our textual and visual results. After presenting our model, we will explain the submitted runs and the obtained results. We will finish by introducing our future work.

## 2 Visual and textual document model

The document model we defined for ImageCLEF 2008 lets us rank documents depending on the query using different methods. Firstly, we explain the key features of our approach to rank documents according to a query using only textual information. Secondly, we describe how we extend the method to handle the visual information. Finally, we present our method for combining textual and visual results.

## 2.1 Textual representation model

As in the vector space model introduced by Salton et al. [7], we represent a document  $d_i$  as a vector of weights  $w_{i,j}$ . Each  $w_{i,j}$  corresponds to the importance of the term  $t_j$  in the document  $d_i$  computed by multiplying  $tf_{i,j}$  and  $idf_j$ , where  $tf_{i,j}$  is the term frequency that characterizes the frequency of the term  $t_j$  in the document  $d_i$ . The  $idf_j$  is the inverse document frequency that quantifies the importance of the term  $t_j$  over the corpus of documents.  $w_{i,j}$  is high when the term  $t_j$  is frequent in the document  $d_i$  but rare in the others. We use  $tf_{i,j}$  and  $idf_j$  defined in the Okapi formula by Robertson et al [6] by :

$$tf_{i,j} = \frac{k_1 n_{i,j}}{n_{i,j} + k_2 (1 - b + b \frac{|d_i|}{d_{avg}})}$$

where  $n_{i,j}$  is the occurrence of the term  $t_j$  in the document  $d_i$ ,  $|d_i|$  the size of the document  $d_i$  and  $d_{avg}$  the average size of all documents in the corpus and  $k_1$ ,  $k_2$  and  $b$  are three constants.

$$idf_j = \log \frac{|D| - |\{d_i | u_j \in d_i\}| + 0.5}{|\{d_i | t_j \in d_i\}| + 0.5}$$

where  $|D|$  is the size of the corpus and  $|\{d_i | t_j \in d_i\}|$  the number of documents where the term  $t_j$  occurs at least one time.

If we consider a query  $q_k$  as a short document, we can represent it as a vector of weights. A score is then computed between the query  $q_k$  and a document  $d_i$  as shown in table 1. The main difference between  $score^1$  and  $score^2$  is the representation of the query. In the first score, the weight of the query is defined by its  $tf$  only while in the second score the weight is equal to  $tf.idf$ . Note that for  $tf_{k,j}$ ,  $b = 0$  because  $|d_k|$  and  $|d_{avg}|$  are not defined for a query.

| Scoring  | Parameters                               | Parameters                        |
|--|--|-----------------------------------|
|  | $tf_{i,j}$                               | $tf_{k,j}$                        |
| $score^1(q_k, d_i) = \sum_{t_j \in q_k} tf_{i,j} idf_j tf_{k,j}$       | $k_1 = 2.2$<br>$k_2 = 1.2$<br>$b = 0.75$ | $k_1 = 8$<br>$k_2 = 7$<br>$b = 0$ |
| $score^2(q_k, d_i) = \sum_{t_j \in q_k} tf_{i,j} idf_j tf_{k,j} idf_j$ | $k_1 = 1$<br>$k_2 = 1$<br>$b = 0.5$      | $k_1 = 1$<br>$k_2 = 1$<br>$b = 0$ |

**Table 1.** Scoring equations and their default parameters[8].

Different sources of text are available. The legend provided with images is often very short and sometimes useless: for example, when the text deals with the copyright of the image or when it gives details about the user who uploaded the image. In order to gain information, we aim at using the original text extracted

from the wikipedia documents in which images appear. We consider a text fragment around the image. The size of the window is tuned using wikipediaMM 2008 collection as a training collection. We add this text to the legend of the image and we index both the added text and the original legend. The indexing is performed with the Lemur software[8].

## 2.2 Visual representation model

In order to combine the visual information with the textual one, we also represent images as a vector of weights. Provided we are able to extract visual words from images, it is possible to use the *tf.idf* formula in the same way as in the textual model. It is therefore necessary to create a visual vocabulary  $V = \{v_1, \dots, v_j, \dots, v_{|V|}\}$  as in [2]. For that purpose, we use 3 different descriptors. The first one (*meanstd*) is the same as in [5]. Each image is partitioned into 16x16 cells. Each cell is described by a six dimensional vector which corresponds to the mean and the standard deviation of  $\frac{R}{R+G+B}$ ,  $\frac{G}{R+G+B}$  and  $\frac{R+G+B}{3*255}$  where  $R$ ,  $G$  and  $B$  are the red, green and blue components of the cell. The second (named *sift<sub>1</sub>*) and third (named *sift<sub>2</sub>*) descriptors are based on the well known sift descriptor [3]. The *sift<sub>1</sub>* firstly detects regions of interest using the MSER method as in [4] while the *sift<sub>2</sub>* one uses a regular partitioning as in the *meanstd* descriptor.

For each of our 3 descriptors, we apply a  $k$ -means algorithm [1] to obtain a vocabulary of 10'000 visual terms. Each visual term represents a cluster of feature vectors.

Then, each image can be represented using a vector of visual terms. Local features are first calculated using one of the 3 descriptors. Then visual terms are determined by seeking, for each feature, the closest visual term (according to the euclidian distance) in the corresponding visual vocabulary.

In the same way as for textual words, the weight of each visual term is computed using a *tf.idf* approach.

## 2.3 Combination

In order to combine textual and visual results we use two different methods. The first one is a simple intersection between the results obtained with the textual query and with the visual one. The second one corresponds to a linear combination between the textual and the visual scores.

$$score(q_k, d_i) = \alpha score_V(q_k, d_i) + (1 - \alpha) score_T(q_k, d_i)$$

The  $\alpha$  parameter lets us add more or less visual information. We calculate its optimal value using the queries from the ImageCLEFwiki 2008 track as a training set.

### 3 Experiments

Using the model described in the previous section, we present runs submitted to ImageCLEFwiki and the results we obtained.

#### 3.1 Submitted runs

| run_id         | run                          | score                     | text          | image                    | combination         |
|----------------|------------------------------|---------------------------|---------------|--------------------------|---------------------|
| <i>LaHC_1</i>  | LaHC_TXT_okapi               | <i>score</i> <sup>1</sup> | <i>legend</i> | -                        | -                   |
| <i>LaHC_2</i>  | LaHC_TXT_tfidf               | <i>score</i> <sup>2</sup> | <i>legend</i> | -                        | -                   |
| <i>LaHC_3</i>  | run_inter_TXT_IMG_Meanstd    | <i>score</i> <sup>2</sup> | <i>legend</i> | <i>meanstd</i>           | intersection        |
| <i>LaHC_4</i>  | run_inter_TXT_IMG_Sift       | <i>score</i> <sup>2</sup> | <i>legend</i> | <i>sift</i> <sub>1</sub> | intersection        |
| <i>LaHC_5</i>  | run_TXTIMG_Meanstd_0.015     | <i>score</i> <sup>2</sup> | <i>legend</i> | <i>meanstd</i>           | $\alpha=0.015$      |
| <i>LaHC_6</i>  | run_TXTIMG_Meanstd_0.025     | <i>score</i> <sup>2</sup> | <i>legend</i> | <i>meanstd</i>           | $\alpha=0.025$      |
| <i>LaHC_7</i>  | run_TXTIMG_Sift_0.012        | <i>score</i> <sup>2</sup> | <i>legend</i> | <i>sift</i> <sub>1</sub> | $\alpha=0.012$      |
| <i>LaHC_8</i>  | run_inter_TXT_IMG_Siftdense  | <i>score</i> <sup>2</sup> | <i>legend</i> | <i>sift</i> <sub>2</sub> | <i>intersection</i> |
| <i>LaHC_9</i>  | run_TXT_100_3_1_5            | <i>score</i> <sup>2</sup> | 100 char      | <i>legend</i>            | -                   |
| <i>LaHC_10</i> | run_TXT_50_3_1_5             | <i>score</i> <sup>2</sup> | 50 char       | <i>legend</i>            | -                   |
| <i>LaHC_11</i> | run_TXTIMG_100_3_1_5_meanstd | <i>score</i> <sup>2</sup> | 100 char      | <i>meanstd</i>           | $\alpha=0.025$      |
| <i>LaHC_12</i> | run_TXTIMG_50_3_1_5_meanstd  | <i>score</i> <sup>2</sup> | 50 char       | <i>meanstd</i>           | $\alpha=0.025$      |
| <i>LaHC_13</i> | run_TXTIMG_Siftdense_0.084   | <i>score</i> <sup>2</sup> | <i>legend</i> | <i>sift</i> <sub>2</sub> | $\alpha=0.084$      |

- *meanstd*: regular partitioning + color descriptor
- *sift*<sub>1</sub>: MSER detector + sift descriptor
- *sift*<sub>2</sub>: regular partitioning + sift descriptor
- *legend*: text of the image document
- *n char*: size (*n* characters) of the text window around the image in the original wikipedia documents

**Table 2.** Presentation of the runs

All the runs are entirely automatic and are summarized on table 2. We define a baseline, *LaHC\_1*, that corresponds to a pure text model. It uses only textual terms for the query and scoring of documents. We calculate the *score*<sup>1</sup> for each image using terms of the textual content. The image name or bounding characters are not considered. We do not use neither feedback nor query expansion. Since *score*<sup>1</sup> is applied, the query terms are weighted with their frequency *tf*.

Using only the text, we perform 3 other runs: the *LaHC\_2* is the same as the baseline except that the query is represented by its *tf.idf* rather than its *tf*. The *LaHC\_9* and *LaHC\_10* are two other text only runs that make use of the bounding text around the image in the original wikipedia document. The *LaHC\_9* adds 100 characters before and after the image while the *LaHC\_10* adds 50 characters.

All other runs exploit both the textual and the visual information of documents. The *LaHC.3*, *LaHC.4* and *LaHC.8* are obtained after an intersection of the text only query results (*LaHC.2*) and the image query using the *meanstd*, the *sift<sub>1</sub>* and the *sift<sub>2</sub>* descriptors. The other runs are obtained from a linear combination of the textual and the visual scores. *LaHC.5*, *LaHC.6*, *LaHC.7* and *LaHC.13* use the textual scores of *LaHC.2* and the visual scores of a visual descriptor (*meanstd*, *sift<sub>1</sub>* and *sift<sub>2</sub>*). *LaHC.11* and *LaHC.12* combine the textual scores of *LaHC.9* and *LaHC.10* with the visual scores of the *meanstd* descriptor.

## 4 Results

All the obtained results are summarized in Table 3. On the whole results, our team ranks 2nd on 8 participants<sup>1</sup>. As we can see, the best results are obtained when we combine the image bounding text and the *meanstd* descriptor. We could have obtained better results if we had combined the image bounding text and the *sift<sub>2</sub>* descriptor. Indeed, comparing results of text-image runs *LaHC.6*, *LaHC.7* and *LaHC.13*, we can notice that the best visual descriptor is *sift<sub>2</sub>*, followed by *meanstd* and *sift<sub>1</sub>*. The last three results obtained after an intersection are the worst results, but if we compare them in term of precision, they are the best ones. Indeed, one document over six is relevant. As a rule, combining the textual information with the visual one always improves the results and return more relevant documents which is very encouraging.

| rank | run                          | map    | num_ret | num_rel_ret |
|------|------------------------------|--------|---------|-------------|
| 5    | run_TXTIMG_100.3.1.5_meanstd | 0.2178 | 44993   | 1213        |
| 6    | run_TXTIMG_50.3.1.5_meanstd  | 0.2148 | 44993   | 1218        |
| 14   | run_TXTIMG_Siftdense_0.084   | 0.1903 | 44993   | 1212        |
| 15   | run_TXT_100.3.1.5            | 0.1890 | 38004   | 1205        |
| 16   | run_TXT_50.3.1.5             | 0.1880 | 37041   | 1198        |
| 20   | run_TXTIMG_Meanstd_0.025     | 0.1845 | 44993   | 1208        |
| 21   | run_TXTIMG_Sift_0.012        | 0.1807 | 44995   | 1200        |
| 24   | run_TXTIMG_Meanstd_0.015     | 0.1792 | 44993   | 1213        |
| 33   | LaHC_TXT_tfidf               | 0.1667 | 35611   | 1192        |
| 44   | LaHC_TXT_okapi               | 0.1432 | 35611   | 1164        |
| 52   | run_inter_TXT_IMG_Siftdense  | 0.0365 | 619     | 142         |
| 53   | run_inter_TXT_IMG_Meanstd    | 0.0338 | 574     | 76          |
| 54   | run_inter_TXT_IMG_Sift       | 0.0321 | 637     | 120         |

**Table 3.** Presentation of the results

<sup>1</sup> <http://imageclef.org/2009/wikiMM-results>

## 5 Conclusion

In this article we proposed improvements to our multimedia model we introduced in [5]. The first one was to use image bounding text extracted from the original documents, the second was to use sift based image descriptors for the visual part and the third one was to add a text/image combination approach. A series of thirteen runs was submitted using the ImageCLEFwiki 2009 collection. The first analysis of the results allowed to make the three following remarks. It's better to use the image bounding text than the legend only. The sift descriptor is better than our previous color descriptor provided it is calculated on a regular partitioning. The text-image combination is a winning strategy which can be implemented by linear combination of textual and visual scores.

For future work, we aim to combine the textual information with more than just one visual descriptor. Moreover as the visual information importance depends on the query, we also plan to learn a different  $\alpha$  parameter for each query.

## 6 Acknowledgements

This work was partly supported by the LIMA project<sup>2</sup> and the Web Intelligence project<sup>3</sup> which are 2 Rhône-Alpes region projects of ISLE cluster<sup>4</sup>.

## References

1. Halil Bisgin. Parallel clustering algorithms with application to climatology. Technical report, Informatics Institute, Istanbul Technical University, Turkey., 2008.
2. Frédéric Jurie and Bill Triggs. Creating efficient codebooks for visual recognition. In *International Conference on Computer Vision*, 2005.
3. David G. Lowe. Object recognition from local scale-invariant features. In *Proceedings of the International Conference on Computer Vision ICCV*, Corfu, pages 1150–1157, 1999.
4. K. Mikolajczyk, T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir, and L. Van Gool. A comparison of affine region detectors. *International Journal of Computer Vision*, 65(1-2):43–72, 2005.
5. Christophe Moulin, Cecile Barat, Mathias Gry, Christophe Ducottet, and Christine Largeton. Ujm at imageclefwiki 2008. In *ImageCLEF 2008*, 2008.
6. Stephen E. Robertson, Steve Walker, Micheline Hancock-Beaulieu, Aaron Gull, and Marianna Lau. Okapi at trec-3. In *Text REtrieval Conference*, pages 21–30, 1994.
7. Gerard Salton, A. Wong, and C. S. Yang. A vector space model for automatic indexing. *Communications of the ACM*, 18(11):613–620, 1975.
8. C. Zhai. Notes on the lemur tfidf model. Technical report, Carnegie Mellon University, <http://www.lemurproject.com>, 2001.

---

<sup>2</sup> LIMA project: <http://liris.cnrs.fr/lima/>

<sup>3</sup> WI project: <http://www.web-intelligence-rhone-alpes.org/>

<sup>4</sup> ISLE cluster: <http://ksup-gu.grenet.fr/isle>