

# Using medGIFT and easyIR for the ImageCLEF 2005 evaluation tasks

Henning Müller, Antoine Geissbühler, Johan Marty, Christian Lovis, Patrick Ruch  
University and University Hospitals of Geneva, Service of Medical Informatics,  
24 Rue Micheli-du-Crest, CH-1211 Geneva 4, Switzerland  
henning.mueller@sim.hcuge.ch

## Abstract

This article describes the use of the *medGIFT* retrieval system for three of the four ImageCLEF 2005 retrieval tasks. We participated in the ad-hoc retrieval task that was similar to the 2004 ad-hoc task, the new medical retrieval task that required much more semantic analysis of the textual annotation than in 2004 and the new automatic annotation task. The techniques used in 2005 are fairly similar to the 2004 techniques for the two retrieval tasks. For the automatic annotation task, scripts were optimised to allow classification with a retrieval system. Unfortunately, an error in the text retrieval system corrupted part of our runs and led to relatively bad results for all runs including text. This error should be fixed before the final proceedings are printed, so correct figures are expected for this.

All retrieval results rely heavily on two retrieval systems: for visual retrieval we use the GNU Image Finding Tool (*GIFT*), and for textual retrieval the EasyIR retrieval system. For the ad-hoc retrieval task, two runs were submitted with different configurations of grey levels and the Gabor filters. No textual retrieval was attempted, but only purely visual retrieval, resulting in generally lower scores than text retrieval. For the medical retrieval task, visual retrieval was performed with several configurations of Gabor filters and grey level and color quantisations as well as several variations of combining text and visual features. Unfortunately, all these runs are broken as the textual retrieval results are almost random. Due to a lack of resources no relevance feedback runs were submitted, which is where *medGIFT* performed best in 2004. For the classification task, a retrieval with the image to classify was performed and the first  $N = 1, 5, 10$  resulting images were used to calculate scores for the classes by simply adding up the score of the  $N$ -images for each class. No machine learning was performed on the data of the known classes, so the results are surprisingly good and were only topped by systems with sophisticated learning strategies optimised for the used data set.

## Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: H.3.1 Content Analysis and Indexing; H.3.3 Information Search and Retrieval; H.3.4 Systems and Software; H.3.7 Digital Libraries

## General Terms

Image retrieval, evaluation, visual retrieval

## Keywords

GIFT, easyIR, visual/textual retrieval, image retrieval

# 1 Introduction

Image retrieval is an increasingly important domain in the field of information retrieval. Compared to text retrieval little is known about how to search for images, although it has been an extremely active domain as well in the field of computer vision as in information retrieval [8, 12, 13, 16]. Benchmarks such as ImageCLEF [2, 3] allow us to actually evaluate our algorithms compared to other systems and deliver us an insight into the techniques that perform well and those that do not perform as good. Thus, new developments can be directed towards these goals and techniques of other well-performing systems can be adapted to our needs.

In 2005, the ad-hoc retrieval task created topics that were better adapted for visual systems using the same database as in 2004. The tasks made available contained three images, so more visual information. We submitted two configurations of our system to this task using visual information only.

The medical retrieval task was performed on a much larger database than in 2004 containing a total of more than 50.000 images [4]. The annotation was also more varied, ranging from a few words in a very structured form to completely unstructured paragraphs. This made it hard to preprocess any of the information, so finally only free-text retrieval was used for our results submission including all XML tags. Also, the tasks were much harder and mainly semantic query tasks, which made the retrieval by visual means more difficult. Due to a lack of resources we could only submit partial results that did not include any relevance feedback or automatic query expansion.

The automatic annotation task was very interesting and challenging at the same time [6]. We did not take into account any of the training data and simply used the retrieval system *GIFT* and a nearest neighbour technique to classify the results. Still, the results were surprisingly good (6th best submission, 2nd best group) and when taking into account the learning data using an approach as described in [10], these results are expected to get better.

ImageCLEF gave us the opportunity to compare our system with other techniques which is invaluable and will provide us with directions for future research.

## 2 Basic Technologies Used

For our ImageCLEF participation, we aim at combining content-based retrieval of images with cross-language retrieval applied on the textual annotation of the images. Based on the results from last year (2004), we used parameters that were expected to lead to good results, plus some new combinations.

### 2.1 Image Retrieval

The technology used for the content-based retrieval of images is mainly taken from the *Viper*<sup>1</sup> project of the University of Geneva. Much information about this system is available [14]. Outcome of the *Viper* project is the GNU Image Finding Tool, *GIFT*<sup>2</sup>. This software tool is open source and can in consequence also be used by other participants of ImageCLEF. A ranked list of visually similar images for every query topic was made available for participants and will serve as a baseline to measure the quality of submissions. Demonstration versions with a web-accessible interface of *GIFT* were also made available for participants to query visually as with feedback in an interactive way as not everybody can be expected to install an entire Linux tool for such a benchmark to use *GIFT*. The feature sets that are used by *GIFT* are:

- Local color features at different scales by partitioning the images successively into four equally sized regions (four times) and taking the mode color of each region as a descriptor;

---

<sup>1</sup><http://viper.unige.ch/>

<sup>2</sup><http://www.gnu.org/software/gift/>

- global color features in the form of a color histogram, compared by a simple histogram intersection;
- local texture features by partitioning the image and applying Gabor filters in various scales and directions. Gabor responses are quantised into 10 strengths;
- global texture features represented as a simple histogram of responses of the local Gabor filters in various directions and scales.

A particularity of *GIFT* is that it uses many techniques well-known from text retrieval. Visual features are quantised and the feature space is very similar to the distribution of words in texts, corresponding roughly to a Zipf distribution. A simple *tf/idf* weighting is used and the query weights are normalised by the results of the query itself. The histogram features are compared based on a histogram intersection [15].

The medical version of the *GIFT* is called *medGIFT*<sup>3</sup> [9]. It is also accessible as open source and adaptations concern mainly visual features and the user interface that shows the diagnosis on screen and is linked with a radiologic teaching file so the MD can not only browse images but also get the textual data and other images of the same case. Grey levels play a more important role for medical images and their numbers are raised, especially for relevance feedback (RF) queries. The number of the Gabor filter responses also has an impact on the performance and these are changed with respect to directions and scales. We used in total 4, 8 and 16 grey levels and for the Gabor filters we used 4 and 8 directions. Other techniques in *medGIFT* such as a pre-treatment of images [7] were not used for this competition due to a lack of resources.

## 2.2 Textual Search

The basic granularity of the Casimage and MIR collections is the case. A case gathers a textual report, and a set of images. For the PathoPic and Peir databases annotation exists for every image. The queries contain one to three images and text in three languages. We used all languages as a single query and also indexed all documents in a single index. Case-based annotation was expanded to all images of the case after the retrieval step, so for us the final unit of retrieval is the image.

### 2.2.1 Indexes

Textual experiments were conducted with the easyIR engine<sup>4</sup>. As a single report is able to contain written parts in several languages mixed, it would have been necessary to detect the boundaries of each language segment. Ideally, French, German and English textual segments would be stored in different indexes. Each index could have been translated into the other language using a general translation method, or more appropriately using a domain-adapted method [11]. However, such a complex architecture would require to store different segments of the same document in separate indexes. Considering the lack of data to tune the system, we decided to index all collections using a unique index using an English stemmer, For simplicity reasons, the XML tags were also indexed and not separately treated.

### 2.2.2 Weighting Schema

We chose a generally good weighting schema of the term frequency - inverse document frequency family. Following weighting convention of the SMART engine, cf. Table 1, we used *atc-ltn* parameters, with  $\alpha = \beta = 0.5$  in the augmented term frequency.

---

<sup>3</sup><http://www.sim.hcuge.ch/medgift/>

<sup>4</sup><http://lithwww.epfl.ch/~ruch/softs/softs.html>

| Term Frequency             |  |
|----------------------------|--|
| First Letter               | $f(tf)$  |
| n (natural)                | $tf$   |
| l (logarithmic)            | $1 + \log(tf)$   |
| a (augmented)              | $\alpha + \beta \times (\frac{tf}{max(tf)})$ , where $\alpha = 1 - \beta$ and $0 < \alpha < 1$ |
| Inverse Document Frequency |  |
| Second Letter              | $f(\frac{1}{df})$  |
| n(no)                      | 1  |
| t(full)                    | $\log(\frac{N}{df})$   |
| Normalisation              |  |
| Third Letter               | $f(length)$  |
| n(no)                      | 1  |
| c(cosine)                  | $\sqrt{\sum_{i=1}^t w_{i,j}^2} \times \sqrt{\sum_{j=1}^t w_{j,q}^2}$                           |

Table 1: Usual *tf-idf* weight; for the cosine normalisation factor, the formula is given for Euclidean space:  $w_{i,j}$  is the document term weight,  $w_{j,q}$  is the query term weight.

### 2.3 Combining the Two

Combinations of visual and textual features for retrieval are rather scarce in the literature [5], so many of the mechanism and fine tuning of the combinations will still need more work, especially when the optimisation is based on the actual query. For the visual query we used all images that are present for a query, including one query containing negative feedback. For the text part, the text of all three languages was used as a combined query together with the combined index that includes the documents in all languages. Results list of the first 1000 documents were taken into account for both the visual and the textual search. Both result lists were normalised to deliver results within the range [0; 1]. The visual result is normalised by the result of the query itself whereas the text was normalised by the document with the highest score. This leads to visual results that are usually slightly lower than the textual results.

To combine the two lists, two different methods were chosen. The first one simply combines the list with different percentages for visual and textual results (textual= 50, 33, 25, 10%). In a second form of combination the list of the first 1000 visual results was taken, and then, all those that were in the first 200 textual documents were multiplied with N-times the value of the textual results.

## 3 The ad hoc retrieval task

For the ad-hoc retrieval task we submitted results using fairly similar techniques as those in 2004. The 2005 topics were actually more adapted to the possibilities of visual retrieval systems as more visual attributes were taken into account for the topic creation. Still, textual retrieval stays very necessary for good results. It is not so much a problem of the queries but rather a problem of the database containing mostly grey or brown scale images of varying quality where automatic treatment such as color indexing is difficult. This should change in 2006 with a new database using mostly consumer pictures of vacation destinations. Such a database could be better analysed automatically using the available color information

We used the *GIFT* system in two configurations, once the normal *GIFT* engine with 4 grey levels and the full HSV space using the Gabor filter responses in four directions and at three scales. The second configuration took into account 8 grey levels as the 2004 results for 16 grey levels were actually much worse than expected. We also raised the number of directions of the Gabor filters to 8 instead of four. The results of the basic *GIFT* system were made available to all participants and used by several. Surprisingly the results of the basic *GIFT* system remain the best in the test with a MAP of 0.0829, being at the same time the best purely visual system participating.

The system with eight grey levels and eight directions for the Gabor filters performed slightly worse and a MAP of 0.0819 was reached. Other visual systems performed slightly lower. The best mono-lingual text systems performed at a MAP of 0.41. Several text retrieval systems performed worse than the visual system for a variety of languages.

## 4 The automatic annotation task

We were new to the automatic annotation task as almost everyone and had mainly used our system for retrieval, so far. Due to a lack of resources no optimisation using the available training data was performed. Still, the tf/idf weighting is automatically weighting rare features higher which leads to a discriminative analysis.

As techniques we performed a query with each of the 1000 images to classify and took into account the first  $N = 1, 5, 10$  retrieval results. For each of these results images from the training set the correct class was determined and this class was thus augmented with the similarity score of the image. The class with the highest final score became automatically the final class selected for the image. For retrieval we used three different settings of the features using 4, 8, and 16 grey levels. The runs with 8 and 16 grey levels also had eight directions of the Gabor filters for indexation. Best results obtained in the competition were from the Aachen groups (best run at 12.6% error rate) that have been working on very similar data for several years, now.

The best results for our system were retrieved when using 5NN and eight grey levels (error rate 20.6%), and the next best results using 5NN and 16 grey levels (20.9). Interestingly, the worst results were obtained with 5NN and 4 grey levels (22.1). Using 10NN led to slightly worse results (21.3) and 1NN was rather in the middle (4 grey levels 21.8; 8 grey levels: 21.1; 16 grey levels 21.7).

As a result we can say that all results are extremely close together 20.6-22.1 %, so the differences do not seem statistically significant. 5NN seems to be the best but this might also be linked to the fact that some classes have a very small population and 10NN would simply retrieve too many image of other classes to be competitive. 8 levels of grey and 8 directions of the Gabor filters seem to perform best, but the differences are still very small.

In the future it is planned to train the system with the available training data using the algorithm described in [10]. This technique is similar to the market basket analysis [1]. A proper strategy for the training needs to be developed to especially help smaller classes to be well classified. These classes cause normally most of the classification problems.

## 5 The medical retrieval task

Unfortunately, our textual retrieval results submitted contained an indexation error, and so the textual results were almost random. We have not identified the error yet, but hope to have it found before the final proceedings are printed. Thus the only textual run that we submitted had only a MAP of 0.0226, whereas the best textual retrieval systems were at 0.2084 (IPAL/I2R). Due to a limitation of resources, we were not able to submit relevance feedback runs, which is the discipline where *GIFT* usually is strongest. The best feedback system was OHSU with a map of 0.2116 for only textual retrieval.

The best visual system is I2R with a map of 0.1455. Our *GIFT* retrieval system was made available to participants and was widely used. Again, the basic *GIFT* system obtained the best results among the various combinations in feature space (map 0.0941), with only i2r having actually better results. The second indexation using 8 grey levels and eight directions of the Gabor filters performs slightly worse at 0.0872.

For mixed textual/visual retrieval, the best results were obtained by IPAL/I2R with map 0.2821. Our best result in this category is using 10% textual part and 90% visual part and obtains 0.0981. These results should be much better when using a properly indexed text base. The following results were obtained for other combinations: 20% visual: 0.0934, 25%: 0.0929, 33%:

0.0834, 50%: 0.044. When using eight grey levels and 8 Gabor directions: 10% visual: 0.0891, 20%: 0.084, 33%: 0.075, 50%: 0.0407. The results could lead to the assumption that visual retrieval is better than textual retrieval in our case, but this holds only true because of our indexation error. We will try to fix the error and deliver proper results as soon as possible to have a correct comparison with the other groups.

A second combination technique that we applied used as a basis the results from textual retrieval and then added the visual retrieval results multiplied with a factor  $N = 2, 3, 4$  to the first 1000 results of textual retrieval. This strategy proved fruitful in 2004 the other way round by taking first the visual results and then augmenting only the first  $N=1000$  results. The results for the main *GIFT* system were: 3 times visual: 0.0471, 4 times visual 0.0458, 2 times visual 0.0358. For the system with 8 grey levels, the respective results are: 3 times visual 0.0436, 4 times visual 0.0431, 2 times visual 0.0237. A reverse order of taking the visual results first and then augment the textually similar would have led to better results in this case but when having correct results for text as well as for visual retrieval, this needs to be proven.

We cannot really deduct extremely much of our current submission as several errors prevented better results.

## 6 Conclusions

Although we did not have any resources for an optimised submission we still learned from the 2005 tasks that the *GIFT* system delivers a good baseline for image retrieval and that it is widely usable for a large number of tasks and different images.

More detailed results show that the ad-hoc task is hard for visual retrieval even with a more visually-friendly set of queries as the image set does not contain enough color information or clear objects, which is crucial for fully visual information retrieval.

The automatic annotation or classification task proved that our system delivers good results even without learning and shows that information retrieval can also be used well for document classification. When taking into account the available training data these results will surely improve significantly.

From the medical retrieval task not much can be deduced for now as we need to work on our textual indexation and retrieval to find the error responsible for the mediocre results. Still, we can say that *GIFT* is well suited and among the best systems for general visual retrieval. It will need to be analysed which features were used by other systems, especially the few runs that performed better.

For next year we will definitely have to take into account the available training data and we hope as well to use more complex algorithms for example to extract objects from the medical images and limit retrieval to these objects. Another strong point of *GIFT* is the good relevance feedback and this can surely improve results significantly as well. Already the fact to have a similar databases for two years in a row would help as such large databases need a large time to be indexed and require human resources for optimisation as well.

## 7 Acknowledgements

Part of this research was supported by the Swiss National Science Foundation with grant 632-066041.

## References

- [1] Rakesh Agrawal and Ramakrishnan Srikant. Fast algorithms for mining association rules. In *Proceedings of the 20th VLDB Conference*, pages 487–499, Santiago, Chile, September 12–15 1994.

- [2] Paul Clough, Henning Müller, and Mark Sanderson. Overview of the CLEF cross-language image retrieval track (ImageCLEF) 2004. In Carol Peters, Paul D. Clough, Gareth J. F. Jones, Julio Gonzalo, M. Kluck, and B. Magnini, editors, *Multilingual Information Access for Text, Speech and Images: Result of the fifth CLEF evaluation campaign*, Lecture Notes in Computer Science, Bath, England, 2005. Springer-Verlag.
- [3] Paul Clough, Mark Sanderson, and Henning Müller. A proposal for the CLEF cross language image retrieval track (ImageCLEF) 2004. In *The Challenge of Image and Video Retrieval (CIVR 2004)*, Dublin, Ireland, July 2004. Springer LNCS 3115.
- [4] William Hersh, Henning Müller, Paul Gorman, and Jeffery Jensen. Task analysis for evaluating image retrieval systems in the ImageCLEF biomedical image retrieval task. In *Slice of Life conference on Multimedia in Medical Education (SOL 2005)*, Portland, OR, USA, June 2005.
- [5] Marco La Cascia, Saratendu Sethi, and Stan Sclaroff. Combining textual and visual cues for content-based image retrieval on the world wide web. In *IEEE Workshop on Content-based Access of Image and Video Libraries (CBAIVL'98)*, Santa Barbara, CA, USA, June 21 1998.
- [6] Thomas M. Lehmann, Mark O. Güld, Thomas Deselaers, Henning Schubert, Klaus Spitzer, Hermann Ney, and Berthold B. Wein. Automatic categorization of medical images for content-based retrieval and data mining. *Computerized Medical Imaging and Graphics*, 29:143–155, 2005.
- [7] Henning Müller, Joris Heuberger, and Antoine Geissbuhler. Logo and text removal for medical image retrieval. In *Springer Informatik aktuell: Proceedings of the Workshop Bildverarbeitung für die Medizin*, Heidelberg, Germany, March 2005.
- [8] Henning Müller, Nicolas Michoux, David Bandon, and Antoine Geissbuhler. A review of content-based image retrieval systems in medicine – clinical benefits and future directions. *International Journal of Medical Informatics*, 73:1–23, 2004.
- [9] Henning Müller, Antoine Rosset, Jean-Paul Vallée, and Antoine Geissbuhler. Integrating content-based visual access methods into a medical case database. In *Proceedings of the Medical Informatics Europe Conference (MIE 2003)*, St. Malo, France, May 2003.
- [10] Henning Müller, David McG. Squire, and Thierry Pun. Learning from user behavior in image retrieval: Application of the market basket analysis. *International Journal of Computer Vision*, 56(1–2):65–77, 2004. (Special Issue on Content-Based Image Retrieval).
- [11] Patrick Ruch. Query translation by text categorization. In *Proceedings of the conference on Computational Linguistics (COLING 2004)*, Geneva, Switzerland, August 2004.
- [12] Yong Rui, Thomas S. Huang, Michael Ortega, and Sharad Mehrotra. Relevance feedback: A power tool for interactive content-based image retrieval. *IEEE Transactions on Circuits and Systems for Video Technology*, 8(5):644–655, September 1998. (Special Issue on Segmentation, Description, and Retrieval of Video Content).
- [13] Arnold W. M. Smeulders, Marcel Worring, Simone Santini, Armarnath Gupta, and Ramesh Jain. Content-based image retrieval at the end of the early years. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22 No 12:1349–1380, 2000.
- [14] David McG. Squire, Wolfgang Müller, Henning Müller, and Thierry Pun. Content-based query of image databases: inspirations from text retrieval. *Pattern Recognition Letters (Selected Papers from The 11th Scandinavian Conference on Image Analysis SCIA '99)*, 21(13–14):1193–1198, 2000. B.K. Ersboll, P. Johansen, Eds.
- [15] Michael J. Swain and Dana H. Ballard. Color indexing. *International Journal of Computer Vision*, 7(1):11–32, 1991.

- [16] Hemant D. Tagare, C. Jaffe, and James Duncan. Medical image databases: A content-based retrieval approach. *Journal of the American Medical Informatics Association*, 4(3):184–198, 1997.