

Crowdsourcing for Medical Image Classification

Alba G. Seco de Herrera, Antonio Foncubierta-Rodríguez,
Dimitrios Markonis, Roger Schaer, Henning Müller

University of Applied Sciences Western Switzerland (HES-SO),
Techno-Pôle 3 · CH-3960 Sierre
alba.garcia@hevs.ch

Abstract. To help managing the large amount of biomedical images produced, image information retrieval tools have been developed to help accessing the right information at the right moment. To provide a test bed for image retrieval evaluation the ImageCLEFmed benchmark proposes a biomedical classification task that focuses on determining the image modality of figures from biomedical journal articles automatically.

In the training data for this machine learning task some classes have many more images than others and thus a few classes are not well represented, which is a challenge for automatic image classification. To address this problem, an automatic training set expansion was first proposed. To improve the accuracy of the automatic training set expansion, a manual verification of the training set is done using the Crowdsourcing platform Crowdflower. This platform allows using external persons and pay for the crowdsourcing or use personal contacts free of charge. Crowdsourcing requires strict quality control or using trusted persons but it can quickly give access to a large number of judges and thus improve many machine learning tasks. Results show that the manual annotation of a large amount of biomedical images carried out in this work can help at image classification.

Keywords. Medical Retrieval; Modality Classification; Crowdsourcing; ImageCLEF;

1 Introduction

Images are produced in hospitals in ever-increasing numbers [1] and provide crucial information for diagnosis, treatment planning and other tasks. Besides in clinical settings, images are also made available via biomedical publications. The biomedical open access literature of PubMed Central alone contained almost 2 million images in 2014. This creates a need for searching in the immense collection of images in institutions and on the World Wide Web, making the data accessible for reuse. Many tools have been developed for these tasks over the past 20 years [2].

Retrieval and classification of medical images have been explored to get additional information for reading and interpretation of medical cases [3] when open questions remain and thus help clinicians in their daily work. Although text queries are commonly used, the visual information of the images can enrich the search. Thanks to benchmarks such as ImageCLEF¹ or Visceral² [4] the retrieval and classification algorithms have been further studied and compared with sometimes

¹ <http://imageclef.org/>

² <http://www.visceral.eu/>

more than 20 research groups participating. In particular, ImageCLEFmed [5] has been proposing several retrieval and classification tasks since 2005. The image and case-based retrieval tasks and the modality classification task are using articles from the biomedical open access literature. The goal of the retrieval tasks is to retrieve images/cases that are similar to a given image/case description. The image modality classification (for modalities such as X-ray, ultrasound, computed tomography, etc.) is used as one of the most important filters to limit the search results in existing systems. Such filtering can improve the precision of the search [6] and reduce the search space [7]. ImageCLEFmed also proposes a medical image classification task based on a proposed hierarchy including medical modalities and other images types appearing in the biomedical literature.

In ImageCLEF 2013 a training set consisting of ~2900 images was distributed to participants and classification methods were evaluated with a test set of ~2600 images. Both sets were obtained from a subset of PubMed Central³ of more than 300,000 images. In the training set some of the image categories were represented by only few images. Therefore, a training set expansion strategy was applied to our multimodal (visual and textual based) classification approach to improve the accuracy precision (from 69.63% to 71.87%) [8].

Crowdsourcing has recently emerged as a tool in bioinformatics for solving large volume of simple human tasks [9]. In this article we propose to use crowdsourcing for two tasks: to verify the automatically detected modality of ~17'000 images and to reclassify the images identified as wrongly classified. Each of these tasks can be solved in a short amount of time (a few seconds) by users familiar with medical images. A short tutorial is also given in the crowdsourcing platform to explain the task and allow quality control. Crowdsourcing was recently used for image annotation in medical imaging, e.g. for evaluation of medical pictograms [10] or for retinal fundus photography classification [11]. Results shows that the manual annotation can improve automatic classification tasks.

The remainder of the paper is organized as follows. Section 2 describes the retrieval and classification tasks as well as the description of the crowdsourcing performance. Section 3 presents the results and Section 4 concludes the paper.

2 Methods

In this section the ImageCLEFmed tasks used in this work are presented. The details of the crowdsourcing performance are also explained.

2.1 Medical ImageCLEF tasks

The ImageCLEFmed benchmark proposes a standard test bed for medical image retrieval that allows researchers to compare their approaches on large and varied data sets including manually generated ground truth [2]. The image-based retrieval task aims to retrieve images for a precise information need expressed through text and example images. On the other hand, the case-based task aims to retrieve cases that are similar to the query case and are useful in differential diagnosis.

Using the modality information of the images can help in the retrieval process to focus on one modality or to remove non clinical images entirely, thus improving the retrieval performance

³ <http://www.ncbi.nlm.nih.gov/pmc/>

[12]. The goal of the ImageCLEFmed modality classification task [5] is to classify the images into medical modalities and other images types, such as Computed Tomography, X-ray or general graphs using the modality hierarchy shown in Figure 1. The work presented in this paper aims to improve the modality classification accuracy to integrate it into the medical retrieval system to enhance and filter its results.

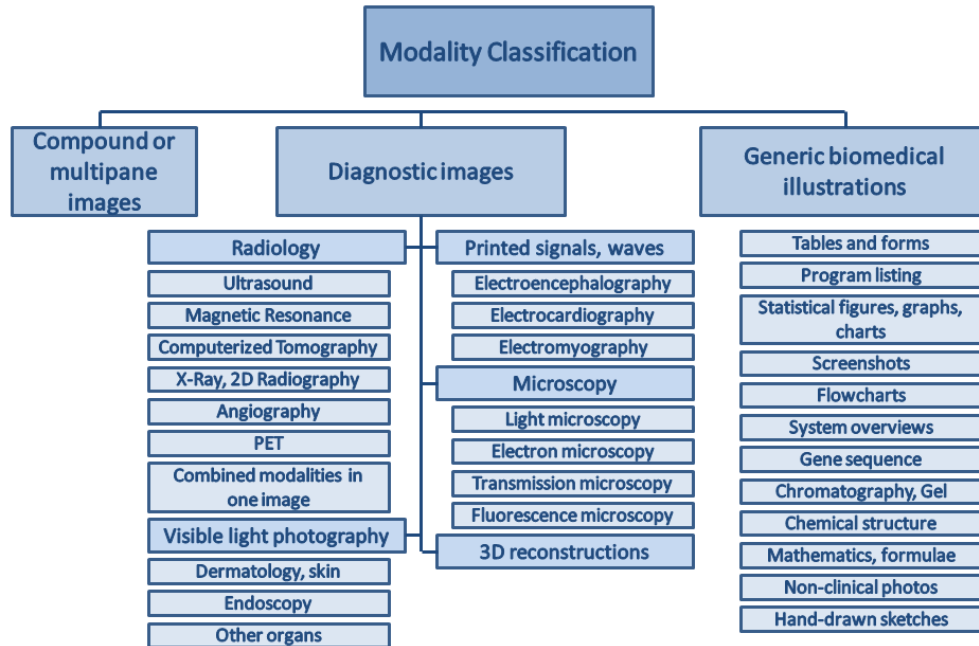


Figure 1. The image class hierarchy used for image classification.

2.2 Training set expansion

Previous work [13] describes the baseline used for the automatic image modality classification. It consists of a multimodal approach based on multiple visual descriptors and a Lucene⁴ baseline using text information. This approach achieves almost 69% of accuracy.

To improve the classification accuracy a training set expansion strategy was applied to better represent all image categories [8]. For this purpose, the dataset of ImageCLEF 2013 medical image retrieval task was used. Each image from the original training set was queried against this dataset and the top results were assigned the class labels of the query. Results that were retrieved by multiple queries belonging to different classes were discarded. This automatic labeling resulted into a larger but “noisy” training set.

2.3 Crowdsourcing

Continuing the work in [14], the Crowdfunder⁵ platform was chosen since it provides an internal interface to be used by a known set of experts. For our experiments, eight experts in the medical imaging domain participated in the crowdsourcing job.

Given the training set from the modality classification task with ~2900 images, an expanded set of images was automatically classified as described in Section 2.2. The internal crowdsourcing

⁴ <http://lucene.apache.org/>

⁵ <http://www.crowdfunder.com/>

interface was used to verify the automatically assigned class for each of the 17002 images of the new set.

A first crowdsourcing task was set up to verify the given tag but, since a large amount of images are compound or multipane images (about 50% of the figures in the biomedical open access literature [15]), an option to correctly define this class of images was added. Therefore, each image was presented with a key question formulated as follow (See Figure 2):

“Does the figure correspond to the category?:

- Yes, perfect classification
- No, compound image
- No, wrong category
- Not sure”

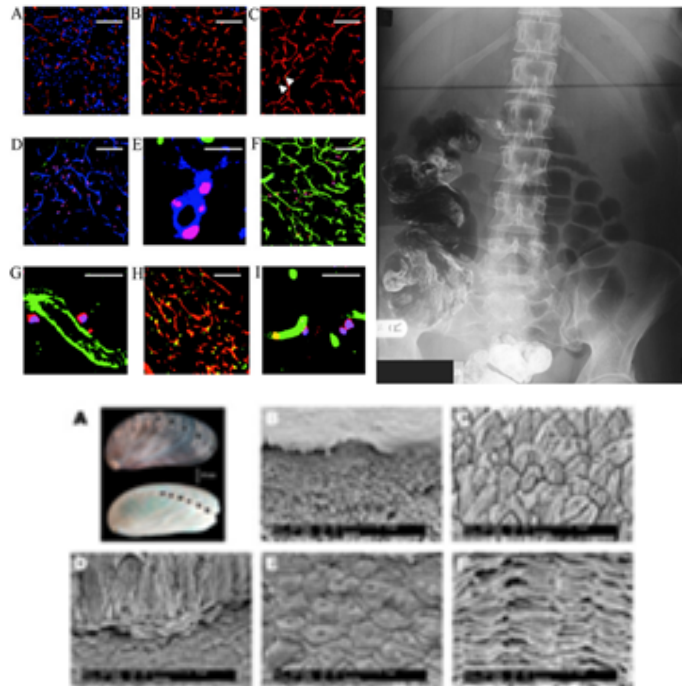


Figure 2. Images automatically classify as “Compound”, “X-ray” and “Electron microscopy” respectively. Crowdsourcing was used to verify this image modality classes.

Using an iterative approach as in [14] images that were wrongly classified by the automated training set expansion were manually reclassified in a second crowdsourcing iteration. The same procedure was applied to the “not sure” category.

In this case the user has to decide to which class each of images belongs across the hierarchy presented in Figure 1. The task was therefore presented in a hierarchical structure where a broad class is first asked and then the subclass. (See Figure 3).

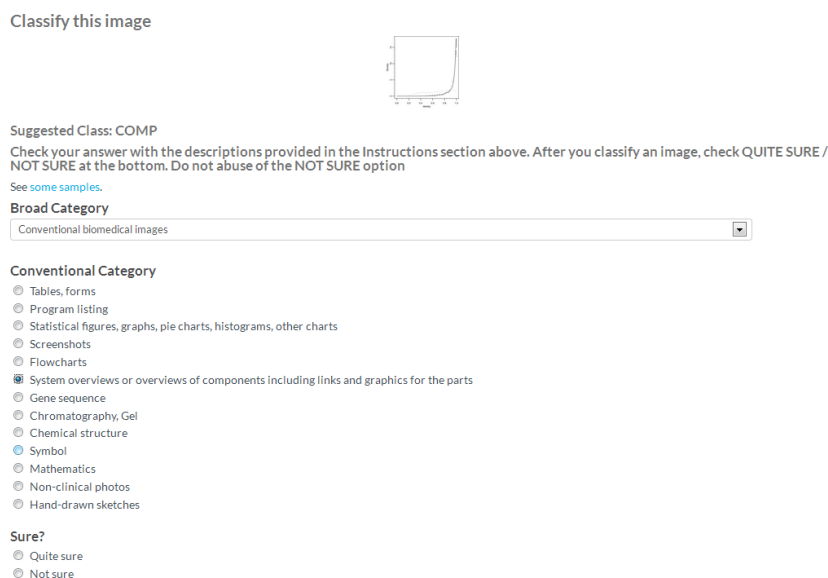


Figure 3: Screenshots of the crowdsourcing interface for image modality classification.

Since this was a more difficult task results, each of the image was classified by two users. A third user judged the images in case of disagreement between the first two answers.

3 Results

The crowdsourcing in this experiment was done with an internal team to limit errors in the classification. A total of eight experts in the medical imaging domain participated in the task. In the past external experts were used but strong quality control is necessary in this case but on the other hand tasks can be finished extremely quickly.

In the first step, 50% of the images were verified by crowdsourcing to augment the training set and automatically classify the remaining images. The results of the crowdsourcing task show that the automatic classification achieved an accuracy of 60% for this additional data set. Thanks to the first question in the platform 21% of the images were reclassified as compound figures during the same crowdsourcing job. Almost 20 % of the images then had to be reclassified manually (see Figure 4).

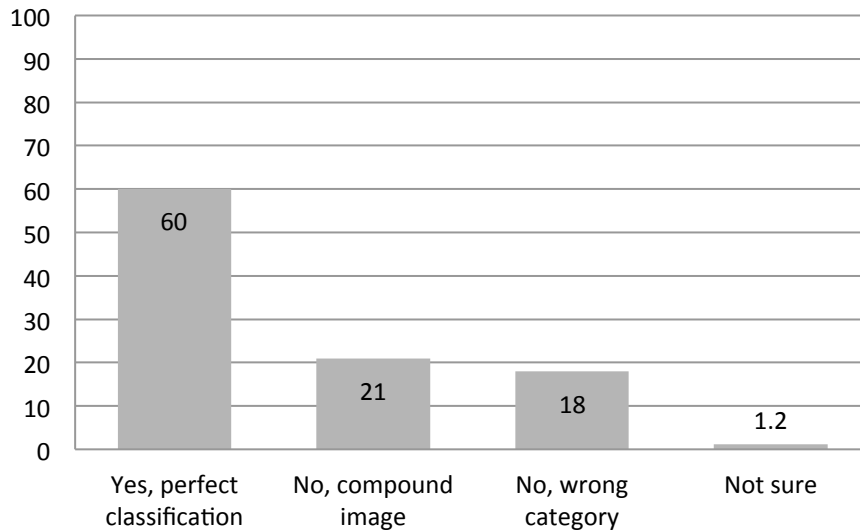


Figure 4. Each bar represents the distribution of each of the answer in the verification crowdsourcing task.

In the second part of the experiment, the correctly classified images and the images classified as compound were added to the initial training set. The new training set containing more than 9000 images was used to automatically reclassify the non-labeled images (images tagged as “wrong category” or “not sure”). A total of ~1600 images automatically reclassified and then verified via crowdsourcing. 16% of the previously wrongly classified images were now well labeled after the automatically reclassification. Figure 5 shows some images that were correctly reclassified.

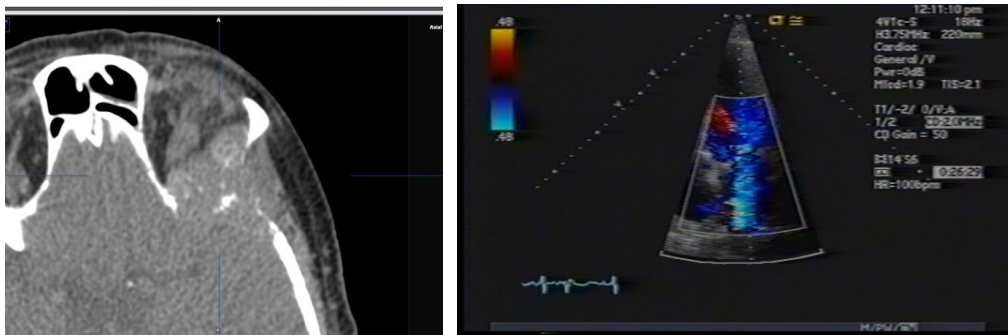


Figure 5. Images correctly reclassified after the training set expansion verification.

In general, the second crowdsourcing task was more difficult for the experts. More knowledge about the modality classes was necessary and indeed the classes were not always easy to identify. Figure 6 shows some examples of images incorrectly labeled which experts found also difficult to classify. Often full size versions of the images were necessary to clearly determine the modality.

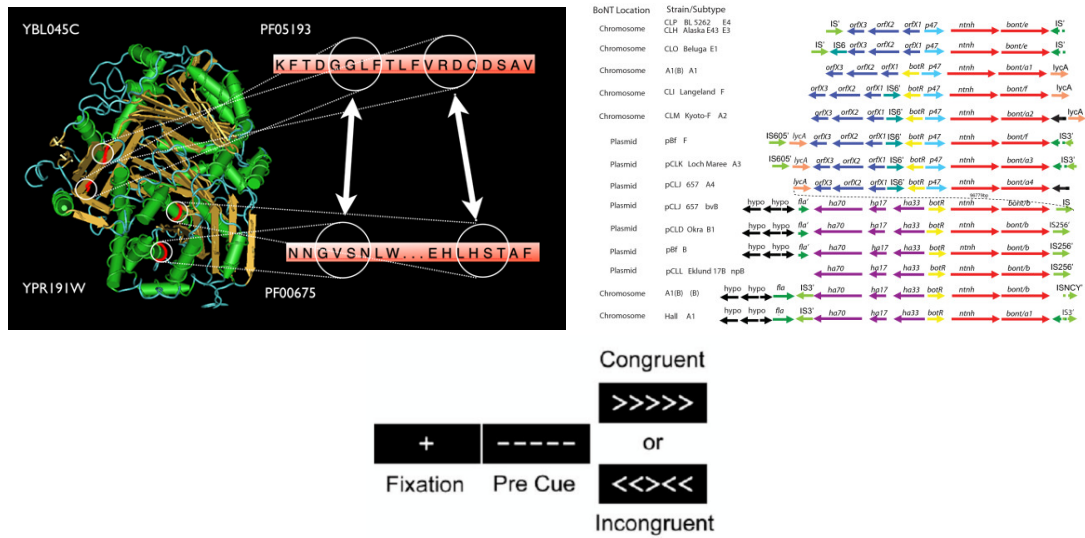


Figure 6. Images incorrectly classified automatically but that were also difficult to classify manually.

Experiments showed that the automatic expansion of the training set improves the accuracy of the ImageCLEFmed modality classification task from ~69% to ~72%. Even better accuracy is expected using the new manually verified and larger training sets.

4 Conclusions

The goal of this work was to use crowdsourcing for improving the quality of an automatic modality classification task that uses visual information of the images and the text of the figure captions. Increasing the size of the training set showed to improve the quality of the automatic classification. Manual correction of such a noisy training set can also significantly improve performance and a crowdsourcing platform can help to ease the process with a simple environment that can be used free of charge with known persons but can also be used with stricter quality control using a larger number of people participating in these platforms for a very limited cost.

The iterative nature of the shown task will be continued to progressively generate a large and discriminative training set so all images of PubMed Central can be automatically classified and then made accessible for retrieval tasks.

Acknowledgments

The research leading to these results has received funding from the European Union's Seventh Framework Program under grant agreement 257528 (KHRESMOI).

References

- [1] C. Akgül, D. Rubin, S. Napel, C. Beaulieu, H. Greenspan and B. Acar, "Content--Based Image Retrieval in Radiology: Current Status and Future," *Digital Imaging*, pp. 208-222, 2011.
- [2] J. Kalpathy-Cramer, A. García Seco de Herrera, D. Demner-Fushman, S. Antani, S. Bedrick and H. Müller, "Evaluating Performance of Biomedical Image Retrieval Systems- an Overview of the Medical Image Retrieval task at ImageCLEF 2004-2014," *Computerized Medical Imaging and Graphics*, 2014.
- [3] E. Uwimana and . M. E. Ruiz, "Integrating an automatic classification method into the medical image retrieval process," in *AMIA Annual Symposium Proceedings*, 2008.
- [4] G. Langs, H. Müller, . B. H. Menze and . A. Hanbury, "VISCERAL: Towards Large Data in Medical Imaging - Challenges and Directions," in *MCBR-CDS MICCAI workshop*, 2013.
- [5] A. García Seco de Herrera, J. Kalpathy-Cramer, D. Demner Fushman, S. Antani and H. Müller, "Overview of the ImageCLEF 2013 medical tasks," in *CLEF 2013 (Cross Language Evaluation Forum)*, Valencia, Spain, 2013.
- [6] M. M. Rahman, D. You, M. S. Simpson, S. K. Antani, D. Demner-Fushman and G. R. Thoma, "Multimodal Biomedical Image Retrieval using Hierarchical Classification and Modality Fusion," *International Journal of Multimedia Information Retrieval*, vol. 23, pp. 159-173, 2013.
- [7] J. Kalpathy-Cramer and W. Hersh, "Multimodal medical image retrieval: image categorization to improve search precision," in *International Conference on Multimedia Information Retrieval*, New York, USA, 2010.
- [8] D. Markonis, A. García Seco de Herrera and H. Müller, "Semi-Supervised Learning for Medical Image," in *Conference and Labs of the Evaluation Forum (CLEF)*, Sheffield, UK, Submitted.
- [9] B. M. Good and A. I. Su, "Crowdsourcing for Bioinformatics," *Bioinformatics*, vol. 29, no. 16, pp. 1925-1933, 2013.
- [10] B. Yu, M. Willis, P. Sun and J. Wang, "Crowdsourcing Participatory Evaluation of Medical Pictograms using Amazon Mechanical Turk," *Journal of Medical Internet Research*, vol. 15, no. 6, 2013.
- [11] D. Mitry, T. Peto, S. Hayat, J. E. Morgan, K.-T. Khaw and P. J. Foster, "Crowdsourcing as a Novel Technique for Retinal Fundus Photography Classification: Analysis of Images in the EPIC Norfolk Cohort on Behalf of the UKBiobank Eye and Vision Consortium," *PLoS One*, vol. 8, no. 8, 2013.
- [12] A. García Seco de Herrera and H. Müller, "Fusion Techniques in Biomedical Information Retrieval," in *Information Fusion in Computer Vision for Concept Recognition*, 2014, pp. 209-228.

- [13] A. García Seco de Herrera, D. Markonis, R. Schaer and I. Eggel, "The medGIFT Group in ImageCLEFmed 2013," in *CLEF 2013 (Cross Language Evaluation Forum)*, Valencia Spain, 2013.
- [14] A. Foncubierta-Rodríguez and H. Müller, "Ground Truth Generation in Medical Imaging: A Crowdsourcing based Iterative Approach," in *Workshop on Crowdsourcing for Multimedia, ACM Multimedia*, Nara, Japan, 2012.
- [15] A. Chhatkuli, D. Markonis, A. Foncubierta-Rodríguez, F. Meriaudeau and H. Müller, "Separating Compound Figures in Journal Articles to allow for Subfigure Classification," in *SPIE Medical Imaging*, Orlando,FL,USA, 2013.