# Combining Medical Domain Ontological Knowledge and Low-level Image Features for Multimedia Indexing

## Dina Demner-Fushman, Sameer K. Antani, Matthew Simpson, George R. Thoma

Lister Hill National Center for Biomedical Communications,
National Library of Medicine, NIH, Bethesda, MD
{ ddemner, santani, simpsonmatt, gthoma} @mail.nih.gov

**Abstract**

Biomedical images are invaluable in establishing diagnosis, acquiring technical skills, and implementing best practices in many areas of medicine. At present, images needed for instructional purposes or in support of clinical decisions appear in specialized databases and in biomedical articles, and are therefore not easily accessible. Our goal is to automatically annotate images extracted from scientific publications with respect to their usefulness for clinical decision support and instructional purposes, and project the annotations onto images stored in databases by linking images through content-based image similarity. This paper presents an overview of our approach to automatic image indexing, content-based image analysis, and the results of a pilot evaluation of an automatic indexing method based on biomedical terms extracted from snippets of text pertaining to images appearing in scientific biomedical articles.

## 1. Introduction

Essential information is often conveyed in illustrations in biomedical publications. These images can be used to illuminate document summaries and answers to clinical questions, to enrich large image collections with textual information from articles, and for instructional purposes. The problem however is to automatically determine which of the images in an article will best serve each of the aforementioned purposes. Our approach to automatic image indexing is to describe (or annotate) an image at three levels of granularity:

- **coarse**, which addresses
  – image modality,
  – relation to a specific clinical task (image utility),
  – body location;
- **medium**, which provides a more detailed description of the image using existing biomedical domain ontologies;
- **specific**, which provides very detailed descriptions of clinical entities and events in an image using terms that are not included in existing ontologies and often are familiar only to clinicians specializing in a narrow area of medicine.

In this paper, we present a pilot evaluation of medium-level indexing that can be achieved by automatically extracting biomedical terms currently available in the largest biomedical domain ontology, the Unified Medical Language System® (UMLS®) Metathesaurus, from snippets of text pertaining to images in scientific biomedical articles (image captions and relevant discussion in the text). We also provide an overview of our research in coarse- and specific-level image indexing and content-based image analysis.

## 2. Background

In our previous exploration of coarse automatic indexing of images by modality (color image, gray-scale image, graph, graphic illustration, etc.) and image utility (suggested by the Evidence Based Medicine paradigm six elements of a clinical scenario that an image might illustrate), we combined image and textual features in a supervised machine learning approach. Textual features were obtained from the captions to the images and paragraphs of text containing discussion ("mentions") of these images. The text and the images were automatically extracted from the HTML-formatted articles. Text was represented as a bag-of-words or as a set of terms obtained by mapping these captions and mentions to the UMLS Metathesaurus. Texture and color features were computed on the entire image without applying any image segmentation techniques.

Texture features were computed as a 3-level discrete 2-D Daubechies' wavelet transform. The four most dominant colors were computed in the perceptually uniform CIE LUV color space and proved most effective. At this coarse level of granularity, a multi-class SVM classifier trained on a bag-of-words representation of image captions performed better in determining image modality (84.3% ± 2.6% accuracy) than when trained on a combination of textual and image features or features reduced to the domain specific vocabulary. For image utility, however, the combination of image and textual features was better than any single-source feature set achieving 76.6% ± 4.2% accuracy (Demner-Fushman et al., 2007).

Often in biomedical publications, several images are combined into a multi-panel figure. This requires sub-figure separation for image analysis to determine image modality. We therefore developed a two-phase algorithm to detect and separate figure panels using cues from caption text analysis, horizontal and vertical profiles and panel edge information (Antani et al., 2008). Further analysis on each image panel revealed its coarse modality. For instance, using color histogram profiles we could determine with sufficient precision if an image is a color image, an illustration/drawing, or a radiographic image (CT, MRI, x-ray, sonogram, etc.). Detecting image modalities is useful in further image analysis and

sub-categorization. Our efforts in this area resulted in development of a method for detecting text overlays on images, arrows, and other content valuable for indexing images by visual content and correlated text description (Antani et al., 2008).

## 2.1 Prior Work in Content-Based Image Retrieval

Our image analysis and image indexing work stems from an ongoing long-term research and development effort into image understanding and content-based image retrieval (CBIR) of biomedical images. We have worked with a large collection of digitized x-ray images of the spine derived from a nationwide health survey to develop image segmentation techniques for extraction of vertebral shape information important to researchers of osteoarthritis and musculoskeletal diseases. Whole and partial shape similarity techniques, multiple object similarity, multidimensional data indexing, relevance feedback, and Web-based frameworks for CBIR have been explored (Hsu et al., 2007).

Subsequently the research has been expanded into localization and similarity matching of pre-cancerous lesions in the uterine cervix on a data set acquired by the National Cancer Institute (NCI) from a multi-year longitudinal study. For this dataset color, texture, and location methods were studied to enable CBIR of several types of regions of interest (Xue et al., 2007). As both data sets have free-text medical records corresponding to the images, we have explored combined text and image retrieval on this data.

Finally, we have also explored automatic coarse-level image labeling and classification on the ImageCLEF 2005 data set using Semantic Error-Correcting Output Codes (SECC) and achieved an overall error rate of 18.7 using 9,000 training images and 1,000 test images (Yao et al., 2006).

Coarse-level image indexing is not sufficient to describe an image taken from a publication beyond achieving retrieval of a particular modality, utility, and location, for example, *ultrasound images for diagnosis of heart conditions*. We hypothesize that medium-level image annotation will facilitate finding images to illustrate summaries and answers to clinical questions, for example, about *echocardiographic finding* of *mitral annular calcification*. Specific-level indexing will be required to answer detailed questions, such as *What is the efficacy of thick acellular human dermis grafts for posterior and middle lamellae reconstruction?*

## 3. Methods

To automatically achieve medium-level indexing we extracted the image captions and mentions from the article text and processed the text using MetaMap, a tool that maps biomedical text to the UMLS (Aronson, 2001). The indexing terms were extracted from the MetaMap machine output, which provides comprehensive information about the mappings of phrases found in the text to the UMLS concepts. The following information was retained: the concept unique identifier (CUI) and semantic type, the preferred UMLS name for the concept, and the offset and length of the substring that was mapped to the concept.

To enable content experts to evaluate the quality of the extracted indexing terms we developed a Web-based evaluation and annotation interface (see Figures 1 and 2). This interface displays an image, bibliographic information about the article from which the image was extracted, and two tabs for annotation and evaluation. The first tab shown in Figure 1 is used for coarse-level image annotation through selecting pre-defined indexing terms for modality, utility and body location. The second tab (Figure 2) serves two purposes:

1. Evaluation of the automatically extracted indexing terms for medium-level indexing;
2. Manual annotation of the image with specific terms, more fine-grained than currently available in the UMLS (specific-level indexing), such as *thick acellular human dermis graft*. Parts of this term can be mapped to the UMLS, but even the closest existing term *Acellular Dermal Replacement* cannot be mapped to the specific term using existing tools.

The purpose of the manual annotation is to identify such missing terms and establish their ontological relations. The results of manual annotation will be used for development and evaluation of automatic indexing methods on all three levels of granularity.

The indexing terms and ontological information extracted from the MetaMap output (Figure 2 top) were evaluated on two axes:

1. Usefulness in image indexing, evaluated on a binary scale.
2. Relevance to the image, evaluated on a five-point scale, ranging from an *exact match* to *unrelated*.

An identified term might not be useful for indexing if it is too broad, too narrow, or unrelated to the image. An unrelated term might be extracted for two reasons:

1. A term might be extracted from the caption text verbatim, but the senses of the term available in the UMLS are not relevant to the image. For example, the string *apex* identified in the caption *Thrombus in left ventricular apex* maps through synonymy to the UMLS concepts:
   - APEX1 gene
   - APEX1 protein, human
   - Highest

The UMLS Metathesaurus does not contain the term *ventricular apex*; and mapping to the correct sense *Cardiac apex* is not possible using strict matching, because the set of synonyms for the *Cardiac apex* concept does not include the term *apex*.
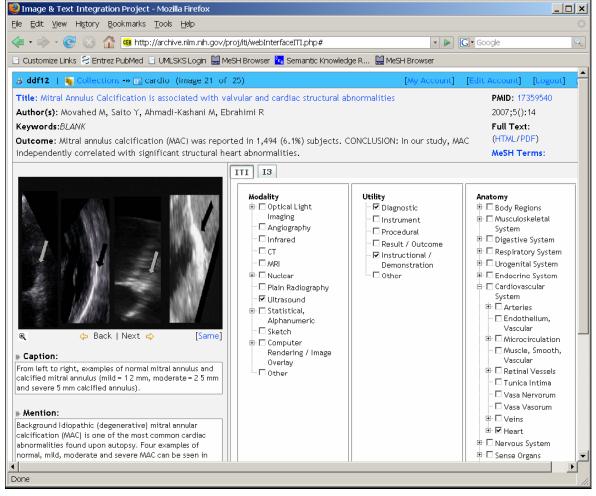
**Figure 2: A Web-based application for image indexing annotation and evaluation. Coarse-level annotation categories.**
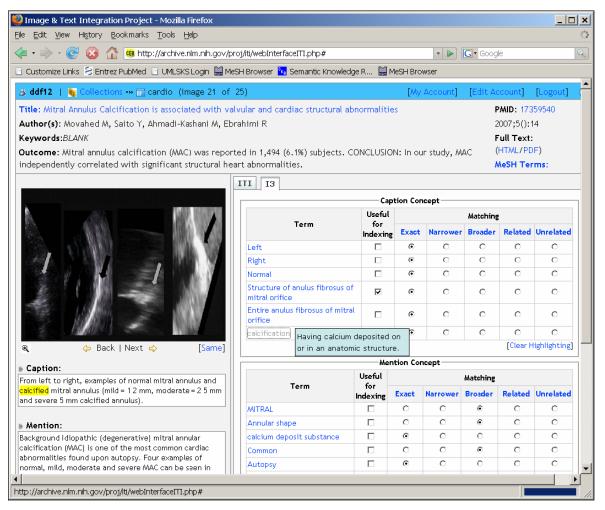


**Figure 1: A Web-based application for image indexing annotation and evaluation. Medium-level indexing evaluation.**

2. A substring identified in the text could be matched to a wrong term in the UMLS Metathesaurus because it is an acronym or abbreviation. For example, the term *LV* identified in the caption an initial increase of *LV filling pressure* is synonymous with:
   - Latvia
   - Leucovorin Calcium
   - Liposome Vesicle

The UMLS Metathesaurus does not contain the expansion of *LV* to *left ventricular* expected in the context of cardiovascular imaging. The assumption that only this sense of the term is expected in the context of cardiovascular imaging is based on the observation that the term is not expanded anywhere in the paper containing the image.

Our interface tool assists the evaluators in determining the sense of the extracted terms through the UMLS definitions which are displayed by positioning the computer mouse over the suggested index term. The tool retrieves the UMLS definitions using the extracted unique concept identifiers. Assistance for determining the origin of an extracted term is provided through highlighting the substring that was mapped to a term in the caption or mention text upon clicking on the suggested index term.

The evaluation interface was used by five physicians and one medical imaging specialist who manually assigned missing specific terms, and evaluated the quality of medium-level indexing terms. The indexing terms were automatically extracted from captions and descriptions of 50 images randomly selected for each evaluator from all images published in *BMC Annals of Facial and Plastic Surgery* and *European Journal of Cardiovascular Imaging* during 2006 and 2007. Their judgments were analyzed to answer the following questions:
1. Do captions and mentions of the image in the text provide information beyond indexing terms assigned by NLM indexers to the papers containing those images?
2. Is the extracted text sufficient for image annotation?
3. Is our extraction method satisfactory?

The first question was answered by intersecting the extracted terms evaluated as useful for imaging with the indexing terms assigned to the papers by NLM indexers and extracted from the bibliographic citations to the papers. These citations in XML format were retrieved using PubMed/MEDLINE®.

The second question was answered by intersecting the additionally assigned terms with the extracted text and with the full-text paper.

The extraction method was evaluated using recall and precision computed for each evaluator as follows: The desired index terms D for the images are the set of extracted terms evaluated as useful for indexing combined with the indexing terms added by the evaluator, A is the set of all suggested indexing terms, and within A there is a set of terms evaluated as useful for indexing C. Precision P and recall R are:
$$P = |C|/|A|$$
$$R = |C|/|D|$$
Precision and recall were computed for each evaluator, and then averaged.

# 4. Results

The six evaluators scored 4, 006 concepts (3, 281 of which were unique) pertaining to 186 unique images extracted from 109 papers. Table 1 presents the average numbers of concepts per image evaluated and found useful for indexing by each evaluator. The majority of the terms rated useful for indexing were also rated as an exact match.

**Table 1: Average number of concepts per image.** Evaluators trained in medical informatics are marked with an asterisk.

| Specialty | Indexing Terms | | |
|---|---|---|---|
| | evaluated | useful | %useful |
| family physician* | 19.26 | 2.38 | 12.4% |
| cardiologist* | 17.80 | 2.02 | 11.4% |
| plastic surgeon* | 17.89 | 1.80 | 10.1% |
| internist* | 17.55 | 2.18 | 12.4% |
| general surgeon | 19.98 | 1.50 | 7.5% |
| medical imaging | 14.46 | 1.40 | 9.9% |
| Mean ± CI | 17.83±2.0 | 1.89±0.4 | 10.6±2.0% |

The 349 exact matches constitute 77.4% of the terms marked as useful for indexing. The remaining 102 selected indexing terms were rated primarily as being broader than an exact description of the image would warrant.

## 4.1 Indexing terms assigned to the article and image annotation

Overall, the evaluators rated 451 extracted terms as useful for indexing and submitted 255 additional indexing terms.

**Table 2: Match between indexing terms assigned to images and papers.** Evaluators trained in medical informatics are marked with an asterisk.

| Specialty | MeSH Terms | | |
|---|---|---|---|
| | extracted | added | %used |
| family physician* | 33.0% | 34.9% | 11.5% |
| cardiologist* | 39.8% | 48.7% | 20.5% |
| plastic surgeon* | 46.9% | 41.2% | 11.1% |
| internist* | 25.0% | 25.7% | 11.7% |
| general surgeon | 33.3% | –– | 7.1% |
| medical imaging | 28.8% | –– | 5.3% |
| Mean ± CI (%) | 34.5±8.2 | 25.1±21.9 | 11.2±5.5 |

Table 2 presents the percentages of terms assigned by the evaluators that match terms assigned by NLM indexers (MeSH terms) to the papers containing the images. In addition, the %used column of the table shows the proportion of the MeSH terms assigned to the paper that

were deemed useful in annotating images.

## 4.2 Locating additional terms in the text

For three of the 255 indexing terms added by the evaluators no image-related text was extracted. Of the remaining 252 added terms, 75 were extracted verbatim from the caption text and 11 from the discussion of the image in the text. Another 139 added terms were generated using captions and mentions through:

- extracting strings with gaps, for example, extracting *Preoperative photograph* from *Preoperative and postoperative photographs*;
- paraphrasing, for example, deriving *elderly* from *89-year old*;
- summarizing, for example, the following mention of the image: *a mobile, left-sided, nasal dorsal implant with tip ptosis, erythema, and swelling of the left nasal vestibule* as *implantation complications*;
- generalizing based on the figure and the caption, for example, *ultrasound*; *surgical method*; or *transthoracic echocardiography*.

The remaining 27 terms were found in the paper title, abstract, and MeSH terms assigned to the paper. Of the 255 additionally assigned terms 103 were subsequently mapped to the UMLS concepts.

## 4.3 Extraction accuracy

The design of the extraction evaluation was recall oriented. All extracted terms were given to the evaluators without any filtering to have enough training examples for learning term selection in the future. Recall and precision achieved by this baseline extraction method are shown in Table 3.

**Table 3: Evaluation of the baseline extraction method.**
Evaluators trained in medical informatics are marked with an asterisk.

| Specialty | Recall | Precision | F-score |
|---|---|---|---|
| family physician* | 0.723 | 0.124 | 0.211 |
| cardiologist* | 0.447 | 0.114 | 0.181 |
| plastic surgeon* | 0.827 | 0.101 | 0.179 |
| internist* | 0.565 | 0.124 | 0.204 |
| general surgeon | 0.333 | 0.075 | 0.122 |
| medical imaging | 0.917 | 0.099 | 0.179 |
| Average | 0.635 | 0.106 | 0.182 |

## 5. Discussion

The results of this baseline pilot evaluation are encouraging. Similarly to Declerck and Alcantara (2006) who identified the title, caption, and abstract of a Web document among the text regions possibly relevant to image annotation, we found captions, mentions, abstracts and titles of scientific publications to provide sufficient information for image annotation. Although the information was easily recognized by the evaluators, on average, only 64% of the desirable indexing terms could be found using the existing extraction methods and

ontologies. More sophisticated mapping algorithms are needed to extract another 15% of the terms, and more complex natural language processing and ontology expansion are needed to identify the remaining terms.

The pilot evaluation clearly indicates that although there is some correlation between the MeSH terms assigned to a paper and image annotation, only a small proportion of the MeSH terms could be used to describe an image, and additional indexing terms have to be extracted from the text.

The variations in the annotation results among the annotators could be partially attributed to the underspecified image annotation rules. The small number of the images annotated by more than one evaluator does not allow computing inter-annotator agreement scores, but there are indications that the differences could be reduced by better defined rules. For example, in one case, two evaluators marked the extracted term *Hypertrophic Cardiomyopathy* as useful, but only one of them also rated *Echocardiography* as a useful term. Had the instructions clearly stated that if a term belongs to the coarse-level annotation, it should not be used for the medium-level description, the discrepancy might have been avoided. We plan to develop a set of specific rules that describe the appropriate terminology, annotation precision, etc. as described in (Grubinger et al., 2006).

## 6. Future work

In the next phase, we will focus on the improvement of the evaluation/annotation interface; improvement of the coarse-level controlled vocabularies; selection of the extracted terms to be suggested as indexing terms; improvement of term extraction, and expansion of the test collection. The implementation of some of the improvements to the interface and coarse-level vocabularies suggested by the evaluators is already underway. Figures 3 presents the changes to the coarse-level annotation tab implemented after the pilot evaluation. The changes involve a better layout, a search function for controlled vocabulary terms for coarse level anatomy annotation, and a new teaching quality annotation axis.

## 7. Acknowledgements

## 8. References

Sameer K. Antani, Dina Demner-Fushman, Jiang Li, Balaji V. Srinivasan, and George R. Thoma. 2008. Exploring use of images in clinical articles for decision support in Evidence-BasedMedicine. *In Proceedings of the 20th SPIE/IS&T Electronic Imaging Conference*, pages 1–10.
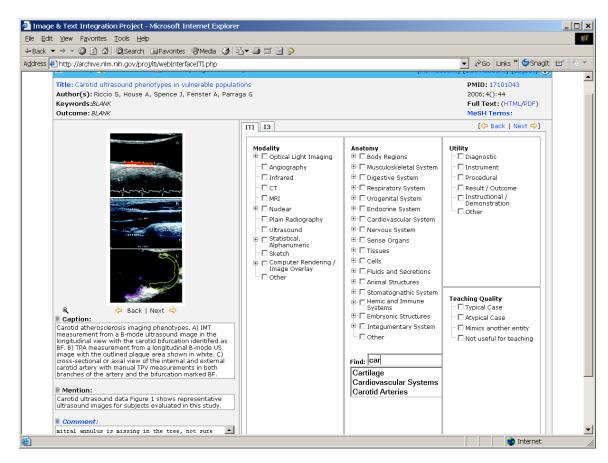
**Figure 3: Modified evaluation interface following the evaluators' feedback.**

Alan R. Aronson. 2001. Effective mapping of biomedicaltext to the UMLS Metathesaurus: The MetaMap program. *In Proceeding of the 2001 Annual Symposium of the American Medical Informatics Association (AMIA 2001)*, pages 17–21.

Thierry Declerck and Manuel Alcantara. 2006. Semantic analysis of text regions surrounding images in Web documents. *In OntoImage 2006 Workshop on Language Resources for Content-based Image Retrieval*, pages 9–12.

Dina Demner-Fushman, Sameer K. Antani, and George R. Thoma. 2007. Automatically finding images for clinical decision support. *In Proceedings of the IEEE Workshop on Data Mining in Medicine (DMMed '07)*, pages 139–144.

Michael Grubinger, Paul Clough, Henning Müller, and

Thomas Deselaers. 2006. The IAPR TC-12 benchmark: A new evaluation resource for visual information systems. *In OntoImage 2006 Workshop on Language Resources for Content-based Image Retrieval*, pages 13–23.

William Hsu, L. Rodney Long, and Sameer K. Antani. 2007. SPIRS: A framework for content-based image retrieval from large biomedical databases. *In Proceedings of the Medinfo Congress*, pages 188–192.

Zhiyun Xue, Sameer K. Antani, L. Rodney Long, Jose Jeronimo, and George R. Thoma. 2007. Investigating CBIR techniques for cervicographic images. *In Proceedings of the 2007 Annual Symposium of the American Medical Information Association (AMIA 2007)*, pages 826–830.

Jian Yao, Sameer K. Antani, L. Rodney Long, George R. Thoma, and Zhongfei Zhang. 2006. Automatic medical image annotation and retrieval using SECC. *In Proceedings of the 19th International Symposium on Computer- Based Medical Systems (CBMS 2006)*, pages 820–825.