

The CLEF Cross Language Image Retrieval Track (ImageCLEF) 2004

Paul Clough¹, Mark Sanderson¹ and Henning Müller²

¹Department of Information Studies, University of Sheffield, Regent Court, 211 Portobello Street, Sheffield, S1 4DP, UK.
{p.d.clough,m.sanderson}@sheffield.ac.uk

²University Hospitals of Geneva, Division of Medical Informatics, 21 rue Micheli-du-Crest, CH-1211 Geneva 4, Switzerland.
henning.mueller@dim.hcuge.ch

Abstract. In this paper we describe ImageCLEF¹, the cross language image retrieval track of the Cross Language Evaluation Forum (CLEF³). We instigated and ran a pilot experiment in 2003 where participants submitted entries for an ad hoc bilingual image retrieval task on a collection of historic photographs from St. Andrews University Library. This was designed to simulate the situation in which users would express their search request in natural language but require visual documents in return. For 2004 we have extended the tasks to include a medical image retrieval task and a user-centred evaluation.

1 Introduction

A great deal of research is currently underway in the field of Cross Language Information Retrieval (CLIR) where documents written in one language are retrieved by a query written in another (see, e.g. [11] and [16]). One can consider CLIR as basically a combination of machine translation (MT) and traditional monolingual information retrieval (IR). Most CLIR research has focused on locating and exploiting translation resources with which the user's search requests or target documents (or both) are translated into the same language. Campaigns such as the Cross Language Evaluation Forum (CLEF) [16] and the Text REtrieval Conference (TREC) [20] multilingual track have helped encourage and promote international research, as well as create standardised resources for CLIR evaluation.

However, one area of CLIR research which has received less attention is image retrieval. In collections such as historic or stock-photographic archives, medical case notes and art/history collections, images are accompanied by some kind of text (e.g. metadata or captions) semantically related to the image [2][12]. Images can then be retrieved using standard IR methods based on textual queries. However, retrieval from an image collection offers distinct characteristics from one in which the document to

¹ ImageCLEF: <http://ir.shef.ac.uk/imageclef2004/>

³ CLEF: <http://www.clef-campaign.org>

be retrieved is natural language text [1][10]. For example, the way in which a query is formulated, the method used for retrieval (e.g. based on low-level features derived from an image, or associated text), the types of query, how relevance is assessed, the involvement of the user during the search process, and fundamental cognitive differences between the interpretation of visual versus textual media. Methods of image retrieval are typically based on visual content⁴ (e.g. colour, shape, spatial layout and texture), or by text/metadata associated with the image (see, e.g. Smeulders et al. [18] and Goodrum [10]).

For those organisations managing image repositories in which text is associated with images (e.g. on-line art galleries), one way to exploit these is by enabling multilingual access to them. To promote research in this area we instigated ImageCLEF [5] as part of the CLEF campaign. We felt this contribution would address an important and timely problem not dealt with by existing cross language evaluation. We envisage ImageCLEF will appeal to both commercial and academic research communities including: cross language information retrieval, image retrieval, and user interaction. The main aims of the ImageCLEF campaign are: (1) to promote and initiate international research for CL image retrieval, (2) to further our understanding of the relationships between CL texts and images for IR, and (3) to create a set of useful standardised resources for CL image retrieval to scientific communities. The paper is divided into the following: in section 2 we describe the ImageCLEF 2003 test collection for an ad hoc retrieval task, in section 3 we describe tasks offered in ImageCLEF 2004 and finally in section 4 we summarise the contents of this paper and provide some ideas for future work in cross language image retrieval.

2 Building a Test Collection for Multilingual Image Retrieval

Evaluation of retrieval systems is either system-focused, e.g. comparative performance between systems or user-centered, e.g. a task-based user study. For many years IR evaluation has been dominated by comparative evaluation of systems in a competitive environment. The design of a standardised resource for IR evaluation was first proposed over 30 years ago by Cleverdon [4] and has since been used in major IR conferences such as TREC [20], CLEF [16] and NTCIR [3]. Over the years the creation of a standard test environment has proven invaluable for the design and evaluation of practical retrieval systems both within and outside a competitive environment. The main components of a TREC-style test collection are: (1) document collection, (2) topics, and (3) relevance assessments.

In TREC, NTCIR and CLEF, participants are given test collection data and topics and asked to submit their entries. A subset, chosen by the organisers, is used to create document pools, one for each topic. Domain experts (assessors) are then asked to judge which documents in the pool are relevant or not. Document pools are created because in large collections it is infeasible to judge every single document for relevance. These assessments are then used to assess the performance of submitted systems. User-centred evaluation is important to assess the overall success of a retrieval system which takes into account other factors other than just system performance, e.g. the design of the user interface and system speed (Dunlop argues this in [7]). A num-

⁴ These are called Content-Based Information Retrieval (CBIR) systems.

ber of researchers have highlighted the advantages of user-centred evaluation, particularly in image retrieval systems (see, e.g. [10], [14] and [7]). One of the main aims of ImageCLEF is to provide both the CLIR and image retrieval communities a number of useful resources (datasets and relevance assessments) to facilitate and promote further research in multilingual image retrieval.

Calls for a TREC-style evaluation for image retrieval systems have been suggested [10][15][19], although Forsyth [9] argues that the evaluation of CBIR systems at the moment is useless because systems are too bad (hence the interest in combining both textual and visual features). We are unaware of existing test collections for CL image retrieval, although evaluation resources do exist to evaluate specific image retrieval tasks, e.g. journalism [13] and CBIR systems, e.g. Benchathlon⁵. One of the largest obstacles in creating a test collection for public use is securing a suitable collection of images for which copyright permission is agreed. This has been a major factor influencing the datasets used in the ImageCLEF campaigns. The ImageCLEF test collection provides a unique contribution to publicly available test collections and complements existing evaluation resources.

2.1 The Existing ImageCLEF Test Collection

Because CL image retrieval encompasses at least two research areas: (1) image retrieval and (2) CLIR, building a comprehensive and suitable test collection is a tall order. Therefore, in 2003 we organised a pilot experiment at CLEF with the following aim: given a multilingual statement describing a user need, find as many relevant images as possible. More formally the task was a bilingual ad hoc retrieval task in which a static collection was searched using previously unseen topics.

The retrieval task was designed to simulate the situation in which a user expresses their need in a language different from the collection, requiring a visual document to fulfil their search request (e.g. searching an on-line art gallery or stock photographic collection). For this retrieval task query translation is the preferred method of bridging the language gap as translating the collection would be both time and resource expensive and less likely in practice. Participants were not constrained in their use of retrieval method, enabling either text or content-based searches (or a combination of both). As a retrieval task there are several challenges other than translation which include: (1) captions typically short in length, (2) images of varying content and quality, (3) bridging the gap between colloquial and domain-specific language used in the captions and cross language queries, and (4) queries short in length thereby providing little context for translation.

The dataset used consisted 28,133 historic photographs from the library at St Andrews University [17]. All images are accompanied by a caption consisting of 8 distinct fields which can be used individually or collectively to facilitate image retrieval (see Fig. 1). The 28,133 captions consist of 44,085 terms and 1,348,474 word occurrences; the maximum caption length is 316 words, but on average 48 words in length. All captions are written in British English and contain colloquial expressions and historical terms. Approximately 81% of captions contain text in all fields, the rest generally without the description field. In most cases the image description is a grammati-

⁵ <http://www.benchathlon.net/>

cal sentence of around 15 words. The majority of images (82%) are black and white, although colour images are also present.

Record ID: JV-A.000460
Short title: The Fountain, Alexandria.
Long title: Alexandria. The Fountain.
Location: Dunbartonshire, Scotland
Description: Street junction with large ornate fountain with columns, surrounded by rails and lamp posts at corners; houses and shops.
Date: Registered 17 July 1934
Photographer: J Valentine & Co
Categories: [columns unclassified][street lamps - ornate][electric street lighting][shepherds & shepherdesses][streetscapes][shops]
Notes: JV-A460 jf/mb



Fig. 1. An example image and caption (see: <http://www-library.st-andrews.ac.uk>).

We generated fifty representative search requests in English (called *topics*) and translated them into 6 different languages: Dutch, Spanish, German, French, Italian and Chinese (provided by the National Taiwan University or NTU). In TREC, CLEF and NTCIR final topics are chosen from a pool of suggestions generated by searchers familiar with the domain of the document collection. Frequently searched subject areas in the St Andrews were identified by analysing log files generated from accesses to a web search engine used by the library. Based on these subject areas we created queries that would test the capabilities of both a translation and image retrieval system, e.g. pictures of specific objects versus pictures containing actions, broad versus narrow concepts, topics containing proper names, compound words, abbreviations, morphological variants and idioms. Each topic consisted of a short title, a longer narrative describing the search request and an exemplar relevant image. For ImageCLEF 2003 only topic titles were translated due to limited resources available to us.

2.2 Relevance Assessments and Evaluation

What turns a set of documents and queries into a test collection are the relevance judgments, manual assessments of which documents are relevant or not for each topic. Judging whether an image is relevant or not is highly subjective (e.g. due to knowledge of the topics or domain, different interpretations of the same document, and searching experience), therefore to minimise this two assessors judged each topic.

We adopted the pooling method as used in TREC, CLEF and NTCIR where a set of candidate documents is created (called the *pool*) by merging together the results of the top n documents from the ranked lists provided by participants. This assumes that highly ranked documents from each entry will contain relevant documents. Ideally, ranked lists should come from a diverse range of systems to ensure maximal coverage. We also supplemented the pooling method with manual interactive searches (also known as *interactive search and judge* or ISJ) to ensure good quality pools (as used in NTCIR). We found assessors were able to judge the relevance of images very quickly

(especially eliminating non-relevant ones) enabling *all* ImageCLEF submissions to be used in creation of the pools (compared to a subset of runs for text-based assessment). One of the authors familiar with the collection assessed all fifty topics to provide a “gold” set of judgments; in addition, ten assessors from the University of Sheffield judged five topics each to provide a second judgment for each topic using a custom-built assessment tool.

Images were judged relevant if *any* part of the image was deemed relevant. Primary judgment was made on the image, but assessors also consulted the image captions. Assessors were asked to judge the relevance of images using a ternary scheme: relevant, partially relevant and not relevant to deal with potential uncertainty in the assessor's judgment (i.e. it is possible to determine that the image is relevant, but less certain whether it exactly fulfils the need described by the topic). Unlike other test collections we provided four sets of relevance assessments (called *qrels*) - strict/relaxed union/intersection - with which to assess system performance based on the overlap of relevant images between assessors and whether the relevance sets include images judged as partially relevant or not. These are further described in [5]. The strict relevance set can be contrasted with a high-precision task; the relaxed set providing an assessment that promotes higher recall.

2.3 Results and Lessons Learned

Four groups entered ImageCLEF 2003: Sheffield University, NTU, University of Surrey and Daedalus, a Spanish R&D organization. All participants used text-based retrieval methods with no content-based image analysis. Results from ImageCLEF have shown that in general CL image retrieval using query translation can achieve relatively high performance for the suggested bilingual search task. However, we found retrieval performance to vary dramatically across both language and topic. The highest result was obtained for French (78% of monolingual); the worst for Chinese (51% of monolingual) indicating there is still room for improvement. In particular, enhancement to deal with poor retrieval caused by translation errors is required. Results from ImageCLEF showed: for Chinese retrieval transliteration of proper names was beneficial, and for other languages thesaurus-based query expansion improved performance. ImageCLEF was effective at attracting new research groups to CLEF and this year is advertised as an entry-level CLIR task.

Based on our experiences from last year we have made the following changes to the ImageCLEF track: (1) to offer greater diversity we have added a medical retrieval task, (2) to promote ImageCLEF as an entry-level CLIR task we are offering topics in 12 languages rather than 6, (3) to encourage participants to exploit visual features we have setup public access to a default CBIR system, (4) due to ambiguity in relevance assessments we have selected more specific topics including queries refined by photographer, location and date (general queries such as “mountain scenery” retrieved too many images and were too laborious to assess), and (5) we are using relevance assessors familiar with the collection (this includes native English speakers who are familiar with colloquial English/Scottish terms, e.g. “perambulator”).

3 The ImageCLEF 2004 track

3.1 The Bilingual Ad Hoc Retrieval Task

A bilingual ad hoc task similar to that run in 2003 is being offered to participants to enable further experiments on the St. Andrews dataset and determine whether improvements can be made on last year's results. Experiments will compare: (1) different methods of query translation (e.g. dictionary-lookup versus MT), (2) query expansion (e.g. global versus local methods), (3) the use of text-based and CBIR methods used either separately or combined, (4) different retrieval models, (5) different indexing methods (e.g. indexing all or some fields) and (6) manual vs. automatic relevance feedback.

A new set of 25 topics has been produced in the same manner as before (decide on general topics and then refine). However, in addition to using St. Andrews query logs, we also used subject areas supplied by staff from St. Andrews' library. Topic refinement is based on the query categorisation scheme suggested by Armitage et al. [1] for picture archives and designed to test a range of different CL and image search parameters. Topics have been translated into the previous languages, plus Japanese, Danish, Russian, Finnish, Swedish and Arabic. One non-intentional but interesting "feature" of translated topics in ImageCLEF 2003 was the introduction of translation errors, e.g. spelling mistakes and erroneous diacritics, resulting in low retrieval performance for some topics. These problems are not addressed by existing CLEF tasks. We will provide two sets of topics: one set will contain spelling errors; the other will be checked and free of such errors.

3.2 The Medical Image Retrieval Task

To offer participants a different domain/scenario and encourage the use of CBIR system we have introduced a task based on medical retrieval. In the ad hoc task it is the query which is multilingual; in the medical retrieval task the document collection is multilingual presenting different CLIR challenges.



Fig. 2. Example images from the CasImage dataset (<http://www.casimage.com/>)

In general, medical practitioners are unsatisfied with retrieving images by text and the implicit knowledge stored in the images plus attached text is rarely used. As a diagnostic aid, being able to search a database of images with a new example would enable them to obtain more evidence. The goal of this task is to investigate the use of

CBIR and text-based retrieval systems for this kind of medical retrieval task. The task is being run by University Hospitals of Geneva who are supplying the medical data, topics and relevance judgments. The medical task is this: given an example image, find similar images which will be helpful in confirming the initial diagnosis. Because the initial retrieval has to be visual, we expect the case notes to be useful in finding additional similar images complementary to CBIR. We also aim to evaluate whether relevance feedback can improve performance, compare relevance feedback using either image/text or both, and whether images alone can be used for pseudo relevance feedback.

The dataset (CasImage) consists of 8,751 anonymised medical images, e.g. scans, and x-rays (see Fig. 2). The majority of images are associated with *case notes*, a written description of a previous diagnosis for an illness the image identifies. Case notes consist of several fields including: a diagnosis, a description, clinical presentation, keywords and title. The task is multilingual because case notes are mixed language written in either English or French. Not all case notes have entries for each field and the text itself reflects real clinical data in that it contains mixed-case text, spelling errors, erroneous French accents and un-grammatical sentences. In the dataset there are 2,078 cases to be exploited during retrieval (e.g. query expansion).

Currently 25 example images (topics) have been chosen as representative from the dataset. A set of ground truths for each topic has already been identified by domain experts based on the CBIR system developed by the third author⁶ and these will form part of the document pools created from participant's entries. Pools will be formed in a manner similar to the ad hoc task and medical practitioners will help judge the relevance of the pools after final submissions. In this task images are judged using a binary relevant or not relevant judgement and assessments will be used to evaluate participant's entries. This retrieval task offers a number of challenges including: (1) combining text and content-based methods of retrieval after an initial visual search, (2) dealing with domain-specific medical terminology, (3) case notes of varying quality in more than one language (i.e. a mixed language index), and (4) the high cost of returning non-relevant images (i.e. mis-diagnosis) which is always inevitable when using visual-only search methods.

3.3. The Interactive Retrieval Task

Campaigns such as iCLEF⁸ have shown the value of user-centred evaluation for CLIR and CL image retrieval would seem to be a rich source for user-centred experiments. Past research has shown that the search activities of a user in an image retrieval system vary between searching for specific images and browsing the image collection (see, e.g. [10] and [6]). For a CL image retrieval system, the issue is how best the system can support the user's search in locating relevant images as quickly, easily and accurately as possible. User-centered evaluation in a variety of contexts and domains will help us determine how CL image retrieval systems can best help users to: (1) formulate their queries (e.g. whether text or visual queries alone are best or can be

⁶ See <http://vipier.unige.ch/> for a list of publications about the VIPER CBIR system.

⁸ See <http://terral.lsi.uned.es/iCLEF/> for information about iCLEF.

used in combination), (2) refine the search request - query reformulation will depend on the outcome of the system and could involve refinements using textual and/or visual features, (3) browse the collection, and (4) identify relevant images (e.g. what additional information would help the user judge the relevance of an image and how best is this displayed).

Cox et al. [6] suggest three classes of image search: (1) target or known-item search (i.e. find a specific image), (2) category search (e.g. “find pictures of the Eiffel Tower”) and (3) open-ended browsing (i.e. wandering through the collection). They argue that the target search encompasses the other categories of search; it is simple for the user to perform and has clear measures of effectiveness. The goal for the user in such a task is given an image to find it again from the collection. Unlike being given a textual topic description, the user must interpret the given image and generate suitable query terms in a given language (different from the document collection). The scenario models the situation in which a user searches with a specific image in mind (perhaps they have seen it before) but without knowing key information thereby requiring them to describe the image instead, e.g. searches for a familiar painting whose title and painter are unknown. This task will use the St. Andrews dataset and our experimental setup will follow the guidelines for user-centred experiments as suggested by iCLEF. This task will be undertaken with collaboration from iCLEF organisers to ensure a consistency in CLEF methodologies. Participants are asked to follow the experimental setup but can perform whatever experiments they like.

A minimum of 8 users and 8 topics are required for this task. Users are given 10/15 minutes to find each image using only CL queries. Topics are general enough so that people unfamiliar with the collection can still perform the searches. Captions must also be translated into this language before being displayed (if at all) to the user. The aim of this experiment will be to observe users search habits and to determine what kind of interface best supports *query refinement*. For example the user is shown a picture of an arched bridge but starts with the query “bridge”. By finding similar images and maybe using keywords from their captions, the user refines the query until the relevant image is found. Query. Topics and systems will be presented to the user in combinations following a *Latin-square* design to ensure user/topic and system/topic interactions are minimised. Qualitative performance measures is captured using questionnaires provided by us, and quantitative measures include: whether the given image is found or not, the time taken to find the image, the number of images viewed before finding the image and number of user interactions required.

4 Conclusions and Future Work

In this paper we have discussed our proposal for three cross-language image retrieval tasks as part of the ImageCLEF campaign. The tasks vary across domain, scenario, where CLIR is used, whether content-based image retrieval is required and whether the task is system or user-centered. Results from ImageCLEF 2003 have shown CL image retrieval to be a success, but large improvements can still be obtained for some languages (e.g. Chinese). Our aim is to promote CL image retrieval and provide a standardised set of resources in the form of test collections (i.e. a collection, topics and relevance assessments) which can be used in further CL image retrieval experi-

ments. In future work we plan to expand the collections and tasks offered in Image-CLEF. In particular we would like to offer collections with non-English captions provide a Web-based image retrieval task and offer further image retrieval tasks, e.g. aspectual retrieval.

References

1. Armitage, L.H. and Enser, P.: Analysis of User Need in Image Archives. In *Journal of Information Science* **Vol. 23(4)** (1997) 287-299
2. Chen, F., Gargi, U., Niles, L. and Schütze, H.: Multi-Modal Browsing of Images in Web Documents. In *Proceedings of SPIE Doc. Recognition and Retrieval VI* (1999) 122-133
3. Chen, K., Chen, H., Kando, N., Kuriyama, K., Lee, S. and Myaeng, S.: Overview of CLIR Task, Third NTCIR Workshop, Japan (2002)
4. Cleverdon, C.W.: The Cranfield Tests on Index Language Devices. In: K. Spark-Jones and P. Willett (eds), *Readings in Information Retrieval*, Morgan Kaufmann (1997) 47-59
5. Clough, P. and Sanderson, M.: The CLEF 2003 Cross Language Image Retrieval Track. In *Proceedings of the Cross Language Evaluation Forum (CLEF) Workshop, Norway* (2003)
6. Cox, I.J., Miller, M.L., Omohundro, M. and Yianilos, P.N.: Target Testing and the PicHunter Bayesian Multimedia Retrieval System. In *Proceedings of Advanced Digital Libraries (ADL'96) Forum*, Washington D.C. (1996)
7. Dunlop, M.: Reflections on MIRA: Interactive Evaluation in Information Retrieval. In *Journal of the American Society for Information Science* **Vol. 51(14)** (2000) 126-1274
8. Flank, S.: Cross language Multimedia Information Retrieval. In *Proceedings of Applied Natural Language Processing and the North American Chapter of the Association for Computational Linguistics* (2000)
9. Forsyth, D.A.: Benchmarks for Storage and Retrieval in Multimedia Databases. In *Proceedings of SPIE International Society for Optical Engineering* **Vol. 4676** (2001) 240-247
10. Goodrum, A.A.: Image Information Retrieval: An Overview of Current Research. In *Informing Science* **Vol. 3(2)** (2000) 63-66
11. Grefenstette, G.: *Cross language Information Retrieval*. Kluwer Academic Publishers, Norwell, MA, USA. (1998).
12. Harmandas, V., Sanderson, M. and Dunlop, M.D.: Image Retrieval by Hypertext Links. In *Proceedings of the 20th ACM SIGIR conference* (1997) 296-303
13. Markkula, M. and Sormunen, E.: Searching for Photos – Journalist's Practices in Pictorial IR. In *Proceedings of Conference on Image Retrieval (CIR'98)* (1998)
14. McDonald, S., Lai, T.S., Tait, J.: Evaluating a Content Based Image Retrieval System. In *Proceedings of SIGIR'01* (2001) 232-240
15. Müller, H., Müller, W., McG. Squire, D., Marchand-Maillet, S. and Pun, T.: Performance Evaluation in Content-Based Image Retrieval: Overview and Proposals. In *Pattern Recognition Letters* **Vol. 22(5)** (2001) 593-601
16. Peters, C. and Braschler, M.: Cross Language System Evaluation: The CLEF Campaigns. In *Journal of the American Soc. for Inf. Sci. and Tech.* **Vol. 52(12)** (2001) 1067-1072
17. Reid, N.: The Photographic Collections in St Andrews University Library. In *Scottish Archives* **Vol. 5** (1999) 83-90
18. Smeulders, A.W.M, Worring, M. Santini, S. Gupta, A. and Jain, R.: Content-Based Image Retrieval at the End of the Early Years. In *IEEE Transactions on Pattern Analysis and Machine Intelligence* **Vol. 22(12)** (2000) 1349-1380
19. Smith, J.R.: Image Retrieval Evaluation. In *Proceedings of the IEEE Workshop of Content-Based Access to Image and Video Databases* (2001) 112-113
20. Voorhees, E.M. and Harman, D.: Overview of TREC 2001, In *NIST Special Publication 500-250: Proceedings of TREC2001*, NIST. (2001)