

# Northumbria Research Link

Citation: Brown, Jean, Mulholland, Richard, Graham, Margaret, Riley, Jon, Vassilev, Vassil, Eakins, John and Furness, Karen (2002) When images work faster than words: The integration of content-based image retrieval with the Northumbria Watermark Archive. *Restaurator*, 23 (3). pp. 187-203. ISSN 0034-5806

Published by: De Gruyter

URL: <http://www.degruyter.com/view/j/rest.2002.23.issue...>  
<<http://www.degruyter.com/view/j/rest.2002.23.issue-3/rest.2002.187/rest.2002.187.xml>>

This version was downloaded from Northumbria Research Link: <http://nrl.northumbria.ac.uk/19394/>

Northumbria University has developed Northumbria Research Link (NRL) to enable users to access the University's research output. Copyright © and moral rights for items on NRL are retained by the individual author(s) and/or other copyright owners. Single copies of full items can be reproduced, displayed or performed, and given to third parties in any format or medium for personal research or study, educational, or not-for-profit purposes without prior permission or charge, provided the authors, title and full bibliographic details are given, as well as a hyperlink and/or URL to the original metadata page. The content must not be changed in any way. Full items must not be sold commercially in any format or medium without formal permission of the copyright holder. The full policy is available online: <http://nrl.northumbria.ac.uk/policies.html>

This document may differ from the final, published version of the research and has been made available online in accordance with publisher policies. To read and/or cite from the published version of the research, please visit the publisher's website (a subscription may be required.)



**Northumbria**  
**University**  
NEWCASTLE

# When Images Work Faster than Words

## The Integration of Content-Based Image Retrieval with the Northumbria Watermark Archive

by A. JEAN E. BROWN, RICHARD MULHOLLAND,  
MARGARET GRAHAM, JON RILEY, VASSIL VASSILEV,  
JOHN EAKINS & KAREN FURNESS

### INTRODUCTION

The Conservation of Fine Art program (School of Humanities) and the Institute for Image Data Research (IIDR), at the University of Northumbria at Newcastle upon Tyne, have been working on a collaborative research project to improve the accessibility of watermark images and related information and in doing so provide a fresh, interactive approach to paper research. The 18 month research project, funded by the Arts and Humanities Research Board (AHRB), has included the evaluation and development of a range of image capture techniques as well as the development of a Content Based Image Retrieval (CBIR) system – SHREW (SHape REtrieval of Watermarks). These have both been integrated into a prototype dynamic web-based archive of watermarks, i.e. The Northumbria Watermark Archive. In the longer term the archive will provide access to information on: paper history and technology; artistic usage; bibliographic usage; conservation criteria as well as suppliers of historic papers etc.

Although the research carried out has been applied to watermarks the approach has the potential to be applied to any numerically large collection of images with distinctive features of colour, shape or texture i.e. coins, architectural features, picture frame profiles, hallmarks, Japanese artists stamps etc. The establishment of an electronic archive undoubtedly offers many advantages to those working with large collections of images in terms of accessibility, however the implementation of a new storage medium and search engine can also introduce difficulties not anticipated. This paper will discuss some of these issues.

#### COLLABORATION BETWEEN DISCIPLINES

In order to develop a digitised database with a CBIR facility for a collection it is generally necessary for two fields of expertise to collaborate i.e. the collection managers and the computer programmers. The original concept for the research may appear relatively straightforward with a clear demarcation between the work to be carried out on image capture, and that to be developed by the computer programmers. However it quickly becomes apparent that the two parts of the team have to develop a greater understanding of their respective disciplines and in particular the terminology that is used. If this is not done there is a strong likelihood that the resulting system will not function in a manner that is compatible with the work of the collection's managers. When dealing with colleagues from the same discipline it is easy to forget that the ease and precision of communication is highly dependent on a shared knowledge and terminology that has been developed over many years. A collaborative project has to develop an effective system of communication over a much shorter period and this should be an early priority in order to maximise the chances of success.

#### PROJECT MANAGEMENT

The research project described here has had two distinct and parallel strands, which were complementary but inter-related. In one strand, we were testing and evaluating methods of extracting the watermarks from the paper and then digitising them for further processing. In the other, we were undertaking software development for describing automatically the shapes of the watermarks. The latter was an iterative process of development, evaluation and refinement. Both parts required appropriate project management to achieve project goals within the set timescale and resources.

It is not the purpose of this paper to discuss in depth issues to do with project management. Suffice to say that any digital image archive requires careful and detailed planning. While it may seem that digitisation of the source material is the most important aspect of any archive creation, in reality, it represents a minor part of the creation of the archive, and more importantly a small percentage of the overall cost. Other issues can often have a severe impact on resources and time. For example:

Copyright clearance: An extremely complex issue that remains largely unclear for electronic archives. Some issues relevant to the UK are presented in Appendix A.

Training and familiarisation: If unfamiliar software or hardware is to be used for the first time, the cost in time and resources of training and familiarisation

must be considered at the beginning of the project. However in our experience this was not required due to the breadth of experience available within the collaborative team.

**Metadata issues:** In any archive it is important that the appropriate standards be used (dealt with in greater depth later in this paper). In addition it is essential that an early decision is made with regard to exactly how much information/metadata is to be input to each record in the archive and whether this is feasible within the timescale available for the project.

**Hardware requirements:** Realistic specifications must be sought, to avoid constant upgrades. Buying new machines may be more cost effective than upgrading older models. It is important to consider that manipulating and storing digital images takes up a large amount of disc space, processor power, and memory (RAM). As a rough guide, most images require 2 or 3 times their uncompressed file size in RAM. Additionally, a large good-quality monitor for image work can be extremely useful, but expensive, and additional items such as PCI (Peripheral Component Interconnect) or AGP (Accelerated Graphics Port) graphics memory cards may be necessary.

**Software Requirements:** Many institutions underestimate the cost of software licenses, especially for image-based software. Legally, a single-user licence permits that software package to be used on one computer only. For a project involving several researchers, several licences or an expensive site licence may be required. A single user licence for Adobe Photoshop® 7.0 currently retails at about £769 (including VAT) in the UK (= €1201 or US\$1148\*)

**Backup and virus protection issues:** Back up facilities must be taken into consideration for any electronic archive. Whether CD/DVD/Zip/Tape or off-site solutions are used, a regular backup procedure implemented at the planning stage is vitally important. Virus protection software should also be purchased installed, and regularly updated. Regular updates for the latest viruses are available for download from the manufacturers website.

**Cross-platform interoperability:** If the archive is to be published on CD or on the web, it should be ensured that those with different operating systems (Macintosh/Windows/Unix), and software applications/versions will be able to use it.

**Obsolescence of hardware and data migration:** It is important that when the funding comes to an end, or when the project is finished, that the archive be maintained, so that it can be accessed in the future. This can mean planning

\* Exchange rate of June 2002.

'data migration' to current storage devices and file formats or depositing the archive with an outside institution. Currently the Arts and Humanities Data Service (AHDS) (<http://www.ahds.ac.uk/>) provide a deposit facility for UK archives related to arts and the humanities. This service is free of charge at the time of writing and the user retains all copyright while AHDS maintains the archive.

## COPYING AND DIGITISING OBJECTS

There are a number of principles that would apply to any collection that is to be copied and digitised. It is essential that the systems recommended for copying the originals, followed by digitising, are relatively straightforward to carry out and not prohibitively expensive. If the system is too complicated, mistakes may be frequent and the resulting images will be less acceptable to the CBIR software. Access to the appropriate information in the database will therefore be more problematic. If the processes are too costly the number of collections able to adopt, or contribute to, the system will be reduced and the collaborative research potential of linking a large number of similar collections will be lost.

A detailed description of the research carried out into recording the watermarks found in sheets of paper can be found in a recent publication.<sup>1</sup> Although the image capture techniques that are described in the publication were selected specifically for a paper-based collection it includes some useful guidelines on how to achieve high quality digital images using a range of different film types and developing techniques. These are essential for the image processing routines involved in CBIR. It was quite clear from the work carried out that a quality digital camera and flatbed scanner could offer a valuable, combined approach to copying and digitising, and this might well be the most effective approach for many collections. However, we also found that alternative approaches were required in situations where the transmitted light necessary for this procedure was not effective i.e. when media obscured the appearance of the watermark or when the paper had been attached to a thick secondary support. This situation might well be more frequent when collections are made from impermeable materials for which transmitted light is inappropriate.

The equipment selected for the digitisation process must provide images of a quality suitable for use with the CBIR system. The specification and price of digitising equipment will inevitably change in the future therefore details of equipment used during our research would have to be considered alongside up-dated, or new, models. Details of the digitising equipment used can be found in the Appendix B.

## CONTENT BASED IMAGE RETRIEVAL

CBIR is a computer based technique for retrieving images from a large collection on the basis of features (such as colour, texture, and shape) that can be automatically extracted from the images themselves. The features used for retrieval can be either primitive or semantic, but the extraction process must be predominantly automatic. Retrieval of images by manually assigned keywords is definitely not CBIR as the term is generally understood – even if the keyword is ‘describe image content’.

CBIR has potential and actual application in engineering, medical imaging, design, journalism, law enforcement, and more recently museology and art history<sup>2</sup>.

Previous research has been carried out into the application of CBIR with watermark images at the University of Geneva<sup>3</sup>. Their system, SWIC (Search Watermarks by Image Content), uses either manually-assigned codes, such as those terms used by the early filigranologist C.M Briquet, the standard descriptors of the IPH (International Association of Paper Historians), or automatically extracted features such as the number, shape and relative locations of image regions.

The IIDR at the University of Northumbria at Newcastle upon Tyne has been carrying out research into shape retrieval for a number of years, particularly in the area of trademarks. Their prototype system – ARTISAN (Automatic Retrieval of Trademark Images by Shape ANalysis) – for abstract trademark images<sup>4</sup> was initially developed in collaboration with the UK Patent Office.

Whilst there are obvious stylistic differences between modern trademark images and historical watermarks, their similarities from an image processing point of view are quite striking. Both are monochrome images made up of a number of individual components and both rely on shape elements (rather than colour or texture) to give them visual impact and distinctiveness. However, digital images of watermarks, whilst similar to trademarks in many ways, offer more complex issues:

Trademarks are generally black and white (binary level) images. Watermark images tend to be greyscale (8-bit) images.

Scale and orientation has little importance in trademark images. However, they are important to the successful identification of watermarks.

Watermark images are encountered in a variety of very different formats (tracings, radiographic images, rubbings, photographs etc.)

Watermark images often suffer from interference due to background noise. This can be due to the chain and laid lines found in the paper, inconsistencies in paper density, or from the presence of a design on the paper surface.

## *When Images Work Faster than Words*

Based on this existing knowledge, the SHREW system has been developed to automatically retrieve watermark images based on shape. SHREW is a modular system, which processes the digital watermark images by extracting the most important shape features and storing them in a database. A new watermark can then be matched against those stored in the database and the images judged most similar to the query shape will then be retrieved in rank order. The stages in the processing of the watermark images are as follows:

**Image clean up:** Extracting the initial shape from a typical watermark image can be problematic, due to the background noise and uneven image intensity distribution. A large proportion of the images are poor quality. Drawing heavily on the Institutes' experience in shape description, various image-processing techniques are carried out in order to enhance the watermark.

**Segmentation:** The watermark image shape is segmented in order to isolate the watermark so that edge detection can be performed.

**Extraction:** Potentially useful shapes are then extracted from the segmented image.

**Assignment of measures:** The shapes are assigned various shape measures and further processing is carried out in order to remove unwanted shapes and close broken lines.

**Construction of shape database:** The remaining shape measures are then added to the database of shape vectors.

Several clean-up operations are performed to reduce the number of usable boundaries for inclusion in the shape-database from which retrieval can subsequently be performed.

For convenient use of the CBIR facility, the SHREW engine is integrated with the database of watermarks, available on the Web. Users, having searched the database and retrieved a set of thumbnails of watermarks matching the search criteria, can then perform a search for similar images based on selecting a particular watermark as a query image. SHREW returns images ranked in order of similarity to the query, from which the user can make further selections. Fig. 1 is giving an example.

### NORTHUMBRIA WATERMARKS ARCHIVE

The Northumbria Watermarks Archive consists of a database of watermark images together with associated and related data. The CBIR facility, SHREW, is linked to the database and is enacted by a user electing to do a search for similar watermark images based on a query image.



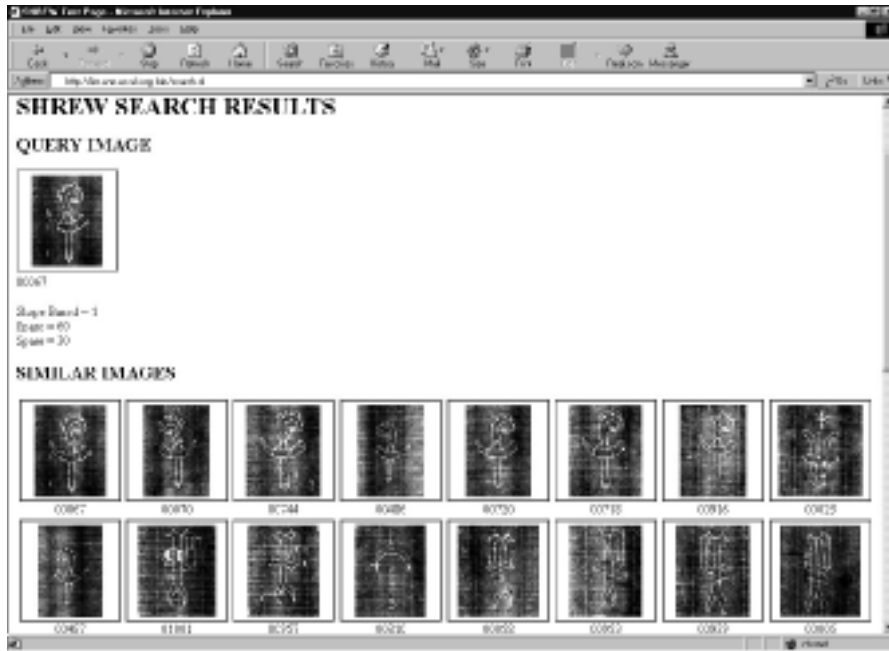


Fig.1: An example of retrieval results from SHREW.

The database that underpins the Archive has been written using public domain software, including:

A fast database for the storage and retrieval of information documents and images using MySQL (<http://www.mysql.org>)

Access over the web via the Apache web server (<http://www.apache.org>) equipped with Tomcat Java server (<http://jakarta.apache.org>)

The front end of the system has been developed as a Java 2-compliant web application consisting of standard servlets and Java Server pages (<http://java.sun.com>)

The integration of the archive with the SHREW content-based image retrieval software was implemented in Perl (<http://cpan.org>)

The overall schema of the system is shown in Fig. 2. The database can be accessed via any standard web browser without additional plugins. It is hoped that if the prototype database proves to be successful, a larger more stable database, based on ORACLE software, will be produced for the Archive. This will provide the structure for a web based archive offering greater flexibility and functionality

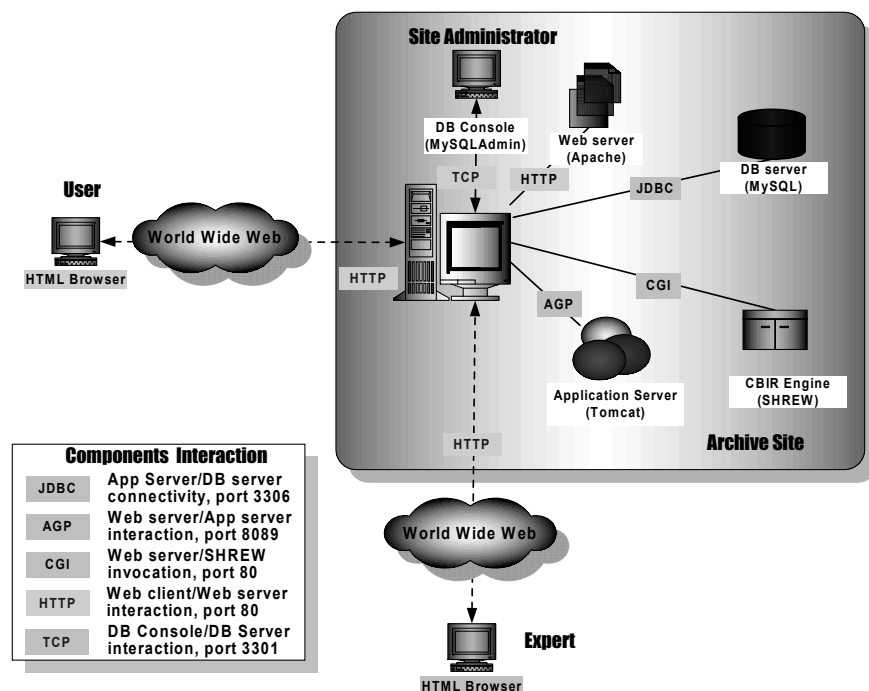


Fig.2: General Architecture of Northumbria Watermark Archive

and the inclusion of more extensive conservation, technical and art historical information.

The process of building the Archive had to solve several technical and organizational problems, common to all public repositories of such kind, such as:

- The cost for implementing and maintenance of an entirely public archive,
- The need to control the process of making additions or alterations to the archive.
- The balance between the graphical content and the documentation of the watermarks.

Firstly, the software used for the archive is not just public domain software, it is also a popular software package that requires minimal effort in its maintenance. MySQL and Postgres were both considered, but the decision was taken to go for the former because of its speed, reliability and wide support of languages. Similar factors determined the choice of Perl as a main integration language.

Secondly, although the archive is in the public domain, the information workflow for uploading images to the repository is not entirely open. Initially a fully

open system was foreseen – including automatic uploading, image processing and repository registration. However, such a workflow can easily lead to discrepancies within an archive and consequently major problems in maintaining it. This problem has not yet been resolved.

The third problem came from the lack of adequate models for the archive. This is discussed further in the next section.

In addition the system has been designed to maintain a full history of any changes that have been made to either the images or the data that have been recorded. Since it is hoped that other collections will adopt a similar format, provision has been made to record any information related to the images themselves and the processing methods that have been used i.e. metadata.

Since our database will combine data and images from more than one source it was necessary to incorporate a system that would provide clarity of ownership. This was one of the early issues raised by potential collaborators, who were anxious that their collection would be subsumed and lost within a system belonging to someone else. As a result a facility has been incorporated that allows the input of all information and images to be tracked to its original source or owner.

The Northumbria Watermark Archive was intended from the beginning to provide dynamic access to those interested in uploading images or information to add to the archive. As mentioned above, work has still to be completed as to which approach to this facility will be the most feasible.

#### THE USE OF STANDARDISED TERMINOLOGY WITHIN A DATABASE

The use of CBIR does not completely exclude the use of key words and descriptors. In fact, providing the descriptive terminology is reasonably comprehensive, it can play a very important role in speeding up access to the various categories, or fields, within which the collection has been organised. One of the major drawbacks to accessing and cross-referencing the major printed publications on watermarks was that they had all used arbitrary descriptive terminology and different cataloguing systems. As a result it was difficult to access information in one publication, let alone cross-reference between publications. In order to maximise access to a new database, it should facilitate links with similar collections thereby accommodating research, collaboration and contribution opportunities beyond the confines of one collection. This can generally be achieved if the registration is carried out using a standardized, descriptive terminology that is common to, or can be adopted by, the majority of collections involved.

In order to achieve this, it is necessary to investigate the standards that have been used to describe similar collections and consider their suitability not only for a digitised database but also for the type of access that you intend to provide for end users. Fortunately, within the museum and library community, there are a wealth of existing standards (eg. CIDOC, SPECTRUM, VAN EYCK and especially the Dublin Core Metadata set) to choose from. In some cases existing standards are insufficiently comprehensive or, alternatively, too detailed. Many standards were not developed for use with an electronic archive and are not immediately compatible. Initially our intention had been to implement the standard of the International Association of Paper Historians (IPH) in 1992, (updated in 1997). This represented the only existing standard for registration of historic papers and watermarks. Although the standard was comprehensive it was difficult to apply to an electronic archive as it failed to address several fundamental principles of database design and many of the classification schemes were not relevant to digital records:

The standard sometimes included more than one type of data within a single field, which overlooked one of the fundamental principles of database design.

Some fields contained multiple discrete pieces of data, which could lead to problems regarding searching, alphabetizing, calculating etc.

Limitations of early database software, led the original designers of the standard to represent field names with numerical codes rather than full names. This was confusing and is no longer a necessity, since the facility to include full text field names is present in all modern database software.

There was no method for those submitting information to record the rationale behind their judgements.

There was no method for tracking changes to the information provided.

There was no 'other' option in many of the fields for image description. As a result, it was impossible to enter a description or search for an extraordinary item not specifically described in the existing standard.

The standard did not distinguish between raw data and scholarly or professional judgements. This made data entry a tedious process: *If reporting papers becomes tedious and time consuming, requiring scholars to do work that goes beyond their research and publication objectives, they will be less likely to take the time to do it<sup>5</sup>.*

In order to adopt the IPH standard it was necessary to modify and develop certain aspects of the original format. Allison<sup>5</sup> has re-organised the IPH standard according to accepted entity-relationship database principles and the abstract

coding, used to identify the various fields, was replaced with more detailed and intuitive field names. This has resulted in a system that is much more user-friendly and compatible with the needs of the end user. The University of Northumbria has further enhanced the standard in order to maximise the potential benefits to as wide a range of end users as possible i.e. conservators, paper and art historians, forensic scientists etc., as well as those with a more general, but non-specific, interest in the collection. As a result additional conservation and technical metadata have been included and thumbnail images have been provided in order to facilitate browsing and fast viewing over the Web.

#### EVALUATION OF THE ARCHIVE

The Northumbria Watermark Archive has been made available via the Web as a prototype system. Whilst evaluation of the retrieval performance of SHREW was built into the research project, an evaluation of the Archive in terms of its effectiveness and usability was not planned within the project.

The research team are currently seeking further funding to develop the Archive into a major resource for those interested in historic papers and watermarks. In addition, further software development is needed, not just to improve its shape extraction and shape matching sub-routines, but also to extend SHREW's capabilities to handle complex and poor quality images and the positioning of chain and laid lines.

In order to inform the future development of the Archive we are preparing an online questionnaire, in order to obtain feedback on the interface, ease of use and the value of the information provided for users, as well as their views about the CBIR facility and the retrieval effectiveness of SHREW.

#### CONCLUSION

The developmental stages of establishing a database combined with CBIR will inevitably require a period of bedding in. However it is hoped that our recent experiences will help others avoid some of the issues described, or at least anticipate them well in advance. There is no doubt that the integration of the CBIR system with a data base can provide an invaluable research tool for a range of different types of collection, providing unprecedented access for research as well as teaching and learning.

APPENDICES

*Appendix A: Copyright Issues and Intellectual Property*

Copyright for digital resources remains an extremely complicated and unsettled issue and varies from one country to another. It can be absolutely crucial to the success or failure of a digital project. The most important issue regarding copyright is that in the European Union (EU), the creator, except when stated otherwise, automatically retains copyright. There is no official copyright register, and no application is necessary. Copyright begins when the material is created.

All digital content is bound by copyright laws for literary works (even images and video) as they are constructed from binary code. In the EU, copyright is held for 70 years after the death of the creator, or 70 years from publication if retained by the publisher. Specific laws holding copyright for 50 years from their creation, or for the full term of the materials involved also bind databases.

Copyright must be sought from the holder for all digital content. Copyright is usually held by the author/creator of the material, but this is not always the case. Permission must be sought from the author/copyright holder in all cases. If it is impossible to contact the copyright holder, or if there is no response, copyright clearance cannot be given. Although the EU directives of the Berne Convention and the Universal Copyright Convention have harmonised certain aspects of European copyright laws, each country will have its own laws to which applicants must adhere

The creator of an image, i.e. the photographer, automatically retains the copyright, unless stated otherwise by the institution, even when permission has been granted for the use of a photograph of an object. Even when an artist has been dead for 70 years, there can still be contractual objections to overcome before photography and/or dissemination is agreed. Additionally, if text or image is re-published in a new addition, copyright is held on the typography/design of that publication, not the original text or image.

For digital material published on the Internet, the issue is more complex. All material on the web is automatically protected by copyright. For clearance purposes, it may be that the copyright laws of the country in which your account is provided must be observed. However, as yet there has been no international agreement on this. In the EU, however, it may be assumed that most countries will be signatories to the Berne and Universal Copyright Conventions, and these laws will apply. As a rule, it is illegal to download material in the UK, which, if you had made an analogue copy, would have infringed copyright.

Databases intended for educational use, or with limited dissemination and access (such as an institutional intranet) may be given clearance more easily than those intended for publication on the web. Not to obtain clearance for material published on the web, could prove to be a very expensive mistake, as doing so may infringe *both* copyright and dissemination laws.

### *Appendix B: Digitising equipment*

The value of a good copy table for the creation of digital images cannot be over-estimated. The IFF copy stand that we utilised, complete with inbuilt tungsten halogen light box, allowed easy copying of flat objects with a digital camera. The intensity of the light produced by the light box allowed excellent transmitted light images of paper objects to be captured.

A comparison of resolution figures can be misleading when selecting a digital camera. Manufacturers often quote in terms of interpolated resolution rather than optical resolution. (Interpolation being the process where the value(s) between pixels are estimated by software). It is essential that optical resolution be used to define the resolution of a camera, as the quality of the software used to create the interpolation varies widely.

Area Array cameras use a fixed grid of vertical and horizontal cells on the CCD (Charged Coupled Device), which replaces the film in the camera. During image capture the entire array is exposed at once and are therefore suitable for moving images and flash photography. The price of these cameras ranges from budget to high-end 35mm and substitute backs for medium format cameras. These cameras tend to use built-in storage and are relatively easy to use. Digital SLR cameras such as the Nikon D1 or Kodak DCS760 offer extremely good image quality, portability and the flexibility of interchangeable lenses. Mid-range cameras (under £1000 / US\$1494 / €1561\*) such as the Nikon Coolpix series are more compact, and can offer comparable results.

High end linear and area scanning array cameras are available as stand alone models or scanning backs for studio cameras with high quality lenses. These offer exceptional quality images with extremely high optical resolutions and subsequently large image files. The image is not usually stored in the camera, but sent to a dedicated computer. This means that portability is limited and as capture times are lengthy – often several minutes per exposure – and not suitable for moving images or flash photography. Suitable high quality lighting is

\* Exchange rate of June 2002.

## *When Images Work Faster than Words*

usually required. Additionally, the high image quality means that extremely large file sizes are created (often 50–100 MB+) requiring storage solutions. Scanning cameras require more skill to use and are considerably expensive to purchase (£12,000 / \$ 17,929 / €18,731\* and upwards).

A good quality flatbed scanner with transparency unit (TPU) and high optical resolution (1600x3200 dpi and higher) is invaluable for digitising large format negatives and radiographs. For 35mm slides and negatives, the quality obtained from a dedicated slide scanner is much higher (2,700–4000 dpi), and is to be preferred. Again optical rather than interpolated resolution should be quoted. For connectivity, SCSI and Firewire® (Apple) are preferred for their fast transfer speed. However, this necessitates the purchase of an additional PCI cards. Most mid-range scanners utilise USB, widely used on both PC and Macintosh platforms, and only marginally slower than SCSI.

The specification of a PC usually needs to be higher when processing high quality (large) image files. The current version of Adobe Photoshop® (7.0) requires *at least*: Windows® 98 (or Me/XP/2000/NT4.0®), Pentium® II or similar, 128Mb RAM, 165Mb available hard disc space. Or for Apple Macintosh systems: PowerPC® G3/G4, MacOS9.0 or OSX10.1, 256Mb RAM, 240Mb available hard disk space. A SVGA monitor capable of resolutions of 1024x768pixels and an Adobe Postscript compatible printer are recommended.

### *Appendix C: Abbreviations*

AGP: Accelerated Graphics Port

AHDS: Arts and Humanities Data Service

AHRB: Arts and Humanities Research Board

ARTISAN: Automatic Retrieval of Trademark Images by Shape ANalysis

CBIR: Content Based Image Retrieval

CCD: Charged Coupled Device

IIDR: Institute for Image Data Research

IPH: International Association of Paper Historians

PCI: Peripheral Component Interconnect /Interface

SCSI: Small Computer Systems Interface

SHREW: SHape REtrieval of Watermarks

SVGA: Super Video Graphics Array

SWIC: Search Watermarks by Image Content

TPU: TransParency Unit

USB: Universal Serial Bus

\* Exchange rate of June 2002.



SUMMARIES

*When images work faster than words – The integration of content-based image retrieval with the Northumbria Watermark Archive*

Information on the manufacture, history, provenance, identification, care and conservation of paper-based artwork/objects is disparate and not always readily available. The Northumbria Watermark Archive will incorporate such material into a database, which will be made freely available on the Internet providing an invaluable resource for conservation, research and education.

The efficiency of a database is highly dependant on its search mechanism. Text based mechanisms are frequently ineffective when a range of descriptive terminologies might be used i.e. when describing images or translating from foreign languages. In such cases a Content Based Image Retrieval (CBIR) system can be more effective. Watermarks provide paper with unique visual identification characteristics and have been used to provide a point of entry to the archive that is more efficient and effective than a text based search mechanism.

The research carried out has the potential to be applied to any numerically large collection of images with distinctive features of colour, shape or texture i.e. coins, architectural features, picture frame profiles, hallmarks, Japanese artists stamps etc. Although the establishment of an electronic archive incorporating a CBIR system can undoubtedly improve access to large collections of images and related data, the development is rarely trouble free. This paper discusses some of the issues that must be considered i.e. collaboration between disciplines; project management; copying and digitising objects; content based image retrieval; the Northumbria Watermark Archive; the use of standardised terminology within a database as well as copyright issues.

*Lorsque les images sont plus rapides que les mots – L'Intégration de CBIR avec les Archives filigranées de l'Université de Northumbria*

Les informations relatives à la fabrication, à l'histoire, à l'origine, à l'identification, à l'entretien et à la conservation d'objets d'art sur papier sont très disparates et souvent difficilement accessibles. Les Archives filigranées de l'Université de Northumbria ont l'intention d'incorporer du matériel d'information dans une banque de données qui sera librement accessible sur Internet fournissant ainsi des ressources inestimables pour la conservation, la recherche et la formation.

L'efficacité d'une banque de données dépend fortement de ses mécanismes de recherche. Les mécanismes basés sur des mots sont fréquemment inefficaces lorsqu'un éventail de terminologies descriptives est utilisé, p. ex. dans la description des images ou dans la traduction à partir de langues étrangères. Dans de tels cas un système CBIR (Content-Based Image Retrieval) peut être plus efficace. Les filigranes confèrent au papier des caractéristiques uniques d'identification visuelle ; ils ont été utilisés comme moyen d'accès aux Archives beaucoup plus efficace qu'un mécanisme de recherche basé sur des mots.

Le système de recherche présenté ici est susceptible d'être appliqué à toutes grandes collections numériques d'images présentant des caractéristiques différentes de couleur, de forme ou de texture, comme p.ex. des pièces de monnaie, des plans de construction, des profils des cadres de tableaux, des poinçons, des estampes d'artistes japonais etc. Bien que l'installation d'un système électronique d'archives intégrant un système de recherche CBIR pourrait indubitablement faciliter l'accès à de grandes collections d'images et aux données s'y rapportant sa réalisation n'en

## *When Images Work Faster than Words*

reste pas moins liée à des problèmes. Cet article énonce certaines des questions qui doivent être prises en compte comme p. ex. la collaboration entre les différentes disciplines concernées, le management du projet, la copie et la numérisation, CBIR, les archives filigranées, la terminologie de la banque de données, les droits d'auteur.

### *Wenn Bilder schneller sind als Worte – Die Verbindung von CBIR mit dem Wasserzeichenarchiv der Universität von Northumbria*

Informationen zur Herstellung, Geschichte, Herkunft, Identifizierung, Aufbewahrung und Konservierung von Kunst auf Papier sind weit verstreut und nicht immer leicht verfügbar. Das Wasserzeichenarchiv der Universität Northumbria versucht entsprechendes Material in eine Datenbank einzubringen, die über Internet frei verfügbar und so ein unschätzbares Hilfsmittel für Konservierung, Forschung, Ausbildung sein soll.

Der Wert einer Datenbank ist in hohem Maße von ihren Suchmechanismen abhängig. Verbale Suchmechanismen sind oft dann wenig effektiv, wenn nicht eindeutig festgelegte, sondern beschreibende und/oder aus anderen Sprachen übertragene Suchkriterien verwendet werden. In diesen Fällen kann ein CBIR-System (Content Based Image Retrieval) effektiver sein. Wasserzeichen sind für Papier ein einzigartiges visuelles Identifikationsmerkmal; das System wurde eingesetzt als besseres Mittel des Zugangs zu dem Wasserzeichenarchiv als ihn ein auf Text basierendes Suchsystem hätte bieten können. Das hier vorgestellte System ist verwendbar für jede aus zahlreichen Einheiten bestehende Sammlung von Bildern, die sich in Farbe und/oder Form unterscheiden, wie z.B. Münzen, Bauzeichnungen, Bilderrahmenprofile, Echtheits- und Künstlerstempel, u.a.m.

Eine maschinenlesbares Archivs mit CBIR-Suchsystem kann ganz sicher den Zugang zu großen Sammlungen von Bildern und ähnlichem Material erleichtern; es zu erstellen ist aber nicht frei von Problemen. Es werden einige dabei zu beachtende Gesichtspunkte diskutiert, wie z.B. Zusammenarbeit der verschiedenen betroffenen Bereiche, Projektplanung, Kopieren und Digitalisieren, CBIR, das Wasserzeichenarchiv, Datenbankterminologie, Urheberrecht.

## REFERENCES

1. Brown, A.J.E., & R. Mulholland: *The Northumbria watermark archive: Using microfocus X-Radiography and other techniques to create a digital watermark database*. Preprints. Works of Art on Paper: Techniques and Conservation, IIC 19<sup>th</sup> International Congress, Baltimore 2002.
2. e.g.: Ward, A. A., M.E. Graham, K.J. Riley & N. Eliot: *Collage and content-based image retrieval: Collaboration for enhanced services for the London Guildhall Library*. Presented at Museums and the Web 2001, Seattle, 15–17 March 2001.
3. Rauber, C., P. Tschudin & T. Pun: *Retrieval of images from a library of watermarks for ancient paper identification*. Proceedings of EVA 97, Elektronische Bildverarbeitung und Kunst, Kultur, Historie. Gesellschaft zur Förderung angewandter Informatik e.V. 1997.
4. Eakins, J. P., J.M. Boardman & M. E. Graham: *Similarity Retrieval of Trademark Images*. IEEE Multimedia, April–June 1998: 53–63.

5. Allison, R.: (1997) *Critique of the IPH standard with proposals for a www distributed database system*. <http://www.bates.edu/Faculty/wmarchive/wm-initiative/iph-commentary.html>. 1997, accessed Jan. 2002.

A. Jean E. Brown<sup>\*</sup>, Senior Lecturer  
Richard Mulholland, Research Assistant  
MA Conservation of Fine Art Programme  
University of Northumbria  
Newcastle-upon-Tyne NE1 8ST  
England  
Tel: +44 (0)191 227 3331  
Fax: +44 (0)191 227 3250  
E-mail: [jean.brown@unn.ac.uk](mailto:jean.brown@unn.ac.uk)

Margaret Graham, Principal Lecturer  
Vassil Vassilev, Senior Lecturer  
School of Computing and Mathematics  
University of Northumbria  
Newcastle-upon-Tyne NE1 8ST  
England

Jon Riley, Software Engineer  
John Eakins, Professor, Director  
Institute for Image Data Research  
University of Northumbria  
Newcastle-upon-Tyne NE1 8ST  
England

Karen Furness, Senior Research Administrator  
Institute For Image Data Research  
D107 Ellison Building  
University of Northumbria  
Newcastle-upon-Tyne NE1 8ST  
Tel: +44 (0) 191 227 4646  
Fax: +44 (0) 191 227 4637  
<http://www.unn.ac.uk/iidr/>

\* Author to whom correspondence should be addressed.