

# **EROS: An Open Source Multilingual Research System for Image Content Retrieval dedicated to Conservation-Restoration exchange between Cultural Institutions**

**Christian Lahanier<sup>1</sup>, Geneviève Aitken<sup>1</sup>, Jiro Shindo<sup>2</sup>,  
Ruven Pillay<sup>2</sup>, Kirk Martinez<sup>3</sup> and Paul Lewis<sup>3</sup>**

1 Centre de Recherche et de Restauration des Musées de France  
6, rue des Pyramides  
75041 Paris Cedex 01  
tel : +33 1 4020 5871  
fax : +33 1 4703 3246  
[christian.lahanier@culture.fr](mailto:christian.lahanier@culture.fr)

2 Digital Publishing Japan  
3 Rue de Medicis  
Paris 75006  
France  
tel: +33 (0) 1 4326 7834  
fax: +33 (0) 1 4326 7834  
[shindo@dp-j.com](mailto:shindo@dp-j.com)  
[ruven@dp-j.com](mailto:ruven@dp-j.com)

3 Intelligence, Agents and Multimedia Research Group  
Dept. of Electronic and Computer Science  
The University of Southampton  
Southampton SO17 1BJ  
tel: +44 23 80 59 44 91  
fax: +44 23 80 59 28 65  
[km@ecs.soton.ac.uk](mailto:km@ecs.soton.ac.uk)  
[phl@ecs.soton.ac.uk](mailto:phl@ecs.soton.ac.uk)

## **Introduction**

In 1990, an EU-funded project entitled NARCISSE (Network of Art Research Computer Image SystemS in Europe) was launched to build a multilingual database to manage museum laboratory documentation relating to painting materials (1), (2) et (3).

This project had been in preparation since 1985 at the French national museum laboratory by a consortium of museum laboratory scientists from across Europe, to fulfil their needs concerning conservation-restoration research, classification, information retrieval and exchange of data techniques.

It was decided to design a system with a well-defined and multilingual vocabulary to describe:

- the works of art by means of historical and museological criteria,
- the technical data relating to the photographic archives (photographic and X rays films),
- the restoration and study reports with information on painting techniques, ageing processes and restoration procedures (4).

The initial computer system was developed by the Sully-Group firm in 1990 and updated in 1996 with a Web client-server system (5). This has been used internally to manage and consult a huge bank of high definition digital images of 15 000 paintings and 35 000 objects. The Intranet access limited the end-users to the useless civil servants of the French Ministry of Culture during the five year period of digitisation of the scientific photographic-archives.

Currently, more than 150 000 photographic and X Ray films on glass or plastic have been scanned at high definitions of up to 6 000 dots x 8 000 lines.

### 1.The C2RMF Installation

The C2RMF installation consists of hardware donated by Hewlett-Packard as part of it's philanthropy program:

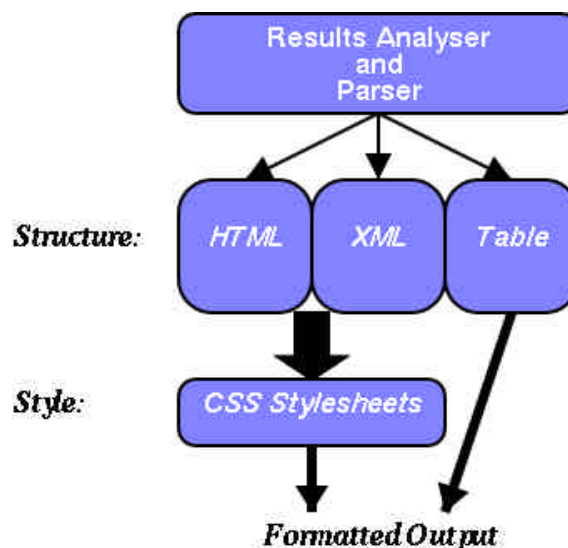
- \* 2 Hewlett-Packard Ultra-SCSI RAID5 dual 700MHz serving a total of 3 TB.
- \* 1 Hewlett-Packard quad 550MHz server
- \* Private 100Mbps switch linking them together

The web server, PHP engine and IIP tile server are installed on the quad-server. The images and databases are on each of the two dual processor Raid servers.

#### 1.1 The development of an advanced open source multilingual archive system

Ten years after the NARCISSE project, a new generation of database system for museum research laboratories is required to face the challenges of the 21<sup>st</sup> century.

- The new system is entirely open source and available for anyone to use, being based on powerful and industry-leading software such as Linux™ OS, Apache, MySQL and PHP.
- Web-based and platform independent, allowing for use internally through an intranet, externally through an extranet and via the internet for general public



- Flexible, allowing for easy customisation for individual needs.
- Fully and transparently multilingual – searches on the data can be performed on all the information in any language.
- Standards compliant – the use of XML etc. allows complex data interaction and analysis to take place both within the server and by the client.
- Distributed – systems in different institutions can be consulted simultaneously and the results aggregated.
- Integrated colour calibrated multi-resolution image viewing based on the ACOHIR system for both flat and 3D objects.
- Extensible and able to evolve over time – allowing for extra modules, for example, for image content-control and secure remote printing (DPJ's VFZ modules).

#### 1.2 The Database Architecture

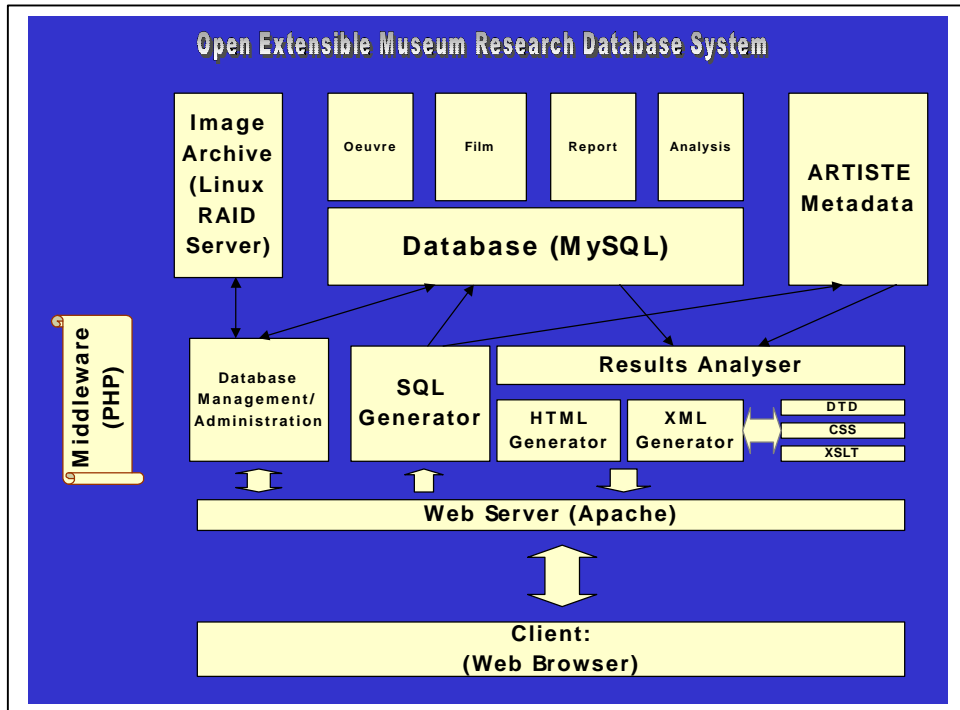


fig 1 Scheme of the system architecture

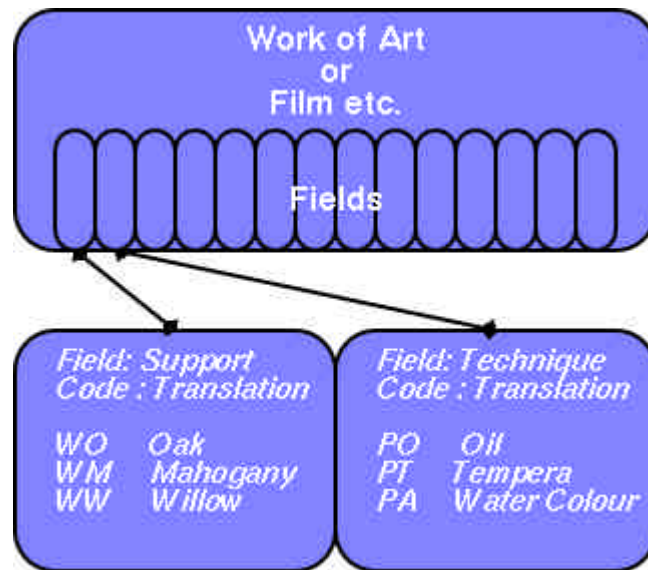
## 1.3 The Database Components

### 1.3.1 Thesaurus

A key feature of the system is the multilingual thesaurus capability. This thesaurus is organised as a set of hierarchical dictionaries for each translatable field and for each available language. The data within the main database is stored in a compact language independent format as short codes. This has the benefit of simplifying the kinds of searches required and the amount of data required to be stored in the main database. This produces a considerable improvement in efficiency and a speeding up of database access. When the search results are presented, the information in the main database is translated via the thesaurus system into the appropriate language.

The thesaurus is not only capable of handling a full lexical hierarchy, but also of handling synonyms and complex character sets, such as Japanese and Chinese, which can be managed via Unicode (UTF-8) encoding.

In addition, the system is able to handle multiple entries within a single field. For example, in the field for support, paintings are sometimes constructed out of several different materials. This information needs to be stored in an efficient and flexible manner. In the EROS database, multiple entries can be stored in a pseudo-mathematical form using the predefined codes. For a support consisting of a mixture of oak and mahogany panel, the resulting entry in the main database would be WO+WM, where + is the character used to separate terms. The thesaurus system, in this case, parses the string and is able to translate each individual term within it. This technique also allows for complex relationships to be created. By using a set of separator characters, various relationships can be defined in a mathematical way. For example, + could represent a mixture, while > could represent a transposition of the support from one material to another.



### 1.3.2 Query Interface

Queries are performed via a web browser based interface. Screens for simple or advanced queries can be easily created and the fields to be viewed customised by the system administrator. In addition, date or numeric size fields can be searched by specifying a range of dates or sizes between which searches are performed.

Users are able to select the working language and the domain of interest as well as the number of results returned and whether film and report results be shown.

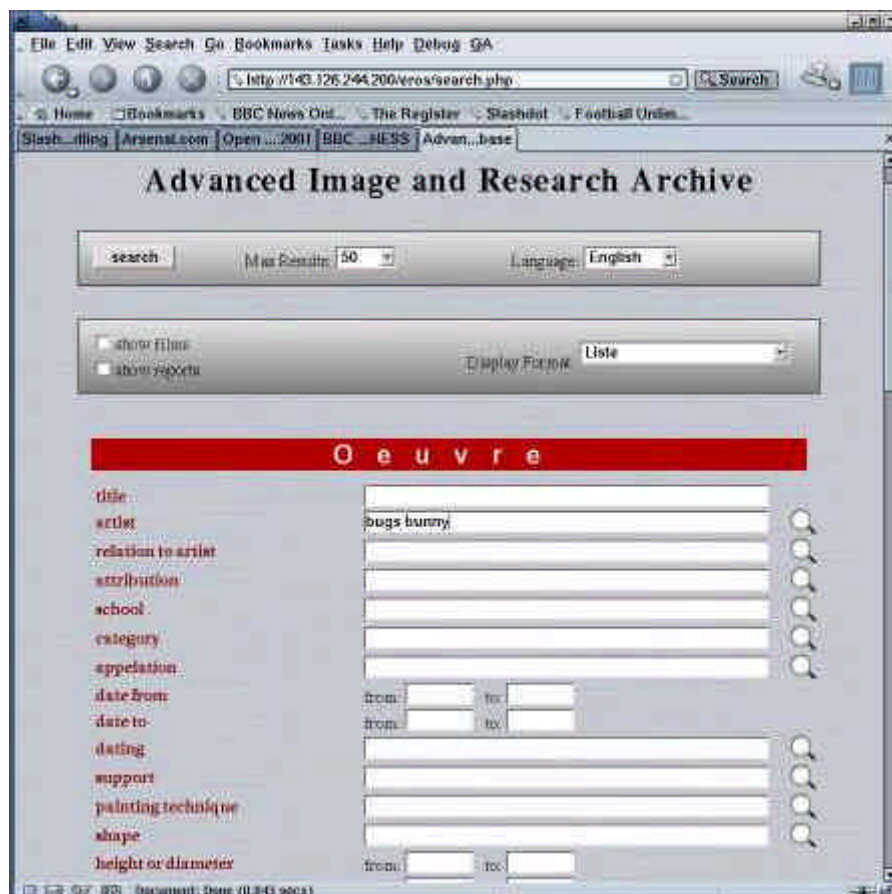
The interface is divided in three parts :

- 1- museological, historical and material information related the works (39)
- 2- technical, or management information related to the photographic document (19)
- 3- technical or management information related to the reports (13)

Where applicable, the user can choose vocabulary from a list of relevant terms classified alphabetically, or can type something directly in. A glossary in nine languages can be consulted on line. Fulltext searches can be made within each field.

The user is also able to specify the display or output format of the results eg. HTML, XML, plain text, formatted tabular, list of images, graphical, statistical analyses etc. It is also capable of label printing for photographic film management.

The paintings are classified by their work number. After a query, the number of matching paintings is shown, together with any matching associated films and reports.



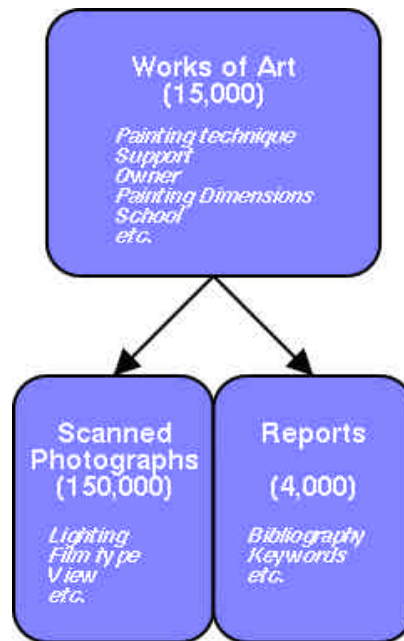
### 1.3.3 Data Entry

Data entry is also performed via a web browser interface similar to that used for querying. Users who enter data need to log on. Write, modification and suppression rights can be assigned and controlled by the system administrator for each user. Information such as name of the user and date of the entry are automatically filled in by the system.

The upper number used in the database as a work number, a film number or a report number appears in a window facing the field line. To maintain the integrity of the data being entered into the system, the controlled lists of relevant vocabulary within the thesaurus are used for each translatable field. When entering data via the web interface, users are required to choose the appropriate term from a menu of this vocabulary. These terms are, of course, available in any language required. When a data is saved into the main database, the vocabulary is translated into the language independent code representation. Where applicable, multiple terms can be combined into the pseudo-mathematical relational format.

### 1.4 Transfer of the NARCISSE database to the EROS system

The transfer of the old database to the new system was performed between July and October 2001. This was a considerable undertaking involving the correction of incorrectly stored data, the normalisation of this data and the creation of controlled lists of relevant terms for each field within the three databases. Scripts were written to automate the transfer and processing of the data efficiently. The database content has furthermore, been split into domains (type of work of art) to simplify the management of the thesaurus.



#### 1.4.1 Reports management

There are currently around 4000 study and restoration reports in the database. These had previously existed only in HTML format. Using tools supplied by the W3C (6), these have been cleaned up and converted into XHTML. All new reports are stored directly as XML and are structured using a custom Document Type Definition (DTD).

The use of XML enables complex structured content to be created. Such content can be analysed or searched in more meaningful and complex ways than was previously possible. Furthermore, the use of Extensible Stylesheet Transforms (XSLT) allows these documents to be delivered in any format eg. HTML, PDF, simple text etc. and with dynamically configurable styles of presentation.

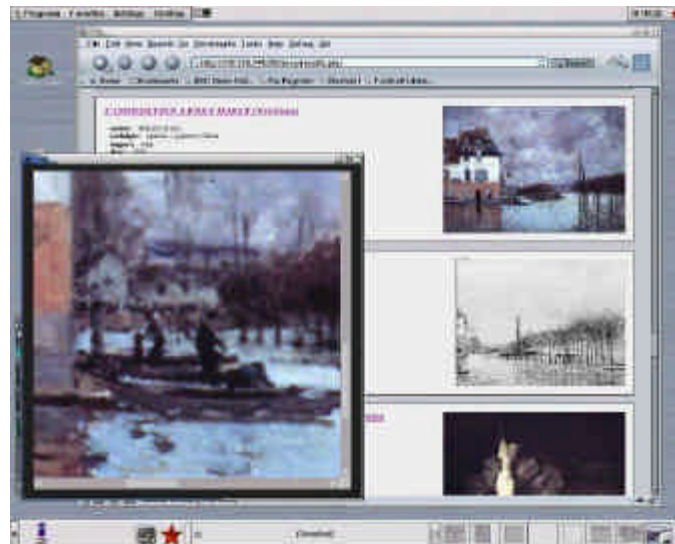
#### 1.5 High Resolution Image Viewing

Another key component of the database system is the capability to remotely view colour calibrated high resolution images of both 2D paintings and 3D objects.

Each image is stored as both a JPEG thumbnail for rapid previewing and in tiled pyramidal TIFF format for high-resolution viewing. A java applet permits multi-resolution viewing in conjunction with a tile server. This viewing system is based on the Internet Imaging Protocol and was initially developed for 2D images as part of the VISEUM (7) (8) project, then extended to handle 3D objects during the ACOHIR (9) project. The viewer works by requesting only the tiles at the appropriate resolution required for viewing a particular part of the image. The requested tiles are then dynamically JPEG encoded by the server and sent to the applet. In this way, images of any size can be viewed quickly across the internet. The 3D objects consist of numerous 2D images taken from various angles around the object. By allowing the viewer to change the angle of view, a 3D object can be navigated around.

The C2RMF image digital archive consists of JPEG encoded tiled pyramidal TIFF images which had been converted from the proprietary "Scopyr" format in use previously and of lossless "zip" encoded tiled pyramidal TIFF images for all new images entered into the system.



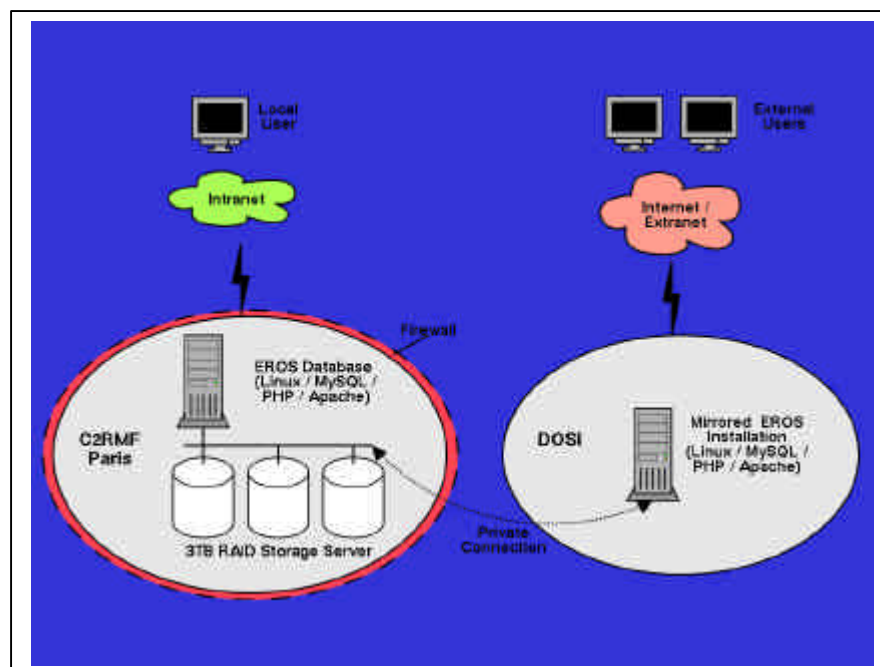


New functionality to be added to the viewer include :

- a scale associated to the image showing the size of a detail,
- multi-windowing to compare various images of the same work or images of different works,
- contrast, brightness, rotation, inversion, superposition, resizing, printing, etc.

### 1.6 Web Site Mirroring

The database has been mirrored at a publicly accessible site at the DOSI (Organisation and Information System Department) to permit link-ups between distributed research databases. The system limits the access to data according to users rights to indoor users (INTRANET), outdoor users (EXTRANET) and to the Web users.



## **1.7 Multilingual Access**

### **1.7.1 Multilingual Glossary Integration**

The multilingual vocabulary (Catalan, Danish, English, French, German, Italian, Portuguese, Spanish) set up for the indexing of the scientific and restoration reports in the frame of NARCISSE has been up dated. It constitutes the list of terms in the reports database to select works of art, images or reports.

More, the definitions of this vocabulary are accessible on line

### **1.7.2 The Database Translation**

The controlled lists of terms as well as the free text information fields (such as the titles) have been translated from French to English, to Portuguese and to Japanese. Unicode module have been integrated for the Asiatic language screens.

## **1.7 VFZoom Module, a "LossLess, Single Source, Controllable, Multi-Use Image Archive"**

The VFZ format is unique and its key features are:

- 100% lossless image format with small file size and 100% lossless reversibility from the VFZ format to common image formats such as TIFF, JPEG etc.;
- able to give any individual access to or use of images perfect control and management, by which the copyright is protected;
- resizable with high speed at the user side from 5 - 1200% of the original image size without loss of quality; original colour information inside RGB channels being exactly retained without artificial colour generation when resizing images;
- 6 quality level structure of an image file to be streamed through network with minimal file size;
- in-built meta data structure;

VFZ technology affords an image holder the opportunity to create a true digital "Lossless, One Source, Controllable, Multi-Use Image Archive", and provides its internal and external end-users with a multitude of flexible functions such as storage, distribution, re-sizing and printing of digital images.

The VFZ format overcomes the currently popular TIFF format with

- a file size reduction for storage without any loss of image quality,
- 100% reversibility from VFZ to TIFF (or several other image formats) without any loss of image quality
- resizability from 5% through 1200% of the original image without loss of quality. The

In addition to these benefits for the image holder who holds images in TIFF, VFZ also provides end-users of digital images, who might have dealt with images in other image formats such as JPEG, with various new uses of lossless images. With VFZ, end-users can freely crop images, resize them from 5% through to 1200% very quickly, save them for storage or printing. Furthermore, in the VFZ format, the image data is stored in six different quality levels within a single VFZ file, and can be transmitted to end-users at any single level so that only the desired quality level and optimal data size are distributed through network. End-users can, at high speed and without bothersome increases in server burden, ascend quality levels as desired utilising the patented streaming function of VFZ technology. Cropping the desired portion of the image and streaming the desired level of quality thus ! allows end-users to access, import

Also VFZ technology includes within the format the capability to store extra information about the images in the form of meta-data. This data can be easily incorporated from or



linked to a contents information database, allowing quick access through a third party search engine to particular images in the archive.

## **2 Image Content Recognition**

Content-based image retrieval is a challenging and active research area with the potential to provide powerful tools for image searching. Although many techniques have been described in the research literature, the capabilities of current content matching systems are still basic general purpose approaches although some powerful applications specific methods can be developed. General techniques based on such features as colour distribution, texture, outline shape and spatial colour distribution have been popular in the research literature and in content based retrieval systems. Within the European ARTISTE project (An Integrated Art Analysis and Navigation Environment) (programme IST1 1999-1 n°11978, from January 2000 to June 2002) people are aiming to develop more targeted solutions for specific image-based queries. The digitisation of the photographic archives at the C2RMF allows us to sample and group together images with similar characteristics thereby providing the reference material for testing such image content recognition software.

### **2.1 Image Content Recognition of Paintings**

Four areas of application have been determined :

#### **2.1.1. The Documentation**

**2.1.1.1** Localisation of a sub-image in an overview which were photographed at different periods of time and in various experimental conditions. The sub-image can be taken during restoration showing lacuna. In fact, the co-ordinates of a slide (sub-image in the overview) will replace the commentary made for each slide in the database.

**2.1.1.2** Identification of a lost photography without any registration number by comparing it to the image metadata.

**2.1.1.3** Iconography characterisation is one the most complex subject. Some shape recognition seems to work such as portrait, landscape, still life, or themes such as crucifixion, or virgin and child.

**2.1.1.4** Label, inscriptions, or stamp characterisation on the back of the work of art.

#### **- 2.1.2. Conservation Restoration (automatic indexing, feature identification)**

**2.1.2.1** making techniques of supports : planks of a panel, canvas texture, stretcher structure and shape,

**2.1.2.2** restoration techniques: cradling, non homogeneous varnish, cheville, butterfly.

**2.1.2.2** ageing: lacuna, cracks, repaints, relining, pliures,

**2.1.2.4** painting techniques : knife, brush, incised decoration of gold leaf,

#### **- 2.1.3. The History of Art (indexing)**

**2.1.3.1** style of the artist : writing and colour.

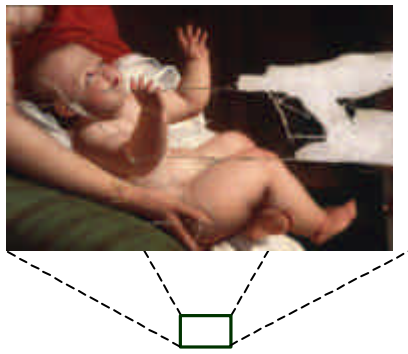
#### **- 2.1.4. Diffusion**

**2.1.4.1** similarity of colour (jade, painting) or of decoration (textile, coin).

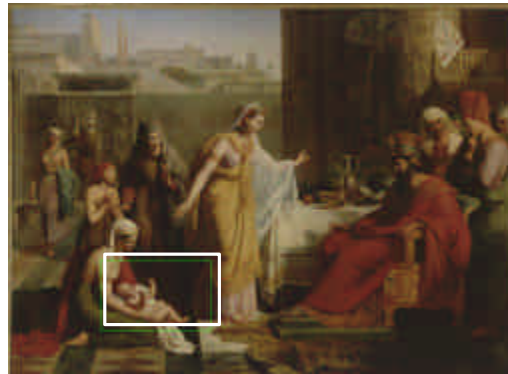
## **2.2 Sub-Image Location in an Overview**

### **2.2.1 Sub-Image Matching and the M-CCV Method**

A robust technique has been developed to retrieve the overview of a painting using a query image which represents all or part of this painting. The query image may have been captured at a different resolution or may have been distorted or degraded. For example, the query image may be a part of an image captured prior to restoration of a painting and the target image may represent the painting after restoration (10). This technique called the Multi-scale Colour Coherence Vector (M-CCV) method uses the colour coherence vector (CCV) (11).



Query



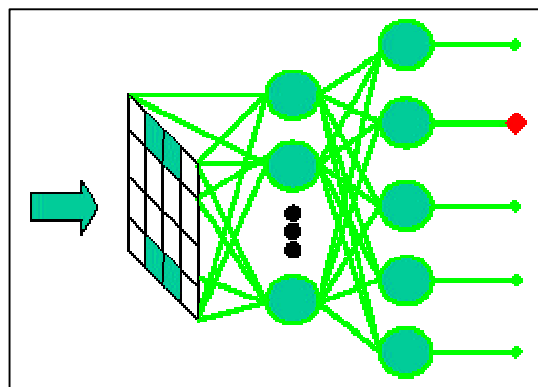
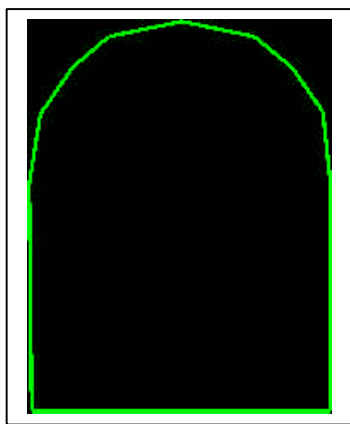
Retrieved result

Figure 1. a query image captured before restoration which is a fragment of Moses presented to Pharaoh painted by Victor Orsel in 1830, photography after restoration conserved at the Fine Art museum, in Lyon inv A 145 3.26x4.31m

### 2.3 Automatic Stretcher Shape Identification

One of the Artiste user-requirements was the identification of painting shapes in order to provide richer statistics for searches. This is useful for restricting areas of interest and avoiding backgrounds etc. In addition, classification of the stretcher type can be carried out in terms of its shape, its size and the nature and weaving of the original canvas for approximately dating of the canvas manufacture.

This is carried out using trained neural networks capable of classifying the shape of the outline for example: oval, circle, etc. This automatic stretcher characterisation of the can then be stored in the database.



### 2.4 Cradling Automatic Recognition

Cradling is used to reinforce aged wood panels and protect them against the environmental effects of humidity and temperature. They are applied fixed or partly flexibly at regular interval onto the back of the painting.

#### 2.4.1 Problematic

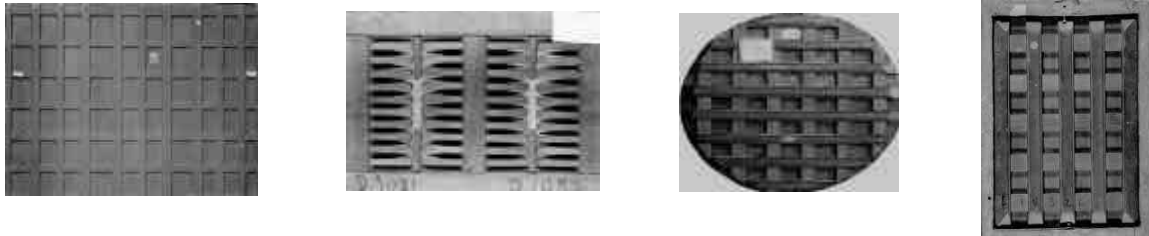
The Painting database at the C2RMF contains 15 000 works and 140 000 digital images. Automatic recognition of a crossed structure on the back of paintings on wood permits the indexing of cradling into the database.

## 2.4.2 Methodology and Results

Eight phases have been set up :

- 1- the sampling of photography of the reverse of painting on wood : 2 775 images
- 2- the analysis (application of the image recognition module) of the samples : 435 cradling identified and 2 340 rejected
- 3- the validation control of the selected and rejected samples (level of errors)
- 4- cluster analysis of the measured parameters (number of vertical and horizontal strokes), width of the strokes, inter-stroke width,
- 5- indexing test through metadata (sub-group identification),
- 6- validation control of the cradling sub-group (by restorers),
- 7- factorial analysis to find correlation between the cradling sub-groups and historical, museological and material information related to the painting,
- 8- interpretation of the correlation.

Cradling examples:



Thanks are given to :

- The Mission of Research and Technology of the French Ministry of Culture for its support, - - Danabalou Mohanassoundaram, engineer at the DOSI for the building of the mirror site at the DOSI and the conversion of the reports to XHTML.

- Rhéda Mamoudi, DESS of Computer Science, who has developed the cradle recognition module.

- Mrs Maria Guerra who translated the Portuguese translation of the Painting database.

- Mrs Junko Koga, from for the Japanese translation of the database.

- Mrs Moen, historian of art, Diplômée de l'Ecole du Louvre, for the translation of the glossary of scientific terms and their definition in Japanese.

Conservation-restoration partners in the ICOM-CC Documentation Group are invited to use and contribute to the collective development of the new database system and to participate in the development of additional functionality such as image content recognition or image processing and statistical analysis modules.

## Bibliography

- (1) Christian Lahanier, David Meili, Michel Aubert "Art and Science" a multilingual CD Rom Intelligent Multimedia Information Retrieval SystemS and Management, Rockefeller University, New York USA, October 11-13 1994
- (2) CD Rom NARCISSE *Glossaire multilangue*. - Lisboa, Arquivos Nacionais- Torre do Tombo, 1993 278p
- (3) Séminaire NARCISSE, *Actes*, Arquivos Nacionais - Torre do Tombo, 1993, 98p
- (4) NARCISSE *Système documentaire des peintures et enluminures*; Lisboa, .Arquivos Nacionais/ Torre do Tombo, nov. 1993, 353p

- (5) Ch. Lahanier, G. Aitken et M. Aubert NARCISSE: une bonne résolution pour l'étude des peintures Techné n°2 1995 pp 178-190
- (6) [www.w3c.org](http://www.w3c.org)
- (7) D. Saunders and J. Cupitt and R. Pillay and K. Martinez (1999) Maintaining colour accuracy in images transferred across the Internet. In : L. MacDonald and M.R. Luo (eds.) Colour Imaging - Vision and Technology. John Wiley, p.215-231.
- (8) K. Martinez (1997) Networking high quality images of art.
- (9) ACOHIR, <http://www.iam.ecs.soton.ac.uk/projects/achir/index.html>
- (10) Chan, S. Martinez, K. Lewis, P. Lahanier, C. Stevenson, J. (2001) "Handling sub-image queries in content-based retrieval of high resolution art images", ICHIM2001 Conference, September, Milan, Italy.
- (11) Greg Pass, Ramin Zabih, and Justin Miller. Comparing Images Using Color Coherence Vectors. In MultiMedia, pages 65-73. ACM, 1996.