

This document describes briefly the most relevant prototypes and demonstrators developed by DELOS partners, which might be the basis for advanced tools and services provided by future Digital libraries.

It also describes briefly the scope of the on going cooperation between DELOS and The European Library for a joint testing and evaluation of some of those technologies.

#### Tools and Services from DELOS partners

OpenDLib

OSIRIS/ISIS: Management and Search of Multimedia Collections

MILOS: a multimedia content management system for multimedia digital library applications

DAFFODIL

A Prototype P2P DL Network

Eurovision: a text-based cross-language image retrieval system.

Ciquest: organising images using concept hierarchies

eaSim and FLOSS

RoleMiner

Sightseeing4U

xSMART

SOMLib-based Interfaces to Digital Libraries

DARE

Annotating web accessible DLs with the MADCOW toolbar

ImageLab VideoBrowser

Content Based Image Retrieval

Content Based Retrieval of 3D Models

Video Semantic Adaptation

Pictorially Enriched Ontology Annotator

Accademia della Crusca (Crusca Academy)

COMISM (COMmunity Interface Semantics Model)

GraphOnto: Ontology Editing and Ontology Mapping, Multimedia MPEG7 and Ontology-Based Metadata Definitions

OntoNL: A Natural Language Interface Generator to Knowledge Repositories

The LUPA Index, A Demonstrator of The Referential Integrity Problem in Large RDF Networks

The ITeM Recommender System

DOMINUS (DOcument Management INtelligent Universal System)

Linking Paper and Digital

The cooperation between DELOS and The European Library

## **Tools and Services from DELOS partners**

### ***OpenDLib***

*Donatella Castelli, CNR-ISTI (Italy)*

OpenDLib is a digital library management system developed at the Istituto di Scienza e Tecnologie dell'Informazione "A. Faedo", of the Italian National Research Council, Pisa, Italy (<http://www.opendlib.com>). It consists of a federation of services that can be customised to satisfy the requirements of a target user community. This federation can be expanded at any time by adding other community specific services. The entire set of services can be managed and hosted either by a single or by a multitude of distributed organisations that collaborate on the maintenance of the shared digital library, each according to their own computational and human resources. An orthogonal system facility enables different user groups to define their own virtual view of the shared digital library, tailored to the specific needs and policies of the group. Several digital libraries based on this system have been built.

### **Information Space**

OpenDLib can handle a wide variety of document types with different formats, media, languages and structures. The same OpenDLib library can maintain, for example, a collection of journals and conference proceedings consisting of articles; a collection of theses in different languages, organised in chapters and sections; a collection of videos, structured into sequences and shots, and a collection of other documents represented only by the set of their bibliographic records. OpenDLib can also manage new types of documents that have no physical counterpart, such composite documents consisting of the slides, video and audio recordings of a lecture, a seminar or courses. It can also maintain multiple editions, versions, and manifestations of the same document, each described by one or more metadata records in different formats.

The documents of an OpenDLib library are organised in a set of virtual collections, each characterised by its own access policy. Authorised people can define new collections dynamically by specifying definition criteria. In the same digital library, for example, it is possible to maintain a collection of grey literature accessible to all users and a collection of historical images accessible only to a specific group of researchers. Each collection is automatically updated whenever a new document matching the definition criteria is published in the library.

### **Functionality**

The basic release of OpenDLib provides services to support the submission, description, indexing, search, browsing, retrieval, access, preservation and visualization of documents. Documents can be submitted as files in a chosen format or as URLs to documents stored elsewhere. They can be described using one or more metadata formats. The search service offers different search options: text free or fielded (with fields selected from a variety of known metadata formats); with or without relevance feedback. Documents retrieved can be navigated over all their editions, versions, structures, metadata and formats. All the above services can be customised according to several

dimensions such as, for example, metadata formats, controlled vocabularies, and browsable fields.

OpenDLib also provides other digital library specific services, such as the control of access policies on documents, and the management of “user-shelves” able to maintain document versions, result-sets, session results, and other information. In addition, a number of administration functions are also given to support preservation of documents, document reviewing process, introduction of new collections, and handling of users and user group profiles.

### ***OSIRIS/ISIS: Management and Search of Multimedia Collections***

*Heiko Schuldt, University of Basel (Switzerland)*

OSIRIS (Open Service Infrastructure for Reliable and Integrated process Support) is a new middleware supporting flow control of processes, and can be an essential component for digital libraries. It supports processes (compound services) as a means for combining existing services and executing them under certain (transactional) guarantees. In a digital library, such services can be used, for instance, for defining search processes and for the maintenance of replica and indexes for the various data repositories of a digital library.

ISIS (Interactive SIMilarity Search) on top of OSIRIS is a collection of services for efficient searches of multimedia collections. It includes meta-data search in combination with efficient and effective content-based similarity search on images, text, audio, and video. It has been developed in the context of ETHWorld, the virtual campus of ETH Zurich.

### ***MILOS: a multimedia content management system for multimedia digital library applications***

*Giuseppe Amato, ISTI-CNR (Italy)*

MILOS is a Multimedia Content Management system that offers functionality specialised to support multimedia digital library applications. It can be seen as the equivalent of a “database system” for document intensive applications (like digital library applications). Building a digital library application with MILOS, to manage and give access to existing or new corpora, implies a reduced effort given that several critical functionalities are already provided by MILOS (<http://milos.isti.cnr.it/>).

MILOS offers three basic functionalities.

#### **1) Management of arbitrary XML encoded metadata:**

MILOS can manage simultaneously heterogeneous metadata. Any arbitrary XML encoded metadata can be handled by MILOS. Advanced search functionalities are provided by means of declarative queries on metadata, in addition to text search queries (in textual elements) and similarity search queries (for instance in MPEG7 visual

descriptors elements). Similarity search is particularly useful when dealing with multimedia documents.

### **2) Transparent management of document storage strategies:**

Documents can be stored using arbitrary strategies depending on the specific requirements of the applications (different access protocols and different locations, e.g. video servers, web servers, file systems, etc.). With MILOS storage strategies are transparent to the application.

### **3) Metadata mapping:**

Metadata schemas used by the applications can be different from metadata actually stored in MILOS. MILOS translates requests made by applications into correct requests to the managed metadata. This, for instance, allows different applications, using different metadata, to access the same MILOS installation, or to use legacy digital library applications with newer data and metadata corpora.

### **Demonstration:**

Several different digital library applications have already been built, with minimal effort, on top of MILOS. These digital library applications have been populated by using material from other existing digital libraries without any adaptation of the corresponding data and metadata. The data sets used in these prototype applications consist of documents and metadata of very different nature:

- the *Reuters*, which is a text intensive data set consisting of newspaper news, (<http://about.reuters.com/researchandstandards/corpus/>);
- the *ACM Sigmod Record*, (<http://www.acm.org/sigs/sigmod/record/xml/>);
- the *DBLP* data sets, which contain heavily structured metadata, (<http://www.informatik.uni-trier.de/%7Eley/db/>);
- the *ECHO* data set, which contains video material and heavily structured metadata, (<http://pc-erato2.iei.pi.cnr.it/echo/>).

Another application (the *MILOS photo book*) provides on-line management and sharing of personal pictures (and their MPEG7 descriptions) and allows image search by using text search on descriptions and visual similarity search on the images.

## **DAFFODIL**

*Claus-Peter Klas, University of Duisburg-Essen (Germany)*

Daffodil is a search system for digital libraries aiming at strategic support during the information search process. From a user point of view this strategic support is mainly implemented by high-level search functions, so-called stratagems, which provide functionality beyond today's digital libraries. Through the tight integration of stratagems and with the federation of heterogeneous digital libraries, Daffodil reaches high effects of synergy for information and services. These effects provide high-quality metadata for the searcher through an intuitively controllable user interface. DAFFODIL supports the complete DL life cycle by providing a personal library, which allows for storing all kinds of DL objects, organized in folders, and allowing for both in-line and out-of-line annotations. Asynchronous collaboration is enabled through sharing of folders and

annotations, whereas a chat tool and a whiteboard tool allow for synchronous collaboration.

André Schaefer; Matthias Jordan; Claus-Peter Klas; Norbert Fuhr (2005).  
*Active Support For Query Formulation in Virtual Digital Libraries: A case study with DAFFODIL*. In: Proc. ECDL 2005.

Claus-Peter Klas; Norbert Fuhr; André Schaefer (2004).  
*Evaluating Strategic Support for Information Access in the DAFFODIL System*. In: Proc. ECDL 2004.

Sascha Kriewel; Claus-Peter Klas; André Schaefer; Norbert Fuhr (2004).  
*Daffodil - Strategic Support for User-Oriented Access to Heterogeneous Digital Libraries*. D-Lib Magazine 10(6).

## **A Prototype P2P DL Network**

*Vassilis Christophides, FORTH-ICS (Greece)*

The main problem in implementing a Digital Library as a federation of independent repositories resides in the different annotating information (e.g., metadata formats, classification schemes) that are maintained in the nodes of a DL network, complicated by the fact that those nodes may join or leave the network at their own will. The interaction with multiple DL nodes to support integrated access despite their heterogeneity is presently beyond the traditional information integration technologies, which impose restrictions on representation and communication languages used at both the semantic and the structural levels. The aim of this prototype is to investigate the peer-to-peer (P2P) resource-sharing paradigm for large-scale distributed Digital Libraries (DL), leading to the so-called P2P DLs. The objective is to support decentralized sharing of data and services in a network of autonomous and heterogeneous DL nodes. A P2P DL keeps a balance between the efficiency provided by a centralized architecture and the autonomy and decentralized sharing/management of data provided by P2P architecture. The prototype provides the following functionality:

**DL network formulation.** We assume that every DL node stores its data in an RDBMS. In order to join the P2P network, a DL node should select one of the RDFS schemas provided by the network, or create an RDFS schema by applying certain kind of schema operations (like union, intersection, difference, selection) on the available RDFS schemas. A tool assists the user to instantiate the RDFS classes by retrieving tuples from the database of the newly joined node. The RDFS schemas provided are views of a given global schema. A DL node can leave the network at its own will. The rest of the network is not affected at all, except that data from that particular DL node is not available any more.

**DL querying.** A user in a DL node may pose queries in the form of RDF triples. A query initiated by a node is sent to its neighboring nodes. Each one of those nodes sends the query to its neighbors, and so on. The query is evaluated at each node and results are sent back to the first node. Since DL nodes maintain different local schemas, the query is reformulated before its evaluation on a DL node to match its local schema. The reformulation is performed using information from the global schema.

### ***Eurovision: a text-based cross-language image retrieval system.***

*Paul Clough, University of Sheffield (UK)*

This system demonstrates how multilingual access has been provided to a digital library collection of historic photographs. Like many collections, the images are accompanied by captions enabling text-based access. The system makes use of online translation tools to translate the user's search request, the user interface and text associated with results into the user's search language. Although a basic interface, this demonstrates what can be done with a machine translation system, as well as where further improvements could be made.

Clough, P. and Sanderson, M., User Experiments with the Eurovision Cross-Language Image Retrieval System, In Journal of the American Society for Information Science and Technology (JASIST), In Print (expected publication 2006).

### ***Ciquest: organising images using concept hierarchies***

*Paul Clough, University of Sheffield (UK)*

Organising a set of documents automatically based upon a set of categories (or concepts) derived from the documents themselves is an obviously appealing goal for Information Retrieval systems: it requires little or no manual intervention (e.g. deciding on thematic categories) and, like unsupervised classification, depends on natural divisions in the data rather than pre-assigned categories (i.e. requires no training data). Concept hierarchy generation is one such method: it automatically associates terms extracted from a document set and organises them into a hierarchy, each term representing a group of documents. A prototype application has been developed, which uses concept hierarchies to organise the results of searching a set of historic photographs. There is also a multilingual version of this, using free online translation tools showing how, like Eurovision, a naive cross-language system can be easily constructed with minimal effort.

Clough, P., Joho, H. and Sanderson, M. (2005), Automatically Organising Images using Concept Hierarchies, Workshop held at the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Workshop: Multimedia Information Retrieval, August 15-19, 2005, in Salvador, Brazil.

Petrelli, D. and Clough, P. (2005) Concept Hierarchy across Languages in Text-Based Image Retrieval: A User Evaluation, In the working notes of the CLEF workshop, Vienna, Austria, 21-23 September 2005.

### ***eaSim and FLOSS***

*Ulrike Steffens, OFFIS (Germany)*

Peer-to-peer architectures are one candidate approach for the realization of distributed digital library architectures. They are of particular interest because they allow for a direct mapping of organizational units to peers within a network. Queries against a peer-to-peer DL are routed along the different peers. eaSim is a Java-based discrete-event system simulator for organisation-oriented peer-to-peer-architectures. The simulator allows comprehensive simulations of query evaluation on the physical layer (physical network), the virtual layer (peer-to-peer-network) and the organisational layer (e.g. social or digital library structures). eaSim is also able to generate lookup queries and peer-to-peer topologies for the layers. Different metrics are used to analyse the performance of peer-to-peer lookup and routing algorithms. The graphical user interface visualizes these metrics at run time to demonstrate the behaviour of the algorithms. The tool /yEd Graph Editor/ can be used to visualize the different layer topologies.

For the organisational layer, a FLOSS (Free/Libre and Open Source Software) structure generator has been developed, which allows the generation of FLOSS structures (graphs) with thousands of developers, projects and their relationships. These structures can be visualized via the /yEd Graph Editor/.

## **RoleMiner**

*Ulrike Steffens, OFFIS (Germany)*

User roles are an important control instrument since they enable the structured assignment of rights existing in organizations and communities. As an example in the context of digital libraries, roles can be used to administrate access rights to protected spaces, collections and resources. The top-down-creation of a role concept within an organization or even in a community from scratch, however, can turn into an obstacle. To overcome this problem, we introduce a tool that makes use of access rights already existing in digital libraries. These rights are analyzed and visualized in order to substantiate the creation of role concepts by operational reality. The tool is based on a data mining algorithm tailored for role mining. In addition to the administration of access rights we plan to extend the tool by detecting typical user roles like "library expert", "novice user" or "interested cultural history". User roles detected by the tool can then be used for better personalization in different digital library services.

## **Sightseeing4U**

*Ansgar Scherp, OFFIS (Germany)*

Sightseeing4U is a generic personalized city guide application employing our MM4U component framework for creating personalized multimedia presentations. The generic tourist guide is applicable to develop personalized tourist guides for arbitrary cities, both for Desktop PCs and mobile devices. The concrete demonstrator we developed for our home town Oldenburg in Northern Germany considers the pedestrian zone and comprises video and image material of about 50 sights. The demonstrator is developed for Desktop PCs as well as PDAs. It supports personalization in respect of the user's interests, e.g., churches, museums, and theatres, and preferences such as the favorite language.

Depending on the specific sightseeing interests the proper sights are automatically selected for the user. This is realized by category matching of the user's interests with the meta data associated to the sights.

## **xSMART**

*Ansgar Scherp, OFFIS (Germany)*

In recent years, many highly sophisticated multimedia authoring tools have been developed. Most of these systems show a limited integration of the context of the targeted user community. The xSMART system (Context-aware Smart Multimedia Authoring Tool) provides a semi-automatic authoring tool that integrates the targeted user context into the different authoring steps and exploits this context to guide the author through the content authoring process. The design of xSMART allows its extension and customization to the requirements of a specific domain through the use of domain-specific wizards. These wizards realize the user interface that best meets the domain-specific requirements and effectively supports the domain experts in creating their content targeted at a specific user context.

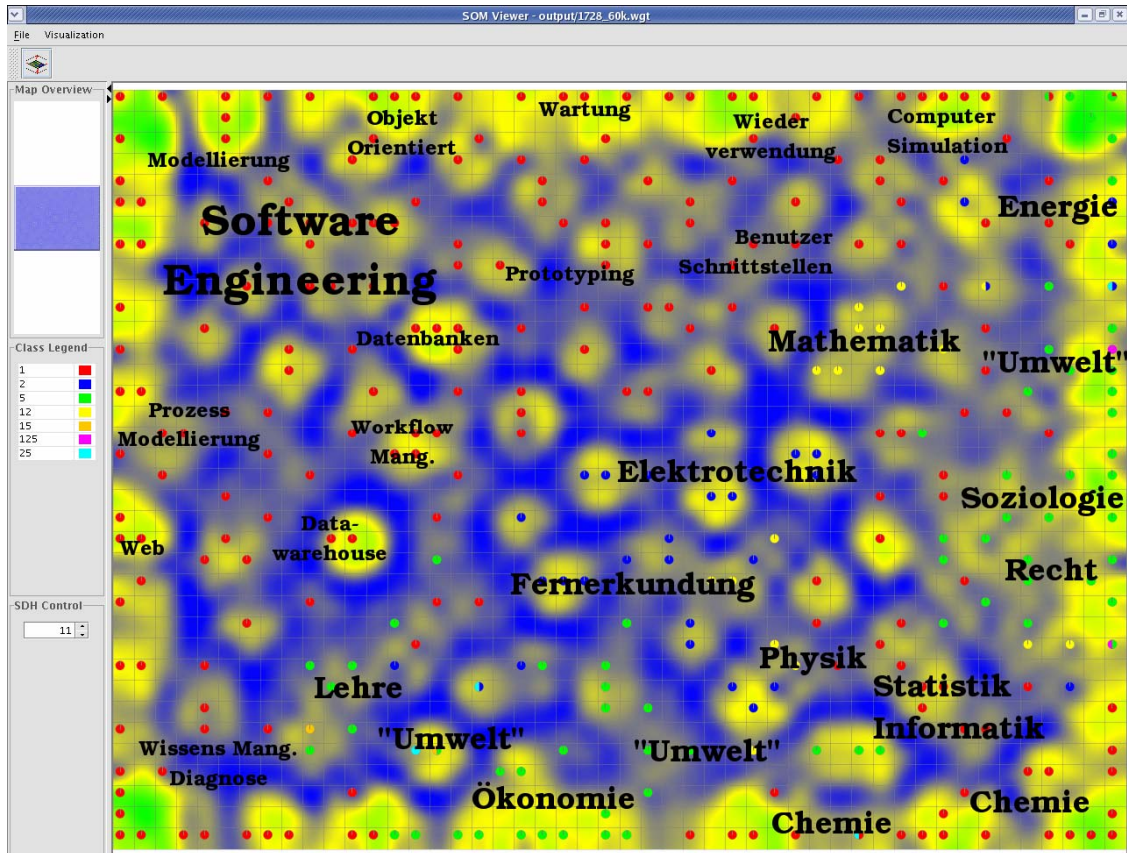
## **SOMLib-based Interfaces to Digital Libraries**

*Andreas Rauber, Technical University of Vienna (Austria)*

The SOMLib Digital Library Project aims at the development of a digital library system supporting intuitive, user friendly browsing of document collections by combining the benefits of conventional library organization with the possibilities offered by digital collections. It is based on the self-organizing map (SOM), a popular unsupervised neural network, used to organize documents by content. Different visualization methods, including map-based metaphors of islands, or bookshelf metaphors, address different user needs.

SOMLib clusters documents by semantic similarity. Different methods for labeling the resulting maps, ranging from keyword selection methods to simple NLP tools, allow users to orient themselves on the map (see picture for an example). The system serves as an additional access method to large document collections, complementing traditional metadata or full-text search. It can be incorporated into basically any traditional DL system as a novel, exiting interface, allowing users to orient themselves in the collection, getting an overview of their holdings. Current prototype studies are in progress for e.g. the Austrian Research Documentation, or an integration into the Greenstone DL system, and may also be used for audio digital libraries, organizing music by sound similarity based n frequency spectra analysis, both on desktop as well as portable/mobile devices.





SOMLib view of the Austrian Research Documentation (subset)

## **DARE**

*Giuseppe Santucci, Università di Roma La Sapienza (Italy)*

The DARE (Drawing Adequate REpresentation) system permits to visually analyze large amounts of data, either exploring single data points' values or interacting with OLAP cubes to discover aggregate values. Data browsing is implemented using up to 6 visual attributes, i.e., the x, y and z axes, color, size, and shape. The OLAP visualization handles 1D, 2D, and 3D visual cubes showing, through the size and the color of each cube element, two summary values (textual reports are available as well). Additional operations like drill down, roll up, and slice and dice are provided. Moreover, the user can switch between elementary and aggregate data at any time, visualizing different data set attributes. The DARE system has been used as a test bed for the IEEE InfoVis Contest 2005 and obtained interesting results [Ber05].

[Ber05] Enrico Bertini, Luigi Dell'Aquila, Giuseppe Santucci Discovering USA technology trends with DARE and SpringView IEEE Symposium on Information Visualization, 2005. InfoVis Contest 2005.

## **Annotating web accessible DLs with the MADCOW toolbar**

*Emanuele Panizzi, Università di Roma La Sapienza (Italy)*

MADCOW is a browser toolbar enabling users to associate a note to any interesting portion of a web page. In fact, while navigating the web with a standard browser, the user can select a portion of text, an image or any other multimedia content in the page and click on the 'create-annotation' button in the MADCOW toolbar. A dialog window pops up, where the user can add his comment, attach any multimedia file of his choice, and chose one out of the 9 annotation types (comment, question, announcement, explanation, integration, etc.)

The web note is automatically saved in the MADCOW server, and can be retrieved in any subsequent browsing session. When the user accesses again the annotated web page, an icon next to the annotated part is shown, and clicking on the icon the annotation is opened in a new browser window. The web note can be accessed also by other users, provided the author declared it as public in the 'create-annotation' window. This fosters cooperative work, as it becomes possible to start discussions about any topic, binding the web notes to the object of discussion. An annotation can be annotated in turn, allowing threads in the discussion.

## **ImageLab VideoBrowser**

*Costantino Grana, Università degli Studi di Modena e Reggio Emilia (Italy)*

An MPEG-1 video browser has been developed, with shot detection, key frame extraction, transition length characterization and MPEG-7 compliant output. The prototype allows different shot detection techniques to be used in order to compare their results. A second step allows sub-shot automatic clip creation. These two processes provide a hierarchical description of the video, which may be employed for searching and indexing in Digital Libraries. Another module deals with Pictorially Enriched Ontologies, a tool for the creation of a set of visual prototypes, which enrich a textual taxonomy, in order to provide descriptive examples for a specific class. The outcome may be linked to the video annotation and allow for automatic mapping of different scenes to classes defined in the ontology. Videos annotated with this technique allow for further search, based on their contents. The system is able to create the ontology (automatic semantic prototype selection based on annotated video classes) and to use it for further annotation of other videos. Searching facilities are not covered.

## **Content Based Image Retrieval**

*Alberto Del Bimbo, University of Florence (Italy)*

An image retrieval system has been developed, which supports content based queries according to features of color and shape.

The system provides to the user a graphical interface that enables queries by using the "query by example" paradigm. In particular, for color-based searches, the user can draw regions, colorize and locate them in order to find images with a predefined arrangement

of color patches. For shape-based retrieval, the user can sketch the profile of a shape he/she is looking for. Similarity by global color is also allowed by letting the user select an example image and asking the system to find similar images.

## **Content Based Retrieval of 3D Models**

*Alberto Del Bimbo, University of Florence (Italy)*

Beside image and video databases, archives of 3D models have recently gained increasing attention for a number of reasons: advancements in 3D hardware and software technologies – in particular for acquisition, authoring and display – their ever increasing availability at affordable costs, and the establishment of open standards for 3D data interchange (e.g. VRML, X3D). Thanks to the availability of technologies for their acquisition, 3D models are being employed in a wide range of application domains, including medicine, computer aided design and engineering, and cultural heritage. In this framework the development of techniques to enable retrieval by content of 3D models assumes an ever-increasing relevance.

A major difficulty for the development of a system for retrieval by content of 3D objects relates to the need to capture the twofold nature by which 3D objects are experienced by humans: view based and structure based. On the one hand, 3D objects can be perceived through multiple 2D views, but on the other hand 3D objects can also be examined by analyzing their 3D structure. This twofold nature is not separable and indeed, our perception of 3D object similarity is purely view based in some cases, purely structural in other cases and a combination of both in the general case.

This prototype presents a new solution combining the advantages of view-based and structure-based approaches to description and matching of 3D objects. The new solution relies on Spin Image signatures and clustering to achieve an effective, yet efficient representation of 3D object content.

## **Video Semantic Adaptation**

*Alberto Del Bimbo, University of Florence (Italy)*

This prototype shows an approach for multimedia access for video presentations on mobile systems. The main idea is to adapt the video content according to the preferences of users, in terms of quality and cost. Video is annotated off-line, extracting highlights and interesting objects, and users may specify the combination of events and objects, providing an index expressing a level of interest in their combination. Then video is automatically adapted according to these preferences, compressing more what is less interesting for the user, and trying to keep a good quality for the most interesting parts.

## **Pictorially Enriched Ontology Annotator**

*Alberto Del Bimbo, University of Florence (Italy)*

This prototype eases the semi-automatic annotation of videos using pictorially enriched ontologies. A typical way to perform video annotation requires to classify video elements

(e.g. events and objects) according to some pre-defined ontology of the video content domain.

Ontologies are defined by establishing relationships between linguistic terms that specify domain concepts at different abstraction levels. However, although linguistic terms are appropriate to distinguish event and object categories, they are inadequate when they must describe specific or complex patterns of events or video entities. In these cases, pattern specifications can be better expressed by using visual prototypes, either images or video clips that capture the essence of the event or the entity. Enhanced ontologies, that include both visual and linguistic concepts, can be useful to support video annotation up to the level of detail of pattern specification.

### ***Accademia della Crusca (Crusca Academy)***

*Alberto Del Bimbo, University of Florence (Italy)*

From the XVII century to nowadays, the Italian "Accademia della Crusca" (<http://www.accademiadellacrusca.it/>) has built five vocabularies, which contain and represent the origin and the evolution of the Italian language during the last four centuries. These vocabularies are artworks in itself, and in the last few years the Accademia has fully transcribed the first four vocabularies and annotated them using the standard XML TEI format (<http://www.tei-c.org/>). This standard format is exhaustive with respect of everything that can appear in a vocabulary, including entries, homographs, abbreviations, foreign words and so on. This prototype, built in cooperation with the Accademia, performs a complete parsing of the transcriptions, creates the inverse indexes and executes simple and advanced search on the indexes. The web application implements several advanced search features such as: search of very short terms (less than 3 letters), multi-biased search, search of word roots and punctuation.

### ***COMISM (COMmunity Interface Semantics Model)***

*Stavros. Christodoulakis, TUC/MUSIC (Greece)*

This prototype implements a generic model and a service oriented framework for Virtual Community Support in Digital Libraries. Virtual Communities consist of people that use the information technology to meet with other people with whom they share some common interests and/or objectives. People join a virtual community in order to satisfy some objectives and/or interests that they have, which are related with the interests and objectives of the virtual community. COMISM captures the most essential functionality needed in various applications of virtual communities, including digital libraries. The model provides generic support for deeper aspects of communities like the personalization of the community information, the support for groups of individuals to satisfy their interests and objectives and to obtain stronger relationships with selected members of the community.

The model has been implemented using a service-oriented architecture to ensure easy integration with existing digital library frameworks. The implementation provides area-separation between communities. For example, a global community may correspond to the whole world, and sub-communities may correspond to continents, countries, cities etc. A member of a virtual community has the ability to create a topic of discussion in an

area that s/he is registered and discuss with the other members of the specific area. S/he can also publish objects of his interest and make them available to others. Information that still evolves, like unfinished conversations, can be stored as well, such as messages, announcements, reviews and member's information objects. A virtual community membership also includes a personal file, where personal info and data can be stored, and which can be private, or public, if the user wants to share it. The personal file wraps up the user's feeling that he belongs to a community.

On top of the service-oriented implementation framework, COMISM provides a web-based user interface to demonstrate the core functionalities of the framework. This web interface could be used as a stand-alone application although the full power of the framework comes from its ability of enhancing other systems, especially digital libraries, with generic virtual community support.

### ***GraphOnto: Ontology Editing and Ontology Mapping, Multimedia MPEG7 and Ontology-Based Metadata Definitions***

*Stavros. Christodoulakis, TUC/MUSIC (Greece)*

GraphOnto is a component that facilitates the generation and population of both standard-based and domain-based ontologies, and their use in multimedia information system components. In the multimedia domain, an Upper Ontology, which captures the MPEG-7 MDS, is utilized and OWL (imported and/or interactively created) domain ontologies extend the upper ontology with domain knowledge. Imported ontologies are parsed so that graphical ontology browsing and editing interfaces are automatically generated. The ontologies are used to guide the metadata definition in a standardized manner. GraphOnto is used in a number of subprojects of the DELOS NoE.

The simultaneous management of multiple ontologies is supported in GraphOnto as well as ontology mappings and constraint checking. In addition, core ontologies are distinguished and receive special treatment. Personalization of the GraphOnto interface on the application, task and user levels is provided.

OWL/RDF metadata are defined using GraphOnto, based on OWL ontologies. The metadata produced may be exported in MPEG-7 compliant syntax, provided that, the ontologies on which metadata definition is based, are integrated with the core ontology capturing the MPEG-7 MDS.

The GraphOnto component may be used either as a standalone ontology and metadata editor or as a semantic indexing tool for multimedia content, working on top of an MPEG-7 compliant Knowledge Base. In particular, the produced MPEG-7 compliant metadata may be stored either in simple files or directly in the MPEG-7 compliant Knowledge Base. In addition, metadata items stored in the Knowledge Base may be reused during metadata definition (e.g. existing metadata items representing soccer players may be reused in a soccer game description).

### ***OntoNL: A Natural Language Interface Generator to Knowledge Repositories***

*Stavros. Christodoulakis, TUC/MUSIC (Greece)*

The objective of the Natural Language Interface Generator (OntoNL) Framework is to provide principles, methodologies and software for the automation of the construction of natural language interfaces to knowledge repositories. These interfaces include capabilities for declaration and manipulation of new knowledge, as well as querying, filtering and ontology driven interaction formulation.

The OntoNL framework is able to address uniformly a range of problems in sentence analysis each of which traditionally would have required a separate computational mechanism. In particular a single architecture can handle both syntactic and semantic ambiguities, can handle ambiguity at both a general and a domain specific environment, and can consult user profiles to personalize the disambiguation

The OntoNL framework makes use of OWL rich vocabulary by using upper and domain ontologies. The semantic search is especially useful in applications where the user searches for concept instances of the model and not for “arbitrary” data. That is, usually the keywords in the query denote one or more concepts. Given an OWL ontology, weights are assigned to links based on certain properties of the ontology, so that they measure the strength of the relation. In this way we can identify related concepts in the ontology to the ones retrieved by the user’s request. The OntoNL framework uses also User Profiles to guide the semantic search in the domain ontology and to rank the results in a way a user meets his preferences.

A Natural Language Interface Generator can perform a complete syntactic analysis of any kind of sentence and a semantic analysis based on a word ontology. The syntactic analysis is based on a methodology that contains a part-of-speech tagging procedure, a grammatical relation annotator, a noun compound analysis component and a synonyms and sense discovery procedure.

The task of POS-tagging is to assign part of speech tags to words reflecting their syntactic category. In our system we adopted a maximum entropy approach, because it allows the inclusion of diverse sources of information without causing fragmentation and without necessarily assuming independence between the predictors. The grammatical relation annotator is responsible of locating the syntactic structure of any sentence and assigns the grammatical relations between head and modifiers, like the subject, the object, the complements, etc. This way we can further help a question answering system that uses the proposed generator for producing a natural language interface. We have implemented a grammatical relation annotation scheme in which each sentence is marked up with a set of grammatical relations, specifying the syntactic dependency which holds between each head and its dependent(s).

The component that is responsible for the noun compound bracketing procedure cooperates with the grammatical relation annotator to provide the correct syntactic structures. Our approach for handling noun compounds is that we use a method to expand n-grams into all morphological forms by the use of morphological tools. The training corpus we use is the domain ontologies used for every different application. This may lead to the conclusion that the test set is very limited in comparison to a linguistic corpus, but it is more specific to the needs of the application. We are interested in the particular needs of the user based on a specific domain. By combining the use of domain ontologies and the WordNet as the training corpus and by taking into account the hyponyms and synonyms of the nouns that constitute the n-gram we maintain all the information needed for the correct bracketing. Since we use the noun compound bracketing methodology to be more accurate when dealing with the user request’s ambiguities we use as a test set the

noun compounds that may appear in the exact user request and as a training corpus the total of domain ontologies used.

We obtain the synonyms and the corresponding senses from a complete word ontology. The output of this methodology is a language model that can be further enhanced with semantic information coming from the domain of the platform to be used.

## ***The LUPA Index, A Demonstrator of The Referential Integrity Problem in Large RDF Networks***

*Martin Doerr, ICS-FORTH (Greece)*

A problem in the large-scale integration of complementary information is that, very often, relevant queries can only be answered by retrieving the end-points of a data path through multiple sources. For instance, a question about the distribution of Greek names on Roman tombstones in Britain requires a source relating the names with the inscriptions, another relating the inscriptions with the stone, another the stone with the place of finding, and another the place with a spot on the map. This is quite different, and by far more challenging, than aggregating all documents about the same topic.

A prototype application has been developed, in the field of iconographic and epigraphic aspects of Roman stone monuments, by manually defining a mapping of each source schema to the CIDOC CRM model (ISO/FDIS21127), transforming automatically the data and storing it in an RDF knowledge base. The main objective is to provide procedures for information integration and global querying over all the contents of the complementary resources. The demonstrator currently contains integrated data from about a hundred thousand records from four distinct data sources.

It can be shown how suitable queries on the integrated data can retrieve candidates for potential duplicates. Depending on the confidence level, duplicates can be removed either automatically or manually. Conflict resolution can be combined with subsequent semiautomatic detection and resolution of further duplicates. This problem is characteristic of all attempts to integrate metadata into large knowledge networks with integrated knowledge management. It requires a transformation into a global schema (here the CIDOC CRM), and a major investment in identifying common items with different identifiers in different sources. It requires a systematic approach and continuous, well-designed procedures improving data quality in parallel with the equally continuous update.

## ***The ITeM Recommender System***

*Giovanni Semeraro, Università degli Studi di Bari (Italy)*

Personalization is an important method for Digital Libraries to take a more active role in dynamically tailoring its information and service offer to individuals in order to meet better their needs. Algorithms designed to support users in retrieving relevant information base their computations concerning relevance on so-called *user profiles*, in which representations of the user interests are maintained. This prototype exploits *supervised*

*machine learning* techniques to induce user profiles from text documents in order to make the access to digital libraries a personalized experience. *Item Recommender* learns *semantic profiles* that capture central concepts representing the interests of the users from documents they deemed relevant. In semantic profiles, keywords are replaced with their meanings, corresponding to synsets (SYNONIM SETS) as defined in the WordNet lexical database.

*Item Recommender* is able to recommend items by learning from both their textual descriptions and the ratings given by users. The system implements the naïve Bayes classifier, an increasingly popular algorithm in text classification applications and is able to classify items as interesting or uninteresting for a particular user by exploiting a probabilistic model learned from training examples. The final outcome of the learning process is a probabilistic model used to classify a new instance as interesting or uninteresting. The model is used as a personal profile including those concepts (synsets) that turn out to be most indicative of the user's preferences, according to the value of the parameters of the model.

### ***DOMINUS (DOcument Management INtelligent Universal System)***

*Floriana Esposito, Università degli Studi di Bari (Italy)*

**DOMINUS** is a system for automated electronic documents processing characterized by the intensive exploitation of intelligent techniques in all the steps involved from document acquisition to document indexing for categorization and information retrieval purposes. It can currently deal with documents in standard formats, such as PostScript (PS) or Portable Document Format (PDF). In order to perform text extraction and indexing from the original PDF/PS document DOMINUS executes a number of steps, which can be summarized as follows: layout analysis of the incoming document; layout correction, if needed; classification; of the document type; understanding of the layout components; text extraction and indexing.

### ***Linking Paper and Digital***

*Moira Norrie, ETH Zurich (Switzerland)*

Despite the emergence of digital technologies, paper persists as a fundamental resource for many human activities. Nowadays, documents tend to be created, stored and distributed electronically, but paper continues to be a preferred medium for many reading and writing activities. Paper is cheap, light, mobile, easily annotated in various ways and supports forms of collaboration difficult to mimic in digital worlds.

We are involved in a number of projects investigating how new technologies can turn paper into an interactive medium. Users can select links on paper using a special pen and activate a range of digital services such as playing a video clip, displaying a web page or retrieving requested information which is then given to the user through generated speech. Interaction can be mixed with capture, enabling users' annotations and drawings to be converted to digital form and also linked back to their location within a paper document.

Several prototype applications are already available.



1) *PaperPoint*: a system to control PowerPoint presentations from a printed overview of the slides

2) *EdFest*: a mobile tourist information system for visitors to the Edinburgh Fringe Festival based on set of interactive paper documents including an event brochure, a map and a bookmark. A text-to-speech engine is used to deliver information via a voice channel. Users can also add reviews of events through hand-written comments in the back of the brochure, in a separate notebook or on post-its. We use OCR software to convert these to text and then the text-to-speech engine to read back these comments to the user for confirmation and also to enable other users to access them from their brochure.

3) *Natural History Museum*:

In conjunction with King's College London, we carried out a user study at the Natural History Museum in London based on an interactive paper brochure that children used to write information about exhibits during their visit round the museum and then later used in the special investigate room of the museum to access associated digital materials and also a play a game.

We are currently developing other applications – including ones based on cross-media educational materials for schoolchildren.

## The cooperation between DELOS and The European Library

DELOS has started a series of actions specifically intended to transfer to some user communities the results achieved and the prototypes developed by the DELOS members. The focus in JPA3 is on the library and cultural heritage communities, through cooperation with the TEL Office (The European Library) and with MICHAEL (representing the community of the Ministries of Cultural Assets). It is also planned to continue the participation in joint activities and events with ELAG (European Library Automation Group) and to start participation in LIBER (Ligue de Bibliothèques Européennes de Recherche) and LIDA (Libraries in the Digital Age) events. Other opportunities for cooperation with museums and archives will be sought and considered during JPA3. Presently, a number of activities aiming at transferring digital library functionality developed by DELOS members to the TEL system are in plan and will be carried out jointly with the TEL Technical Working Group. The following four major areas of common interest have been identified.

*The reference model.* It will be possible to capitalize on the expertise of the TEL participants to bring into the reference model activities valuable input for its refinement; at the same time it is expected that the conceptual framework resulting from the reference model will provide valuable input to TEL for extensions and refinement of the architecture of The European Library system.

*Multilingual capabilities.* TEL users should be able to access and search the library in their own (or preferred) language, retrieve documents in other languages and have the results presented in an interpretable fashion (e.g. possibly with a summary of the contents in their chosen language). One of the first activities will be a feasibility study aimed at producing guidelines and strategy for the attainment of this long-term ambitious goal.

*Personalization.* A first step will be the development of personalization guidelines to identify those services more suitable for personalization, followed by a second step for integration and testing of existing prototype software and development of possible improvements over those services.

*User interface and Visualization.* The first step will be an evaluation of the existing TEL user interface, as well as the exploration of additional services, especially for supporting query formulation, collection navigation and results visualization. The activities in this area will focus on four topics: evaluation, support for query formulation, virtual collections and navigation, presentation/visualization of results.

For each area of interest a specific new Task has been defined in Workpackage 8, and the new Tasks are: T8.2 (Validation and refinement of the reference model through interaction with TEL), T8.3 (Multi-Lingual Information Access in TEL), T8.4 (Personalization capabilities in TEL), T8.5 (User interface design for TEL, navigation and visualization services).

Given the renewed interest of the European Commission for the cultural sector, including the launching of a European Digital Library, it is essential for DELOS to participate in these efforts by providing not only scientific contributions, but providing also tools and prototypes to be used in real-world applications. DELOS has started a series of actions

specifically intended to transfer to some user communities the results achieved and the prototypes developed by the DELOS members. Presently the focus is on the library and cultural heritage communities, through cooperation with the office of The European Library (representing the National libraries of all the European countries), and with MICHAEL (representing the community of the Ministries of Cultural Assets). It is also planned to continue the participation in joint activities and events with ELAG, the European Library Automation Group. The cooperation with TEL has been structured into four “technology transfer” Tasks, and will focus primarily on the integration of DELOS-provided software and prototypes into the existing European Library system, which presently provides access to the collections of 15 National Libraries in Europe, with 8 more libraries expected to be accessible by the end of 2006.

#### **Task A – Validation and refinement of the digital library reference model through interaction with TEL**

An important task in the process of defining the DELOS digital library reference model is its validation against the requirements of concrete digital libraries and existing digital library management systems. The expertise of the TEL participants (especially the members of the TEL Technical Working Group TWG) will bring into the reference model activities valuable input for its refinement; and similarly the conceptual framework resulting from the reference model will provide valuable input to TEL for the extension/refinement of the architecture of The European Library.

TEL has a business process-oriented view on requirements for digital libraries. The DELOS view is bottom-up, trying to identify basic digital library services. A major outcome of the validation of the reference model will be an attempt to map the business process-oriented view to the reference model to basic digital library services.

The cooperation will capitalize also on the results of another on-going effort of DELOS aimed at integrating different digital library services (provided by the DELOS partners) into an existing middleware platform (OSIRIS Work-flow platform). The resulting “integrated prototype” should provide a concrete example of an implementation of the reference architecture, and be a test-bed for digital library functionality. Currently, TEL runs a metadata repository over the different national libraries giving links to documents in the partner organizations. The services and the infrastructure provided by the integrated prototype might be used by TEL for implementing value-added services of The European Library.

In order to obtain results based on real-world data, DELOS will harvest the available TEL metadata and make use of them in the integrated prototype, possibly also adding links to the TEL metadata in the services implemented within the integrated prototype. On the basis of those results, TEL will be able to evaluate the possibility of testing, integrating or adding to the TEL architecture some of the services developed within the integrated prototype.

#### **Task B – Multi-Lingual Information Access in TEL**

The ultimate goal of this task is to provide the capability to users of TEL to access and search the library in their own (or preferred) language, retrieve documents in other languages and have the results presented in an interpretable fashion (e.g. possibly with a summary of the contents in their chosen language). The problem is complex and many

factors are involved. These include: the number of languages involved, the current heterogeneous setup of TEL, the lexical tools and resources needed.

- *Number of Languages.* The number of different languages represented in TEL constitutes a major hurdle for Multi-Lingual Information Access (MLIA), as ideally it should be possible to launch a query in any one of the national languages of the TEL collections and retrieve relevant material in any one of the collections. Possible approaches to the problem might be the use of multilingual ontologies, metadata and subject authority data, similar document search.
- *Heterogeneous set-up.* A major problem is represented by the heterogeneous set up of TEL, as it is not clear whether the existing infrastructure is able to accept a cross-language query result.
- *Resources Needed.* Any cross-language strategy implies the acquisition and development of appropriate lexical tools and linguistic resources such as stemmers, morphologies, bilingual dictionaries, etc.

The implementation of MLIA in TEL is an ambitious task that can be considered a medium/long-term goal, to be achieved through a series of intermediate steps. The first step is the establishment of a joint DELOS/TEL working group to perform a feasibility study aimed at producing the following output:

- Guidelines as to how the TEL infrastructure should be adapted to be ready for the requirements of multilingual access and output, given the current state of play to be able to access the data at the national libraries. Recommendations as to how the libraries should be delivering their data may form part of these guidelines.
- Guidelines for the preparation of multilingual textual resources, which will be included into TEL (either digitized from existing text or borne digital)
- Definition of strategies that should be adopted by TEL with respect to enabling TEL users to search in their preferred language and retrieve documents in other languages
- Identification of the most promising possible implementation directions: linking metadata, similar document search, etc.

### **Task C – Personalization capabilities in TEL**

The final objective of this task is to produce guidelines and prototype software for new added-value services of interest to the final users, initially selecting those services that present a lower-risk of failure when personalized. A first short-term objective is addressing the development of personalization guidelines, and a second medium-term objective will address integration and testing of existing prototype software and development of possible improvements.

The initial personalization topics to explore are listed below. The exploration will start with a study of the existing access logs from TEL, in order to perform an analysis and categorization of context, to derive specifications for new types of logged data and suggestions on the design of innovative personalized services:

- query expansion (e.g. given a query related to “stars”, distinguish between a hobby astronomer versus a cosmology researcher)
- profile building
- notification about new material based on profiles

- recommendations based on profile similarity
- annotation sharing based on profiles
- provision of added-value links and/or service (e.g. OpenURL, etc.), based on preferences and rights of user or organization.

The first expected result is a comprehensive report containing suggestions for added-value personalized services based on their potential and user-perceived relevance. The report will also provide suggestions about data to be logged for a better evaluation of the site use, which would be the base for evaluating the personalisation possibilities of the site. It is then expected to develop an initial prototype of a toolkit for log analysis and personalization services, which will lead to a refined toolkit and a final report containing agreed guidelines for personalized services.

**Task D – User interface design for TEL, navigation and visualization services.** This task addresses the overall design of the TEL user interface, as well as the exploration of additional services, especially for supporting query formulation, collection navigation and results visualization. The activities in the short-medium term will focus on four topics:

*Evaluation.* As part of an on-going task in DELOS a comparative evaluation between the current TEL interface and an appropriate variant of the DAFFODIL framework is being performed. This evaluation will follow both an analytical and an empirical approach. The goal of the analytical evaluation is to assess the functional similarities and differences between the two systems. The goal of the empirical evaluation is to evaluate how well each tool supports the users needs.

The analytical evaluation will consider the usability, functions for search, browsing and result display, and the feedback/help functions of both. For this purpose, the methodology and the questionnaires developed by DELOS will be used. The empirical evaluation will be based on a user-centered and qualitative approach. Its focus is on the users experience with the tools, considering user characteristics, preferences and strategies, the types of activities/tasks users perform, and the environment in which the search tool is used. The evaluation must take into consideration the current necessary portal nature of the site.

*Support for query formulation.* The DAFFODIL interface already provides functions that help the user in formulating better queries. Most basic, a built-in spell checker will flag search terms not contained in the dictionary, and will propose correct variants. For advanced query formulations, a syntax checker will point out syntactically incorrect formulations. Finally, there is a 'related term' tool that proposes (statistically) similar terms for any of the query terms entered. The comparative evaluation will show to what extent these tools are useful for the TEL users, and then possible integration into the TEL system will be evaluated.

*Virtual collections and navigation.* It is planned to provide and test an add-on service for building virtual digital collections starting from a set of real ones. For the definition of the virtual collections, the service should distinguish between expert users (like technicians and librarians) and end users. The service relies on automatic batch techniques of indexing, clustering, and classification of existing collections, allowing a visual navigation of their content, for an easier definition of the virtual collections. The benefits expected are that the user can be presented with cross collection views, can deal with smaller set of more relevant data and therefore queries can be processed in a faster way.

*Presentation/visualization of results.* The DAFFODIL system already provides functions for relevance ranking or quick filtering of results, as well as extracting attributes like author names or frequent terms from the result set. In addition, it is planned to provide and test an add-on service that allows end users to interact with the query result in a more effective way. Several techniques can be used in this context: real time indexing, cluster-gather algorithms, smart use of relevance factor, information visualization techniques. One of the objectives of this activity is to understand which of those techniques are more useful in the TEL environment, in addition to motivate the end users to explore more large datasets and have “more fun” while exploring the digital library. Also expert users like librarians could use this service to better visualize (and hence understand) the results produced by the navigation service described above, in order to define more easily customized views for end users.