

## 15. Project Information

Institution : University of Cape Town  
Short title : High Performance Digital Libraries on-Demand  
Area : Information and Communication Technology  
Sub-area : Human-Information Interaction  
Project start year : 2007  
Project end year : 2010  
Budget start year : 2008  
Budget end year : 2010  
Descriptive title : High Performance Digital Libraries using Transparent  
Grid-based Resources on-Demand  
Chosen years : 3

## 16. Problem Identification

Systems to manage information (digital libraries/repositories, learning management systems, etc.) provide key technologies for the storage, preservation and dissemination of knowledge in its various guises e.g., research documents, theses and dissertations, cultural heritage documents. Each of these is currently the subject of communities working furiously to increase the amount of data and information available to ever-larger numbers of students and scholars. But, as the sheer quantity of information and its use increases, the standard tools for managing such information do not cope well and there is a need for new techniques to continue to provide useful and innovative services to users. Building on prior work and the current state-of-the-art in high performance computing, this study thus proposes to look into how information management systems can be made scalable by transparently consuming grid-like resources, efficiently and on-demand, especially to suit the resource-constrained environments of developing countries.

## 17. Rationale and Motivation

Digital library systems (DLses) are rapidly growing in popularity as the technology matures and also because of the advocacy of groups such as the Open Access and Electronic Thesis and Dissertation communities. The effect of this popularity is that there are now more accessible collections, growing at relatively high rates - Lyman and Varian (2003) estimated 5 exabytes of new digital information in 2002 alone! Newer buzzwords such as Web 2.0 (O Reilly, 2005) mask an emerging trend where users produce far more information than in the past, exemplified by the rapidly growing popularity of user-oriented information exchanges spaces such as MySpace, Facebook and YouTube.

There is a need for tools to manage these large and growing collections and meta-collections and make them accessible to the relevant audiences. However, these tools are not readily available and popular DL systems do not always scale appropriately (Haedstrom, 2003; Imafouo, 2006). While much research has gone into the scalability of Web-delivered DL content (see, for example, (Andresen, et al., 1996)), access to services is only one dimension of the management tasks, which typically also include internal data processing for classification, preservation-related manipulation and ingest procedures. Commercial service providers such as Google provide some scalable services but these are proprietary and normally restricted to some types of services. These service providers, however, recognise the crucial need for scalability and digital library tools also need to adopt this as an underlying philosophy.

Digital library tools need to be accessible to users and managers of collections of varying sizes. On the one hand, it is difficult to predict how large a collection may become; and on the other hand, it is not cost-effective to reserve resources (such as storage space) in advance. In the spirit of high performance computing initiatives, a scalable system should therefore support a range of collection sizes and various levels of service provision on-demand.

In recent years some developments in the DL community at large have investigated this need

for greater scalability. One of the most prominent efforts is the Storage Resource Broker (Moore, et al., 2000) from San Diego Supercomputing Centre, which mediates access to large-scale heterogeneous storage. This technology can be used as an underlying layer for some repository systems, notably DSpace. An alternative, the Fedora repository system (Staples, et al., 2003), claims to support 10 million digital objects. In the experimental arena, the DILIGENT (2007) project, with EU funding, is integrating digital library technology into the EGEE European Grid project. While most of these efforts help with scalability, there is still much work to be done and only some aspects are addressed in each project. Most importantly, scalability of systems in developing countries is a somewhat different problem because of the exaggerated scarcity of storage, computation and networking resources.

Prior work done by the applicant focused recently on flexibility of digital library systems based on component architectures. As component architectures naturally lend themselves to distributed systems, these components were extended (in a 2005-2006 study) to support migration and replication of basic services. Then, on a protocol level, harvesting of metadata (according to the Open Archives Initiative Protocol for Metadata Harvesting) was tested in parallel systems with different architectures and significantly greater efficiency of resource use was demonstrated. Current related efforts include:

- How to re-architect the underlying Web server to support generic mobile components.
- How to build information retrieval systems on clusters and grids that make most efficient use of limited resources.
- How to make management of Grid-based systems easier for end-users.

These earlier and current studies have served to illustrate the possibilities that are available and demonstrate the limitations of various approaches in South Africa. Cluster computing, for example, is not as popular as it is abroad - here there are fewer and smaller clusters dedicated to particular projects, with only the Centre for High Performance Computing building a system on the scale of peer international efforts.

These projects and future endeavours are thus very strongly influenced by local conditions in South Africa. Typical information management systems do not have a cheap or endless supply of storage, networking or computational resources. As such, there is a desperate need to make efficient use of what resources do exist. Also, it is necessary that systems are able to start small and grow over time. The Open Access community in Southern Africa is growing with support from funding agencies such as OSI. Simultaneously, the Mellon Foundation, through its Aluka initiative, is assisting with the digitization of cultural heritage documents. The applicant is working closely with some of these projects to develop their collections and services - recognising that current solutions do not scale, there is a need to look into scalability issues at the same time.

Research into scalable on-demand information management satisfies multiple themes of the research agenda defined by NRF in its ICT programme. In terms of the first theme of Software Development and Integration, scalable information systems require a radical rethink in terms of core software architecture for Internet-based applications.

In terms of Telecommunications and Networking, this work will contribute to most of the identified research themes. Protocol design is a key part of developing distributed systems that underpin Grid technologies and form the basis of many modern distributed DL systems. Local conditions of poor network connectivity impact on such experimental protocol and system design to encourage minimalist approaches and robustness of algorithms. At the same time, network management must be instituted at a high level to manage these scarce resources. Lastly, this work will contribute to the areas of distributed systems and services since the system under examination contains network-distributed services. Investigations into aspects such as scalability and robustness of systems generalise to other distributed systems, while services within the context of digital libraries are driven by user needs in a networked environment so include the full gamut of possibilities.

Scalable digital libraries as networked information systems fall within the research theme of Human-Information Interaction, especially in the use, processing and retrieval of information. The focus of scalable systems is mostly on the processing of information as an enabler for subsequent retrieval and use, while each of the latter activities also requires elements of scalability.

Finally, from a Developmental angle, digital libraries fundamentally change the landscape of access to information and therefore affect the quality of life for all. This work, by designing scalable transparent digital libraries, will make them more accessible to archivists working in conditions with limited resources. This, in return, makes information more readily available to the ordinary student, teacher or researcher (within the constraints of basic Internet access), levelling the playing fields by removing access to information as a barrier to learning.

#### References:

Andresen, Daniel, Tao Yang, Omar Egecioglu, Oscar H. Ibarra, and Terence R. Smith (1996), Scalability Issues for High Performance Digital Libraries on the World Wide Web, Technical Report 1996-03, Department of Computer Science, University of California Santa Barbara, March 1996.

Diligent (2006) A Digital Library Infrastructure on Grid Enabled Technology. Website <http://www.diligentproject.org/>

Haedstrom, Margaret (2003), Research Challenges in Digital Archiving and Long term Preservation, NSF Post Digital Library Futures Workshop, 15-17 June 2003, Cape Cod.

Available [http://www.sis.pitt.edu/dlwkshop/paper\\_hedstrom.html](http://www.sis.pitt.edu/dlwkshop/paper_hedstrom.html)

Imafouo, Amlie (2006), A Scalability Survey in IR and DL, TCDL Bulletin, Volume 2, Issue 2. Available <http://www.ieeetcdl.org/Bulletin/v2n2/imafouo/imafouo.html>

Lyman, Peter, and Hal R. Varian (2003), How Much Information 2003?, University of California. Available <http://www2.sims.berkeley.edu/research/projects/howmuch-info-2003/index.htm>

Moore, R., C. Baru, A. Rajasekar, B. Ludascher, R. Marciano, M. Wan, W. Schroder and A. Gupta (2000), Collection-Based Persistent Digital Archives Parts 1 and 2, D-Lib Magazine, April/March 2000. Available <http://www.dlib.org/dlib/march00/moore/03moore-pt1.html> and <http://www.dlib.org/dlib/april00/moore/04moore-pt2.html>

O'Reilly, Tim (2005), What is Web 2.0: Design Patterns and Business Models for the Next Generation of Software, 30 September 2005. Available <http://www.oreillynet.com/pub/a/oreilly/tim/news/2005/09/30/what-is-web-20.html>

Staples, Thornton, Ross Wayland and Sandra Payette (2003), The Fedora Project: An Open-source Digital Object Repository System, D-Lib Magazine, April 2003. Available <http://www.dlib.org/dlib/april03/staples/04staples.html>

## 18. Research Aims

The aim of this research is to develop techniques for building scalable digital information management systems based on efficient and on-demand use of generic grid-based technologies.

Specifically, the following questions may be investigated (subject to sufficiently many students with appropriate skills joining the team):

- a) Can we migrate a typical DL architecture to a Grid system such that remote resources are transparently brought into service when needed?
- b) How do we recast standard services such as information retrieval to operate on small to terabyte-scale information stores using transparent Grid resources? (currently in progress)
- c) How do we allocate, distribute and schedule resources for maximal efficiency of data transfer algorithms such as OAI-PMH?
- d) Can we layer a typical DL architecture over a volunteer computing paradigm such as BOINC?
- e) How do we support the use of generic Grid systems by higher-level users and services? How do we design management tools that will ease the adoption of Grid technology for DL and possibly other applications? (currently in progress)

In summary, this research aims to look at various aspects that will affect the adoption of grid technology for digital archives in resource-constrained environments. The expected outcome of this work is a set of proven guidelines and experimental tools to move the digital library community closer to an ideal of simple, flexible, scalable and robust digital library architectures, building on an underlying Grid fabric.

## 19. Workplan: Research Activities

Student support is a primary aim of this project. As such, the project will be conducted as a series of related sub-projects, with collaboration on the interfaces and external dependencies.

These can be enumerated as follows:

Development of a high-level interface to manage Grid resources. This supporting activity makes Grid technology more accessible to users and developers without detailed knowledge of Grid technology. This work is currently being addressed by one MSc student, under supervision of the principle investigator. Time-frame: 2007-2008

Refactorisation of services to operate in a resource-constrained environment. An MSc student is currently working on the information retrieval aspect, under supervision of the principle investigator. It is envisaged that a second student may work on additional aspects of the project related to other services such as browsing and visualization. Time-frame: 2007 to 2009

Design and implementation of a framework for transparent Grid-based DL systems, to be conducted by an MSc student under the supervision of the principle investigator. Time-frame: 2008-2009

Data transfer algorithms for optimal use of networking resources with dynamic and large data stores. Experiments have been done by the principle investigator in 2006, continuing in 2007 and possibly 2008. These will be extended by an MSc student in 2009-2010.

The use of volunteer computing is a focused project that is suitable for an MSc student, under supervision of the principle investigator. Time-frame: 2009-2010.

It is intended that student assistants will help with the building of prototypes, as part of their research training in (pre-MSc) Honours degrees. This has worked effectively in the past and has helped to encourage Honours students to enroll for further postgraduates degrees.

Some equipment for this work is already available and is being acquired during the course of 2007. Additional machines will be needed for students and to bolster server and local computational resources in 2008 and beyond.

## 20. Workplan: Research Approaches/Methods/Techniques

The methodology for this work has two identifiable aspects: those elements applicable to all of the sub-questions that will be addressed; and elements that relate to each individual major issue. Thus, these will be treated as separate sections:

### 1. Overall research design

#### 1.1. Data collection

As an enabling mechanism for the other aspects of this project, data collections need to be established in the form of local archives. These will be obtained by working with researchers to archive their publications in subject repositories recommended by their individual communities e.g., the applicant worked with the law faculty at UCT in 2005 to set up <http://lawspace.law.uct.ac.za/>, which archives law-related documents. In addition, institutional repositories for documents such as post-prints and electronic theses and dissertations are becoming popular and it is proposed that these will be buttressed by additional components to make them interoperable with our and other efforts. In all of these cases, the work involves consultation with staff and the use of standard or emerging tools in the digital library community - some software development in related projects helps to further enhance tools for local use. The applicant is a key member of the joint NRF-CHELSEA working group on a SA national ETD project as well as a key member of the Sivulile Open Access working group. Both efforts are centred on building digital collections and providing access to them. The collections built over time will serve a two-fold purpose: the data can be used by experimental systems and the researchers - who would have had exposure to digital libraries - can provide a base group for evaluation of systems. Essentially, this nurtures a tradition of digital library use and understanding, which is lacking in South Africa and must be developed to support research and development efforts.

As a second front to data gathering, the OAI Protocol for Metadata Harvesting can be used to obtain data from remote sources. The OAI maintains a registry of archives currently sharing their collections of metadata and/or data freely with external entities. This includes notable organisations such as the US Library of Congress and collections such as OCLC WorldCat - a collection of more than 4 million pointers to dissertations. Such data, or subsets thereof, can readily be obtained as and when necessary, over the Internet, in order to supplement and complement local collections and to expose local audiences to international resources and vice versa.

#### 1.2. Specialized equipment and infrastructure

There is no unconventional equipment required. Instead, medium- to low-end PC servers using open source software are becoming the norm in both the digital library and high performance computing communities and these are all that is needed. A 14-node computing cluster has already been set up and used in previous projects and this will be utilised for some experiments in this study.

As much of the work in this project is designed to be conducted by postgraduate students working independently, commodity workstations are required for each of those researchers, who begin working on the project in 2008 and 2009. Some servers may need to be replaced and/or upgraded and there probably will be maintenance costs to replace some of the cluster

nodes over time.

Also needed for all servers and workstations is a local network infrastructure and a realistically fast external Internet connection. External research collaboration and data collection from outside sources may require a stable and broadband connection.

### 1.3. Systems Integration and Dissemination

While the research question has been divided into sub-questions aimed to define and encourage student participation in the larger project, integration among all of these parts is essential to prove the viability of the overall approach. This integration will be led by the principal investigator, to oversee the various parts and maintain a consistent and workable model and toolset as time progresses.

Dissemination is usually considered to be a separate activity post-experimentation, but in the context of networked digital libraries, it is a primary need since collaborators frequently provide support for experiments with remote systems. The principal investigator will encourage dissemination and actively seek collaborators to help with experimental validation of the approach taken by this project.

## 2. Specific issues

### 2.1. High-level interface to manage Grid resources

Grid computing is considered by many to be a complex technology that is not simple to interface with or develop for. Can an interface be developed to support easier adoption and use of grids? How does this relate to Web 2.0 technologies and newer modes of browser-centric service-oriented systems? Can such an interface be generalised to work with any or all Grid middleware systems? Can enhanced usability and abstraction promote use of Grids? To address these questions, this project includes investigation into current interfaces and interface technology and implementation and evaluation of a browser-based layer over and around Grid technology.

This work is being carried out in conjunction with one Masters student in 2007-2008, initially funded by a 2006-2007 NRF grant. Substantial responsibility for the design and evaluation will be delegated to the Masters student, under supervision of the principal investigator.

### 2.2. Refactorisation of services

Digital library services usually provide the ability to search and browse through information. Beyond the obvious services, there also are recommender systems, ratings services, subscriptions, management interfaces and a host of others. Each of these services includes user and/or component interaction and some form of data processing. These activities are sometimes very computationally intensive, more so when data collections grow in size. Thus it is vital that scalable services process data and requests using generic storage and computation resources, such as those found in a Grid. This project thus centres around investigation of how to refactor DL services such that they can be implemented as a layer over a Grid. Questions to be resolved include how to handle updates in data collections, how to most effectively use storage and network resources and the granularity for sub-operations. An information retrieval system prototype is being built and evaluated for its ability to efficiently use available resources.

This work is being carried out in conjunction with one Masters student in 2007-2008, initially funded by a 2006-2007 NRF grant, continuing a 2006 Honours project. A second student will be recruited in 2008 to work on other services. Substantial responsibility for the design and evaluation will be delegated to the Masters students, under supervision of the principal investigator, with Honours students focusing mainly on the implementation of tools.

### 2.3. Transparent Grid-based DL

A transparent Grid-based DL system is one where Grid-based remote resources are brought into service based on changing needs, either as large datasets need to be processed or to deal with flash crowds. Until such events occur, the DL system can be restricted to a smaller computing environment (nominally a single host node). Thus, the system should scale on demand and not lock in resources unnecessarily. Questions to be answered in this study include what support is needed at a lower level to implement such technology, appropriate scheduling and resource reallocation algorithms and algorithms for mapping resources to requests. A prototype framework and reference system will be built to demonstrate how a typical DL system can be layered over a Grid, using remote facilities only when required.

This work will be carried out in conjunction with one Masters student in 2008-2009. Substantial responsibility for the design and evaluation will be delegated to the Masters students, under supervision of the principal investigator, with a possible Honours project focusing mainly on the implementation of tools.

### 2.4. Data transfer optimisation

The proposed high-level Grid architecture will optimise the use of computational resources but data transfer optimisation requires a thorough analysis of storage and network utilisation. The OAI-PMH used to transfer metadata among systems has been shown to be reasonably efficient when the client operates in a parallel manner. This study will extend existing work to investigate how the server also can be made to operate on parallel requests. If the dataset can be processed in parallel by multiple clients connecting to corresponding servers, some level of global resource allocation should result in minimal data transfer overhead. Questions that need to be resolved include how weak/strong data consistency is to be maintained, how to handle updates, global vs. local state management and algorithms to balance ideal resource allocation

with overheads. An enhanced protocol and algorithm will be developed and experiments will be conducted to measure the improvement or degradation under different circumstances on a reference implementation.

This work will be carried out in conjunction with one Masters student in 2009-2010. Substantial responsibility for the design and evaluation will be delegated to the Masters students, under supervision of the principal investigator, with a possible Honours project focusing mainly on the implementation of tools.

#### 2.5. Volunteer computing

Volunteer or cause computing (such as SETI@home) is currently popular because of the ease of adoption and the large installed base. It should be possible to harness such systems for digital library systems. This project will thus investigate the use of a system such as BOINC to form the underlying fabric of a typical DLS. A prototype will be built as a proof-of-concept and evaluated for its performance and stability when conducting typical DL operations.

This work will be carried out in conjunction with one Masters student in 2009-2010. Substantial responsibility for the design and evaluation will be delegated to the Masters students, under supervision of the principal investigator.

### 21. Potential Impact on HR Development

This proposal is specifically designed to encourage the training of students in the principles and techniques at the intersection of modern networked systems and high performance computing. This is a critical skill needed in the IT community to develop and enhance future infrastructure for information sharing. While postgraduate students will benefit directly, their experiments will involve scores of researchers and fellow students who will gain a better understanding and appreciation for high performance digital libraries by exposure to the technology. There are currently 5 MSc students working on related projects in digital library systems at UCT, 2 of whom are funded through NRF, while 3 have external funding. One of these students has just finished his studies and submitted his dissertation for examination. Two students will join this project in 2008 as their work lies at the intersection of this and a previous NRF-funded project. As this project continues to investigate emerging problems, 1 additional MSc student will be recruited and trained in 2008/2009 and 2 more in 2009/2010. In addition, Honours and 3rd year students will assist where possible, as a way of encouraging them to pursue/continue research-led degrees in future. This is an absolutely vital component for promoting research among young students.

This project also aims at promoting, directly and indirectly, collaboration among researchers through the medium of digital libraries. Such collaborations enhance the quality of individual research works, while contributing to a sense of community among researchers.

### 22. Potential Impact on Redress and Equity

Digital libraries, by their very nature, address imbalances in the information arena by making information more accessible to regular users. This is especially relevant in the area of journal publications. These are traditionally expensive to obtain and therefore only easily available to universities that subscribe at high costs to themselves, the state and their students/academics. One of the main aims of digital library research is to make it possible for academics to archive, review and publish work without the need for expensive intermediaries such as publishing houses. Even if such works are already published, simple digital libraries can be used to make the publications available locally, as is allowed by many publishers already (e.g., ACM, Springer). As such archives become commonplace, large data sets and complex services require more computing power. This research ensures that the technology is in place to enable flexible information management systems for those who need them the most.

Eventually, digital libraries are the mechanism to get high quality information about Africa out to the rest of the world, and in the reverse direction, for us to obtain high quality electronic resources. Scalability ensures that storage of and access to these resource collections does not become a stumbling block as digital archiving is widely adopted in Africa.

In terms of redress and equity among students, a concerted effort is always made to train previously disadvantaged students, as evidenced by the list of past and current students.

### 23. Potential Outcomes

This project has many potential outcomes for the digital libraries and high performance computing research communities.

The research community in digital libraries has not concentrated on scalability as a key concern and this needs to be addressed. Some parallel efforts will complement this work, together hopefully leading to a change in the way production systems are built. Similarly, the high performance computing community has not placed much emphasis on digital archiving. Thus any results will foster cross-pollination between the two communities.

This project has implications for the library and archiving community in terms of cost savings and the ability to gradually adopt new technology. Using information management systems that generalise to high performance platforms, it will be possible for archivists to start every project in proof-of-concept mode and scale upon (expected) success or with the arrival of funding or larger data sets. This eliminates the need for large proposals and injection of funds with

unknown returns.

In terms of artefacts, this project will contribute a suite of tools that may be used to design digital libraries or that may be adapted to similar projects in distributed systems. Software packages produced as part of the work also may be used in research environments to make local resources accessible to local and/or international audiences. Finally, the collections developed during this project will serve as testbeds for future experimentation locally and abroad, where they will encourage future international efforts to acknowledge and cater for the needs of countries such as South Africa.

Ultimately, by advancing the science of building information management systems, this project hopes to make information more accessible to users. In South African society, where the cost of information is abnormally high, any move towards acceptance of digital libraries makes it easier for people to access information that would normally be beyond their reach (philosophically and physically). A typical example of this is the Networked Digital Library of Theses and Dissertations (NDLTD), which is making it possible for students and researchers to obtain electronic copies of theses and dissertations online. This was previously only possible through inter-library loan. Now, anybody anywhere in the world can find and get access to a thesis if it is part of a member digital library at a participating institution. NDLTD maintains a union collection of metadata from all over the world and this collection has been steadily growing in size. This is a prime example of a system where scalability needs to be part of the underlying fabric. This research can help to provide answers when scalability questions arise in projects such as NDLTD.

## 24. Progress to Date: Summary

The principle investigator has actively developed and supports a framework for building digital library systems as collections of inter-connected components. These are used in various projects, including NDLTD, and have formed the basis for a number of experiments at UCT and beyond. Experiments to validate the architecture were conducted with various communities and case studies. The results of initial work have been published in conferences, journals and magazines (as listed in the appropriate sections).

During 2005-6, an MSc student (Muammar Omar) worked on migration and replication of components in a small component farm. His dissertation was submitted and is currently being examined.

During 2007, two students began working on the high-level Grid management interface and tera-scale information retrieval subprojects (as listed in previous sections). These students have both developed detailed proposals and are currently working on development of algorithms and prototypes.

Collaboration has been established with the Centre for High Performance Computing (CHPC). The applicant has been part of discussions and workshops related to this activity (CHPC, SCAW and Grid). The applicant is part of the High Performance Computing research group at UCT, as a collaborative venture among staff with varied interests intersecting with HPC techniques. This provides a platform for discussion of technologies and techniques across subdisciplines.

From a data collection perspective, a local archive of research publications has been established and is in a prime position to serve as a future testbed for experimental services. Also, working with the ETD Africa and Sivulile initiatives, the principle investigator is assisting universities in South Africa and Africa in general to set up electronic thesis and dissertation and open access projects - these will serve not only the needs of the local communities but also jettison the African community into the digital library age and provide case studies and testbeds for current and future research. To this end, the applicant has been an invited speaker at advocacy and training workshops, conferences and symposia in Johannesburg (2003, 2005), Addis Ababa (2004), Pretoria (2004, 2005, 2007), Windhoek (2004), Cape Town (2005), Singapore (2006), Stellenbosch (2006) and Maseru (2007).

The NRF-funded Scalable Information Management project runs from 2006-2007. Four MSc students are investigating various aspects of building scalable systems and a number of initial results and changes in the high performance computing research landscape have informed the move to start investigating Grid technologies as soon as possible. Thus, some students will continue their work in the framework of the new project.

## 25. Progress to Date: Research Outputs

The projects denoted Open Digital Libraries, which defined basic architecture and core interfaces, and Flexible Digital Libraries, which investigated higher level integration of components, together form the core underpinning for building generic information management systems. The former was funded by the US-NSF and Mellon foundation and was the subject of the PhD of the applicant. The latter was funded by the NRF in 2004-2005 and both have resulted in various publications and technical reports as listed in the relevant section of the proposal. Some of these publications have a direct bearing on the concepts of scalability and remote management and will therefore be built upon in future work.

Scalable Information Management is a relatively new project undertaken in 2006-7. Initial publications in 2006 are listed in the relevant section. In addition, 2 papers related to scalability have been submitted to an international conference and 3 are in preparation for local and international conferences. One student has submitted his thesis and another is expected to complete his thesis within the course of 2007.