# Digital libraries and society:
# New perspectives on information dissemination

**Ian H. Witten**

Department of Computer Science, University of Waikato, New Zealand

ihw@cs.waikato.ac.nz


Corresponding author:


Ian H. Witten

Dept of Computer Science,

University of Waikato,

Hamilton, New Zealand.

Phone +64 7 838-4246

Fax: +64 7 858-5095

ihw@cs.waikato.ac.nz

# Digital libraries and society:
# New perspectives on information dissemination

## ABSTRACT

Digital libraries are large, organized collections of information objects. Well-designed digital library software has the potential to enable non-specialist people to conceive, assemble, build, and disseminate new information collections. This has great social import because, by democratizing information dissemination, it provides a counterbalance to disturbing commercialization initiatives in the information and entertainment industries. This article reviews trends in today's information environment, introduces digital library technology and explores the use of digital libraries for disseminating humanitarian information in developing countries, a context that is both innovative and socially motivated. We demonstrate how currently available technology empowers users to build and publish information collections. Conventional public libraries are founded on the principle of open access, and extending this to digital libraries presents a challenge to HCI—a challenge that is magnified if open access is extended to those who create library collections too.

**Keywords**: Digital libraries, human-computer interaction, social issues in computing, developing countries, internationalization

## INTRODUCTION TO THE CHAPTERS IN THIS SECTION

Digital libraries are large, organized collections of information objects. Whereas standard library automation systems provide a computerized version of the catalog—a gateway into the treasure-house of information stored in the library—digital libraries incorporate the treasure itself, namely the information objects that constitute the library's collection. Whereas standard libraries are, of necessity, ponderous and substantial institutions, with large buildings and significant funding requirements, even large digital libraries can be lightweight. Whereas standard libraries, whose mandate includes preservation as well as access, are "conservative" by definition, with institutional infrastructure to match, digital libraries are nimble: they emphasize access and evolve rapidly.

The five chapters in this section provide an excellent illustration of the huge variety of interesting issues in digital library research that impacts the Asia Pacific region.

Unlike the New World, where most of the research on technological aspects of digital libraries originates, Asia has an exceptionally rich cultural heritage. This manifests itself in a huge legacy of documents, in various forms—from paper to palm leaves—and in various different conditions. The need to preserve this legacy is particularly pressing in today's world, where political instability is rife and climate change is beginning to have an effect. As recent events have shown, disasters, both manmade and natural, can have a devastating effect on fragile cultural artifacts. S.M. Shafi from the University of Kashmir, India, surveys the many issues involved with digital archiving of medieval manuscripts in Asia.

Libraries are pillars of education, and it is natural to expect that digital libraries will provide new opportunities for innovative educational practices. These will be particularly relevant to the Asia Pacific region because of the huge disparities in access to education between the different communities there. Peer-to-peer learning has always been a crucial factor in personal development, although it is frequently ignored in educational studies. Natalie Lee-San from Nanyang Technological University, Singapore, describes her studies of how digital libraries can provide an innovative, perhaps revolutionary, environment for peer-to-peer learning amongst youths. She touches on many practical issues: gender differences, different learning styles, different levels of media and computer literacy, and age-related differences.

Many economies in the Asia Pacific region are agriculturally-based. Modern agriculture is a knowledge-based activity that can benefit greatly from digital libraries. Mila Ramos from the International Rice Research Institute in the Philippines describes a large-scale digital library system designed to support the growth, nurturing, harvesting, and distribution of that most Asian of staples, rice. This digital library supports an institute whose goal is to improve the well-being of present and future generations of rice farmers and consumers, particularly those with low incomes. The institute's library houses the world's most comprehensive collection of technical literature on rice, and provides a widely-used international reference service. As in many specialized libraries, digital library technology is seen to have special advantages in a world of shrinking library budgets.

Intellectual property issues are a central driving force behind the market in information of which libraries are a part. And the questions become more complex as the nature of today's information shifts from a primarily book-based culture to one that embraces all types of multimedia objects, and large, carefully-curated, collections of such objects. The chapter on "Multimedia digital library as intellectual property" by Hideyasu Sasaki and Yasujshi Kiyoki at

Keio University in Japan clarify the copyright situation as it affects multimedia collections and compilations. They go on to discuss the patentability of particular retrieval mechanisms, an essential component of digital libraries.

The fifth chapter in this section on how digital library research impacts the Asia Pacific region is the present one, on digital libraries in society. Most existing digital library projects, being research-oriented, are predicated on state-of-the-art equipment and interfaces, academic and research institutions, special collections. In contrast, this article argues for universal access: digital library technology can and should be available to everyone, on all platforms, in all countries; and it can and should enable ordinary people to exercise their creative powers to conceive, assemble, build, and disseminate new information collections that are designed not just for western academics but for a wide diversity of different audiences throughout the world. Though less glamorous, this may, in the end, be a more important goal for society. Digital libraries pose an inherent tension between the technologist's desire for advanced solutions that use the latest and greatest hardware and software, and the librarian's desire for wide, cross-platform availability and long-term preservation—as epitomized by the sustained success of paper as a delivery medium. To achieve universal access for both information consumers and collection-builders is really a problem for HCI.

In the next subsection we examine the social need for digital libraries, particularly in developing countries, by briefly sketching some trends in commercial publishing and contrasting them with a growing international perspective of information as a public good. We draw out the implications for the user interface, which is the principal bottleneck in allowing non-specialist people to make public information available in focused collections that are universally usable. Then we introduce digital library technology and illustrate it with a particular example, the Greenstone digital library software, which is designed for a broad user base and is in widespread use in many corners of the world—from Uganda to the US, Kazakhstan to Canada, Nepal to New Zealand. Following that, we review a project that is applying digital library technology to the distribution of humanitarian information in the developing world, a context that is both innovative and socially motivated. Next we discuss issues of universal access and illustrate them with reference to the Greenstone software. We include a brief demonstration of a prototype system that is intended to allow anyone to build and disseminate information collections, and illustrates some human interface challenges that arise when providing necessarily complex functionality to a non-computer-oriented user base. We close with the hope that future digital libraries will find a new role to play in helping to

reduce the social inequity that haunts today's world, both within our own countries and between nations.

## BOOKS, LIBRARIES, AND THE SOCIALLY DISADVANTAGED

Today, the long-standing three-way tension between the commercial interests of publishers, the needs of society and information users, and the social mandate of public libraries, is being pulled and stretched as never before.

First, the very notion of a "book" is evolving in many different directions: books become more interactive; publishers rent content; books are distributed under restrictive conditions that mechanically prohibit sharing. While it would be premature to make specific predictions, it seems likely that these trends will further disadvantage the disadvantaged—particularly those in poorer countries who have yet to benefit from ready access to ordinary books. Second, a huge body of information is becoming freely available on the Internet. Much is of questionable quality, but some is very good indeed. In many cases the information is provided for the "public good" rather than for commercial profit, and the redistribution of such information is likely to be encouraged, rather than prohibited, by those who make it available. Initiatives like UNESCO's "Information for all" program and the upcoming World Summit on the Information Society highlight the importance of public information; they are founded on the belief that information literacy will help alleviate many of the problems confronting human societies. Third, the implications for libraries are mixed. Whereas new controls by publishers over how the content they own may be used presents libraries with significant problems, the ready availability of "public good" information meshes well with library philosophy. A new role is emerging for information professionals who can select material, index it, add appropriate metadata, and redistribute it in added-value form for the good of society. Suitable technological infrastructure is being provided by the open source movement, which is making available high-quality software for repackaging and distribution of information (and not just on computer networks).

We expand on each of these points below, and then summarize the prospects of digital libraries and what we see as the implications for the field of HCI.

### Books

What future has the book in the digital world? The question is a complex one that is being

widely aired (see Lynch, 2001, for a particularly thoughtful and comprehensive discussion). Authors and publishers ask how many copies of a work will be sold if networked digital libraries enable worldwide access to an electronic copy of it. Their nightmare is that the answer is *one*: how many books will be published online if the entire market can be extinguished by the sale of one electronic copy to a public library (Samuelson and Davis, 2000)? To counter this threat, the entertainment industry is promoting new "digital rights management" (DRM) schemes that permit a degree of control over what users can do that goes far beyond the traditional legal bounds of copyright. Indeed, the acronym is more aptly expanded as "digital restrictions management" because it is concerned solely with content owners rights and not at all with user's rights. It is, in effect, a "private governance system in which computer systems regulate which acts users are and are not authorized to perform" (Samuelson, 2003). Anti-circumvention rules are sanctioned by the Digital Millennium Copyright Act (DMCA) in the US (similar legislation is being enacted in other countries). The DMCA has been used, for example, to prosecute a Norwegian teenager for writing software to play a DVD that he had purchased on a computer for which no commercial playback systems exist.

Can DRM be applied to books? The motion picture industry can compel manufacturers to incorporate encryption into their products because it holds key patents on DVD players. Commercial book publishers are promoting e-book readers that, if adopted on a wide scale, would allow the same kind of control to be exerted over reading material. Basic rights that we take for granted (and are legally enshrined in the concept of copyright)—such as the ability to lend a book to a friend, resell it on the second-hand market, keep it indefinitely, continue to use it when your e-book reader breaks down, donate it to charity, preserve it for your grandchildren, copy excerpts without resorting to a handwritten transcription—are in jeopardy. DRM allows such rights to be controlled, monitored, and withdrawn instantly, and DMCA legislation makes it illegal for users to seek redress by taking matters into their own hands. Fortunately, perhaps, lack of standardization and compatibility issues are delaying consumer adoption of e-books.

In the realm of scholarly publishing, digital rights management is more advanced. Academic libraries license access to content in electronic form, often in tandem with purchase of print versions too. They have been able to negotiate reasonable conditions with publishers—probably because they represent the lion's share of the scholarly market. However, the extent of libraries' power in the consumer book market is moot. One can envisage a scenario where publishers establish a system of commercial, pay-per-view, libraries for e-

books and refuse public libraries access to books in a form that can be circulated (Roehl and Varian, 2001, describe an interesting parallel between historical circulating libraries and video rental stores).

These new directions present our society with puzzling challenges, and it would be rash to predict what society's response will be. But one thing is certain: they will surely increase the degree of disenfranchisement of those who do not have access to the technology.

**Public information**

In parallel with publishers' moves to reposition books as technological artifacts with refined and flexible control over how they can be used, an opposing trend has emerged: the ready availability of free information on the Internet. Of course, the world-wide web is an unreliable source of enlightenment, and undiscriminating use is dangerous—and widespread. As early as 1996 complaints arose that the Web's contents are largely unattributed, undated, unannotated, unreliable; information about author and publisher is unavailable or incomplete; far too many resource catalogues ("hubs") are chasing far too few original or non-trivial documents ("authorities") (Ciolek, 1996)—complaints that are very familiar today. But one thing has changed: search engines and other portals have enormously increased our ability to locate information that is at least ostensibly relevant to any given question. Teachers complain bitterly that students view the Web as a replacement for the library, harvesting information indiscriminately to provide answers to assignments that are at best shallow and at worst incoherent and incorrect. One consolation is that the very same search facilities can be used to detect plagiarism.

Nevertheless, the Web abounds with accessible, high-quality information. Many social groups, non-profit societies and charities make it their business to create sites and collect and organize information there. To take a single example at random, a Google search for *diabetes* returns three national diabetes associations (US, Canada, UK) in the top ten hits, and of course many more exist. Each of these sites offers a cornucopia of valuable information on the disease, which is not commercial and provided for the public good. Widespread use is strongly encouraged, and it seems likely that arrangements could be made for re-distribution of the material presented there, particularly it was intended as a not-for-profit service and appropriate acknowledgement was made.

One of the key problems with information distribution via the Web is that it disenfranchises

developing countries. Although the Web does not extend into the homes of the socially disadvantaged in developed countries either, various national programs are working to provide access (such as the Bill and Melinda Gates Foundation grants to public libraries). But network access varies enormously across the world. Whereas in 1998 more than a quarter of the US population were surfing the Internet, the figures for Latin America and the Caribbean was 0.8%, for Sub-Saharan Africa 0.1%, and for South Asia 0.04% (UNDP, 1999). Schools and hospitals in developing countries are poorly connected. Even in South Africa, the best-connected African country, many hospitals and 75% of schools have no telephone line. Universities are better equipped, but even there up to 1,000 people can depend on just one terminal. The Internet "is failing the developing world" (Arunachalam, 1998).

Prompted by this inequity, the importance of information, and particularly public information, is today being highlighted by prominent international bodies. For example, UNESCO's "Information for all" programme was established in 2001 to foster debate on the political, ethical and societal challenges of the emerging global knowledge society and to carry out projects promoting equitable access to information. It reflects a growing awareness that information is playing an increasing role in generating wealth and human capital, and that participation in the "global knowledge society" is essential for social and individual development. Information literacy is described as "a new frontier" by the Director of UNESCO's Information Society Division (Quéau, 2001). The International Telecommunications Union has established a World Summit on the Information Society, held in Geneva in December 2003 and Tunis in 2005, to promote a global discussion of the fundamental changes that are being brought about in our lives by the transformation from an industrial to an information society, and to confront the extreme disparities of access to information between the industrialized countries and the developing world.

**Libraries and their role**

What is the librarian to make of all this? The mandate of today's public libraries, in sharp contrast to that of publishers, is to facilitate the open distribution of knowledge. Librarians strive to enable the free flow of information. Their traditions are liberal, founded on the belief that libraries should serve democracy. To help fulfill their mission as resource centers for citizens, public libraries maintain collections of records, policy statements, government documents, and so on. A recent promotional video from the American Librarian's Association exults that "the library is democracy's place of worship" (ALA, 2002).

Clearly, the impending redefinition of the book as a digital artifact that is licensed rather than sold, tied to a particular replay device, with restrictions that are clearly laid out and mechanically enforced, is an innovation that goes right to the heart of libraries. The changing nature of the book may make it hard, or even impossible, for libraries to fulfill their mandate by providing quality information to readers. And on the other hand, the emergence of a vast storehouse of information on the Internet poses a different kind of conundrum. Librarians, the traditional gatekeepers of knowledge, are in danger of being bypassed, their skills ignored, their advice unsought. Search engines send users straight to the information they require—or so users think—without any need for an intermediary to classify, catalogue, cross-reference, advise on sources.

The ready availability of information on the Internet, and its widespread use, really presents librarians with an opportunity, not a threat. Savvy users realize they need help, which librarians can provide. A good example is Infomine, a cooperative project of the University of California and California State University (amongst others) (Mason *et al.*, 2000). Infomine contains descriptions and links to a wealth of scholarly and educational Internet resources, each of which has been selected and described by a professional academic librarian who is a specialist in the subject and in resource description generally. Participating librarians see this as an important expenditure of effort for their users, a natural evolution of their traditional task of collecting and organizing information in print.

What kind of technical infrastructure is needed to support and promote this kind of work? Open source software is a powerful ally for librarians who wish to extend liberal traditions of information access. These systems make the source code freely available for others to view, modify, and adapt; and the very nature of the licensing agreement prevents the software from being appropriated by proprietary vendors. But the open-source movement is more than just a vehicle for librarians to use: its link with library traditions goes much deeper. Public libraries and open source software both enshrine the same philosophy: to promote learning and understanding through the dissemination of knowledge. Both are pervaded by a sense of community, on the one hand the kind of inter-institutional cooperation exemplified by inter-library loan and on the other teams of designers and programmers that frequently cross national boundaries.

New trends in information access present librarians in developed countries with difficult and conflicting challenges. Meanwhile, however, the situation in the developing world is dire.

Here, traditional publishing and distribution mechanisms have failed tragically. For example, according to the 1999 UN Human Development Report (UNDP, 1999), whereas a US medical library subscribes to about 5,000 journals, the Nairobi University Medical School Library, long regarded as a flagship center in East Africa, last year received just 20 journals (compared with 300 a decade before). In Brazzaville, Congo, the university has only 40 medical books and a dozen journals, all from before 1993, and the library in a large district hospital consisted of a single bookshelf filled mostly with novels.

## Digital libraries and the challenge for HCI

Traditional libraries are substantial institutions that occupy physical space, present a physical appearance, and exhibit tangible physical organization. When standing on the threshold of a large bricks-and-mortar library you gain a sense of presence and permanence that reflects the care taken in building and maintaining the collection inside. Digital libraries, in contrast, are lightweight. But they provide potentially far greater accessibility, which means that they will have even greater social effects. Once created, they can, without significant institutional support, continue to serve users. They can be distributed throughout most of the developed world over the Internet. In developing countries and remote corners of the developed world they can be circulated on removable media—CD-ROM, DVD, or 100 Gb disk units the size of videocassettes—and updated over radio (or in Internet cafés). Issues of copyright pose difficult problems, but they are manageable. For example, there is plenty of non-copyright material, or material whose owners are prepared to donate copyright for socially useful purposes, and trends towards more open access to academic and humanitarian information are visible. Not everyone sees digital rights management and the DMCA as the way forward, and in the longer term publishers, to remain viable, will have to investigate alternative revenue models for the information they own. No wonder international organizations such as the United Nations, along with many smaller non-government organizations (NGOs), are keenly interested in digital library technology.

Advances in digital library technology are radically lowering the bar for the design and production of richly-organized, coherent, focused collections of information. Now, anyone with access to sufficient source material can use public-domain software to build large, fully-searchable, collections the size of traditional personal or institutional libraries—in minutes. Let the minutes stretch to hours and the collection can be polished, organized, branded, distributed. It can include fully-illustrated text, images, video, music. It can present attractively-designed

pages with consistent use of icons. Keywords, key phrases, even acronyms and their definitions, can be extracted—automatically—and used to underpin novel means of access. Let the hours stretch to days and metadata can be manually added that permits further levels of organization. Given access to programming skills, creative new facilities that stretch the imagination can be rapidly integrated into the system.

All this, one might say, can be done with ordinary Web sites: there is no need for digital library technology. However, bitter experience has shown that all but the most rudimentary sites do require significant institutional support—for organization and maintenance. The Web is littered with incomplete, unfinished, unmaintained, out-dated, inconsistently-organized, useless information collections. Just as traditional library cataloging procedures integrate new works into existing collections with minimal overhead so that they immediately become first-class members of the collection, so digital libraries allow new documents to be added completely automatically. In the case of traditional libraries this is done through the small but non-negligible overhead of generating a new catalog entry (one to two hours per book). With ordinary Web sites it requires inserting links manually into index pages and the like, and may involve adding links not only into the new document but also into existing ones that ought to reference it—it's like rewriting the book, and maybe revising all other books in the library too! In contrast, digital libraries bring access structures instantly and effortlessly up to date whenever new documents are added.

The challenge for HCI is to design and build digital library systems that fulfill the potential of digital libraries as a "killer app" for computers in developing countries, which will bring concomitant benefits in almost every other sphere of application. For the information consumer we need access to information that is guaranteed across space, time, and culture. We need flexible distribution mechanisms for documents, and for information objects of all types, that can be accessed on all computer platforms—including the lowliest. We need a choice of distribution over the Web or on removable media such as CD-ROM or DVD. Digital libraries can incorporate flexible presentation that caters to individual differences, such as large-font displays or spoken output for the visually impaired. Libraries are places where information is preserved, not rendered obsolete, and digital libraries must instill confidence that information prepared today can be accessed next week, next decade, next century—regardless of technological changes. An important aspect of digital libraries is their ability to work in local languages, promoting pluralism and reducing the risks of homogeneity. Because language is the vehicle of thought, communication, and cultural identity, this will encourage diversity and

strengthen individual cultures. But there is a long way to go: even Unicode is woefully incomplete in certain areas, such as African languages.

Naturally, today's digital library systems focus principally on the reader: the consumer of the material stored in the library's treasure-house. But digital libraries make a more radical, and perhaps ultimately more important, contribution by empowering ordinary users to conceive, assemble, build, and disseminate new information collections themselves. In principle, modest computing resources are quite sufficient to enable users to build new collections by gathering together material in local files or on the Web (or both); augmenting it with appropriate metadata that supports convenient search and browsing operations; incorporating advanced features like key-phrase extraction, document summarization, and metadata extraction; designing an attractive and functional interface; and publishing the collection on a variety of different media that are suitable for the intended readership. The HCI challenge is to realize this potential for users—such as most librarians—who have a strong understanding of information and its organization, but no more interest in computers than they have in paper-making technology, last millennium's vehicle for information dissemination.
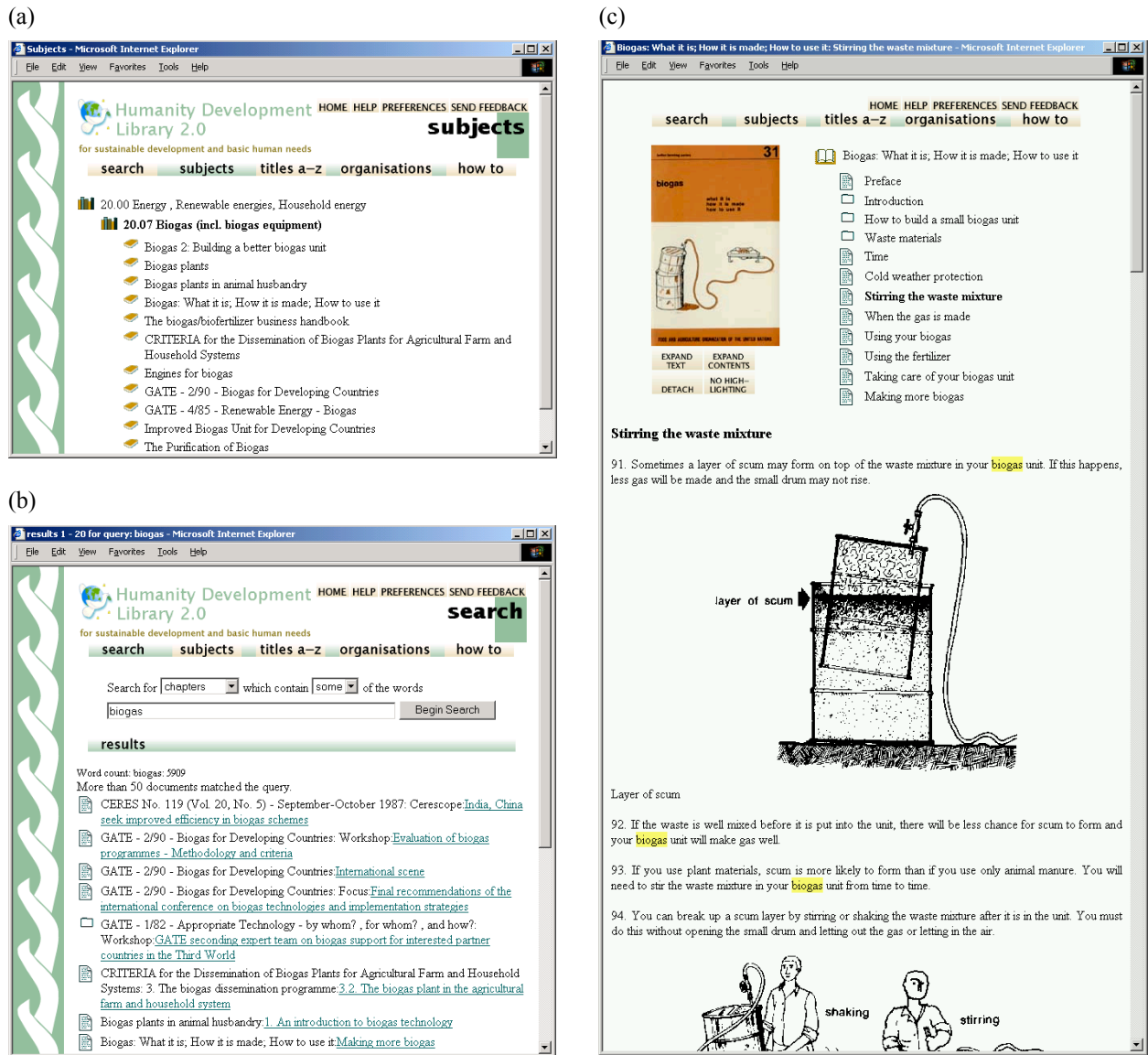
(a)

(c)

(b)

Figure 1 The Humanity Development Library: (a) Searching for *biogas*; (b) Browsing by subject; (c) Reading a document

# WHAT ARE DIGITAL LIBRARIES?

A digital library is an organized collection of information,

> *a focused collection of digital objects, including text, video, and audio, along with methods for access and retrieval, and for selection, organization, and maintenance of the collection.*

> (Witten and Bainbridge, 2002)

This definition deliberately accords equal weight to user (access and retrieval) and librarian

(selection, organization and maintenance). The latter functions are often overlooked by digital library proponents, who often work from a technology perspective rather than from the viewpoint of library or information science, but it is precisely these aspects that allow digital libraries to be used to democratize information dissemination.

As a concrete example, consider the Humanity Development Library, a collection of some 1,200 authoritative books and periodicals, produced by many disparate organizations—UN agencies and other international organizations—on various areas of human development, from agricultural practice to economic policies, from water and sanitation to society and culture, from education to manufacturing, from disaster mitigation to micro-enterprises. It contains 160,000 pages and 30,000 images, which if printed would weigh 340 kg, cost $20,000, and occupy a small library book stack. Instead, it takes the form of a digital library and is distributed on a CD-ROM throughout the developing world at essentially no cost. (It's also on the Web at nzdl.org.)

The Humanity Development Library is produced using the Greenstone software, a freely-distributed open-source project whose aim is to create novel digital library technologies and make them available for others to use. Greenstone digital libraries are arranged in *collections*. A collection comprises several (typically several thousand, or several million) documents, and a library may include several collections, each organized differently. Collections built with Greenstone offer simple but effective searching and browsing facilities based on metadata and the full text of electronic documents. Each collection is individually designed to take advantage of whatever metadata is available. All collections support full-text searching and most provide several different browsing options, although they differ depending on the collection design and the metadata available. Typically you can search for particular words that appear in the text, or within a section of a document, or within a title or section heading. A variety of interfaces exist for browsing collections by *title*, *subject*, *date*, or any other metadata chosen by the collection designer.

Figure 1 shows snapshots of the Humanity Development Library. In Figure 1a documents are being searched for chapters containing the word *biogas*. In Figure 1b the collection is being browsed by subject: by clicking on the bookshelf icons the user has discovered an item under *Section 20, Energy, Renewable energies, Household energy*. Pursuing an interest in energy from biogas, the user selects a book by clicking on its book icon (Figure 1c). All the icons in the screenshots of Figure 1 are clickable. Those at the top of the page return to the library home

page, provide help text, and allow you to set user interface and searching preferences. The navigation bar beneath gives access to the searching and browsing facilities, which differ from one collection to another. This particular collection can be searched by book, chapter, or section, and browsed by subject, title, organization, and "how to" metadata as indicated by the navigation bar.

Documents are presented as Web pages generated by Greenstone from the source material. In Figure 1c the book's cover is displayed as a graphic on the left, and an automatically constructed table of contents appears at the start of the document. The current focus, *Stirring the waste mixture*, is written in bold in the table of contents; its text (including any illustrations) starts further down the page. Incidentally, the material in this collection on building household biogas plans is fascinating, though it is not directly relevant to this article.

Greenstone collections present documents as automatically-generated Web pages. This allows documents in different source formats to be presented in a consistent manner, and lets users view the entire collection with a standard Web browser—no special viewing applications are required. However, the collection maintainer may choose to present the original source document (whether Word, PDF, PostScript, PowerPoint, Excel, a QuickTime movie, an audio file, or whatever) instead of, or as well as, the HTML version, and rely on the user's web browser to select a suitable application to display the document. In general, Greenstone deals well with documents and metadata in a wide variety of different formats.

## USING DIGITAL LIBRARIES TO DISSEMINATE HUMANITARIAN INFORMATION

Digital libraries provide perhaps the first really compelling *raison d'être* for computing technology in the developing world. Priorities in these countries include health, agriculture, nutrition, hygiene, sanitation, and safe drinking water. Though computers *per se* are not a priority, simple, reliable access to practical information relevant to these basic needs certainly is. In an article entitled "The promise of digital libraries in developing countries" Witten *et al*. (2002) mention ten information collections, available on the Web (at nzdl.org) and CD-ROM, from organizations ranging from UN agencies to small NGOs, in which Greenstone is being used to deliver humanitarian and related information in developing countries. For example, the Humanity Development Library described above is a compendium of practical information aimed at helping reduce poverty, increasing human potential, and giving a practical and useful education. Rather than recapitulating the brief summaries of collections that appear in the

above-cited paper, we describe four new ones that have been created recently, and distributed in the same way (Figure 2).

The Researching Education Development library is a project of the Department for International Development (DFID), a British government department responsible for promoting development. Its central focus is a commitment to an internationally agreed target of halving the proportion of people living in extreme poverty by 2015. Associated targets include ensuring universal primary education, gender equality in schooling, and skills development. It works in partnership with other governments and multilateral institutions, with business and the private sector, with civil society and the research community. It has created a CD-ROM library containing many education research papers and other documents. Each one represents a study or piece of commissioned research on some aspect of education and training in developing countries.

The Energy for Sustainable Development library was initiated as part of the outreach phase of the World Energy Assessment, which was initiated jointly by the United Nations Development Programme (UNDP), the United Nations Department of Economic and Social Affairs (UNDESA), and the World Energy Council (WEC), along with funding from the United Nations Foundation. This library contains a broad and valuable collection of 350 documents (26,000 pages) from UNDP, UNDESA, WEC and many other organizations. It includes titles that all these organizations have published on the subjects of energy for sustainable development—technical guidelines, journals and newsletters, case studies, manuals, reports, and other training material. The documents are in English, Spanish and French, and one document has Arabic, Russian and Chinese translations as well.

The UNAIDS Library contains publications in the "Best Practice" collection (including key materials, case studies, technical updates, and points of view) which form a unique resource for those working in planning and practice. It is produced by the Joint United Nations Programme on HIV/AIDS, whose global mission is to lead, strengthen and support a response to the AIDS epidemic that will prevent the spread of HIV, provide care and support for those infected by the disease, reduce the vulnerability of individuals and communities to HIV/AIDS, and alleviate the socioeconomic and human impact of the epidemic.

Figure 2 A selection of recent humanitarian digital library collections on CD-ROM

The Health Library for Disasters is the result of a collaboration between the emergency and disaster programs of the World Health Organization (WHO) and the Pan American Health Organization (PAHO), with the participation of many other organizations: the United Nations High Commissioner for Refugees (UNHCR), the United Nations Children's Fund (UNICEF), the International Strategy for Disaster Reduction (EIRD); the Red Cross Movement (ICRC and IFRC); the SPHERE Project; non-governmental organizations such as OXFAM; and national organizations such as the National Emergency Commission of Costa Rica. It contains more that 300 technical and scientific documents on disaster reduction and public health issues related to emergencies and humanitarian assistance. A follow-up to the Spanish-language Biblioteca Virtual de Desastres discussed by Witten *et al.* (2002), it includes technical guidelines, field guidelines, case studies, emergency kits, manuals, disaster reports, and other training materials.

## UNIVERSAL ACCESS

Universal access to digital libraries presents huge challenges to software engineers and HCI practitioners. The Greenstone digital library software allows us to glimpse some of the issues, although it certainly does not yet effectively address them all. We summarize some technical details in the next subsection, before turning to more interesting questions of access for readers, collection builders, and international users.

### Platforms and distribution

Most digital libraries are accessed over the web, using any web browser. However, in many environments, particularly in developing countries, web access is insufficient and the system must run locally. And if people are to build and control their own libraries, a centralized solution is inadequate: the software must run on their own computers. Thus digital library systems intended for broad access should run on a wide variety of computer systems, particularly low-end ones.

Developed under Linux, the Greenstone server runs on any Windows, Unix, or MacOS/X system. All versions of Windows are supported, from 3.1 up (including 3.1/3.11, 95/98/ME/, NT/2000 and XP). Supporting primitive platforms poses substantial challenges of a rather mundane nature: for example, Microsoft compilers no longer support Windows 3.1 and it is necessary to acquire obsolete versions (e.g. at software auctions). Under Windows, pre-built collections can be viewed on any system with at least 8 Mb RAM, but collections cannot be built under Windows 3.1/3.11—for this at least a Pentium processor is generally required, except for very small collections. The fact that Greenstone does not run on early Macintosh systems is a serious drawback in certain environments (e.g. many schools).

In an international cooperative effort established in August 2000 with UNESCO and the Belgium-based Human Info NGO, Greenstone is being distributed widely in developing countries with the aim of empowering users, particularly in universities, libraries, and other public service institutions, to build their own digital libraries. UNESCO recognizes that digital libraries are radically reforming how information is acquired and disseminated in its partner communities and institutions in the fields of education, science and culture around the world, and particularly in developing countries. Their hope is that this software will encourage the effective deployment of digital libraries to share information and place it in the public domain.

The UNESCO distribution of Greenstone is a CD-ROM that contains the full source code and executable binaries for Windows and Linux, along with all necessary associated software (e.g. Perl for Windows). Full documentation (four PDF manuals) and five demonstration collections are included. The current CD-ROM is trilingual, with complete interfaces, instructions, and documentation in English, French and Spanish. For those with Web access, the same package is also available for download from the Greenstone Web site (greenstone.org), often in a form that is slightly ahead of the CD-ROM version—for example, many other language interfaces are included and full documentation is available in Russian and Kazakh too. Providing accessibility in different languages is more difficult than one might at first realize. As well as the manuals, installation instructions and installation prompts, the licensing agreement, and the readme files have to be translated too.

## Access for readers

Greenstone collections like the Humanity Development Library can be published as standalone collections on removable media such as CD-ROM, or presented on the Web. CD-ROM is a very practical format in developing countries. Any Greenstone collection can be converted into

a self-contained Windows CD-ROM that includes the Greenstone server software itself (in a version that runs right down to Windows 3.1) and an integrated installation package. The installation procedure has been thoroughly honed to ensure that only the most basic of computer skills are needed to install and run a collection under Windows.

Even standalone Greenstone users interact through a Web browser: Netscape is supplied on each CD-ROM for those who do not already have a browser. In standalone mode the software runs locally but incorporates a Web server so that if the system happens to be connected to a network—say a hospital or school intranet—information is available to other machines that may not possess CD drives. This happens automatically: no special configuration is necessary. Another difficult engineering challenge is checking for the existence of a network. While installed network software is easily detected, it is hard to determine non-intrusively whether it is operational (sending oneself a message often results in the user being asked to dial their local Internet service provider). Incorrectly installed or configured software is endemic in developing countries, because computers there are often cast-offs whose software is inappropriate to their present environment, yet system support to rectify the problems is unavailable. It is essential for universal access that such problems are addressed properly and solved satisfactorily without involving the user, even though they are mundane and time-consuming.

Greenstone provides some support for the visually impaired by incorporating a "textual" mode of access that replaces all images by textual prompts. This output is suitable for users with speech synthesizers or other specialized access devices. However, the facility is not well advanced: in particular, we have not yet refined it through usability testing and interface improvement.

**Building new collections**

Effective human development blossoms from empowerment rather than gifting. As the Chinese proverb says, "Give a man a fish and he will eat for a day; teach him to fish and he will eat for the rest of his days." Disseminating information originating in the developed world, like the Humanity Development Library, is a useful activity for developing countries. But a more effective strategy for sustained long-term human development is to disseminate the capability of creating information collections, rather than the collections themselves. This will allow developing countries to participate actively in our information society, rather than observing it from outside. It will stimulate the creation of new industry. And it will help ensure that

(a)                                                           (b)



(c)                                                           (d)
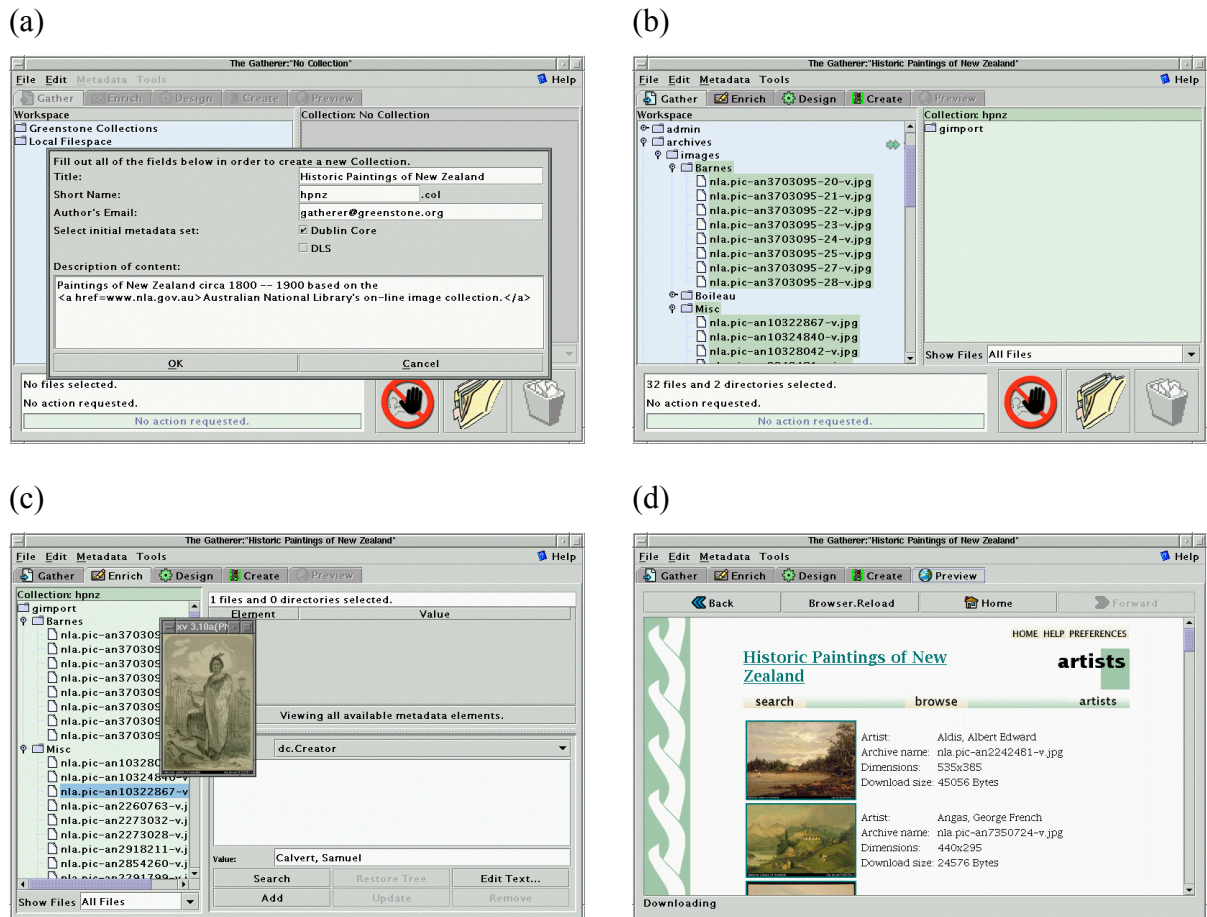


Figure 3 Building a collection with the Gatherer.

intellectual property remains where it belongs, in the hands of those who produce it.

Users whose skills resemble those of librarians rather than computer specialists should be able to build and distribute their own digital library collections. As an initial step in this direction, Greenstone includes an interface called the "Collector" that is intended to help people build their own library collections (Witten *et al*. 2000). Collections may be built and served locally from the user's own web server, or (given appropriate permissions) remotely on a shared digital library host. End users can build new collections styled after existing ones from material on the Web or from their local files—or both, and collections can be updated and new ones brought on-line at any time. The interface, which is intended for non-professional end users, is modeled after widely used commercial software installation packages (such as InstallShield[i]), frequently called software "wizards"—a term we deprecate because of its appeal to mysticism and connotations of utter inexplicability. We chose this interaction style because it simplifies the choices and presents them very clearly.

Figure 3 shows a new Greenstone librarian interface, currently undergoing beta testing by

UNESCO at sites in Argentina, India, Kazakhstan, Mexico, and South Africa, which builds on lessons learned from the Collector. It incorporates a great deal of additional functionality, particularly the ability for users to associate metadata with any item or group of items, and to reuse metadata elements without retyping them. In Figure 3 it is being used to collate a selection of images for a digital library collection, augment these source documents with textual metadata and then build and view the collection. From here, it is a matter of a few further clicks to produce a self-installing CD-ROM version of the collection.

In this illustration the user is developing a digital library collection of historic paintings of New Zealand. The user creates a new collection using the file menu (Figure 3a), and a resulting popup window prompts for some general information about the collection. Once the user has filled out this form the main window becomes active. A series of panels guide the user through the processes required to build the collection. The left-hand pane of the *Gather* panel, shown in Figure 3b, shows the file system and the right-hand one represents the contents of the collection, initially empty, which the user populates by dragging and dropping files. In Figure 3c the user has moved to the *Enrich* panel and is adding textual metadata (the name of the artists) to the selected documents. The next two panels, *Design*, *Create*, help the user structure the collection, control its appearance, and build it. On completion the result is viewed in the *Preview* panel. Figure 3d shows a page from the newly built collection, in which source documents are alphabetically listed by artist. Shown alongside each thumbnail are the artist's name, its catalog number, image dimensions, and its download size. The full-size image is shown by clicking on the thumbnail. The user may skip backwards and forwards through the panels using them to augment and enhance the collection, perhaps adding further source documents, editing metadata values, and altering the collection's appearance.

## Customization

An important component of access is allowing people to control the appearance of the collections they create. Many who build digital libraries want to brand them to ensure that they with an appropriate personal, institutional, or corporate image, and some can only contemplate software solutions that allow them to do so. Although Greenstone comes with the standard appearance shown in Figure 1—a distinctive bar down the side of all pages except those that show documents in the library, a green access bar with yellow buttons, etc.—the interface is highly configurable.

Greenstone creates all pages that appear on the screen on the fly: none are stored in advance.

They are generated using macros, written in a simple language specially designed for the job, that perform textual replacement. One reason is that Greenstone accommodates a large number of different interface languages (see next subsection), and macros help cope with this. All text fragments are couched as macro definitions. To add a new language, just the macro contents need to be translated—no web pages need be reworked. Every page displayed by the system is passed through a macro interpreter that expands all the macros on the page. The interpreter checks a language variable and uses the macro definitions pertaining to it, which loads the page in the appropriate language.

Macros can have parameters. In this case, the parameter is the language variable: it causes the appropriate text fragment to be used for the macro's expansion. If there is no Arabic version for a particular macro, the interpreter will automatically substitute the default version (English). This lets system developers experiment with the interface without having to worry about translating every little bit of new text immediately. Defaulting to English is not ideal—it reflects an Anglo-centric mindset—but it seems better than displaying nothing. (However, if "nothing" were preferred, it would be a simple matter to alter the software to default to the empty language!)

Macros are also used to deal with display variables. Whenever a web page contains information that is not known in advance—like the number of documents returned by a search, or the value of a particular metadata item, or the content of a document page—a macro name is used in the page description. Unlike language macros, these macros are *dynamic*: their content is not stored in advance but generated by the system in accordance with the value of the variable in question.

Users can completely alter the form of the user interface by rewriting the macro files—or even by writing their own web pages which embed the dynamic macros that generate bits of Greenstone output. Figure 4 shows a Greenstone interface that has been heavily customized by Lehigh University Library, Pennsylvania (bridges.lib.lehigh.edu). The standard Greenstone appearance has been completely obliterated in favor of an in-house style, yet all the Greenstone functionality is available. Figure 4 shows a thumbnail of a book cover on the left, and full bibliographic information on the right. Entering a book displays facsimile images of its pages.

**Internationalization**

The international Unicode character set is used throughout Greenstone, and documents in any
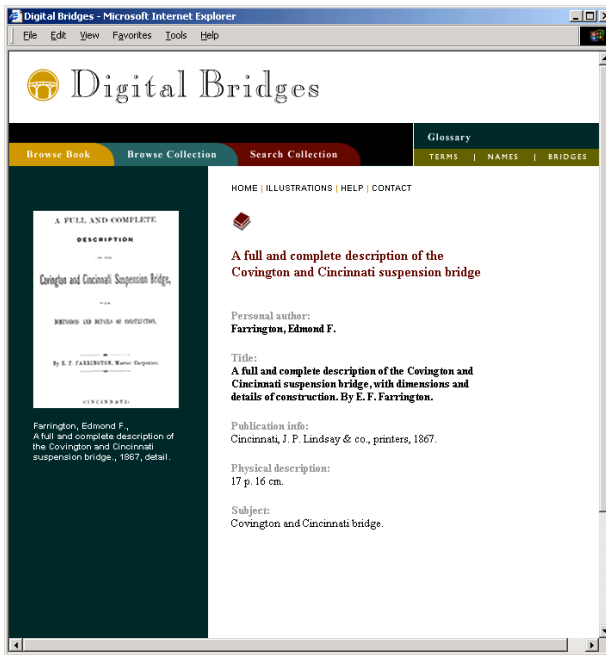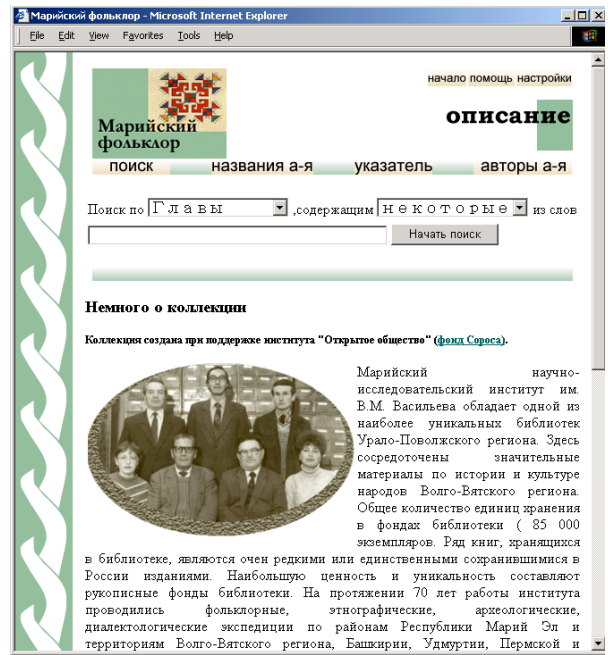
Figure 4 A customized interface to Greenstone



Figure 5 A Russian digital library collection

Unicode-supported language and character encoding can be imported. (In fact, the software can automatically detect the language and encoding of most documents.) Collections of documents in Arabic, Chinese, Cyrillic, English, French, Spanish, German, Hindi, and Maori are publicly available. The New Zealand Digital Library Web site (nzdl.org) hosts many of these, and the Greenstone Web site (greenstone.org) links to sites that contain further examples.

It makes little sense to have a collection whose content is in Chinese or Russian, but whose supporting text—instructions, navigation buttons, labels, images, help text, and so on—are in English. Consequently, the entire Greenstone interface has been translated into a range of languages, and the interface language can be changed by the user as they browse from the *Preferences* page. As noted above, all the language fragments in the interface (and also the contents of language-dependent images) are stored in macro files. These have been translated by Greenstone users in other parts of the world and contributed back to the project. (The same mechanism provides text-only versions of the interface to accommodate visually impaired users.) Figure 5 shows an example: a Russian collection. Currently, interfaces are available in Arabic, Czech, Chinese, Dutch, French, Galician, German, Hebrew, Indonesian, Italian, Kazakh, Maori, Portuguese, Russian, Spanish, Turkish, and English.

Managing the organizational and software complexity of any comprehensive and evolving open source software system presents a significant challenge. However, the challenge is greatly

magnified when the interface is available in different languages, for enhancements to the software and changes to the interface must be faithfully reflected in each language version. No single person knows all interface languages; no single person knows about all modifications to the software—indeed there is likely no overlap at all between those who translate the interface and those who build the software. Currently, Greenstone has about twenty interface languages and there are around 600 linguistic fragments in each interface, ranging from single words like *search*, through short phrases like *search for*, *which contain*, *of the words*, to sentences like *More than ... documents matched the query*, to complete paragraphs like those in the on-line help text. Maintaining the interface in many different languages is a logistic nightmare. The solution adopted by Greenstone is to incorporate a language translation facility, which allows authorized people to update the interface in specified languages. A standard version control system is used to manage software change, and from this the system automatically determines which language fragments need updating and presents them to the human translator.

## CONCLUSIONS

By allowing people to easily create and disseminate large information collections, digital libraries extend the applications of modern technology in socially responsible directions, and counter a possible threat towards the commercialization of information in line with practices developed by the entertainment industry. As far as the developing world is concerned, digital libraries may prove to be a "killer app" for computer technology—that is, an application that makes a sustained market for a promising but under-utilized technology. The World-Wide Web is often described as the Internet's killer app. But the Internet does not really extend to developing countries, and the developing world is missing out on the prodigious amount of basic, everyday human information that the Web provides, and its enormous influence on promoting and internationalizing business opportunities. There is little incentive to make copies of the entire Web available locally because of its vast size, rapid change, and questionable information value per gigabyte. However, it is easy to provide focused information collections on both the Web and, in exactly the same form, on removable media such as CD-ROM, DVD, or bulk disk storage devices—indeed, the Greenstone software described above allows one to create a complete, runnable, self-installing CD-ROM image from a Web collection in just a few mouse clicks.

Public libraries are founded on the principle of universal access, and digital libraries should be too. This provides HCI with enormous practical challenges. Universal access means running on

low-end devices, but one does not want to provide a lowest-common-denominator solution that sacrifices high-end capability where it is available. Universal access means that interfaces should be available in the world's languages, but one does not want the burden of translation to stifle the development of new functionality and features. Universal access means educating users: UNESCO is mounting training courses on building collections with Greenstone in Bangalore, Almaty, Senegal, and Suva, and discussions are underway for Latin America; the Tulane Institute has run courses that use Greenstone collections as a resource in many locations in Africa (e.g. Burkina Faso, Cameroon, Cote d'Ivoire, Democratic Republic of Congo, Ghana, Rwanda, Senegal, Sierra Leone, Togo) and Latin America (e.g. Argentina, Bolivia, Colombia, Ecuador, Guatemala).

Universal access also means that non-textual material should enjoy first-class status in a digital library—perhaps first-class status in "the literature." This has important cultural ramifications. It should be possible to create digital library collections intended for use by people in oral cultures, who may be illiterate or semi-literate. Or people who, though they can read and write their own language, cannot speak or read the language of the digital library. Imagine having access to collections that spring out of the rich cultures of China or Arabia, created by people who grew up in these cultures, without having to learn a new language. More practically—since you, the reader, being culturally privileged, can probably access this kind of information in translation—imagine giving someone in the highlands of Peru, fluent and literate in her native language of Quechua, first-hand access to the information in humanitarian collections such as the Humanity Development Library (currently available only in English and French) or the Biblioteca Virtual de Desastres (until recently available only in Spanish). Opening up digital libraries for the illiterate is a radical and potentially revolutionary benefit of new interface technology.

An important, and liberating, different between digital libraries and conventional ones is that anyone should be able to create their own digital collections. This presents HCI challenges that are difficult yet more conventional: providing non-computer users with access to advanced and complex functionality. Users should be able to collect their own source material, provide their own metadata, design their own collections, and present it through their own interface.

Digital libraries give software engineers and HCI practitioners a golden opportunity to help reverse the negative impact of information technology on developing countries and reduce the various "digital divides" that cleave our world (Norris, 2001)—the "social divide" between the

information rich and the information poor in our own nations, the "democratic divide" between those who do and do not use the panoply of digital resources to engage, mobilize and participate in public life, as well as the "global divide" that reflects the huge disparity in access to information between people in industrialized and developing societies.

## ACKNOWLEDGEMENTS

## REFERENCES

ALA (2002) "Rediscover America @ your library." Video produced by the American Library Association, Chicago, IL. Available from www.ala.org/@yourlibrary/rediscoveramerica.

Arunachalam, S. (1998) "How the Internet is failing the developing world." Presented at *Science Communication in the Next Millennium*, Egypt; June.

Ciolek, T. M. (1996) "The six quests for the electronic grail: Current approaches to information quality in WWW resources." *Review Informatique et Statistique dans les Sciences humaines (RISSH)*, No. 1-4. Centre Informatique de Philosophie et Lettres, Universite de Liege, Belgium. pp. 45-71.

Lynch, C. (2001) "The battle to define the future of the book in the digital world." *First Monday*, Vol. 6, No. 5; June.

Mason, J., Mitchell, S., Mooney, M., Reasoner, L. and Rodriguez, C. (2000) "INFOMINE: Promising directions in virtual library development." *First Monday*, Vol. 5, No. 6; June.

Norris, P. (2001) *Digital divide? Civic engagement, information poverty and the Internet worldwide*. Cambridge University Press, New York

Quéau, P. (2001) "Information literacy: a new frontier." *UNISIST Newsletter*, Vol. 29, No. 2, pp. 3-4.

Roehl, R. and Varian, H.R. (2001) "Circulating libraries and video rental stores." *First Monday*, Vol. 6, No. 5; May.

Samuelson, P. and Davis, R. (2000) "The digital dilemma: a perspective on intellectual property in the information age." Presented at the Telecommunications Policy Research Conference, Alexandria, Virginia; September.

UNDP (1999) *Human development report 1999*. UNDP/Oxford University Press, New York.

Witten, I.H., Loots, M., Trujillo, M.F. and Bainbridge, D. (2002) "The promise of digital libraries in developing countries." *The Electronic Library*, Vol. 20, No. 1, pp. 7–13.

Witten, I.H. and Bainbridge, D. (2002) *How to build a digital library*. Morgan Kaufmann, San Francisco.

---

[i] www.installshield.com