

Chapter 4.21

Digital Library Structure and Software

Cavan McCarthy
Louisiana State University, USA

INTRODUCTION

Digital libraries (DL) can be characterized as the “high end” of the Internet, digital systems which offer significant quantities of organized, selected materials of the type traditionally found in libraries, such as books, journal articles, photographs and similar documents (Schwartz, 2000). They normally offer quality resources based on the collections of well-known institutions, such as major libraries, archives, historical and cultural associations (Love & Feather, 1998). The field of digital libraries is now firmly established as an area of study, with textbooks (Arms, 2000; Chowdhury & Chowdhury, 2003; Lesk, 1997); electronic journals from the US (D-Lib Magazine: <http://www.dlib.org/>) and the UK (Ariadne: <http://www.ariadne.ac.uk/>); even encyclopedia articles (McCarthy, 2004).

BACKGROUND

Digital libraries require appropriate presentation and careful logical organization to make them easily accessible, but arrangements typical of Web systems are inadequate for them. The classic Web structure, where random links can be created between any pair of pages, is not appropriate to highly organized data. The other classic arrangement is the tree or directory structure often found in computerized systems, where the user starts from a “trunk” or “root directory” and goes to a branch, then a subdivision of that branch. This is effective for individual images, but is inadequate for navigating sequential pages, as in a digital library system presenting lengthy texts. Before discussing the different software solutions available, it is useful to review the principle types of digital material currently offered by digital libraries.

DIGITAL LIBRARY MATERIALS

At this time digital library resources can be divided into three categories: images, texts and other resources:

Images

Image access is used for individual visual resources, such as photographs, posters, drawings, etc. The classic procedure uses a series of three types of image. Scanning produces a high-quality archive image, which is then used to generate an access image, for general public use. Finally, a small thumbnail image is produced, for quick reference (Boss, 2001; Lee, 2001). In more detail:

Archive Image

A high-quality image, scanned directly from the original, destined for long-term preservation. Normally an uncompressed TIF (Tagged Image File Format) image is used here; TIFs offer the highest quality images and a resolution of 600 dpi (dots per inch) is standard. As scanning is an expensive operation, which exposes original materials to possible damage, the archive image will be carefully preserved. It must always exist at the system level, but is not necessarily available to the end-user. TIF files occupy significant server space and imply lengthy download times. Another factor is that some DL will want to sell their own hard-copy prints of quality images.

Access Image or Working Image

A quality image, adequate for consultation and serious study by digital library users. This is normally a high-quality JPG (Joint Photographic Experts Group) image, generated from the archival TIF. JPG files are widely used on the Internet and offer quality spatial and color reproduction and a high compression ratio. For DL purposes JPG images will often be generated at a resolu-

tion of 300 dpi; a size of 640x480 pixels is also common.

Thumbnail Image

A small reference image, which gives the user a general idea of the Access image, before downloading that image. Typically a medium to low quality JPG, generated from the Access image, but about one-tenth of its size, and commonly produced at a resolution of 72 dpi. GIF format (Graphic Interchange Format) can also be used for thumbnails (Arizona, 2000; Western, 2003).

Text

Multi-page text documents, such as books, or journal articles require special procedures. Numerous options are possible and the principle alternatives for input, simple text presentations and pagination will be examined in turn; the earliest procedures will be discussed first.

Text Input

Manual keyboarding was originally adopted by Project Gutenberg, the first significant text-oriented digital library (<http://promo.net/pg/>), founded in 1971. This is a laborious process which severely limits productivity, and is now rarely used.

OCR (Optical Character Recognition) software is now routinely used to scan text into digital libraries. OmniPage Pro (<http://www.scansoft.com/omnipage/>) or the Russian software ABBYY (<http://www.abbyy.com/>) are frequently cited in the digital library context. OCR text requires careful revision, because even 99.99% accuracy means that there will be one mistake every couple of pages, but only a person fully conversant with the literature will be able to identify errors at this level. Many digital library texts are older books whose ornate type faces or soiled pages can generate additional OCR errors.

6 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/digital-library-structure-software/29474

Related Content

Evolutionary Computation as a Paradigm for Engineering Emergent Behaviour in Multi-Agent Systems

Robert E. Smith and Claudio Bonacina (2003). *Intelligent Agent Software Engineering* (pp. 118-136).

www.irma-international.org/chapter/evolutionary-computation-paradigm-engineering-emergent/24147

ART-Improving Execution Time for Flash Applications

Ming Ying and James Miller (2011). *International Journal of Systems and Service-Oriented Engineering* (pp. 1-20).

www.irma-international.org/article/art-improving-execution-time-flash/55059

SNI Field Blocking and Internet Censorship

JiYoung Jung, Minwoo Park, Hee Kyoung Shin and Yongtae Shin (2022). *International Journal of Software Innovation* (pp. 1-12).

www.irma-international.org/article/sni-field-blocking-internet-censorship/289601

Building Defect Prediction Models in Practice

Rudolf Ramler, Johannes Himmelbauer and Thomas Natschläger (2014). *Handbook of Research on Emerging Advancements and Technologies in Software Engineering* (pp. 540-565).

www.irma-international.org/chapter/building-defect-prediction-models-in-practice/108635

Towards Test-Driven and Architecture Model-Based Security and Resilience Engineering

Ayda Saidane and Nicolas Guelfi (2014). *Software Design and Development: Concepts, Methodologies, Tools, and Applications* (pp. 2072-2098).

www.irma-international.org/chapter/towards-test-driven-architecture-model/77791