

**Trusted Digital Repositories:
Attributes and Responsibilities**

An RLG-OCLC Report

RLG
Mountain View, CA
May 2002

© 2002 Research Libraries Group

All rights reserved

First published in May 2002

Adobe and Acrobat Reader are registered trademarks of Adobe Systems Incorporated in the United States and/or other countries.

RLG, Inc.

Mountain View, California 94041 USA

www.rlg.org

Executive Summary

In March 2000 RLG and OCLC began a collaboration to establish attributes of a digital repository for research organizations, building on and incorporating the emerging international standard of the *Reference Model for an Open Archival Information System (OAIS)*. A working group was created to reach consensus on the characteristics and responsibilities of trusted digital repositories for large-scale, heterogeneous collections held by cultural organizations. A draft report was issued in August 2001 and in the extended comment period that followed a variety of interested individuals and organizations around the world contributed numerous thoughtful and helpful suggestions that have been incorporated into this final report.

As the Commission on Preservation & Access (CPA)/RLG Task Force on Archiving of Digital Information did in their 1996 report, this working group recognizes the development of national—and, increasingly, international—systems of digital repositories that are or will soon be responsible for the long-term access to the world's social, economic, cultural, and intellectual heritage in digital form. Also like the Task Force, the working group understands that content creators, owners of information, and current and potential users must be able to trust repositories with this responsibility.

A trusted digital repository is one whose mission is to provide reliable, long-term access to managed digital resources to its designated community, now and in the future. In this report the working group has articulated a framework of attributes and responsibilities for trusted, reliable, sustainable digital repositories capable of handling the range of materials held by large and small research institutions. The framework is broad enough to accommodate different situations, architectures, and institutional responsibilities while providing a basis for the expectations of a trusted repository. The critical component will be the ability to prove reliability and trustworthiness over time.

Finally, the report recommends that RLG, OCLC, and other organizations:

1. Develop a process for the certification of digital repositories.
2. Research and create tools to identify the significant attributes of digital materials that must be preserved.
3. Research and develop models for cooperative repository networks and services.
4. Develop systems for the unique, persistent identification of digital objects that expressly support long-term preservation.
5. Investigate and disseminate information about the complex relationship between digital preservation and intellectual property rights.
6. Determine the technical strategies that best provide for continuing access.
7. Define the minimal-level metadata required for long-term management and develop tools to automatically generate and/or extract as much of it as possible.

Just a few years ago, the development of trusted digital repositories seemed far in the future, but today it is an immediate challenge. The expert involvement and community consensus developed during the course of this work suggest that organizations and funding agencies will need to work together in the very near future to address the needs articulated in this report.

RLG/OCLC Working Group on Digital Archive Attributes

Neil Beagrie
Assistant Director, Digital Preservation
Joint Information Systems Committee

Catherine Lupovici
Head, Digital Library Department
Bibliothèque nationale de France

Dr. Marianne Doerr
Director, Center for Digitization
Bayerische Staatsbibliothek

Kelly Russell
Project Manager (1998-2001), Cedars
(CURL Exemplars in Digital Archiving)

Dr. Margaret Hedstrom
Associate Professor, School of Information
University of Michigan

Colin Webb
Director, Preservation Services Branch
National Library of Australia

Maggie Jones
Project Manager (2001-2002), Cedars
(CURL Exemplars in Digital Archiving)

Deborah Woodyard
Digital Preservation Coordinator
The British Library

Anne Kenney
Associate Director
Department of Preservation & Conservation
Cornell University

RLG Liaison
Robin Dale
Program Officer
Member Programs & Initiatives

OCLC Liaison
Meg Bellinger
President, Preservation Resources

Contents

Executive Summary	i
RLG/OCLC Working Group on Digital Archive Attributes	iii
Introduction	1
Intended Audience.....	2
Terminology.....	3
1 Trusted Digital Repositories	5
A Definition.....	5
Some Examples.....	5
Establishing Trust.....	8
<i>Earning the Trust of Designated Communities</i>	9
<i>Trusting Third-Party Providers</i>	9
<i>Trusting Digital Documents</i>	10
2 Attributes of a Trusted Digital Repository	13
Compliance with the <i>Reference Model for an Open Archival Information System (OAIS)</i>	13
Administrative Responsibility	13
Organizational Viability.....	14
Financial Sustainability	14
Technological and Procedural Suitability.....	14
System Security.....	14
Procedural Accountability.....	15
3 Responsibilities of a Trusted Digital Repository	17
High-Level Organizational and Curatorial Responsibilities	17
<i>The Scope of Collections</i>	17
<i>Preservation and Lifecycle Management</i>	18
<i>The Wide Range of Stakeholders</i>	18
<i>Ownership of Material and Other Legal Issues</i>	18
<i>Cost Implications</i>	19
Operational Responsibilities	21
<i>Negotiating for Appropriate Information from Content Providers</i>	21
<i>Obtaining Sufficient Control of the Information</i>	23
<i>Determining the Repository's Designated Community</i>	27
<i>Ensuring the Information to Be Preserved Is Independently Understandable to the Designated Community</i>	27
<i>Following Documented Policies and Procedures</i>	28
<i>Making the Preserved Information Available to the Designated Community</i>	29
<i>Advocating Good Practice in the Creation of Digital Resources</i>	31
4 Certification of Trusted Digital Repositories	33

5 Summary and Recommendations	37
Appendix A: OAIS Technical Overview.....	41
The OAIS Reference Model and Digital Preservation Metadata.....	41
<i>Preservation Description Information (PDI)</i>	42
<i>Representation Information (RI)</i>	42
<i>International Collaboration and Digital Preservation Metadata</i>	43
The OAIS Functional Model	43
<i>Submission and “Pre-Ingest” Activities</i>	44
<i>Ingest</i>	45
<i>Archival Storage</i>	45
<i>Data Management</i>	46
<i>Preservation Planning</i>	46
<i>Archive Administration</i>	47
<i>Access/Dissemination</i>	47
Appendix B: The Evolution of “Trust” in Computing Systems.....	49
Appendix C: Operational Responsibilities Checklist.....	55
Glossary	57
Selected Resources.....	61
Projects	61
Publications	61

Introduction

All research resources need care and attention to survive, but digital research resources need more attention, often much sooner than resources on paper. The inherent fragility of digital materials leaves only a small window of opportunity to address this problem before we start to lose resources on an ever larger scale.

As the traditional custodians of cultural heritage, libraries, archives, and museums are actively addressing methods and strategies for preserving digital materials. The challenge is great: cultural institutions are rapidly creating, converting, and acquiring material in a vast variety of formats, from word-processed documents to still images, and from data sets to electronic records. As collections grow, the needs associated with their maintenance and long-term viability grow too. A necessary outgrowth of this process has been the development of digital archives or digital repositories.

In 1994 the Commission on Preservation & Access (CPA)/RLG Task Force on Archiving of Digital Information began work together to describe and explore the nature of a reliable repository for digital materials. The major findings of the 1996 CPA/RLG report included these key points:

- Long-term preservation of digital information on a scale adequate for the demands of future research and scholarship will require a deep infrastructure capable of supporting a distributed system of digital archives.
- A critical component of the digital archiving infrastructure is the existence of a sufficient number of trusted organizations capable of storing, migrating and providing access to digital collections.
- A process of certification for digital archives is needed to create an overall climate of trust about the prospects of preserving digital information.¹

By 1998 a few large institutions were beginning to grapple with building digital repositories, but most cultural organizations still wanted specific guidance for developing either their own local repositories or trusted third-party digital repository services. A study and report by Dr. Margaret Hedstrom and Sheon Montgomery, *Digital Preservation Needs and Requirements in RLG Member Institutions*, examined the responsibilities of archives, libraries, museums, and other repositories for preserving and providing access to valuable, digital resources. The study found that two-thirds of respondents had assumed responsibility for digital information but 42% of those institutions reported that they lacked the operational and/or technical capacity to mount, read, and access some digital material in their holdings. And because of the nature of the materials held by cultural organizations, there was little surprise that three-fourths of respondents reported that irreplaceable information would be lost if their digital materials were not adequately preserved.² With the responsibility—and in some cases legal

¹ John Garrett and Donald Waters, *Preserving Digital Information: Report of the Task Force on Archiving of Digital Information* (Commission on Preservation and Access and RLG, 1996) www.rlg.org/ArchTF/index.html.

² Margaret Hedstrom and Sheon Montgomery, *Digital Preservation Needs and Requirements in RLG Member Institutions* (RLG, 1998) www.rlg.org/preserv/digpres.html.

obligation—to preserve these materials, institutions were faced with building local digital repositories or identifying trusted third parties that could meet their digital preservation needs. In either case, the task was daunting because community consensus had yet to be reached on a preferred digital archiving infrastructure and on what it meant to be a trusted digital repository.

As the same time as the survey, work on some of the needs articulated by the CPA/RLG Task Force was being advanced by the Consultative Committee for Space Data Systems in its *Reference Model for an Open Archival Information System (OAIS)* and by the many groups and individual institutions that were designing their own digital repository systems. The *Reference Model for an Open Archival Information System* was designed “to create a consensus on what is required for an archive to provide permanent, or indefinite long-term, preservation of digital information.”³ By establishing a common framework of terms and concepts, the Reference Model allowed existing and developing archives to be compared and contrasted, building the stage for a better understanding of the requirements of a full digital archive—not just the technical system requirements, but “an organization of people and systems, that has accepted the responsibility to preserve information and make it available for [its] Designated Community.”⁴

In March 2000 RLG and OCLC began a collaboration to establish attributes of a digital repository for research organizations, building on and incorporating the soon-to-be-international standard of the OAIS Reference Model. RLG and OCLC recognized that, despite emerging OAIS-related projects and initiatives in Europe and Australia, a definition and consensus on the characteristics and responsibilities of a sustainable digital repository for large-scale, heterogeneous collections held by cultural organizations was still needed. They formed a working group comprising digital preservation experts from around the globe who represented key research organizations involved in the long-term maintenance of digital materials (see page iii, RLG/OCLC Working Group on Digital Archive Attributes).

This report describes a framework of attributes and responsibilities for trusted repositories for digital content capable of handling the range of materials held by large and small research institutions. It builds on the foundations laid down in the CPA/RLG report, including its concept of a “deep infrastructure” and on the more recent work on the OAIS Reference Model, which provides a high-level, generic model for the environment, producers, users, data types, and information flows of a digital repository.

Intended Audience

The guidance and recommendations included here are primarily intended for cultural institutions such as libraries, archives, museums, and scholarly publishers and are specifically aimed at those with traditional or legal responsibility for the preservation of our cultural heritage. At the same time, the types of technologies needed for trusted digital preservation should be applicable to any organization interested in the long-term maintenance of and continuing access to digital materials. Although the report highlights some key strategic issues, its main focus is practical so that it can be useful to senior administrators as well as to those implementing digital archiving services. And while this report is concerned with

³ Consultative Committee on Space Data Systems, *Reference Model for an Open Archival Information System (OAIS): CCSDS 650.0-R1.1* (Washington, DC: CCSDS Secretariat, February 16, 2001): iii.

⁴ *Ibid.*, 1-1.

technology, it is not a technical document. It does not assume detailed technical knowledge on the part of the reader, but some basic level of technical understanding and a general awareness of digital preservation and related issues will be necessary.

This report provides guidance relevant to local, regional, national, and, international efforts—digital preservation is not limited by geography. The coordination of digital preservation activities will be of paramount importance to ensure convergent rather than divergent development; future scholarship and research will depend heavily on interoperability and collaboration.

Terminology

Digital preservation interests a range of different communities, each with a distinct vocabulary and local definitions for key terms. While a general glossary is provided (see page 55), it is important to draw attention to the meaning of a few terms within the context of this document.

For the purposes of this report, “digital preservation” is defined as the managed activities necessary for ensuring both the long-term maintenance of a bytestream and continued accessibility of its contents. If discussion pertains specifically to one, a more precise term is used.

The OAIS Reference Model uses “digital archive” to mean the organization responsible for digital preservation; this paper uses “archive” in place of “repository” only when “archive” is taken directly from the OAIS Model.

The phrase “designated community” is taken directly from the OAIS Reference Model and is defined as “an identified group of potential users of the archives’ contents who should be able to understand a particular set of information.” A “designated community” is multifaceted and decisions about what to preserve must take into account not only the needs of current users, but also those of users far into the future. The complex issues related to designated communities are addressed throughout the report.

1 Trusted Digital Repositories

A Definition

A trusted digital repository is one whose mission is to provide reliable, long-term access to managed digital resources to its designated community, now and in the future. Trusted digital repositories may take different forms: some institutions may choose to build local repositories while others may choose to manage the logical and intellectual aspects of a repository while contracting with a third-party provider for its storage and maintenance. Whatever the overall infrastructure, however, to meet expectations all trusted digital repositories must

- accept responsibility for the long-term maintenance of digital resources on behalf of its depositors and for the benefit of current and future users;
- have an organizational system that supports not only long-term viability of the repository, but also the digital information for which it has responsibility;
- demonstrate fiscal responsibility and sustainability;
- design its system(s) in accordance with commonly accepted conventions and standards to ensure the ongoing management, access, and security of materials deposited within it;
- establish methodologies for system evaluation that meet community expectations of trustworthiness;
- be depended upon to carry out its long-term responsibilities to depositors and users openly and explicitly;
- have policies, practices, and performance that can be audited and measured; and
- meet the responsibilities detailed in Section 3 of this paper.

Some Examples

Scenario 1: A national library responsible for ensuring long-term accessibility to large, diverse, and growing collections of digital resources, including online publications, complex multimedia products, the digital output of large imaging programs, and a range of special databases. The community it serves is extremely diverse and may be defined as anyone, anywhere, anytime with access to a contemporary, lowest-common-denominator personal computer.

This repository operates as part of a legal deposit environment and is mandated to provide secure, long-term access to materials it accepts as a part of legal deposit. The producer/creator community may include almost anyone: large commercial publishers that already supply print-based resources to the library; new commercial publishers; individuals engaged in vanity publishing; research networks establishing scholarly journals; government agencies; digitization contractors; the institution's own staff; writers depositing papers, including computer files; etc.

The national library is also establishing collaborative distributed archiving of some classes of digital collections, such as online publications, with regional and higher-education libraries and other memory institutions, as well as publishers. Different partners exercise different levels of archiving function and responsibility over varying periods. Policies specifying the nature of responsibilities in the distributed archiving arrangement are in place for all participating institutions.

The digital repository is being built in-house and implemented both selectively and incrementally. Because of the nature of the legal deposit laws, the library must be able to prove the authenticity of the digital materials and therefore strong authenticity and system security will be an early, mandatory component. The library also recognizes a need to take more comprehensive views of digital resources at specified points in their lifecycle.

Scenario 2: A large university library with a growing collection of digital materials to support teaching and research, including online databases, electronic journals, digitized materials, digital output of university staff and students (e.g., theses and dissertations), digital course materials, and institutional records in electronic form. This repository serves primarily the university's faculty, staff, and students, but secondary users include the wider academic community and local community members who purchase library privileges. The university assumes that researchers will have access to the digital collections on campus via the local area network or via the Internet using the high-speed academic network and a personal computer. Descriptive access to the collections is available through the library management system and all access privileges are controlled by a user authentication system in order to meet licensing and intellectual property rights arrangements.

The producer/creator community for the university library, like the national library, may include almost anyone, from an individual scholar to a large commercial publisher. Locally, the university will have some control over the creation of digital materials, but most come from producers over whom the library exercises little or no influence.

The university computer service is contracted to host the digital repository. The system, currently under review, is based on the OAIS Reference Model.

Scenario 3: A museum with a growing collection of digital materials, including surrogates of museum objects, surrogates created for online exhibitions, and original digital art. The museum serves a very diverse community comprising students, researchers, artists, the general public, and organizations seeking digital material for commercial use.

The producer/creator community for the museum, unlike that for the national and university libraries, is generally individual artists over whom the museum has little control. The museum does have control over how digital surrogates are created at its request but not over the creation of others. The responsibility for archiving the digital materials for the long term—regardless of creator—belongs to the museum alone.

The museum uses a content management system to provide day-to-day access to the digital collections, but the system was never intended to facilitate archival storage. Because the museum lacks technical infrastructure and qualified staff, it will contract with a third-party archiving service so that its materials will be professionally

managed, controlled, and backed up to meet its long-term management responsibilities. The commercial service is OAIS-compliant and can provide services from Ingest through Access, however access to materials in the digital repository—including the rights to retrieve, use, alter, or delete items—is restricted to a limited number of museum staff.

Scenario 4: A virtual digital repository for e-journals created through a distributed system of networked computers. Many institutions participate in this repository system, providing local storage and backup for multiple titles and files. Data management is provided and controlled by sophisticated, open-source software developed through the collaboration of several institutions.

Each e-journal title is stored at at least four different geographic locations, greatly reducing the risk of loss due to individual or even multiple server failures. If a server fails or files are corrupted, the validation software automatically detects and repairs or reloads the affected material to the local server, then reconnects it to the networked system. Access to these preserved files is controlled by license and enforced through user authentication software.

With such open-source software, minimal technical administration, and low-cost hardware and software, even smaller libraries can afford to preserve and provide trusted access to their e-journals.

Scenario 5: A small cultural institution with a fairly large digital library collection, which it is legally bound to preserve though it does not yet have the infrastructure, money, or expertise to build full, local digital archiving capacity.

The institution does have expertise in content description and owns and maintains a robust discovery and delivery system that provides access to the mainly image-based materials. This access system, however, was not designed to store and manage the large master files of these images for the long term. The institution decides to use a commercial, third-party service to support the long-term storage and maintenance of the digital files.

After careful investigation, discussion, and analysis, it is determined that the institution and vendor can work together to provide an OAIS-compliant archival system for the master image files, with access mechanisms limited only to file retrieval. Agreements and contracts are drawn up. The institution submits master files and necessary metadata to the vendor in an agreed-upon Submission Information Package. The vendor takes responsibility for the data/bit management aspects of the storage and backup of the files. Changes to files or procedures are undertaken only upon agreement by both parties.

These are just a few examples of the many different approaches that national agencies, research institutions, and cultural organizations will take to establish digital repositories for all or part of their digital collections. The infrastructure of the institution (a large university or national repository versus a small library, archives, or museum) will be a determining factor in the nature of the overall digital repository system, but another factor will be equally important: the repository's "designated community"—its identified group of potential users—will determine what is deposited (content and format), how the digital information is managed and preserved, and how it is disseminated and accessed. Despite their different organizational models, all digital repositories will need to address the same underlying issues of not only functionality, but of reliability.

Establishing Trust

The 1996 RLG/CPA report made a clear statement about trust in digital archives:

For assuring the longevity of information, perhaps the most important role in the operation of a digital archive is managing the identity, integrity and quality of the archives itself as a trusted source of the cultural record. Users of archived information in electronic form and of archival services relating to that information need to have assurance that a digital archives is what it says it is and that the information stored there is safe for the long term.⁵

Although much has been accomplished globally to move the digital archiving agenda forward in the years since this report, there is as yet no collective agreement on a more exact definition of “trusted archives.”

The archival and computer professions promulgate a host of concepts and terms that lay a foundation for establishing the defining characteristics of dependable digital archiving repositories. Commonly used terms such as “reliable,” “responsible,” “trustworthy,” and “authentic” help to define the nature of the archival enterprise and its myriad relationships with those creating, managing, and using digital objects. Computer scientists worldwide have grappled with definitions and performance measures of trusted military systems for nearly 20 years (see Appendix B, The Evolution of “Trust” in Computing Systems). Likewise, the airlines industry has required trustworthy, responsible, and authentic systems. In the last decade, groundbreaking work by archivists in Australia, North America, and Europe has resulted in fundamentally new approaches and tools that specify the nature and performance of accountable record-keeping systems. And in the past few years digital library experts have contributed their experience to a growing body of literature and applications pertaining to the construction and maintenance of secure systems accommodating large quantities of digital resources. But what is meant by the phrases “trusted archives” or “trusted repository”?

“Trust” is defined by the *Merriam Webster Dictionary* as

assured reliance on the character, ability, strength, or truth of someone or something . . . one in which confidence is placed . . . a charge or duty imposed in faith or confidence or as a condition of some relationship . . . something committed or entrusted to one to be used or cared for in the interest of another.⁶

By this definition, most cultural institutions are already trusted. Libraries, archives, and museums are entrusted with the materials and objects that document our cultural heritage. They are trusted to store these valuable materials. They are trusted to provide access to them in order to document and reveal history as well as to foster the growth of knowledge. They are trusted to preserve these items to the best of their ability for future generations.

Cultural institutions have excelled in preserving large amounts of cultural heritage in the form of physical objects. In many cases, these physical objects or documented, reliable surrogates are available to patrons as “proof” of the institution’s capability to collect and preserve for the long term. But since digital information is less tangible and much more mutable than other materials, trust and reliability may be more difficult to prove.

⁵ Ibid.

⁶ Merriam Webster Web site, www.m-w.com.

By its very nature, digital information can be transitory and difficult to preserve; certainly the traditional methods of preservation are less applicable. The digital landscape—and the digital management landscape—are also quickly evolving with the exponential growth of digital information. No one institution will be able to preserve it all—in fact, it will take many institutions, organizations, businesses, and others to preserve the cultural information that had previously been placed solely in the care of cultural institutions. How can the average public citizen trust that digital material is being preserved? And if the cultural institution cannot create and manage a local digital repository, can this trusted institution trust a third-party provider?

Minimally, three levels of trust apply to the establishment of trusted digital repositories:

1. How cultural institutions earn the trust of their designated communities.
2. How cultural institutions trust third-party providers.
3. How users trust the documents provided to them by a repository

Earning the Trust of Designated Communities

Institutions responsible for the preservation of nondigital material already tend to enjoy a fairly high level of public trust because libraries have reliably preserved a large amount of the human record over time.⁷ While the challenges presented in the preservation of digital information are much different and require new solutions, the public will likely have at least some trust that cultural institutions will succeed based on the success of the past.

Thus far, libraries, archives, and museums have shown they can create and provide access to digital materials. Users now rely on institutions to provide ongoing development of systems that support long-term access to the materials. Over time, institutions will keep the users' trust so long as they sustain reliable access to information.

Trusting Third-Party Providers

Service providers gain the trust of cultural institutions through a combination of proven reliability, fulfillment of contractual responsibilities, and demonstrated sensitivity to community issues. These attributes are easily measured, however institutions are reluctant to engage the third-party digital service providers that have not proven their reliability and without demonstrated experience, the service provider cannot prove reliability. To resolve the tension between a repository's appropriately high standards and the attempts to meet the challenge, a combination of repository attributes and other criteria must be identified to foster interaction and begin to lay the foundation for trust between cultural institutions and third-party providers.

One option may be to identify certain attributes in a third-party service provider that the institution requires of itself. In an interview published in 2000 in *RLG DigiNews*, Kevin Guthrie (president, JSTOR) opened his response to a question about JSTOR's plans for creating a trusted repository in this way:

Well, trust in this context is both very important and rather difficult to define. It is important because the goal is to be able to establish a relationship whereby a library can rely on a third party to provide a service that has been a

⁷ There are controversial exceptions, for example, Nicholson Baker, *Double Fold: Libraries and the Assault on Paper* (New York: Random House, 2000).

core function of a library; that is, archiving. That is no small responsibility and any enterprise that aims to provide such a service is going to have to earn a very high level of trust.

At this early juncture, I don't think there exists a standard definition or "litmus test" of what it will take to become such a trusted archives. I do think that the mission of the enterprise is a fundamental component of the assessment.⁸

Guthrie's statement suggests that libraries, museums, and archives would be more likely to extend a provisional trust to a third-party repository that shares a similar mission (e.g., nonprofit, for scholarly use). A large cultural institution that builds a repository to both meet its own needs and serve as a third-party digital repository may be perceived as inherently trustworthy, at least initially.

In the absence of trusted repositories or reliable, proven practices, a program for certification could provide a basis for trustworthiness. Certification would dictate criteria that must be met and would employ mechanisms for assessment and measurement (see Section 4, Certification). Cultural institutions would have a tool for measuring currently available services, and repositories would have a known set of best practices or standards to meet in order to gain the business of cultural institutions. Certification, achieved on a periodic basis for several years, could resolve tension between the immediate need for trusted archives and the need to develop and prove reliability over time. Once a number of organizations have achieved a solid reputation and gained the trust of the community, the need for a certification process may fade.

Trusting Digital Documents

A user must be able to trust digital documents provided by digital repositories. Authentication of digital information includes the ability to detect change to a digital document. Peter Graham has said:

The great asset of digital information is also its greatest liability: the ease with which an identical copy can be quickly and flawlessly made is paralleled by the ease with which a change may undetectably be made.⁹

Both intentional and accidental changes have ramifications. As Clifford Lynch has noted,

It is very easy to replace an electronic dataset with an updated copy, and . . . the replacement can have wide-reaching effects. The processes of authorship . . . produce different versions which in an electronic environment can easily go into broad circulation; if each draft is not carefully labeled and dated it is difficult to tell which draft one is looking at or whether one has the "final" version of a work.¹⁰

How can a user be certain that the document received is the one requested? How can the document be verified to be the exact item deposited into the digital repository in the past?

⁸ "Developing a Digital Preservation Strategy for JSTOR, an interview with Kevin Guthrie," *RLG DigiNews* 4, no. 4 (August 15, 2000) www.rlg.org/preserv/diginews/diginews4-4.html#feature1.

⁹ Peter Graham, "Issues in Digital Archiving," in *Preservation: Issues and Planning*, eds. Roberta Pilette and Paul Banks (Chicago, IL: American Library Association, 2000): 101.

¹⁰ Clifford Lynch, "Accessibility and Integrity of Networked Information Collections" (Office of Technology Assessment, Congress of the United States, July 5, 1993): 68.

Henry Gladney recently asserted that the data integrity and information trustworthiness of digital documents are somewhat easily addressed: "Proof of digital document authenticity and provenance can be achieved by message authentication codes signed by trusted institutions."¹¹ Mechanisms such as checksums (or similar error-detection techniques) and public key encryption systems already support authentication for many communities, including e-commerce. The technologies underpinning authentication will continue to evolve to meet broad, cross-community demand and the cultural community will be able to take advantage of future development.

If a trusted repository requires an authentication system, then institutions must include in all digital documents the information that interacts with the system. At this time, almost all extant or emerging metadata standards and best practices include a placeholder for authentication information; unfortunately, they do not all require it. This should change.

Recommendation

Investigate and define the minimal-level metadata required to manage digital information for the long term. Develop tools to automatically generate and/or extract as much of the required metadata as possible.

¹¹ Henry Gladney, *Digital Document Quarterly* 1, no. 1 (1Q2002) home.pacbell.net/hgladney/ddq_1_1.htm.

2 Attributes of a Trusted Digital Repository

The attributes of a trusted, reliable digital repository need to be identified. A framework of attributes must accommodate all different situations and institutional responsibilities while providing a basis for expectations of a trusted repository. The following list reflects the emerging expert community's thinking about such attributes:

- Compliance with the *Reference Model for an Open Archival Information System (OAIS)*
- Administrative responsibility
- Organizational viability
- Financial sustainability
- Technological and procedural suitability
- System security
- Procedural accountability

Compliance with the *Reference Model for an Open Archival Information System (OAIS)*

A trusted digital repository will make sure the overall repository system conforms to the OAIS Reference Model. Effective digital archiving services will rely on a shared understanding across the necessary range of stakeholders of what is to be achieved and how it will be done. The Reference Model supplies a common framework, including terminology and concepts, for describing and comparing architectures and operations of digital archives. As well, the OAIS provides both a **functional model**—the specific tasks performed by the repository such as storage or access—and a corresponding **information model** that includes a model for the creation of metadata to support long-term maintenance and access. Organizations and institutions building digital repositories should commit to understanding these models and make sure all aspects of the overall system conform.

Administrative Responsibility

A trusted digital repository will provide evidence that it has a fundamental commitment to implementing the range of community-agreed standards and best practices that affect its operations—particularly those that directly influence its viability and sustainability. Administrative responsibility extends to meeting appropriate national and/or international standards for the physical environment, backup and recovery procedures, and security systems. The trusted repository will meet or exceed community standards for performance and will collect and share data measurements routinely with depositors. It will involve external community experts in validating and/or certifying its processes and procedures on a regular schedule. Written agreements with depositors will address all appropriate aspects of acquisition, maintenance, access, and withdrawal. Further, ongoing risk management and contingency planning will play a routine part of the organization's annual strategic planning

activities. A reliable repository will commit itself to transparency and accountability in all its actions.

Organizational Viability

Organizations choosing to become trusted digital repositories will establish themselves in ways that demonstrate their viability. Their mission statements will reflect a commitment to the long-term retention, management of, and access to digital cultural assets on behalf of depositors and users. Their legal status and standing will be appropriate to the range of responsibilities they are undertaking. Their business practices will be transparent and forthright. Staffing levels and areas of expertise will be appropriate to the work undertaken; further, staff training and professional development opportunities, including conference attendance and participation, will be given priority to ensure the currency of staff skill sets. The repository will continually review its policies and procedures to ensure that appropriate growth can occur and that new processes and procedures are tested for scalability. A formal succession plan or escrow arrangements will be developed in consultation with community experts, depositors, and peer organizations that identifies all relevant content and designates trusted inheritors should the repository cease to exist.

Financial Sustainability

A trusted digital repository should be able to prove its financial sustainability over time. Overall, trusted repositories will adhere to all good business practices and should have a sustainable business plan in place. Normal business and financial fitness should be reviewed at least annually. Standard accounting procedures should be used. Both short- and long-term financial planning cycles should demonstrate an ongoing commitment to a balance of risk, benefit, investment, and expenditure. Operating budgets and reserves should be adequate.

Technological and Procedural Suitability

Community experts currently advocate a range of preservation strategies. A trusted digital repository will consider all relevant options and will communicate openly about the suitability of various strategies. It will ensure that it has in place all appropriate hardware and software for inventory management functions, including all the forms of acquisition, storage, and access it offers. The repository will also have policies and plans for replacing technology as needed. The repository will comply with all relevant standards and best practices, ensuring that staff have adequate expertise to understand and implement them. It will also undergo regular external audits on its system components and performance.

System Security

All systems used in the operation of a trusted digital repository will be designed to assure the security of the digital assets. Policies and practices will meet community requirements, particularly those pertaining to copying processes, required redundancy of data, authentication systems, firewalls, and backup systems. The repository will have written policies and plans for disaster preparedness, response, and recovery, and staff will be trained appropriately. Special attention will be given to processes that address data integrity to avoid

loss of data, detect changes in data, and restore lost or corrupted data. Any detected changes (including loss or corruption and restoration) will be documented and the depositor will be notified both of the changes and any actions taken.

Procedural Accountability

A trusted digital repository is responsible for a range of interrelated tasks and functions (see Section 3, Responsibilities); it will therefore be accountable for all relevant policies and procedures. Repository practices will be documented and made available on request. Monitoring mechanisms that measure and ensure the continued operation of all systems and procedures will be in place. Preservation strategies undertaken (e.g., migration, emulation, etc.) will be recorded and justified in the context of community-wide best practices. Feedback mechanisms will be in place to support the resolution of problems and to negotiate the evolving requirements between the repository, any third-party service providers, and the designated communities.

3 Responsibilities of a Trusted Digital Repository

High-Level Organizational and Curatorial Responsibilities

Research repositories need to understand fully what responsibilities they should assume for the preservation of digital materials. Organizational responsibility must be understood at three basic levels. Organizations must first understand their own local requirements. Second, they need to understand which other organizations might share some of the responsibilities through geography or arrangements such as consortial agreements or shared user communities, disciplines, or format of materials. Third, they need to understand which responsibilities can be shared and how. Assuming that the general model for digital repositories is more or less distributed, its success relies on shared understanding across the federation or network of repositories of their respective duties and roles. Comprehensive coverage within the collections and effective interoperability across repositories will rely on such understandings.

Although a detailed discussion is beyond the scope of this report, a summary of these major factors is useful:¹²

- the scope of collections;
- preservation and lifecycle management;
- the wide range of stakeholders;
- ownership of material and other legal issues; and
- cost implications.

The Scope of Collections

Digital materials for libraries and archives range from simple (e.g., text-based) digital files to complex multimedia and database resources. The sheer variety of digital materials and the role that they play in the collection make development and application of collections policies very challenging. The existence or lack of a physical equivalent influences decisions about whether and how the digital resource is preserved. For materials that have a physical counterpart, preservation decisions take into account considerations such as the condition of the original materials and the reason for digitizing (e.g., for increased access to the materials). Materials that are “born digital” can present more challenging problems because their “being digital” is not only a method of access, it represents their value as an information artifact. For many born-digital resources, effective preservation will rely as much on preservation of the object’s digital characteristics or properties as on preservation of its basic intellectual content. More importantly, when a library or archives digitizes its own collections, it can control decisions about standards, formats, quality control, and documentation. The preservation of materials generated outside may not include this degree of control.

¹² For a more complete discussion of the roles and responsibilities of different stakeholders in the lifecycle of digital materials, see Neil Beagrie and Daniel Greenstein, *A Strategic Policy Framework for Creating and Preserving Digital Collections*. Version 5.0 (Arts and Humanities Data Service Executive, 1998, updated July 2001) ahds.ac.uk/strategic.pdf.

Preservation and Lifecycle Management

Preservation decisions for digital items cannot wait until continued use of the materials has proved they are worth keeping. Postponing preservation decisions can and most often will result in preservation actions that are more complex, more labor intensive, and more costly. A resource can even be held hostage by an obsolete piece of software. It is also important to accept the fact that digital information is more transitory and mutable, so it may not survive benign neglect. Preservation requires active management that begins at the creation of the material and depends on a proactive approach by digital repositories and the cooperation of the stakeholders, including data providers.¹³

The Wide Range of Stakeholders

Content creators, systems developers, custodians, and future users are all potential stakeholders in the preservation of digital materials, and this complicates the determination of responsibilities—who, when, and for how long. Often, those creating digital materials or designing digital content management systems do not take great interest in their long-term preservation. For example, commercial publishers are justifiably interested in the preservation of their materials only as long as they are commercially viable, while libraries and their users are often interested in continued access to materials long after they cease to turn a profit. Similarly, for archives, it is usually when an electronic records management system is designed (well before any records are created) that key decisions are made that affect the long-term preservation of the records themselves. In both cases, decisions about how the materials are handled when created or maintained determine how or whether the repository can preserve them.

Ownership of Material and Other Legal Issues

Responsibility for preservation has traditionally been considered alongside ownership of the materials; that is, the owner of the materials was responsible for determining their life span. However, ownership of digital materials is not often straightforward. While a book can be taken into the collection and set upon the shelf, digital materials are less tangible. For a growing number of digital materials considered “integral” to research collections and archives, access is provided through licensing arrangements—often through a regional or national consortium. Licensing arrangements can apply to either the digital content itself or to software necessary for specific functionality and access to the content. Although the organization may own the right to access material or use the software for a specified period, there is often no guarantee of rights beyond the terms of the license. While commercial publishers are beginning to provide some guarantee of continuing access, most licensing agreements are still perilously vague about how the digital repository will be maintained and how long-term access will be ensured. Reliance solely on creators or producers of digital materials for long-term preservation is potentially risky, not least because digital resources are not generally created or engineered with long-term preservation in mind.

It will be critical in the future for research repositories to work as closely as possible with content creators to ensure that long-term preservation responsibilities are clearly understood and documented in licensing agreements; this is currently being explored by The Andrew W. Mellon Foundation’s e-Journal archiving program.¹⁴ It will require increased cooperation and

¹³ Ibid.

¹⁴ A program funded by The Andrew W. Mellon Foundation designed to plan the development of e-journal repositories meeting specific requirements developed by the Digital Library Federation (DLF), Council for

effective communications with publishers, software suppliers, and other producers to ensure that what is deposited is a copy of the data object in the format most suitable for preserving the materials over the long term. Understanding the important difference between long-term preservation and short-term access—particularly while materials are still commercially viable—is critical. Libraries may require different license arrangements for long-term preservation than for end-user access.

Often, rights that relate to the software and systems used to create the material impinge on its preservation. Very little, if any work, has been done with software vendors to raise awareness about the longevity of their materials in the interests of future scholarship and research.

Digital preservation has even wider legal implications. How preservation infringes on copyright remains unclear. For example, the content creator does not usually own the rights to the software and systems used to create the digital file. This raises legal issues when access or changes to those systems are necessary. In such cases, at best, a repository will need to arrange separate rights clearance for long-term maintenance; at worst, preservation will be compromised because rights clearances for access cannot be obtained. Some work has been done on the establishment of repositories for software to help address these concerns, however the research repository community will need to make an appeal to have this conflict taken into consideration in the creation or renewal of national deposit legislation.

Recommendation

Investigate and disseminate information about the complex relationship between digital preservation and intellectual property rights.

Cost Implications

Although not a great deal is known about the costs of preserving complex digital objects over time, there is an accepted wisdom in the library community that digital preservation will require ongoing resource commitments—potentially more than for traditional materials. Traditional and digital preservation should be compared with caution, because the complex dependencies between long-term maintenance and continuing access make comparison problematic. Indeed, for digital materials that have no analog equivalent, comparison is meaningless. Although it may be too early to compare the costs of digital and traditional preservation meaningfully, one thing is certain: preserving digital materials will require resource commitments over time. While traditional materials have ongoing costs for stable storage environments, digital materials will also require periodic analysis and the application of new technical strategies to ensure continuing access. Digital preservation is also likely to draw on resources longer than traditional preservation does, and it may be the case that different technical strategies (e.g., different types of migration or emulation) will require different costing timeframes and schedules.

In a recent publication entitled *Preservation Management of Digital Materials*, Jones and Beagrie suggest that digital preservation costs are based on four interrelated factors:

Library and Information Resources (CLIR), and Coalition for Networked Information (CNI). Seven major libraries have received grants, including the New York Public Library and the university libraries of Cornell, Harvard, MIT, Pennsylvania, Stanford, and Yale; see www.clir.org/diglib/preserve/ejp.htm.

- The need to actively manage inevitable changes in technology at regular intervals and over a (potentially) infinite time frame.
- The lack of standardization in both the resources themselves and the licensing agreements with publishers and other data producers, making economies of scale difficult to achieve.
- The as yet unresolved means of reliably and accurately rendering certain digital publications so that they do not lose essential information after technology changes.
- That for some time to come digital preservation may be an additional cost on top of the costs for traditional collections, unless cost savings can be realized.¹⁵

Digital technologies and applications shift rapidly; strategies to preserve objects resulting from new approaches must keep pace. It is therefore difficult, if not impossible, to establish concrete costs for all associated activities. Further, the uncertainty of the financial commitment represented by digital preservation makes assuming preservation responsibilities more complex. Work can be done to understand how costs will play out and where savings can be made, but the preservation of digital materials cannot wait for exact information because it may never appear. What will be important is an understanding of where the main costs are likely to fall and how, within existing practices, these can be incorporated to achieve economies of scale. In addition, the ways that other stakeholders (e.g., content providers) can decrease costs (e.g., by changing the way they supply materials for deposit—in specific formats or with better descriptions) should be explored. Repositories that currently provide guidelines for depositors include the Cornell University Library, Arts and Humanities Data Service (AHDS), and the National Library of Australia.¹⁶ Repositories will also need to understand more about the advantages of collaborative approaches. The costs may be easier to absorb across multiple institutions.

Research repositories need to begin work now—“this is likely to be a better strategy than only discussing and studying the problem.”¹⁷ Integrating digital preservation into the everyday management and organization of the library or archives will help ensure that the necessary skills and knowledge are embedded within the organization for the earliest and most effective savings or economies of scale.

Recommendation

Research and develop models for cooperative repository networks and services.

¹⁵ Neil Beagrie and Maggie Jones, *Preservation Management of Digital Materials: A Handbook* (2001) www.jisc.ac.uk/dner/preservation/workbook/.

¹⁶ *Report of the Digital Preservation Policy Working Group on Establishing a Central Depository for Preserving Digital Image Collections*, Version 1.0 (March 2001) www.library.cornell.edu/imls/image%20deposit%20guidelines.pdf; AHDS provides guidelines for each of its service providers in Visual Arts, Performing Arts, Electronic Texts, History, and Archaeology, www.ahds.ac.uk/dephow.htm; National Library of Australia, *Safeguarding Australia's Web Resources: Guidelines for Creators and Publishers* (2000) www.nla.gov.au/guidelines/2000/webresources.html.

¹⁷ Johan Steenbakkens, *The NEDLIB Guidelines: Setting up a Deposit System for Electronic Publications*, NEDLIB Report Series, report 5 (NEDLIB Consortium, 2000), available from www.kb.nl/coop/nedlib/.

Operational Responsibilities

The CPA/RLG report recommended “a dialogue among the appropriate organizations and individuals on the standards, criteria and mechanisms needed to certify repositories of digital information as archives.”¹⁸ The OAIS Reference Model is a useful framework, lending depth and breadth to the definition of a trusted digital repository and its associated attributes. The model is also the framework for identifying the responsibilities of a trusted digital repository.

The following list of responsibilities is taken from work done by the OAIS community to define the principle obligations of an OAIS-compliant repository, whether it is a component of an institution’s overall digital library system or a third-party archiving service. One addition to the list acknowledges the repository’s critical role in the promotion of standards in the area. When an institution and third-party service work together, the institution would take responsibility for selection, curatorial review, preparation/verification of metadata, and determination of value to a designated community while the third-party service would handle system architecture, file management, authentication and validation mechanisms, etc.

A reliable digital repository:

- negotiates for and accepts appropriate information from information producers and rights holders;
- obtains sufficient control of the information provided to support long-term preservation;
- determines, either by itself or with others, the users that make up its designated community, which should be able to understand the information provided;
- ensures that the information to be preserved is “independently understandable” to the designated community; that is, that the community can understand the information without needing the assistance of experts;
- follows documented policies and procedures that ensure the information is preserved against all reasonable contingencies and enables the information to be disseminated as authenticated copies of the original or as traceable to the original;
- makes the preserved information available to the designated community; and
- works closely with the repository’s designated community to advocate the use of good and (where possible) standard practice in the creation of digital resources; this may include an outreach program for potential depositors.

Negotiating for Appropriate Information from Content Providers

This responsibility refers to all transactions between the content providers (including creator, rights holders, etc.) and the repository prior to formal submission into the repository.¹⁹ The nature of these interactions is largely determined by the control or influence the repository

¹⁸ John Garrett and Donald Waters, *Preserving Digital Information: Report of the Task Force on Archiving of Digital Information* (Commission on Preservation and Access and RLG, 1996) www.rlg.org/ArchTF/index.html.

¹⁹ In the OAIS Reference Model, submission of materials into the repository is referred to as the Ingest function; see Appendix A.

has over the creation or submission of the resources. Some repositories require content creators to conform to specific restrictions on the type or format of materials, while others exert little or no influence over the creation of materials. In most cases, some activity to prepare the materials for submission (e.g., creation of metadata and documentation) will be required. In a few cases, “aggressive rescue” or salvage measures will need to be taken. Negotiations would cover:

- Legal issues, involving all negotiations concerning copyright and other rights (e.g., privacy, donor restrictions) and appropriate clearance for long-term maintenance and continuing access, as well as short-term or immediate access in some cases. Separate negotiations may be necessary for current access and for long-term preservation. Long-term preservation may be jeopardized if submission of materials is based solely on negotiations for current access, which can be more restricted to protect the commercial interests of the content provider.
- Preservation metadata based on agreed-on specifications. A repository has to ensure that agreements are in place with content providers about the bibliographic and technical metadata that accompany the submitted materials.²⁰
- Authenticity checks, confirming that the digital materials submitted to the repository are exactly what the content provider intended.
- Record keeping, with adequate documentation for the transactions between the repository and content provider.

Fulfilling this responsibility requires:

- Well-documented and agreed-on policies about what is selected for deposit, including, where appropriate, specific required formats.
- Effective procedures and workflows for obtaining copyright clearance for both short-term and immediate access, as necessary, and preservation.
- A comprehensive metadata specification and agreed-on standards for its implementation. This is critical for federated or networked repositories and includes standards for the provision of rights metadata from content providers and for representing technical and administrative metadata.²¹
- Procedures and systems for ensuring the authenticity of submitted materials.
- Initial assessment of the completeness of the submission.
- Effective record keeping of all transactions, including ongoing relationships, with content providers.

²⁰ In the OAIS Reference Model, bibliographic metadata is referred to as Preservation Description Information (PDI) and technical metadata is referred to as Representation Information (RI); see Appendix A.

²¹ Work already done in this area will provide a useful starting point. ONIX International, an initiative involving key players in the publishing industry, is developing a standard for the exchange of metadata for publishing, including rights metadata (www.editeur.org/onix.html). The National Information Standards Organization has created a *Draft Standard for Trial Use* specifically addressing *Technical Metadata for Digital Still Images* (www.niso.org/committees/committee_au.html) and the OCLC/RLG Preservation Metadata Working Group has developed a comprehensive preservation metadata framework applicable to a broad range of digital preservation activities (www.oclc.org/research/pmwg/).

Obtaining Sufficient Control of the Information

This responsibility refers to activities following the submission or identification of material and involves preparation of the object for storage in the repository, including:

- **Analysis of the digital content:** At this point, the repository must, in consultation with the depositor/rights owner and systems managers, assess the digital object and determine which of its properties are significant for preservation. Rights clearances that have been obtained for copyrighted materials influence this analysis, as do collections management policies and documented collection strengths. This is the moment to apply policies about what formats are acceptable; any necessary migrations should have taken place.²² Once decisions have been made about what will be preserved, this assessment should be automated as much as possible and the content analyzed systematically as material is deposited into the repository.

Assessing an Object's Significant Properties

Recent work has shown that before digital materials are placed into a repository, the owning institution needs to decide what level of preservation is appropriate for each digital object or class of objects. Decisions about an object's "significant properties" influence decisions about the level and method of access as well as the level of preservation metadata required for long-term maintenance. Cedars and the National Library of Australia have explored the concept of significant properties, or an object's "preservable essence."²³ Particularly for repositories with a wide variety of digital materials, how decisions are made about significant properties will be important.

For traditional materials, access and preservation are often one and the same and are generally done by the same organization. A set of papers will, in all likelihood, continue to be readable and therefore useable. However, for digital materials, simply maintaining a bytestream does not necessarily ensure the digital material will be preserved at an acceptable level to both the repository and the users. For digital materials, the level of "access" in the technical sense depends on judgments made by the collection manager and it should reflect the repository's collections management policies.

A digital object's significant properties are not absolute, nor are they static. How a repository determines which properties are significant depends on the nature of the organization, the services it provides, and its role in preservation. The collection manager judges the appropriate levels of preservation and access to fulfill the repository's responsibilities and meet the needs of its stakeholders. Materials deemed to be part of the collection's core might retain all of their original functionality, while more peripheral materials might not include the full complement of these functions or properties. A library may have a choice of formats and associated functionality for a digital object (e.g., PDF, HTML, XML, or SGML for an electronic journal) and the authority to choose what to preserve.

²² A repository may require the depositor to migrate material to an acceptable format for submission, or for accuracy and reliability the repository may choose to perform migrations itself.

²³ The Cedars Project Report, April 1998–March 2001, www.leeds.ac.uk/cedars/OurPublications/cedarsrepmar01exec.html.

However, to ensure the authenticity of an electronic record in a legal context, archives may face more restricted choices; for example, they may not be allowed to make any functional changes to the electronic record. Since collections management policies may change over time, repositories may need to reevaluate and reassess their digital objects' significant properties.

The significant properties of a digital object (i.e., the acceptable level of functionality) dictate the underlying technical form that needs to be documented and supported to ensure preservation of those properties and the amount of metadata, including detailed technical metadata, that must be stored alongside the bytestream to ensure the object is accessible to the agreed-on level. Naturally, decisions about the significant properties have important resource implications: the more significant properties deemed to be necessary, the more associated metadata will be required. The creation and maintenance of the detailed metadata associated with the object's significant properties are critical to the repository's preservation function—the detailed descriptions and the technical information necessary for rendering the bytestream into a meaningful digital object ensure *long-term* preservation.²⁴ How continuing access is provided over time can and should be kept separate, conceptually, from this basic preservation function.

Significant properties—a simple example: A repository decides that the only significant property of an electronic journal published on the Web is the text within the journal, not its layout and formatting. There is no need to store information about the HTML environment, but only to include information about retrieving or rendering an ASCII text file.

Significant properties—a more complex example: An electronic journal that is published on the Web in HTML format includes a database that provides access to the original research data. Although end users currently access the journal in HTML, these pages are created on the fly from SGML. For archiving, the repository takes the SGML files and decides that the significant properties include the hypertext links (internal) as well as any multimedia functions (e.g., sound and video clips) and the functionality of the database, so the object is to be preserved at full functionality. Therefore, the required technical metadata includes robust technical descriptions of the objects, including the SGML DTD and other information about the systems and the software necessary to run the video and sound clips, as well as information about the database (which may or may not use standard SQL), and, finally, the arguably less complex technical metadata about retrieving the text and images.

The significant properties of digital materials need to be determined by policy makers, which a large repository cannot possibly manage object by object. Policies will need to apply to different classes of object and will need to be systematized

²⁴ In the OAIS Reference Model, the detailed pathways that render a preserved bytestream into a meaningful digital object are called Representation Networks; see Appendix A. Representation Information provides the technical means to render the digital objects and also provides detailed human-readable descriptions of the technical environment as necessary. It is the key to the preservation function with OAIS. Cedars implemented the OAIS Model so that Representation Information can take the form of a reference to existing information in the network for other objects already in the archive using the same technology. The Representation Information is therefore not duplicated and the archive saves resources.

and automated. Further, the creation of detailed technical metadata to support significant properties may be beyond the human resources of most research libraries and archives. For this reason, work with software suppliers and systems designers could play a key role: it would benefit many libraries and archives to develop digital repository management systems that provide for the automatic generation of technical metadata for materials submitted to the repository.²⁵ Likewise, repository management systems could then be designed to incorporate standard technical metadata as it is submitted alongside the digital object by the content provider, supporting important collaboration with content providers and creators.²⁶ In this sense, digital repositories can take advantage of the fact that technology is still evolving to influence developments to better accommodate long-term preservation.

Recommendation

Research and create tools to identify the attributes of digital materials that must be preserved.

- **Continuing access arrangements:** A repository needs to choose a strategy for continuing access, which will need to be reevaluated regularly as technology changes. For example, if an object relies on a complex technical environment or uses proprietary technology, an emulation of that environment might be desirable either now or in the future, which affects the level of technical metadata required. Indeed, it may be that the repository stores an emulator, in which case some standards for the development of archival-quality emulators will be necessary.²⁷
- **Verification of metadata:** Any metadata that accompanies the object when it is submitted to the repository must be verified and, as necessary, enhanced to support the object's long-term maintenance as well as continuing access.
- **Unique and persistent identification of materials:** Much work has been done on the need for unique and *persistent* naming. Nowhere is this more relevant than in long-term preservation of digital materials. A repository needs to ensure that an accepted, standard naming convention is in place that identifies its materials uniquely and persistently for use both in and outside the repository. Equally important is a system of reliable linking/resolution services in order to find the uniquely named object, no matter its physical location. In a distributed model, it is particularly critical that the participating organizations agree to standard naming conventions and resolution options.

²⁵ Cedars has done some preliminary work on providing a network of Representation Information as part of the archive's function. During the Web-based ingest process, content providers simply choose the type of digital file they are submitting from a drop-down menu. If the type of file is not included in the drop-down menu, a request is forwarded to the archive administrator and the necessary Representation Information is created. In addition, there may be some scope for applying the work done on "canonicalization" of digital materials to assist in the automatic creation of Representation Information; see Clifford Lynch, "Canonicalization: A Fundamental Tool to Facilitate Preservation and Management of Digital Information," *D-Lib Magazine* 5, no. 9 (September 1999) www.dlib.org/dlib/september99/09lynch.html.

²⁶ NISO (National Information Standards Organization) is working on this specifically for digital still images. Such work might prove applicable to a wider range of digital materials.

Recommendation

Design and develop systems for the unique, persistent identification of digital objects that expressly support long-term preservation.

- **Creation of the Archival Information Package:**²⁸ Digital repositories can store a digital object and its associated metadata in two ways: as a single bytestream or separately. For practical reasons, repositories may prefer to store the digital object within the repository and provide only pointers or references to the associated metadata in other systems, such as bibliographic data stored in the library management system. Such “virtual encapsulation” avoids duplicating metadata but separating a digital object and its metadata may present problems in the future. Some experts feel that long-term preservation may be best served by storing the digital content and all of its relevant metadata as a single file. The Metadata Encoding and Transmission Standard (METS) for encoding descriptive, administrative, and structural metadata addresses this need. Depending on its use, a METS document could be used in the role of Submission Information Package (SIP), Archival Information Package (AIP), or Dissemination Information Package (DIP) (see Appendix A).
- **Authentication and integrity checking:** The repository needs to ensure that mechanisms are in place for verifying the digital object, including all associated metadata. The repository should verify that the digital object can be rendered (or at least traced) from the encapsulation back into its original form as it was submitted to the repository. This should include verifying not only the integrity of the bytestream but also confirming the object’s usability and functionality.
- **Archival storage:** Whether archival storage is centralized or distributed, it relies on a robust and well-documented policy for storage and maintenance and for the expected level of service. For archival storage by a third party, service-level agreements are essential. The policy must include systems for routine integrity checking of the bytestream, once it has been established within the storage facility, redundancy of data storage, and for disaster preparedness, response, and recovery.

Fulfilling this responsibility requires:

- Detailed analysis of an object or class of objects to assess its significant properties. Analysis should be automated as much as possible and informed by the collections management policy, rights clearances, the designated community’s knowledge base, and policy restrictions on specific file formats.
- Verification and creation of bibliographic and technical metadata and documentation to support the long-term preservation of the digital object according to its significant properties and underlying technology or abstract form, with monitoring and updating of metadata as necessary to reflect changes in

²⁷ CAMiLEON is investigating the use of emulators for continuing access to digital materials. Although as yet inconclusive, the work will eventually result in recommendations for appropriate technology for preservation-quality emulators; see the CAMiLEON Web site: www.si.umich.edu/CAMILEON/.

²⁸ See Appendix A.

technology or access arrangements. This involves understanding how strategies for continuing access, such as migration and emulation, influence the creation of preservation metadata.

- A robust system of unique identification.
- A reliable method for encapsulating the digital object with its metadata in the repository.
- A reliable archival storage facility, including an ongoing program of media refreshment; a program of monitoring media; geographically distributed backup systems; routine authenticity and integrity checking of the stored object; disaster preparedness; response and recovery policies and procedures; and security.

Determining the Repository's Designated Community

Preservation takes place for the designated community: whether to preserve an object or class of objects is initially determined by how the repository's designated community values its content. Likewise, the creation of the technical infrastructure to ensure access to the object depends entirely on the community's technical capability or knowledge base, which determines the minimum level of associated technical metadata for long-term access. The knowledge base may not be the users' technical expertise but the technical capability they have either as actual technical knowledge or through access systems.

Traditionally, knowledge of a library's designated community was gleaned through face-to-face interaction. Generally, the user community was assumed to fit within a broadly defined research or academic community. The user communities of digital repositories, however, are much more difficult to discern. Not only must repositories meet the needs of current users, long-term archiving is meant to benefit future generations—perhaps hundreds of years from now. Anticipating these users' needs or their technical means is even more difficult. Any technical infrastructure that meets the needs of current users must not preclude other possible or different uses and users in the future.

Fulfilling this responsibility requires:

- Analysis and documentation of the repository's current designated community as well as the possible needs and modes of access of future users.
- For federated or cooperating repositories, a shared understanding of the designated communities that are to be served.

Ensuring the Information to Be Preserved Is Independently Understandable to the Designated Community

“Independently understandable information” is information that the designated community can understand without the assistance of experts.²⁹ Making digital information independently understandable poses formidable challenges because it is not itself humanly readable at any level—it relies on further digital information to make it meaningful. However, it is possible

²⁹ The term “independently understandable” is defined by OAIS as “a characteristic of information that has sufficient documentation to allow the information to be understood and used by the designated community without having to resort to special resources not widely available, including named individuals.” (Committee for Space Data Systems, *Reference Model for an Open Archival Information System (OAIS) Red Book*, February 2001)

to stipulate that the technical metadata required for rendering binary data into meaningful digital objects correspond to the lowest common level of technical knowledge or capability in the designated community. For example, in the current technical climate, one might assume that any user could use a Web browser and an HTML file.

Access and dissemination of objects from the repository will also need to reflect changes in both the technology and in the knowledge base of the designated community. It may be necessary to provide different migrated versions of objects as technologies change; whether this is also reflected in changes to the digital object and technical metadata will be determined by the organization's policies. Many repositories may change their continuing access methods without changing the stored object itself. In other words, on-the-fly migration may be provided for materials just for access, without the migrated version itself being stored.

Clearly defining the current designated community and its level of technical capability will help limit the resources necessary to support this “lowest-common-denominator” approach. For organizations such as national repositories that have a very loosely defined community with a diverse knowledge base, this could prove very labor intensive.

Fulfilling this responsibility requires:

- Well-maintained and documented technical metadata that is kept aligned with the knowledge base of the designated community and with changing technologies.
- A “technology watch” to manage the risk as technology evolves and to provide continuing access and updated methods of access as necessary, such as new migrations or emulators.

Recommendation

Investigate and determine which technical strategies best provide for continuing access to digital resources.

Following Documented Policies and Procedures

In the past, some organizations may have relied on vague or even unwritten policy for the management of traditional collections. However, to ensure effective and efficient mechanisms for long-term preservation of and continuing access to its digital contents, a repository requires well-documented and widely adopted policies—and well-documented procedures.³⁰ For distributed repositories, this means clearly articulated responsibilities across participating organizations and consortia.

For research repositories, a strategy or policy for preservation of digital materials may necessarily relate more widely to the organization's information strategy as a whole. More particularly, however, a policy for the preservation of digital files needs to sit comfortably within or alongside policies for nondigital content. The link between policy and procedure will also be critical. If the policy of a research repository sets different levels of collection for long-term retention, each level will need a corresponding procedure. Over time, linked policies and procedures will help to reduce costs by supporting automation and scaling. Rather than consider each digital object individually at the point of deposit, procedures will automatically apply, based on the policy for a particular part of the collection.

A Simple Example of Policy and Procedure

The collections policy of the University of Anytown applies these levels of collecting to its subject areas:

- Comprehensive/Research
- Study and Teaching
- Minimal

Their policy also applies these levels of responsibility for long-term preservation:

- Archival (kept forever)
- Served (available for the foreseeable future)
- Mirrored (responsibility taken only for the short term)
- Linked (no long-term responsibility assumed)

Each of the levels has a corresponding technical procedure for digital archiving. For “Archival” materials in a “Comprehensive” collection, all significant properties would be preserved and levels of necessary metadata would be assigned automatically (ideally) at the time of submission. At the other extreme, “Linked” materials from a “Minimal” collection might have very little, if any, preservation.³¹

Fulfilling this responsibility requires:

- Policies for collections development (e.g., selection and retention) that link to technical procedures about how and at what level materials are preserved and how access is provided both short and long term.
- Policies for access control to ensure all parties are protected, including authentication of users and disseminated materials.
- Policies for storage of materials, including service-level agreements with external suppliers.
- Policies that define the repository’s designated community and describe its knowledge base.
- A rigorous system for updating policies and procedure in accordance with changes in technology and in the repository’s designated community.
- Explicit links between these policies and procedures, allowing for easy application across heterogeneous collections.

Making the Preserved Information Available to the Designated Community

Providing access to materials is an integral responsibility of a digital repository, but access must be clearly defined in order for a repository to understand its implications. Immediate access to materials will require different policies, such as license arrangements, and therefore different management than access to materials over time. If materials are only accessible in a particular format to a specific group of users for a designated period, different mechanisms

³⁰ Margaret Hedstrom and Sheon Montgomery, *Digital Preservation Needs and Requirements in RLG Member Institutions* (RLG, 1998) www.rlg.org/preserv/digpres.html.

³¹ This example is based on the Berkeley Digital Library SunSITE’s *Digital Library SunSITE Collection and Preservation Policy* (1996) sunsite.berkeley.edu/Admin/collection.html.

will need to be in place than might be appropriate later. Access arrangements will change in accordance with changes in licenses, law, and technology, and with local resource constraints. Repositories also need to ensure as far as possible that decisions about access when materials are submitted do not limit what might be possible in the future.

- Resource discovery: A repository's users need to find materials. Many libraries and archives provide access through their main catalog. In practice, many objects to be deposited in a repository will arrive with existing—often very rich—bibliographic descriptions such as MARC or Dublin Core, either accompanying the object or available in an existing system.
- Authenticity: The authenticity of digital materials is more complex and potentially troublesome than that of traditional library or archival materials.³² While traditional materials can be physically verified, digital objects have less obvious evidence of authorship, provenance, or even context. For this reason, they give rise to suspicions that will be assuaged only by rigorous mechanisms throughout the repository for ensuring that the digital object is what it purports to be and that it is what was originally deposited into the repository. Authenticity checks are required at all functional levels of the digital repository: At submission, mechanisms must ensure that the object, as received, is what was intended by the content provider.³³ The stored material needs a regular system of integrity checks to ensure the bytestreams are maintained. Physical and system procedures must be maintained and the available mechanisms for access to the original bytestreams must be regularly checked; migrated versions must be verified and available emulators tested. Finally, the information provided to the user—the copy of the bytestream as well as the necessary metadata and rendering software—all requires verification.
- Legal issues: Legal restrictions—licenses and legislation—govern access to materials in a repository, and over time these change. Digital repositories require an infrastructure that can support different access arrangements for different materials and different types of users.
- Pricing: Repositories that govern access with a fee structure require mechanisms for managing electronic commerce.
- User support: Access, particularly to older materials, requires some level of user support. To a great extent, this will be determined by designated community's knowledge base or technical capability.
- Record keeping: As part of the repository's administration function, it may be advisable to keep track of the dissemination of objects out of the repository.

Fulfilling this responsibility requires:

- A system for discovery of resources.
- Appropriate mechanisms for authentication of the digital materials.

³² Charles T. Cullen et al., *Authenticity in a Digital Environment* (CLIR, May 2000) www.clir.org/pubs/reports/pub92/contents.html.

³³ In the OAIS Reference Model, submission is referred to as the Ingest function; see Appendix A.

- Access control mechanisms in accordance with licenses and laws, and an “access rights watch.”
- Mechanisms for managing electronic commerce.
- User support programs.

Advocating Good Practice in the Creation of Digital Resources

If digital repositories advocate standards for the creation of digital materials, they will be able to achieve important economies of scale and reduce costs, as well as ensure the creation of rich digital resources worthy of long-term preservation and digital materials that are more amenable to digital preservation practices and processes. It is critical to involve both the repository’s designated community (e.g., content providers, content creators, and users) and software suppliers and it may require an outreach program for potential depositors. Standards and best practices for the creation, description and necessary metadata are emerging for many types of files and formats.³⁴

It will be important for key players to facilitate this dialogue; organizations such as RLG, OCLC, or IFLA (International Federation of Library Associations and Institutions) are well placed to take the lead at the international level, where this may be more effective.

Fulfilling this responsibility requires effective mechanisms for advocating good practice for content providers.

³⁴ For example, best practices related to *Benchmarking digital reproductions of printed monographs and serials* have recently been released by the Digital Library Federation (www.diglib.org/standards/bmarkfin.htm); *Guides to Quality in Visual Resource Imaging* was released by RLG and DLF in 2000 (www.rlg.org/visguides/); and the Joint Information Systems Committee has released *Working with the Distributed National Electronic Resource (DNER): Standards and Guidelines to Build a National Resource* to facilitate the development of UK digital collections (www.jisc.ac.uk/dner/development/guidance/DNERStandards.html).

4 Certification of Trusted Digital Repositories

At least two viable models for certification of trusted digital repositories are in use and well known within the library and archives community:

- The audit model is applicable to depositories holding government records, especially electronic records. In the US, such depositories must meet guidelines created by legislation or by agencies such as the Department of Defense.³⁵
- The standards model operates in various places in the library and archives community. Two examples are guidelines for producing preservation-quality microfilm and ISO interlibrary lending.³⁶ Institutions involved in these activities adhere to standards established by appropriate agencies. Peer institutions “certify” the product or service by their acceptance and/or use of it.

While both models work well, neither can completely address the range of activities, functions, and responsibilities associated with digital repositories.

In 1999, experts gathered at the Archival Workshop on Ingest, Identification, and Certification Standards (AWIICS) to begin developing standards specifically appropriate to the needs of digital repositories. Leading the discussion on certification, Bruce Ambacher, National Archives and Records Administration, Center for Electronic Records, identified four general approaches to certification: individual, archival program, process, and data.³⁷ These four approaches refine the audit and standards models.

- Individual: Individual, professional certification or accreditation is sometimes referred to as personnel certification. Traditionally, archivists have been certified through a combination of education, work experience, and a competencies examination administered by the Academy of Certified Archivists (ACA). Nothing equivalent exists for electronic archiving or digital repository management.
- Program: Certification of a program or institution can be achieved through a combination of self-evaluation using standardized checklists and criteria and site inspections typical of program accreditation. Three models are the Society of American Archivists Evaluation of Archival Institutions, the HMC Approval from the Historical Manuscripts Commission (UK), and the Museum Assessment Program from the American Association of Museums. For these certifications,

³⁵ Department of Defense, United States, *Design Criteria Standard For Electronic Records Management Software Applications*, jtc.fhu.disa.mil/recmgt/dod50152.doc; National Archives and Records Administration, *NARA Regulations in the Code of Federal Regulations, Regulations in 36 CFR Chapter XII, Subchapter B - Records Management* (2001) www.nara.gov/nara/cfr/subch-b.html; National Archives and Records Administration, *Basic Laws and Authorities of the National Archives and Records Administration* (2000) www.nara.gov/nara/basiclaws.html; US House of Representatives, *Downloadable U.S. Code: Records Management By Federal Agencies*, 44 USC - Chapter 31 (January 2000) uscode.house.gov/title_44.htm.

³⁶ Nancy Elkington, ed., *RLG Preservation Microfilm Handbook* (Mountain View, CA: RLG, 1992); Nancy Elkington, ed., *RLG Archives Microfilming Manual* (Mountain View, CA: RLG, 1994); Interlibrary Loan Protocol Implementers Group (IPIG), *Profile for the ISO ILL Protocol: Version 2.0* (April 2001) www.nlc-bnc.ca/iso/ill/document/ipigwp/profile/ipv2_0.pdf.

³⁷ *Archival Workshop on Ingest, Identification, and Certification Standards (AWIICS): Draft Report* (1999) ssdoo.gsfc.nasa.gov/nost/isoas/awiics/.

areas assessed include legal authority, governing authority, financial resources, staff, facilities, collection development, collection preservation, access, and outreach.

- **Process:** Process certification assesses methods and procedures that can be subjected to quantitative or qualitative guidelines for adherence to internal and external requirements. External standards to evaluate archival and digital repository processes may include the ISO 9000 family, DoD 5015.2-STD, the Public Record Office Standard (Victoria), the Public Record Office (UK), and BSI DISC 0008.³⁸
- **Data:** Data certification addresses the persistence or reliability of data over time and data security. Certification for data persistence would include both internal and external quality control through standards such as ISO 9000:2000 and procedures manuals. It would also include documenting the processes for migrating data, creating and maintaining metadata, updating data or files, and authenticating new copies. The Public Key Certification Policy and Certification Practices Framework addressed data security—because of the e-commerce boom, however, not for digital archiving. This framework, which deals with user authentication and user communication in e-commerce transactions, handles access control issues for repositories and removes the need for additional data security certification.³⁹

The participants in the AWIICS workshop agreed that elements of each of these four processes could form a certification program that provides layers of trust. A layered approach should convey a high degree of confidence that the information a repository disseminates is the same as the information it ingested and preserved, with full documentation for all necessary modifications. A preliminary checklist created at the workshop will serve as a tool for further work on the OAIS standardization.⁴⁰ Both the checklist concept and the certifiable elements envisioned at the workshop provide a base for a certification framework.⁴¹

Recommendation

³⁸ International Organization for Standardization (ISO), *ISO 9000 / ISO 14000*, www.iso.ch/iso/en/iso9000-14000/iso9000/iso9000index.html; Department of Defense, United States, *Design Criteria Standard For Electronic Records Management Software Applications*, jtc.fhu.disa.mil/recmgt/dod50152.doc; Public Record Office Victoria (Australia), *Standard for the Management of Electronic Records, PROS 99/00* (April 2000) www.prov.vic.gov.au/vers/standards/standards.htm; Public Record Office (UK), *Management, Appraisal and Preservation of Electronic Records* (1999)

www.pro.gov.uk/recordsmanagement/eros/guidelines/; British Standards Institution, *Code of Practice for Legal Admissibility and Evidential Weight of Information Stored Electronically*, DISC PD 0008:1999 (1999).

³⁹ Internet Engineering Task Force (IETF), *Internet X.509 Public Key Infrastructure Certificate Policy and Certification Practices Framework* (March 1999) www.ietf.org/rfc/rfc2527.txt.

⁴⁰ Archival Workshop on Ingest, Identification, and Certification Standards (AWIICS), *Certification (Best Practices) Checklist*, ssdc.gsfc.nasa.gov/nost/isoas/awiics/CertifBase.ppt.

⁴¹ The Certification group at the AWIICS workshop recommended that “accreditation of archives is important and should be pursued, but that it can only be accomplished when best practices are in place.” Since the OAIS Reference Model was only in draft at the time, certification activities were delayed until the Reference Model was ready for ISO standardization; Donald Sawyer and Jerry Winkler, “Digital Archive Directions Workshop Extremely Successful,” *NOST Hosts Archiving Workshop* 14, no. 4 (September 1998) nssdc.gsfc.nasa.gov/nssdc_news/sept98/01_j_garrett_0998.html.

Develop a framework and process to support the certification of digital repositories.

5 Summary and Recommendations

RLG and OCLC together asked the Working Group on Digital Archive Attributes to reach consensus on the characteristics and responsibilities of a sustainable digital repository for large-scale, heterogeneous collections held by cultural organizations. The working group has articulated a framework of attributes and responsibilities for trusted repositories for digital content capable of handling the range of materials held by large and small research and cultural institutions. When all are present and met, those attributes and responsibilities identify reliable, sustainable digital repository infrastructures and form the basis for the development of future trusted services.

The framework is broad enough to accommodate different situations, architectures, and institutional responsibilities while providing a basis for the expectations of a trusted repository. A trusted digital repository is more than just an organization responsible for storing and managing digital files. A trusted digital repository is one whose mission is to provide reliable, long-term access to managed digital resources to its designated community, now and in the future. Some institutions may choose to address all functional and informational responsibilities in a complete local system. Others may choose to manage the logical and intellectual aspects of a repository while contracting with a third-party provider for digital file storage and maintenance. A critical component is the overall infrastructure supporting the reliability and sustainability of digital repositories so that organizations and their designated communities can trust that digital resources will be preserved for the long-term.

For organizations wishing to provide digital repository services, building the required trust via reliable, proven practice will take time. Yet the need for trusted services is immediate. Since action needs to be taken now to preserve the already large body of digital materials, a program for certification should be developed to provide a basis for trustworthiness. Certification would dictate criteria that must be met and would employ mechanisms for assessment and measurement. Cultural institutions would have a tool for measuring currently available services, and service providers would have a known set of best practices or standards to meet in order to gain the business of cultural institutions. Certification, achieved on a periodic basis for several years, could resolve the tension between the immediate need for trusted archives and the need to develop and prove reliability over time.

Finally, this working group did not attempt to identify best practices of specific technical mechanisms or strategies directly related to long-term storage and maintenance of digital materials. It is recognized that issues such as authentication, use of persistent identifiers, redundancy, and metadata are all necessary for the long-term viability of digital collections. Further research is needed in these critical areas before best practices can be established and endorsed, but that research was out of scope for this working group.

Additional work must take place to enable the establishment of a network or federation of trusted digital repositories that meet the needs of cultural institutions. The working group has identified and prioritized a series of recommendations and proposes the following as next steps.

Recommendation 1: *Develop a framework and process to support the certification of digital repositories.*

A certification framework and certification process for digital repositories are crucial and their absence has been an impediment to assigning trust. Model processes, including checklists for certification reviews, should be developed incorporating the community-approved attributes of trusted digital repositories, the work of the ISO Archiving Series, and other relevant projects.

Recommendation 2: *Research and create tools to identify the attributes of digital materials that must be preserved.*

Assessing the significant properties of a digital object or of a class of digital objects will be a critical component of an effective long-term digital repository. More work is required on the possible significant properties for different classes of digital materials and how these properties in turn determine the underlying technical form or structure to be preserved and the necessary technical metadata. Tools should be developed to enable creators to identify the significant properties at the point of creation rather than at a later point in history.

Recommendation 3: *Research and develop models for cooperative repository networks and services.*

Archivists and librarians need more thorough understanding of how cooperative digital repositories and repository networks can be implemented and managed, including the use of third-party service providers. Models for the establishment of cooperative archiving services will be useful and necessary, as will be examples of service-level agreements as they apply to digital repositories (e.g., service-level agreements for external suppliers of archival storage).

Recommendation 4: *Design and develop systems for the unique, persistent identification of digital objects that expressly support long-term preservation.*

The pressing need for unique and persistent systems of identification for digital information has supported a great deal of work in this area. However, it is not yet clear that current approaches, if any, are best suited to the purposes of long-term preservation. Systems should be designed—or existing systems modified—to enable long-term maintenance, storage, and access to digital resources.

Recommendation 5: *Investigate and disseminate information about the complex relationship between digital preservation and intellectual property rights.*

Digital preservation activities have legal implications. How preservation infringes upon current copyright laws remains unclear and is a current impediment to large-scale action. Research is needed to identify where current copyright protections inhibit digital preservation and how technical strategies might impinge on copyright laws. Tools should be developed to identify the roles and responsibilities of content creators and organizations preserving information. Work is needed on models for obtaining copyright clearance and models for contracts or agreements between rights owners/producers and archives/libraries.

Recommendation 6: *Investigate and determine which technical strategies best provide for continuing access to digital resources.*

Libraries and archives need more practical experience using both migration and emulation as strategies for continuing access. Guidance should include both technical and legal

(copyright) implications for migration and emulation. Further work is necessary on how specific strategies for continuing access may affect preservation metadata created when materials are submitted to the collection.

Recommendation 7: *Investigate and define the minimal-level metadata required to manage digital information for the long term. Develop tools to automatically generate and/or extract as much of the required metadata as possible.*

Work is needed to define what technical and administrative metadata will be needed for the long-term preservation of digital files. What metadata will be critical to support continued rendering and functionality? What descriptive information will be needed to support long-term semantic interoperability? Metadata elements recording authentication information (checksums, CRC, MNP, etc.) should become mandatory for all digital objects being deposited into trusted digital repositories. Tools should be developed to support the automatic generation and extraction of metadata.

Appendix A: OAIS Technical Overview

A review of the Open Archival Information Systems Reference Model provides insight into the necessary functions of a long-term digital repository. Its adoption by libraries and repositories is increasing, due mainly to continued involvement of the library and archives community in its development. Although work is still to be done on many of the standards necessary for its effective distributed implementation, the Reference Model provides a critical framework for establishing or enhancing digital archiving services.

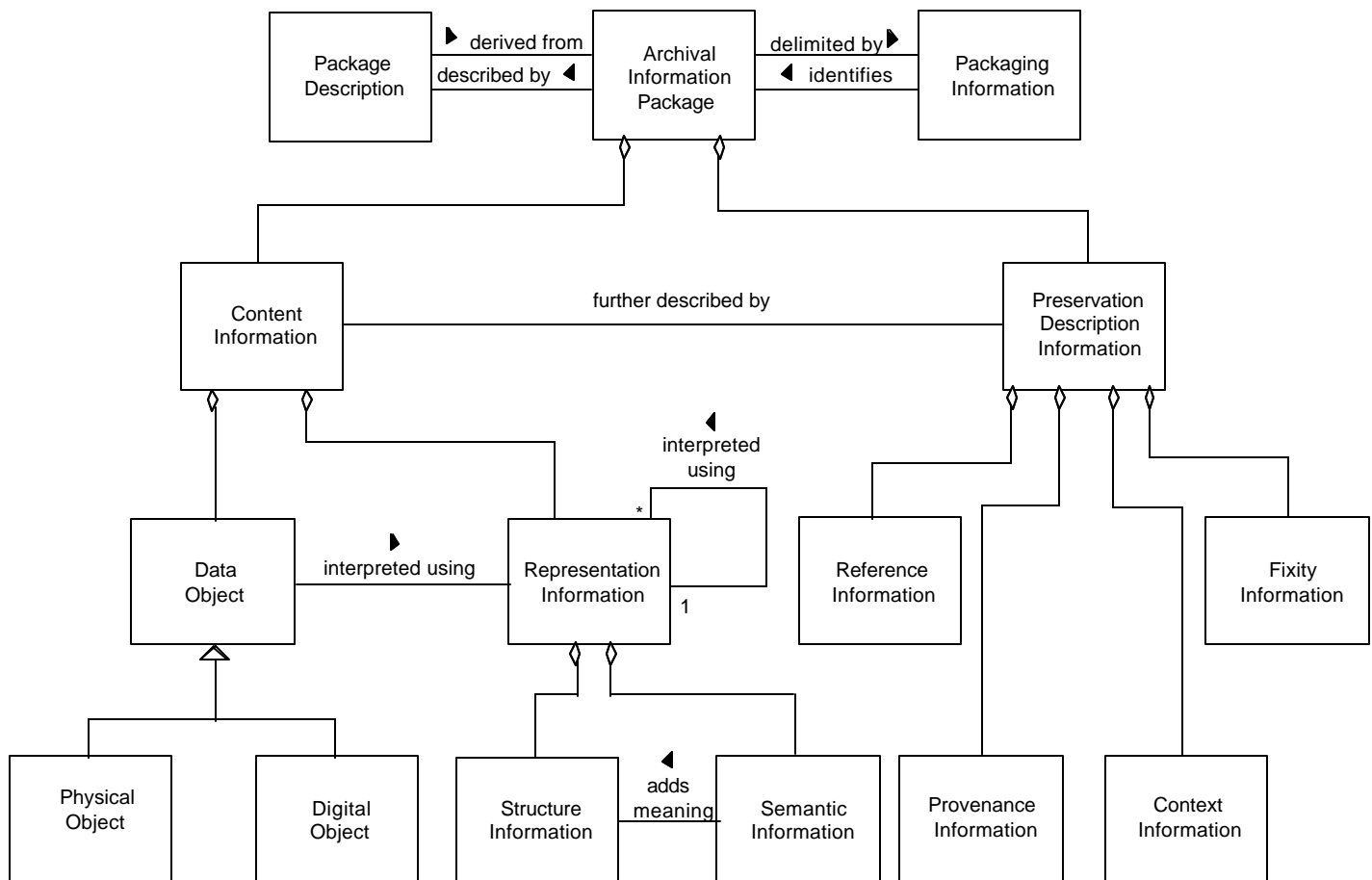


Figure 1. The Archival Information Package (AIP), detailed view⁴²

The OAIS Reference Model and Digital Preservation Metadata

The foundation of an OAIS digital repository is the **Information Package**, which includes both a digital object and the necessary associated metadata. As objects are submitted to the repository, they arrive as a **Submission Information Package (SIP)**, which contains the

⁴² Consultative Committee on Space Data Systems, *Reference Model for an Open Archival Information System (OAIS): Draft Recommendation for Space Data System Standards, CCSDS 650.0-R-1, Red Book, May 1999, www.ccsds.org/RP9905/RP9905.html.*

digital object and any other information the content provider deemed relevant or available. On submission, the SIP is enhanced as necessary and encapsulated as an **Archival Information Package (AIP)**, including the data object and all its associated metadata. The AIP is the cornerstone of the digital repository. When a user requests access to an object, a **Dissemination Information Package (DIP)** is provided, which typically contains a copy of the digital object as well as the necessary metadata and support systems to retrieve and use the digital object.

The original digital object is stored as a bytestream in the AIP, along with the metadata necessary for making that bytestream into a meaningful and useable digital resource. In the future what is known about digital materials will come from the information or metadata stored with them—the binary data is meaningless without some description of what it is and how it works. OAIS provides for two main types of metadata: **Content Information** contains both the digital object (as a bytestream) and all the necessary technical metadata (called **Representation Information** or **RI**). **Preservation Description Information (PDI)** is all of the other descriptive information that, although not critical for actual conversion of the bytestream, is necessary for long-term preservation. The creation and maintenance of the preservation metadata (both the technical RI and the PDI) will most likely represent the bulk of the initial costs and complexities of digital preservation.

Preservation Description Information (PDI)

Preservation Description Information is the metadata that includes traditional bibliographic information as well as more detailed information to ensure the material is effectively preserved. This metadata includes details of ownership, copyright, and other intellectual property rights, such as licensing arrangements and access restrictions. Preservation Description Information is broken down into Provenance, Reference, Fixity, and Context Information.

Representation Information (RI)

Representation Information is the technical metadata for mapping the bytestream into specific data types and formats. Without adequate Representation Information, the bytestream is not retrievable as a meaningful digital object: the Representation Information provides meaning to the bits. In practice, Representation Information involves many different descriptions for a variety of relevant technologies. For example, even a simple Web page that contains graphics requires descriptions of the Web environment (browser, etc.), the text (ASCII standard), and the image files. This is a recursive system—all the Representation Information requires additional Representation Information in order to be understood. Representation Information is therefore likely to involve references to other Representation Information elsewhere in the repository and will take the form of a **Representation Network**. In theory, if the digital object is to remain accessible for the long term, this recursion will stop only when a physical form is encountered, such as a system specification or technical manual. However, in practice, the “end node” of the Representation Network will be the software or hardware that is part of the knowledge base of the designated community, which allows end users to understand the digital object without recourse to the repository or the data creator. It is the level of technology or technical capability that can be assumed to be supported outside the repository itself—a kind of lowest common denominator. These end nodes of the Representation Information will need to be closely monitored. As any technology threatens to become obsolete, Representation Information will need to be generated and stored in the repository to support that technology.

International Collaboration and Digital Preservation Metadata

A number of initiatives have developed preliminary specifications for preservation metadata (both RI and PDI). Although beyond the scope of this report, the work of Harvard University, the National Library of Australia, the NEDLIB and Cedars projects, and others have played a key role.⁴³ RLG and OCLC are now taking forward the work done by these and other organizations to build consensus and develop a standard framework for metadata for the long-term retention of digital materials.⁴⁴

The OAIS Functional Model

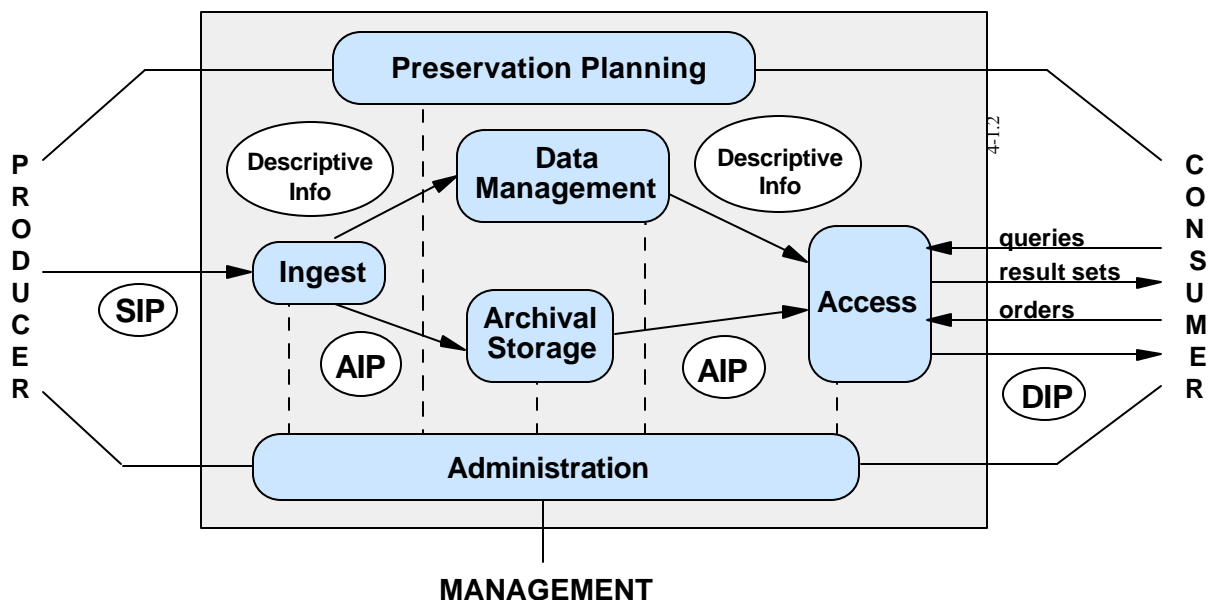


Figure 2. Overview of the Open Archival Information Systems Reference Model

Along with its model for necessary metadata, OAIS includes a comprehensive logical model for a repository's functions. Any scheme for describing a reliable digital repository's attributes will require that a variety of communities share an understanding of its basic functions. The OAIS functions include:

- Submission and “pre-Ingest” activities
- Ingest

⁴³ Harvard University Library, *Digital Repository Services (DRS) User Manual for Data Loading*, available from hul.harvard.edu/ois/systems/drs/doc.html; Harvard University Library, Library Digital Initiative, *Image Reformatting*, hul.harvard.edu/ldi/html/reformatting_image.html; Harvard University Library, *Library Preservation at Harvard: Image Digitization*, preserve.harvard.edu/resources/digital.html; National Library of Australia, *Preservation Metadata for Digital Collections: Exposure Draft (1999)* www.nla.gov.au/preserve/pmeta.html; Catherine Lupovici and Julien Masanes, *Metadata for the Long Term Preservation of Electronic Publications*, NEDLIB Report Series, report 2 (NEDLIB Consortium, 2000), available from www.kb.nl/coop/nedlib/. Cedars Project Team, *Metadata for Digital Preservation: The Cedars Outline Specification (June 2000)* www.leeds.ac.uk/cedars/metadata.html.

⁴⁴The OCLC/RLG Working Group on Preservation Metadata, *Preservation Metadata for Digital Objects: A Review of the State of the Art. A White Paper*. (January 2001) www.oclc.org/research/pmwg/presmeta_wp.pdf.

- Archival Storage
- Data Management
- Preservation Planning
- Archive Administration
- Access/Dissemination⁴⁵

The following descriptions of these functions are based both on the Reference Model itself and on implementations of the Model in a library or archive, mainly by the Cedars and NEDLIB projects⁴⁶.

Submission and “Pre-Ingest” Activities

Before a repository can accept responsibility as a reliable archiving service, management tools must be in place, covered by a well-documented and agreed-on collections policy document. For most libraries and archives, this will be in the form of a collections management/development policy. Many of the same criteria apply to both digital and traditional. However, for digital materials these elements are the most critical:

- Evaluation criteria for assessing potential submissions; that is, selection criteria for digital preservation.
- Collections development strategies and technical strategies for continuing access.
- Collections development procedures, including review procedures pertaining to retaining and deaccessioning materials.

Where appropriate, the repository needs to ensure availability of copyright and other intellectual property rights or privacy/confidentiality information, including licenses, schedules for deposit (where regular updates will be forthcoming), and appropriate documentation, and may even include details of preferred formats and media.

As part of the pre-Ingest activities, the repository also needs to:

- Check any existing deposit schedules to ensure everything expected has been received.
- Assign the digital object’s unique identifier(s), if not already available, and provide labels for the physical artifact.
- Check for viruses and validate the integrity of the digital object and its physical carrier.
- Assess in detail the significant properties of the digital object, such as its look and feel, or functionality.
- Validate or improve the documentation.

⁴⁵The OAIS Reference Model contains a complete glossary of its terms.

⁴⁶On Cedars, see The Cedars Project Report, April 1998–March 2001, www.leeds.ac.uk/cedars/OurPublications/cedarsrepmar01exec.html; Johan Steenbakkens, *The NEDLIB Guidelines: Setting up a Deposit System for Electronic Publications*, NEDLIB Report Series, report 5 (NEDLIB Consortium, 2000); available from www.kb.nl/coop/nedlib/.

- Where appropriate, reformat the digital object according to repository policies.
- Ensure that all necessary metadata for long-term maintenance and continuing access accompanies the object.

Ingest

Once the digital resource has been properly prepared, the Ingest function allows for materials submitted as a Submission Information Package (SIP) to be prepared as Archival Information Packages (AIPs) for storage. At this stage, the Content Information and related PDI are established in the repository. Creation of the AIP is the foundation of the long-term preservation function.

Ingest to the repository on a practical level involves:

- Assignment and/or validation of unique identifier. This identifier must be unique not only within the repository but also within the repository's wider community or the federation of which it is a part.
- Selection and validation of the agreed-on underlying technology or underlying abstract form based on the object's significant properties.⁴⁷
- Transformation of the object as it was submitted, along with its associated metadata, into a bytestream that can be stored on suitable hardware in the repository.
- Establishment of necessary Representation Information.
- Verification of all Preservation Description Information.

In a distributed environment, the Ingest function can be performed across participating repositories. In practice, this requires the adoption of strict policies for the assignment of unique identifiers or use of existing identifiers and a well-documented and fully implemented specification for preservation metadata—both PDI and RI. The outcome of ongoing work with unique identifiers and emerging standards in this area will be important.

Archival Storage

The Archival Storage function is the necessary services for the effective storage and retrieval of AIPs. Such functions include:

- Moving AIPs from Ingest into permanent storage.
- Managing the storage hierarchy.
- Refreshing the storage media.
- Providing all necessary information to allow objects to be disseminated from the repository.

Like the Ingest function, Archival Storage can be centralized or distributed. Whichever model is adopted, it will rely on a robust and well-documented policy that describes how the material is stored and maintained and what level of service is expected. Archival Storage can

⁴⁷ For a discussion of an object's significant properties, see Section 3.

be done either in-house or by a third party; service-level agreements are essential with third parties.

Archival Storage must also include systems for routine integrity checking. The stored bytestream should be checked regularly to assess whether it is truly identical to the original bytestream. Authentication and integrity checking should be a regular activity of the repository's administration.

Another critical component of any storage facility is a robust system for disaster recovery, with these main components:

- An ongoing program of media refreshment, transferring bytestreams onto newer, fresher media. This includes a program of monitoring repository media for possible degradation and subsequent integrity checking on refreshed bytestreams.
- Geographically distributed backup systems—ideally, more than one.

Data Management

Data Management covers all aspects of an OAIS repository and is essential for both long-term preservation and day-to-day administration and use. It represents good record keeping at every stage described by the OAIS functions. The activities in the Data Management component are determined by the policies developed and maintained by the repository's management and administration. Some of the components that the OAIS Reference Model includes in Data Management are:

- Pricing information (if applicable) and access controls.
- Customer profiles.
- Tracking of user requests.
- Security information, including any usernames, passwords, digital certificates—anything used to authenticate users of the repository.
- Statistical information to improve operation.
- Accounting information.

Cedars added these elements to Data Management:

- Records from pre-Ingest negotiations, such as immediate or short-term access arrangements.
- Policies for and monitoring of the allocation of unique identifiers.
- Maintenance of records of holdings for use with finding aids.⁴⁸

Preservation Planning

Preservation Planning is the function responsible for monitoring the OAIS environment and providing recommendations to the repository (through the Archive Administration function) to ensure that materials are accessible to the designated community over the long term. This function detects critical information about shifts in the knowledge base of the designated

⁴⁸The Cedars Project Report, April 1998–March 2001, www.leeds.ac.uk/cedars/OurPublications/cedarsrepmar01exec.html.

community, allowing the repository to accommodate shifts in technology. Preservation Planning includes:

- Monitoring the designated community.
- Monitoring technology.
- Monitoring the significant properties of the repository's contents.
- Developing preservation strategies and standards for continuing access.
- Developing packaging designs and migration or routine transfer plans.

The addition of this function to OAIS in March 2001—thanks particularly to the efforts of NEDLIB and the National Library of Australia—reflects the influence of libraries and archives on the Reference Model. Implementation of the Model by libraries and archives has made clear the need for a function that explicitly allows for the application of processes such as migration or emulation for continuing access.

Archive Administration

Archive Administration is all services needed for day-to-day maintenance of the repository. This includes management of the systems configurations (statistics, etc.) and repository policy development and maintenance. Archive Administration includes:

- Negotiating submissions agreements with content producers and providers.
- Reviewing procedures.
- Maintaining systems configurations for hardware and software.
- Developing and maintaining repository policies and standards, including policies and standards recommended and enhanced by Preservation Planning.
- Providing user support.
- Interacting with management outside of the repository.

To this list, Cedars added:

- Reviewing and maintaining Representation Information and Networks, based on recommendations from Preservation Planning.
- Negotiating user access agreements with service providers or others.
- Communicating with other repositories.⁴⁹

Archive Administration may also include monitoring of legal status and relevant changes in national and possibly international law.

Access/Dissemination

In the OAIS Reference Model, access to archived materials is provided through the Dissemination Information Package (DIP): a copy of the digital object along with the necessary metadata, and software as necessary. In addition to preparation of the DIP,

⁴⁹ Ibid.

Dissemination requires mechanisms for both verifying the integrity of the information in the DIP and for ensuring that users have permission for access to the material.

The DIP differs from the AIP in that it contains only a copy of the digital object (or a new object generated from the AIP such as in a migration) and only the metadata necessary for access to the appropriate level. It probably does not include the full complement of rich, descriptive PDI or RI. The important separation of AIP and DIP allows for materials to be disseminated in different ways and by different methods of continuing access. For example, a database might be available both as a migrated version accessible in the current version of Microsoft Access or with an emulator of its original database package. This separation nicely illustrates the difference between long-term maintenance and continuing access.

Although it would be possible to simply copy the AIP and distribute it to a user, this is probably not practical. Most users would have no need for all the detailed metadata that supports the bytestream long term. The DIP should contain only the metadata necessary for the required level of access. For example, a PDF file may be disseminated as a bytestream with a copy of Adobe Acrobat® Reader®, while within the repository the technical metadata includes much richer information about the PDF format.

Technically, the level of access provided for digital materials depends on the judgment of the archivist and/or the collection manager, based on the objects' significant properties. Strategies such as emulation or migration of some kind are the obvious choices for providing access over time. For many materials, access will change over time and will vary according to the designated community. For example, changes in licensing arrangements or in copyright may require updates to the way access is provided, but not necessarily to the way the bytestream is preserved. Likewise, particular designated communities may have different levels of access to particular materials at different times (e.g., researchers vs. undergraduate students or rights holders vs. users).

Digital materials cannot be considered preserved without meaningful access. However, publishers and other rights holders are often cautious about the preservation of their materials if they think unlimited or unrestricted public access is a necessary precondition. The distinction between access for current users and long-term maintenance and access needs to be explored more deeply and better understood to ensure interests in the one do not jeopardize the other.

Appendix B: The Evolution of “Trust” in Computing Systems

A quick survey of the uses of the terms “reliable” and “trust” since the late 1980s can blaze a path toward clearer thinking about the nature of organizations that will undertake long-term digital preservation services. At the same time, the language and concepts used to express expectations can be used to build toward a common definition that bridges the jargon and experiences of relevant expert communities.

A Guide to Understanding Audit in Trusted Systems was issued by the National Computer Security Center (NCSC) in 1988. This guide was an effort to identify the ways that trusted systems could be assessed and measured by third parties and focused very specifically on systems that processed classified military information. A trusted system had to be one that had built-in audit trails:

A trusted computer system must provide authorized personnel with the ability to audit any action that can potentially cause access to, generation of, or effect the release of classified or sensitive information. The audit data will be selectively acquired based on the auditing needs of a particular installation and/or application. However, there must be sufficient granularity in the audit data to support tracing the auditable events to a specific individual (or process) who has taken the actions or on whose behalf the actions were taken.⁵⁰

In 1992 the NCSC issued *Guidelines for Writing Trusted Facility Manuals*, which defined a trusted computer system as one “that employs sufficient hardware and software assurance measures to allow its use for simultaneous processing of a range of sensitive or classified information.” Two further definitions appearing in the document’s glossary are illuminating:

Trusted Computing Base (TCB): The totality of protection mechanisms within a computer system—including hardware, firmware, and software—the combination of which is responsible for enforcing a security policy. A TCB consists of one or more components that together enforce a unified security policy over a product or system. The ability of a TCB to enforce a security policy correctly depends solely on the mechanisms within the TCB and on the correct input by system administrative personnel of parameters (e.g., a user’s clearance) related to the security policy.

Trusted Path: A mechanism by which a person at a terminal can communicate directly with the TCB. This mechanism can only be activated by the person or the TCB and cannot be imitated by untrusted software.⁵¹

These NCSC documents augment a basic interpretation of “trust” with the concepts of auditability, security, and communication. It should be noted that here communication refers to exchanges between a person and a machine. Later uses of the term denote different forms of exchange.

In 1996 the University of Pittsburgh’s School of Information Sciences developed a set of *Functional Requirements for Evidence in Recordkeeping*, which stipulated:

⁵⁰ *A Guide to Understanding Audit in Trusted Systems* (1998)
www.radium.ncsc.mil/tpep/library/rainbow/NCSC-TG-001-2.html.

⁵¹ National Computer Security Center, *Guidelines for Writing Trusted Facility Manuals*, NCSC-TG-016, Yellow-Green Book (October 1992).

Organizations must comply with the legal and administrative requirements for recordkeeping within the jurisdictions in which they operate, and they must demonstrate awareness of best practices for the industry or business sector to which they belong and the business functions in which they are engaged.⁵²

Compliance and auditability are linked concepts here, establishing a direct connection between performance and assessment. By using the term “conscientious,” the authors affirmed their belief that a repository must be meticulous in order to act responsibly.

Mark Stefik, principal scientist in the information sciences and technology laboratory at the Xerox Palo Alto Research Center, has been writing about trusted systems for some years. In his 1997 *Scientific American* article, “Trusted Systems,” and in an article published that same year in the *Berkeley Technology Law Review*, “Shifting the Possible: How Trusted Systems and Digital Property Rights Challenge Us to Rethink Digital Publishing,” Stefik made a series of points about ways to ensure the security and unchangeability of materials on the Web so as to protect both creators and distributors from unlawful copying by consumers:

Computer scientists recognize trusted systems as those that follow rules that govern terms, conditions, and fees for using digital works.

Knowing what the rules are is a central part of the design of trusted systems.

Necessary security requirements vary based on the type of work being protected, but requirements for the most valuable works should include detection and prevention of tampering.

One trusted system has to be able to recognize another trusted system.

Copying can be permitted if it is strictly controlled in accordance with the creator’s and/or distributor’s and rights and interests (particularly those related to charging fees).

Certification of trusted systems can ensure that they are compliant and can be relied on to follow appropriate rules and instructions.⁵³

Stefik reinforced previously identified concepts of security, communication, and compliance but added certification, copying controls, and following rules to the amalgam.

In her 1998 article, “Building Recordkeeping Systems: Archivists are Not Alone on the Wild Frontier,” Margaret Hedstrom, associate professor in the School of Information, University of Michigan, surveyed the evolving field of electronic and hybrid (that is, involving paper-based and electronic source material) record-keeping systems, suggesting numerous ways that the archival profession could take better advantage of work going on in the computer and digital library communities to develop trusted systems:

Recent research and development efforts within the archival community combined with new methodologies for trusted systems provide archivists with a variety of tools to enhance the integrity, reliability, and usefulness of

⁵² *Functional Requirements for Evidence in Recordkeeping* (1996)
web1.archive.org/web/19981205092442/www.sis.pitt.edu/~nhprc/prog1.html.

⁵³ Mark Stefik, “Trusted Systems,” *Scientific American* (March 1997)
www.sciam.com/0397issue/0397stefik.html;
Mark Stefik, “Shifting the Possible: How Trusted Systems and Digital Property Rights Challenge Us to Rethink Digital Publishing,” *Berkeley Technology Law Journal* 12, no. 1 (spring 1997)
www.law.berkeley.edu/journals/btlj/articles/12_1/Stefik/html/text.html.

electronic recordkeeping systems...Recent research also illustrates that strategies and tactics for electronic recordkeeping rarely involve a simple choice between policy, standards, systems design and implementation. Rather, archivists and records managers need to pursue the right combinations of policies, standards, and system design methodologies that organizations can implement and that offer solutions which are affordable and commensurate with the risks and benefits involved.⁵⁴

Hedstrom further posited that organizations “are seeking trusted recordkeeping systems that follow rules for records creation, maintenance, and preservation at all times.” And that:

The expansion of electronic commerce into personal and retail consumption depends...on the ability of individuals and organizations to communicate and conduct business using trusted systems that are not predicated on prior established relationships or formal contractual agreements.

Hedstrom’s points underscored the need for trusted systems that follow rules and that are surrounded by appropriate combinations of policies and standards so that risk, benefit, and cost are balanced. Communication (between people and businesses) and new tools developed for electronic commerce may well enhance our ability to build effective systems for long-term archiving services.

In a 1998 paper (revised in 2000) describing the Stanford Archival Vault prototype, Stanford computer scientists Brian Cooper, Arturo Crespo, and Hector Garcia-Molina articulated the steps necessary to implement a reliable digital object repository.⁵⁵ Emphasizing replication strategies, the authors argued convincingly that their methodology “provides an extremely reliable storage infrastructure for preserving digital objects, even as hardware, software and organizations evolve.” They identified these key factors:

- Avoidance of erasure (including deletion and overwriting by users) through write-once policies.
- Remote backup agreements that incorporate replication policies by the remote system provider.

A “reliability layer” within the distributed archival repository architecture encompasses a series of functions and mechanisms that the authors believe results in a reliable environment for preserved objects. Some of the functions identified include:

- Detection and restoration of missing/corrupted information.
- Communications among trusted components.
- User security, intellectual property management, query processing.
- Import/export facility to move objects into and out of the store.

Here was another bid for three definitional components: communication, security, and replication. Communication in this instance involves exchanges between parts of a system and also between federated member systems. The authors also introduced two new

⁵⁴ Margaret Hedstrom, “Building Recordkeeping Systems: Archivists are Not Alone on the Wild Frontier,” *Archivaria* (1998): 44-71.

⁵⁵ Brian Cooper, Arturo Crespo, and Hector Garcia-Molina, Hector, “Implementing a Reliable Digital Object Archive,” in *Proceedings of the Fourth European Conference on Research and Development in Digital Libraries (ECDL)* (2000) dbpubs.stanford.edu/pub/2000-28.

concepts: backup policies and avoiding, detecting, and restoring lost or corrupted information.

More recently, Cooper and Garcia-Molina expanded on their thinking about reliable digital object archives and discussed the benefits of creating replication partnerships that function as peer-to-peer trading networks.⁵⁶ In describing how these networks would interwork, the authors outlined the challenge of estimating the reliability of each partner and suggested factors that should be considered:

- Frequency of past failures (loss of data) as a predictor of future failures.
- Use of reliable hardware (disks).
- Presence of successful security measures.
- Reputation (perceived reliability).

Crespo and Garcia-Molina further refined their thinking about reliability and trust by identifying reputation and performance as important factors in determining trustworthiness.

In May 2000 the Digital Library Federation (DLF) proposed a set of Minimum Criteria for an Archival Repository of Digital Scholarly Journals.⁵⁷ Focusing on only one class of digital object (digital scholarly journals), their seven criteria include a mix of definitional and functional requirements. Two criteria refer to the defining characteristics of the organization itself:

Criterion 1. A digital archival repository . . . will be a trusted party that conforms to minimum requirements agreed to by both scholarly publishers and libraries.

Criterion 2 A repository will define its mission with regard to the needs of scholarly publishers and research libraries. It will also be explicit about which scholarly publications it is willing to archive and for whom they are being archived.

The DLF placed trustworthiness at the center of its requirements, although it did not specify exactly how trust should be demonstrated. Further work involving experts, institutions, and publishers in designing trusted digital archives for journal publications, is now funded by The Andrew W. Mellon Foundation.

Trustworthiness here takes on a fuller aspect with the addition of notions about organizational mission, the stipulation of agreements between creators and providers, and the assertion that trusted repositories will openly share information about what they are preserving and for whom.

The US Defense Advanced Research Projects Agency (DARPA) described its Information Assurance and Survivability (IA&S) technologies while referencing a new solicitation for a Composable High Assurance Trusted Systems program in a report issued in March 2001. In defining its overall aspirations, DARPA summarized its thinking about trusted systems in this way:

⁵⁶ Brian Cooper and Hector Garcia-Molina, "Creating trading networks of digital archives," in *1st ACM/IEEE Joint Conference on Digital Libraries* (2001) dbpubs.stanford.edu/pub/2001-23.

⁵⁷ Digital Library Federation, *Minimum Criteria for an Archival Repository of Digital Scholarly Journals*, version 1.2 (May 2000) www.clir.org/diglib/preserve/criteria.htm.

Confidence in future systems must be achieved through system and network-level technologies involving approaches such as layered complementary mechanisms that will be cost-effective and scalable within three to five years. Proposed approaches must demonstrate the ability to support the advanced functionality of future trusted systems while maintaining a high level of confidence in the protection of these systems.⁵⁸

This reference enlarged the body of related terms and concepts by incorporating complementarity, cost-effectiveness, scalability, and protection.

Dr. Audun Jøsang, senior research scientist at Queensland University of Technology, has spent much of the past fifteen years developing models and tools for evaluating trusted systems. One of his most innovative contributions is his *A Metric for Trusted Systems* wherein he proposed a formal model for quantifying subjective beliefs about trustworthiness through the use of evaluative processes.⁵⁹ According to Jøsang, “Trust is a subjective belief [and] trust management for open computer networks . . . includes . . . the factors which influence users’ trust in web sites and e-commerce.” Jøsang’s work involves the development of what he calls Subjective Logic and a trust inference engine based on Subjective Logic “to assist users and organisations to make trust assessments about remote parties on the Internet.”⁶⁰ He suggests that security evaluation is a well-understood method for determining trust in implemented system components and that a successful evaluation leads to the determination of an assurance level that reflects the trustworthiness of a system component.

Jøsang’s work deepens a collective understanding of trust and at the same time provides tools to ensure that *evaluation of system components* is methodical and trustworthy in itself.

⁵⁸ Defense Advanced Research Projects Agency (DARPA), BAA #01-24, *Composable High Assurance Trusted Systems, CBD Reference* (March 2001) www.darpa.mil/ito/Solicitations/CBD_01-24.html.

⁵⁹ A. Jøsang and S.J. Knapskog, *A Metric for Trusted Systems* (1998), available from citeseer.nj.nec.com/129647.html.

⁶⁰ Audun Jøsang, security.dstc.edu.au/staff/ajosang/.

Appendix C: Operational Responsibilities Checklist

Negotiating for and Accepting Information from Content Providers

- Well-documented and agreed-on policies about what is selected for deposit, including, where appropriate, specific required formats.
- Effective procedures and workflows for obtaining copyright clearance for both short-term and immediate access, as necessary, and preservation.
- A comprehensive metadata specification and agreed-on standards for its implementation.
- Procedures and systems for ensuring the authenticity of submitted materials.
- Initial assessment of the completeness of the submission.
- Effective record keeping of all transactions, including ongoing relationships, with content providers.

Obtaining Sufficient Control of the Information

- Detailed analysis of an object or class of objects to assess its significant properties. Analysis should be automated as much as possible and informed by the collections management policy, rights clearances, the designated community's knowledge base, and policy restrictions on specific file formats.
- Verification and creation of bibliographic and technical metadata and documentation to support the long-term preservation of the digital object according to its significant properties and underlying technology or abstract form, with monitoring and updating of metadata as necessary to reflect changes in technology or access arrangements. This involves understanding how strategies for continuing access, such as migration and emulation, influence the creation of preservation metadata.
- A robust system of unique identification.
- A reliable method for encapsulating the digital object with its metadata in the archive.
- A reliable archival storage facility, including an ongoing program of media refreshment; a program of monitoring media; geographically distributed backup systems; routine authenticity and integrity checking of the stored object; disaster preparedness; response and recovery policies and procedures; and security.

Determining the Repository's Designated Community

- Analysis and documentation of the repository's current designated community as well as the possible needs and modes of access of future users.
- For federated or cooperating repositories, a shared understanding of the designated communities that are to be served.

Ensuring the Information to Be Preserved Is Independently Understandable

- Well-maintained and documented technical metadata that is kept aligned with the knowledge base of the designated community and with changing technologies.
- A “technology watch” to manage the risk as technology evolves and to provide continuing access and updated methods of access as necessary, such as new migrations or emulators.

Following Documented Policies and Procedures

- Policies for collections development (e.g., selection and retention) that link to technical procedures about how and at what level materials are preserved and how access is provided both short and long term.
- Policies for access control to ensure all parties are protected, including authentication of users and disseminated materials.
- Policies for storage of materials, including service-level agreements with external suppliers.
- Policies that define the repository’s designated community and describe its knowledge base.
- A rigorous system for updating policies and procedure in accordance with changes in technology and in the repository’s designated community.
- Explicit links between these policies and procedures, allowing for easy application across heterogeneous collections.

Making the Preserved Information Available to the Designated Community

- A system for discovery of resources.
- Appropriate mechanisms for authentication of the digital materials.
- Access control mechanisms in accordance with licenses and laws, and an “access rights watch.”
- Mechanisms for managing electronic commerce.
- User support programs.

Advocating Good Practice in the Creation of Digital Resources

- Effective mechanisms for advocating good practice for content providers.

Glossary

Archival Information Package (AIP): An Information Package, consisting of the Content Information and the associated Preservation Description Information, which is preserved within an OAIS.

Archival Storage: The OAIS entity that contains the services and functions for the storage and retrieval of Archival Information Packages.

Content Information: The set of information that is the primary target for preservation. It is composed of a Data Object and its Representation Information. For example, Content Information could be a single table of numbers representing, and understandable as, temperatures, but excluding the documentation that would explain its history and origin, how it relates to other observations, etc.

Context Information: The information that documents the relationships of the Content Information to its environment, including why the Content Information was created and how it relates to other Content Information.

Data migration: One current strategy for providing continuing access to archived digital materials over time. It is perceived to be the most reliable strategy for continuing access to many types of digital materials because it has been used for years for routine migration of homogeneous digital materials. However, the library community lacks documented practical experience that shows it is a reliable approach for heterogeneous digital collections such as multimedia CD-ROMs or electronic journals, so it has yet to see widespread adoption for these materials. For the purposes of this report, the definition of “data migration” is based on that provided in the CPA/RLG report: “a set of organized tasks designed to achieve the periodic transfer of digital materials from one hardware/software configuration to another, or from one generation of computer technology to a subsequent generation.”⁶¹

Although in some cases data migration is straightforward, it can be complex and may not be reversible, which will have important implications for its effectiveness as a continuing access strategy. For example, in the early 1990s many organizations with files stored in EBCDIC code opted for immediate migration to ASCII despite the fact that the ASCII character set didn't allow for all character information from EBCDIC files to be migrated. This was not a reversible migration and some information was lost. When UNICODE was subsequently introduced, its character set did allow for all EBCDIC characters to be represented.

Organizations that had discarded the original EBCDIC files lost information irretrievably.⁶²

Currently, there is a lack of practical experience to provide guidance on migrating complex digital materials through different data formats. CAMiLEON and Cornell University have some initial results comparing different approaches to data migration.⁶³

⁶¹ John Garrett and Donald Waters, *Preserving Digital Information: Report of the Task Force on Archiving of Digital Information* (Commission on Preservation and Access and RLG, 1996) www.rlg.org/ArchTF/index.html.

⁶² Example adapted from David Holdsworth and Derek M. Sergeant, *A Blueprint for Representation Information in the OAIS Model* (2000) esdis-it.gsfc.nasa.gov/MSST/conf2000/PAPERS/D02PA.PDF.

⁶³ CAMiLEON, part of the JISC/NSF International Digital Libraries Program (DLI2), is a partnership between the University of Michigan and the University of Leeds to explore use of emulation as a strategy for digital

Data object: A digital object (can also be a physical object).

Designated community: An identified group of potential users of the archive's contents who should be able to understand a particular set of information. The designated community may be composed of multiple user communities.

Digital preservation strategy: A digital preservation strategy is a particular technical approach for providing continued access to archived digital materials. At this time, three main strategies are used to keep materials within the repository "fresh" and ensure that they are accessible using current technology: data migration, persistent object transformation, and technology emulation.

Dissemination Information Package (DIP): The Information Package, derived from one or more AIPs, received by the user in response to a request to the repository.

Fixity Information: The information that documents the authentication mechanisms and provides authentication keys to ensure that the Content Information object has not been altered in an undocumented manner. For example, a CRC code for a file or a checksum.

Independently understandable: A characteristic of information that has sufficient documentation to allow the information to be understood by the designated community without having to resort to special resources not widely available, including named individuals.

Information Package: Content Information and associated Preservation Description Information that is needed to aid in the preservation of the Content Information. The Information Package has associated Packaging Information used to delimit and identify the Content Information and Preservation Description Information.

Ingest: The OAIS entity that contains the services and functions that accept Submission Information Packages from producers, prepare Archival Information Packages for storage, and ensure that Archival Information Packages and their supporting Descriptive Information become established within the OAIS-compliant repository.

Knowledge base: A set of information, incorporated by a person or system, that allows that person or system to understand received information.

Long term: A period long enough to raise concern about the effect of changing technologies, including support for new media and data formats, and of a changing user community.

Long-term preservation: The act of maintaining correct and independently understandable information over the long term.

Metadata for digital preservation: The effective management and use of digital resources in a repository will rely on a robust system of resource description—for resource discovery, access, and preservation. Metadata research continues to generate interest worldwide; to date, most activity has focused on metadata for resource discovery (e.g., MARC, Dublin Core, CIDOC). However, there is increasing awareness that reliable digital repositories will depend on the creation and storage of information required to support a chosen preservation strategy, such as migration, emulation, or technology preservation. This information will

preservation, see www.si.umich.edu/CAMILEON/; Gregory W Lawrence et al., *Risk Management of Digital Information: A File Format Investigation* (CLIR, June 2000) www.clir.org/pubs/reports/pub93/contents.html; Paul Wheatley, *Migration—a CAMiLEON discussion paper*, www.leeds.ac.uk/camileon/.

need to describe the data in detail including, broadly speaking, both descriptive and structural metadata. Although these have been defined in different ways, within the context of the OAIS Model, preservation metadata takes two forms:

- Preservation Descriptive Information, which includes general resource description as well as rights management information and descriptions of actions taken for the purposes of preservation.
- Representation Information, which maps the stored data into more meaningful concepts; that is, systems information that renders bits and bytes into a meaningful digital object. For example, the ASCII definition, which maps data (bits) into readable symbols.

Open Archival Information System (OAIS) Reference Model: Developed by the Consultative Committee on Space Data, a conceptual framework and reference tool for defining a digital repository. It provides a model of the environment, functions, and data types for implementing a digital repository. The OAIS is undergoing approval as an ISO standard and its publication as an international standard is expected later this year.

Persistent object transformation: Although not widely used in the library and archives community, a strategy for providing continuing access that has been widely publicized internationally. Persistent object transformation may appear to be migration by another name; in fact, it takes a longer-term approach, focusing not on the object's technical environment or on moving it into current technology, but attempting to define the essential attributes and methods of the object. These attributes and methods are then made explicit within the objects themselves through tagging and/or encapsulation that are independent of the object's current or original technical infrastructure.⁶⁴

Preservation Description Information (PDI): The information that is necessary for adequate preservation of the Content Information; it can be categorized as Provenance, Reference, Fixity, and Context Information.

Producers: The people or systems that provide information to the repository.

Provenance Information: The information that documents the history of the Content Information. This information tells the origin or source of the Content Information, any changes that may have taken place since it was originated, and who has had custody of it since it was originated. Examples of Provenance Information are the principal investigator who recorded the data, and the information concerning its storage, handling, and migration.

Reference Information: The information that identifies and, if necessary, describes one of more mechanisms used to provide assigned identifiers for the Content Information. It also provides identifiers that allow outside systems to refer unambiguously to a particular Content Information; for example, an ISBN.

Rendering software: Software that displays Representation Information of an Information Object in forms understandable to humans.

Repository: An organization that intends to maintain information for access and use.

⁶⁴ Reagan Moore et al., "Collection-Based Persistent Digital Archives—Part 1," *D-Lib Magazine* 6, no. 3 (March 2000) www.dlib.org/dlib/march00/moore/03moore-pt1.html, and Reagan Moore et al., "Collection-Based Persistent Digital Archives—Part 2," *D-Lib Magazine* 6, no. 4 (April 2000) www.dlib.org/dlib/april00/moore/04moore-pt2.html.

Representation Information: The information that maps a data object into more meaningful concepts. For example, the ASCII definition that describes how a sequence of bits (i.e., data object) is mapped into a symbol.

Representation Network: The full set of Representation Information that describes the meaning of a digital object. This can apply to a single data object or to an entire repository.

Significant properties: The technical and other characteristics of the digital object that the depositor and/or repository agree to be most important for preservation over time. For digital materials, simply maintaining a bytestream does not ensure the digital object will be preserved at a level acceptable to the repository and its users. A digital object's significant properties are not assumed to be absolute; repositories will make judgments that fulfill their preservation responsibilities and meet the needs of their user communities and the wishes of the depositor. Significant properties may apply to a single digital object or to an entire class of digital objects within a repository. This term was first coined, defined, and employed by Cedars.

Submission Agreement: An agreement reached between an OAIS-compliant repository and the producer that specifies a data model for a data submission. This data model identifies format/contents and the logical constructs used by the producer and how they are represented on each media delivery or in a telecommunication session.

Submission Information Package (SIP): An Information Package that is delivered by the producer to the repository for use in the construction of one or more AIPs.

Technology emulation: A strategy for continuing access to digital materials that mimics or re-creates the digital object's original technical environment using current technology. Access to the object relies on a copy of the original bytestream (as deposited) and an emulation of its original operating environment. Emulation can take place at either the hardware or software level. Emulation may be particularly useful for preserving the "look and feel" of the object; however, like migration, there is little if any practical experience in applying emulation in a production environment. Preliminary work by IBM and others suggests that emulation may be the most practical (and in some cases the only) strategy for some more complex or esoteric digital materials, particularly for preserving the original functionality.⁶⁵ Recent findings of the CAMiLEON project challenge the perceived wisdom that migration and emulation are two completely different approaches: preliminary work with obsolete computer systems suggests that, as we build a more complete understanding of the various methods for migrating data, it will become clear that emulation and migration represent related approaches on a graduated scale.⁶⁶ It may be that skepticism about emulation will be revealed as a lack of understanding of the complexities of data migration.

⁶⁵ Raymond A. Lorie, *The Long Term Preservation of Digital Information* (New York: ACM Press, 2001) doi.acm.org/10.1145/379437.379726. For more on different approaches to using emulation for digital preservation, see David Holdsworth and Paul Wheatley, *Emulation, Preservation and Abstraction*, www.leeds.ac.uk/camileon/; and Jeff Rothenberg, *Avoiding Technological Quicksand: Finding a Viable Technical Foundation for Digital Preservation* (CLIR, January 1999) www.clir.org/pubs/reports/rothenberg/contents.html.

⁶⁶ Paul Wheatley, *Migration—a CAMiLEON discussion paper*, www.leeds.ac.uk/camileon/.

Selected Resources

Projects

CAMiLEON (Creative Archiving at Michigan and Leeds: Emulating the Old on the New): www.si.umich.edu/CAMILEON/

Cedars (CURL Exemplars in Digital Archives): www.leeds.ac.uk/cedars/

NEDLIB (Networked European Deposit Library): www.kb.nl/coop/nedlib/

PANDORA (Preserving and Accessing Networked Documentary Resources of Australia)
Project: pandora.nla.gov.au/index.html

Preserving Access to Digital Information (PADI): www.nla.gov.au/padi/

Publications

Beagrie, Neil and Daniel Greenstein, *A Strategic Policy Framework for Creating and Preservation Digital Collections*, E:Lib Supporting Study P3 (London: Library and Information Technology Centre, 1998) ahds.ac.uk/strategic.htm.

Consultative Committee for Space Data Systems (CCSDS), *Reference Model for an Open Archival Information System (OAIS)* (July 2001) www.ccsds.org/documents/pdf/CCSDS-650.0-R-2.pdf.

Garrett, John and Donald Waters, *Preserving Digital Information: Report of the Task Force on Archiving of Digital Information*. (Washington, DC: Commission on Preservation and Access, and Mountain View, CA: RLG, 1996) www.rlg.org/ArchTF/index.html.

Gatenby, Pam "Digital Archiving: Developing Policy and Best Practice Guidelines at the National Library of Australia" (paper presented at An Interactive Workshop sponsored by ICSTI and ICSU Press, January 2000) www.icsti.org/2000workshop/gatenby.html.

Hedstrom, Margaret and Sheon Montgomery, *Digital Preservation Needs and Requirements in RLG Member Institutions* (Mountain View, CA: RLG, 1998) www.rlg.org/preserv/digpres.html.

Jones, Maggie and Neil Beagrie, *Preservation Management of Digital Materials Workbook* (London: Re:source, 2000) www.jisc.ac.uk/dner/preservation/workbook/.

OCLC/RLG Working Group on Preservation Metadata, *Preservation Metadata for Digital Objects: A Review of the State of the Art. A White Paper* (January 2001) www.oclc.org/research/pmwg/presmeta_wp.pdf.

Preservation 2000: An International Conference on the Preservation and Long Term Accessibility of Digital Materials, Conference Papers (2000) www.rlg.org/events/pres-2000/prespapers.html.

Public Record Office, *Management, Appraisal and Preservation of Electronic Records* (Kew: Public Record Office, 1999) www.pro.gov.uk/recordsmanagement/eros/guidelines/.

RLG DigiNews, a bimonthly newsletter: www.rlg.org/preserv/diginews.

Rothenberg, Jeff, *Avoiding Technological Quicksand: Finding a Viable Technical Foundation for Digital Preservation* (Washington, DC: Council on Library and Information Resources, 1999)
www.clir.org/pubs/reports/rothenberg/contents.html.