

**Erik Mitchell**

---

**Metadata  
Standards and  
Web Services  
in Libraries,  
Archives,  
and Museums**

**An Active Learning Resource**

---





**Metadata Standards  
and Web Services in Libraries,  
Archives, and Museums**



# **METADATA STANDARDS AND WEB SERVICES IN LIBRARIES, ARCHIVES, AND MUSEUMS**

*An Active Learning Resource*

Erik Mitchell



An Imprint of ABC-CLIO, LLC  
Santa Barbara, California • Denver, Colorado

Copyright © 2015 by Erik Mitchell

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, except for the inclusion of brief quotations in a review, without prior permission in writing from the publisher.

**Library of Congress Cataloging-in-Publication Data**

Mitchell, Erik T., 1972–

Metadata standards and Web services in libraries, archives, and museums : an active learning resource / Erik Mitchell.

pages cm

Includes bibliographical references and index.

ISBN 978-1-61069-449-0 (paperback) — ISBN 978-1-61069-450-6 (ebook) 1. Metadata.

2. Information organization. 3. Web services—Library applications. 4. Museums—Information technology. 5. Archives—Information technology. I. Title.

Z666.7.M58 2015

025.3—dc23 2015015545

ISBN: 978-1-61069-449-0

EISBN: 978-1-61069-450-6

19 18 17 16 15 1 2 3 4 5

This book is also available on the World Wide Web as an eBook.

Visit [www.abc-clio.com](http://www.abc-clio.com) for details.

Libraries Unlimited

An Imprint of ABC-CLIO, LLC

ABC-CLIO, LLC

130 Cremona Drive, P.O. Box 1911

Santa Barbara, California 93116-1911

This book is printed on acid-free paper (∞)

Manufactured in the United States of America

# CONTENTS

Chapter 1	Introduction to the World of Digital Information Organization	1
	<i>What This Book Is About</i>	1
	<i>Metadata and Its Roles in Everyday Life</i>	3
	<i>What Is Information?</i>	4
	<i>The Information Lifecycle</i>	6
	<i>"Aboutness" and the Role of Classification in Information</i>	9
	<i>The Internet and Its Impact on Information Use</i>	10
	<i>The Purpose of This Book and What You Can Expect</i>	13
	<i>Notes</i>	14
Chapter 2	Information Systems as Boundary Objects	15
	<i>Information Seeking Behavior</i>	15
	<i>Process-Based Models</i>	17
	<i>Cognitive and Affective Models</i>	18
	<i>Information Seeking in Context</i>	19
	<i>Connections between Information Seeking and Technology</i>	20
	<i>Understanding Document Structure</i>	21
	<i>Digital Document Case Study: HTML</i>	22
	<i>Document Structure Overview</i>	24
	<i>HTML in Context: Web Browsers, Web Servers, and the HTTP Protocol</i>	26
	<i>Verifying Adherence to the HTML Schema and Serialization Standards</i>	27
	<i>Conclusion</i>	28
	<i>Notes</i>	29

Chapter 3	Design of Information Systems	31
	<i>Information Services in Libraries: The Integrated Library System and MARC</i>	31
	<i>Components of Information Service Design</i>	36
	<i>Exploring the Connection between Document Structure and Use</i>	38
	<i>The Model–View–Controller Paradigm</i>	47
	<i>The Document Object Model (DOM)</i>	49
	<i>JavaScript</i>	51
	<i>Conclusion</i>	57
	<i>Notes</i>	57
Chapter 4	Information Organization Models	59
	<i>A Broad Model of Organized Information</i>	60
	<i>Cataloging Principles</i>	65
	<i>Content Rules</i>	73
	<i>General Data Models</i>	83
	<i>Conclusion</i>	95
	<i>Notes</i>	96
Chapter 5	Metadata Standards, Contents, and Values	99
	<i>The Relationship between Metadata, Resources, and Information Systems</i>	99
	<i>Metadata Schemas and Their Roles in Information Systems</i>	103
	<i>Vocabularies</i>	103
	<i>Taxonomies</i>	104
	<i>Thesauri</i>	105
	<i>Ontologies</i>	105
	<i>Types of Metadata Schemas</i>	105
	<i>Cataloging Principles and Metadata Schemas</i>	107
	<i>Engaging in Metadata Creation—Selecting and Applying Schema</i>	107
	<i>The Structure of a Metadata Schema</i>	108
	<i>Exploration of Specific Metadata Schemas, Vocabularies, Thesauri, Taxonomies, and Ontologies</i>	117
	<i>Understanding Types of Information Organization Structures</i>	117
	<i>Metadata Schemas Common in Libraries, Archives, and Museums</i>	117
	<i>Creating Metadata</i>	127
	<i>Contents and Values</i>	129
	<i>Relationships between Metadata Schemas, Vocabularies, Thesauri, Taxonomies, and Ontologies</i>	130
	<i>RDF and Linked Data</i>	132
	<i>RDF Records Serialized as XML</i>	133
	<i>RDF Schema (RDFS)</i>	133
	<i>RDF Data Structures</i>	134



<b>Contents</b>	<b>vii</b>
<i>Example of RDF in Use</i>	136
<i>RDF-based Databases</i>	138
<i>RDF and Vocabularies</i>	139
<i>Conclusion</i>	140
<i>Notes</i>	141
Chapter 6 <i>Serialization</i>	143
<i>Serialization and Exchange Formats</i>	143
<i>History and Future Directions of Commonly         Used Serialization Formats</i>	155
<i>Serialization and Available Technology</i>	171
<i>Linked Data in LAM Institutions</i>	173
<i>Metadata Exchange and Web Services</i>	175
<i>Conclusion</i>	184
<i>Notes</i>	185
Chapter 7 <i>Creating, Using, and Evaluating Metadata in Digital Information Systems</i>	187
<i>Quality Control</i>	188
<i>Digital Libraries</i>	205
<i>Examples of Large-scale Digital Information Systems</i>	209
<i>Conclusion</i>	215
<i>Notes</i>	215
Chapter 8 <i>Using Metadata to Create Information Services</i>	217
<i>Review of Programming Foundations</i>	217
<i>The Extensible Stylesheet Language Family</i>	225
<i>Conclusion</i>	248
<i>Notes</i>	248
Chapter 9 <i>Future Trends in Information Systems, Metadata, and Information Use</i>	251
<i>Library Systems and Metadata Are Increasingly Open</i>	252
<i>LAM Information Resources and Metadata Are         Becoming Increasingly Networked</i>	253
<i>Demonstrating Value Will Be a Key Challenge for         LAM Institutions around Metadata and         Information System Design</i>	254
<i>LAM Institutions Must Be Creative to Address         Metadata Quality Issues</i>	255
<i>Cloud Computing Will Continue to Influence         Technology, Metadata, and Information Services</i>	258
<i>Researchers Need New Metadata and Information         Gathering Platforms</i>	260

<i>Technology Capabilities Lag Behind Innovation in</i>	
<i>Data Storage and Sharing</i>	261
<i>Whither LAM Metadata?</i>	262
<i>LAM Communities Will See Continued Change in</i>	
<i>User Needs and Information Seeking</i>	263
<i>Library and Information Science Draws on</i>	
<i>Expertise in Multiple Domains</i>	264
<i>Notes</i>	265
 <i>Bibliography</i>	 267
<i>Index</i>	273

## **Chapter 1**

# **INTRODUCTION TO THE WORLD OF DIGITAL INFORMATION ORGANIZATION**

### **WHAT THIS BOOK IS ABOUT**

The field of library and information science (LIS) is concerned with the roles, uses, and impact of information in people, across communities and throughout history. Librarianship and other information careers have built a shared understanding of what information is, how it is created and used, and what impact it can have on individuals and communities. In addition, for much of the 20th century, libraries, archives, museums (LAMs) and similar cultural heritage institutions have shared a consistent view of information management, albeit not necessarily consistent systems and standards. Nevertheless, this shared vision has enabled LAM institutions to move quickly into new fields such as information retrieval, information systems analysis, and data curation as the practice of information creation, management, and use changed. Along the way, LAM institutions have had to redefine their information systems to meet shifting technologies and user needs and have had to develop new systems and standards to accommodate digital information resources and networked information communities. To cite a familiar example, library catalogs, initially developed to help users find and gain access to resources, have been re-invented multiple times as information technology changed. Published first on scrolls, later in books, and more recently on catalog cards, in microfiche, and ultimately as web-based applications, library catalogs have sought to provide an important service to library users. In the last decade, however, the needs of information communities have changed as have the formats of the information materials they work with. The development of the Internet into a user-initiated publishing platform and a network of resources turned search engines by Yahoo and Google into world-wide catalogs of information that allowed users to find information in this *de facto* virtual library created by what is generally called the World Wide Web, or simply *the web*. The evolution of the web and the information

services around it meant that issues that drove the creation of LAM catalogs, such as metadata creation and management, user-focused information needs analysis, and the creation and management of digital documents, became important areas of scholarship and development outside of the LIS domain. It also meant that traditional approaches to description needed to be reconsidered for increasingly complex resources.

The growth of prominence of both information and information technology in our daily lives has also meant that librarians, archivists and other information professionals need to be aware of the digital literacies that support information creation, management, and use. In some cases this means being able to design, build, manage, and facilitate use of digital information tools and resources. Every day we use websites that effortlessly allow us to search, browse, select, and use information. We most likely all have our preferred web search engine that prioritizes results for us and may even show us selected e-mails or file results from our cloud-stored personal documents. The simplicity of these sites masks the complexity underneath them. Designing effective information services involves a complex understanding of user needs, technology tools, document and metadata standards, and modes of information interaction. It means understanding the relationship between web design, information organization, and information technology issues. Underneath each of these services there is a world of technology and information-seeking/interaction theory that guides how we interact with the service as well as allowing us understand what information the service can make available to us. These interactions are supported by information systems that store documents, create and manage representations of these documents for search and retrieval, and provide discovery and interaction methods that make the service usable. Information system designers should design systems by thinking about potential users who have information needs, questions, and abilities that guide their use of the system. Between these systems that store and serve resources and a user's information need we have an interaction guided by information-seeking behaviors, human-computer interactions, and personal information-management practices. Many of these systems and services run on computers that are designed to serve thousands or hundreds of thousands of users at the same time; however, the tools and technologies are not so complex that they cannot run on almost any laptop or desktop computer. These two elements combined, the software and the hardware, could be considered to be an "information system."

This book explores the world of information organization from the perspective of digital information service design, seeking ways to understand the shift of information management activities in an expanding universe of digital information. In order to do this it combines technical and conceptual content to help the reader better understand what information organization tasks are and how information systems affect how we use information. The book will focus upon the theories and mechanics of information creation, its use, and reuse in a digital age. The goal is to guide the reader through the domains of information organization, information technology, and information interaction and to explore these domains as the three building blocks of information systems and services. In exploring information

system design, this book explores in some detail the connection between metadata and web services, demonstrating how metadata and web services are created, used, and reused.

## METADATA AND ITS ROLES IN EVERYDAY LIFE

The simplest definition of *metadata*, and perhaps the easiest to grasp conceptually, is “data about data.” Yet, that simple definition fails to reveal what is, in fact, a very complicated and sometimes controversial area of study and practice. Metadata is created every time we add context to information or data, and those who create metadata are often unaware that they are creating it. In the days of the card catalog metadata was something that was created from the title page of a book and captured on a 3x5 card. In our increasingly digital society metadata is generated with every phone call, with every social media post, with every saved document on your computer. The word metadata famously surfaced in 2013 when news of the National Security Agency’s collection of information about every phone call made in the United States emerged. It is likely that before the president of the United States gave details about the gathering of metadata for these phone calls people would not have even considered that the simple act of pushing that “send” button on their mobile phone leaves a record of the number called, the number calling, the duration of the call, and the location where the call was made. Perhaps more importantly, the term *metadata* likely would not have entered the popular vernacular.<sup>1</sup>

In the digital world, metadata surrounds every document, every interaction, and every communication. For example, a modern digital camera will attach metadata to the image file it produces using the method provided by the image format’s standard. Such information can include a timestamp, geolocation information if the camera has GPS support, image compression technique used, F-stop setting, aperture setting, and whether the camera’s flash was used.<sup>2</sup> If a person is somewhat old fashioned and buys music CDs instead of downloading them from Amazon or Apple’s iTunes Store, the process of *ripping* the CD—that is, reading each audio track and compressing it to an MP3, AAC, or Ogg Vorbis format—necessarily requires the creation and attachment of metadata to the new audio file. The metadata is most often obtained by the querying of Gracenote’s Compact Disk Database (CDDDB) by the ripping program.<sup>3</sup> When the ripping program is, for example, Apple’s iTunes, the compressed tracks are stored in the user’s iTunes library and the metadata is stored in the library’s database. However, the metadata is also stored in the file itself, making it accessible to a portable music player. The act of creating a new written document with office software such as Microsoft Office also offers an opportunity for inserting metadata about the document—information such as the author’s name, contact information, document title, subject, keywords, and comments.

Although metadata has become ubiquitous and is often created using automated techniques, the use of metadata creation remain the same as they were in a pre-digital age. For the university student of twenty years ago, conducting

research would require searching the library's catalog as well as journal indexes—both collections of metadata, in this case, data about books and journals in the stacks and the information within. The catalogs and indexes are metadata about the holdings of a library's traditional printed materials and the articles contained within the journals in the stacks. Indeed, a library catalog may be the first example of the creation of metadata, produced by the librarian of the first Sumerian clay tablet archive, perhaps nothing more than a list of holdings and the shelf location of each tablet. University students of today have a vastly wider range of information resources at their disposal, yet the problem remains the same: What information is available, where is it located, and what sources of metadata are available to be searched in order to find it? And, of course, when we use a computer to access a file on some sort of storage device such as a hard disk or a thumb drive, we are using metadata in the form of the directory path and file name so that the program, through the underlying operating system, can determine the file's location on that device and open it. One common thread running through the above scenarios is the need for standards.

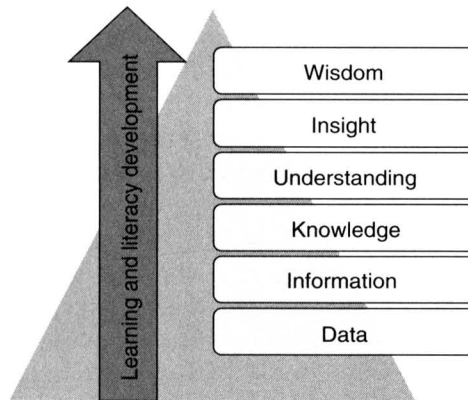
## WHAT IS INFORMATION?

Information is both a very concrete concept—for example, the contents of a printed book—and an ephemeral concept, such as the process of communicating to others, your internal understanding of data, and your memory of what others said to you. Information theory has been a central part of the work of 20th-century information theorists<sup>4</sup> and is a common thread through many works on the evolution of computers and digital information environments. Several detailed histories of information and information technology discuss the roles of key individuals, technologies, and movements in the evolution of information across human history and consider the thread of information as it transforms societies.<sup>5</sup> An overriding theme in these histories is an interest in how *information objects*—including books, scrolls, web pages, and other forms of physical and digital documents—are created, organized, and used. Each of these histories also explores other types of information, including that which forms the foundation of internalized knowledge—information that is collectively shared and understood in communities, and information as an element of a cultural value system. Wright,<sup>6</sup> for example, suggests that documents need not be understood to have impact in communities. Legal documents, for example, are rarely fully understood outside the small communities that create and make use of them and generally contain a jargon not often used by ordinary people; yet these documents and the information contained within them govern interpersonal relationships, business relationships, and have real-world legal outcomes.

The field of information science views the abstract concept of information through a multifaceted lens. Michael Buckland's seminal work exploring the facets of information<sup>7</sup> studied information in three forms: objects, memory, and communication. Buckland's work in this area dates back to the 1990s, when the world

Figure 1.1.

## The Data-to-Wisdom Hierarchy



of information was much different but the underlying theory is as relevant in the digital world as it was when applied to physical information objects.

While Buckland's article views information through a form and formats lens, another popular view of information is as a building block in the continuum from data to wisdom. Figure 1.1 shows an adapted view of this hierarchy that is popularly represented in education literature. This view of information places "information" as a stepping-stone on a hierarchy of increasingly contextualized information. From an information science perspective, knowledge, insight, and wisdom are derivative products of the state of being informed, and data is simply information without context, for example, using a number without specifying any units.

### The Data-to-Wisdom Hierarchy

In the education field, questions about the origin and nature of information focus on differentiating information from data, knowledge, and wisdom, and on aligning the concept of learning with different cognitive outcomes (e.g., remembering, understanding, analysis). This hierarchy was first described by Benjamin Bloom as *Bloom's Taxonomy*<sup>8</sup> and has enjoyed a number of updates to reflect current thinking in educational practice. For example, the initial Bloom's taxonomy viewed evaluation as the highest level of learning and differentiated synthesis from analysis. In updated versions, the highest level of learning is represented as "creating," with "evaluating" serving the penultimate role.<sup>9</sup>

Bloom's Taxonomy can be a useful tool in evaluating user information engagement and learning. In addition to defining cognitive approach to understanding the hierarchy of information engagement (i.e., to remember and/or understand it), information use (i.e., to apply, analyze, and/or evaluate it), and information creation, Bloom's taxonomy discusses affective (e.g., emotional) and psychomotor

(e.g., skills and abilities) characteristics of individuals. These areas map to similar focus areas in the information community, such as Kuhlthau's exploration of the connection between information-seeking activities and affective impact. Her Information Seeking Process (ISP) model is built around the cognitive and affective experience of uncertainty and the effect of increased confidence on the information seeker as the search and the collection of wanted information progress.<sup>10</sup> Kuhlthau and Bloom's works are very much directed at exploring Buckland's cognitive and communicative forms of information. Both Bloom's Taxonomy, which puts structure around the amorphous cognitive state of being informed, and Kuhlthau's synthesis of cognitive, affective, and communicative roles into a holistic view of information seeking help sequence the activities of becoming informed.

Information systems play an important role in these interactions, often serving as a way to search for and filter through information, as a structure to put around information (e.g., a classification system), or a storage mechanism for keeping track of complicated data. In fact, information objects have a strong symbiotic relationship with people and communities of practice. An information object's meaning is reflected in the language used to create it, in the structural cues and prompts used to represent key information (e.g., title, author, tables, pictures, graphs), and in the organizational concepts used to deliver information. These structures serve as a shared, common framework through which information is shared and preserved, elevating the details of the document from data, to information, to wisdom.

## THE INFORMATION LIFECYCLE

When information is recorded and becomes an object to be managed, it is important to consider the steps through which this resource progresses. In LIS this model is often referred to as an information lifecycle. The development and refinement of these models has been a subject of great interest and study. As a result, several lifecycles have been identified, including information creation lifecycles, curation lifecycles, and even information-seeking lifecycles. The lifecycle model fits the concept of information in part because of the connection to Buckland's views of information being a thing, a process, and knowledge, which agree with the notion that information, is created, managed, accessed, and preserved.

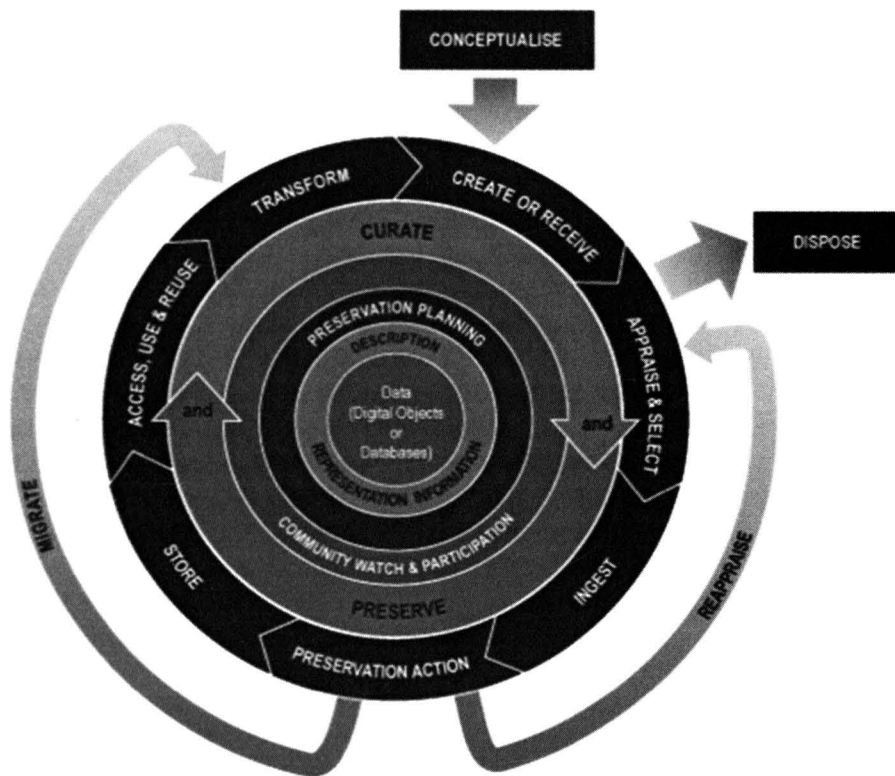
Lifecycle models describe the movement of information through a process (as a thing), emphasize the need to support dissemination and use (again, as a process), and result in the creation of new information about what the next steps users of the information will take as part of their information-seeking activity (as knowledge).

While each of these lifecycle models features different activities, they all share a number of commonalities in that each of these models tends to follow the "plan → acquire → manage → provide access → analysis" approach.<sup>11</sup>

Each of these activities ought not to be thought of as being performed sequentially for a given object or collection but rather as a process of its own that may have multiple step dependencies and may be active at different stages during the information object management life cycle. For example, it is reasonable to expect that



Figure 1.2.

Digital Curation Centre Curation Lifecycle (<http://www.dcc.ac.uk>)

data gathering and analysis related to the information object will occur throughout the object's lifecycle.

In addition, it should not be assumed that activities such as discovery, storage, and preservation are bounded or one-time events. Each of these activities may involve one or more information systems, will likely involve multiple information organization activities and standards, and will likely require analysis of different users' information needs and their respective literacies.

While each of these models treats information organizational activities as part of a much larger system, the actual processes of organization, description, and meta-data creation have their own processes and cycles. For example, the preservation processes in these models often follow the Open Archival Information System (OAIS) model.<sup>12</sup> Each step in a lifecycle model has corresponding standards relating to the digital documents, document management, and document use. Furthermore, a close alignment exists between the steps in an information lifecycle and the standards that are useful in guiding a resource through the process. These standards are drawn from a number of communities and standards bodies, and all of them represent best practice in information management in their respective

Figure 1.3.  
Library and Archives Canada Records Management Model (<http://www.bac-lac.gc.ca>)

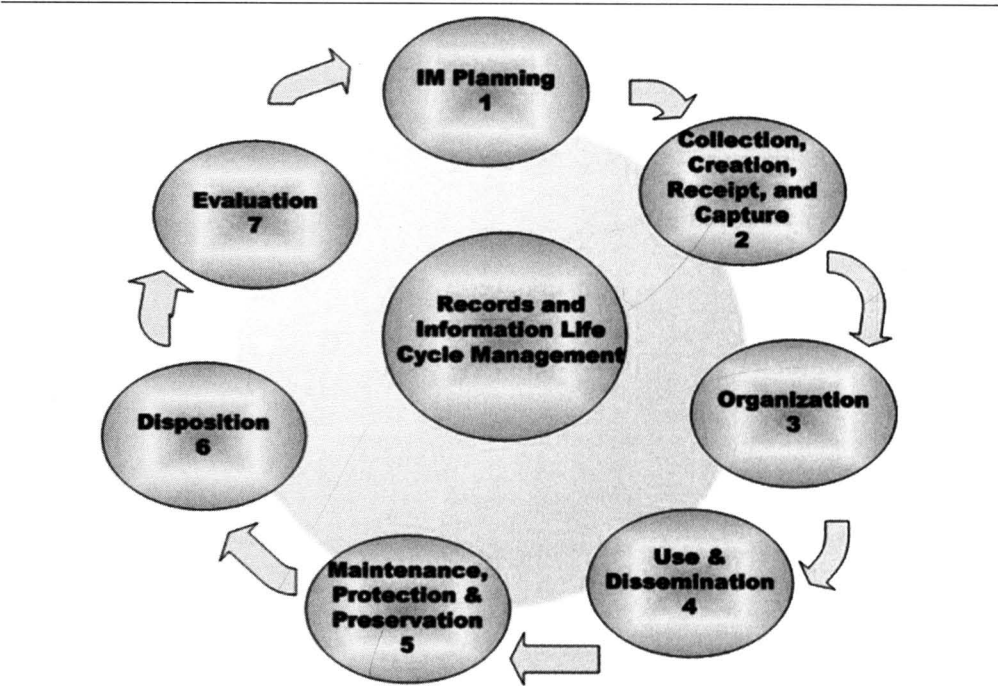
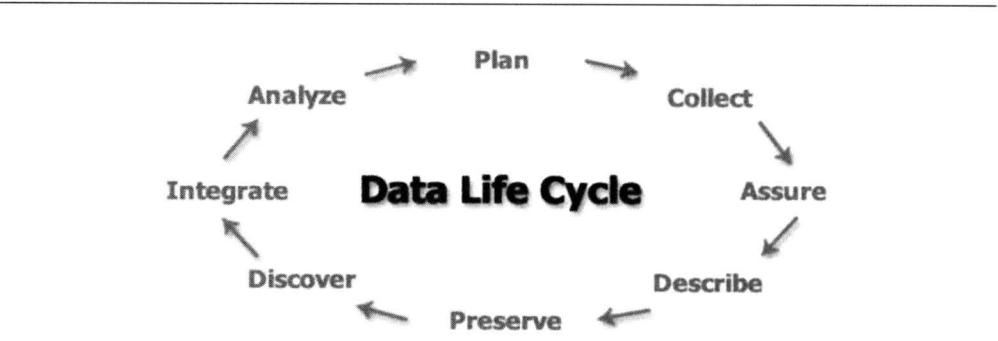


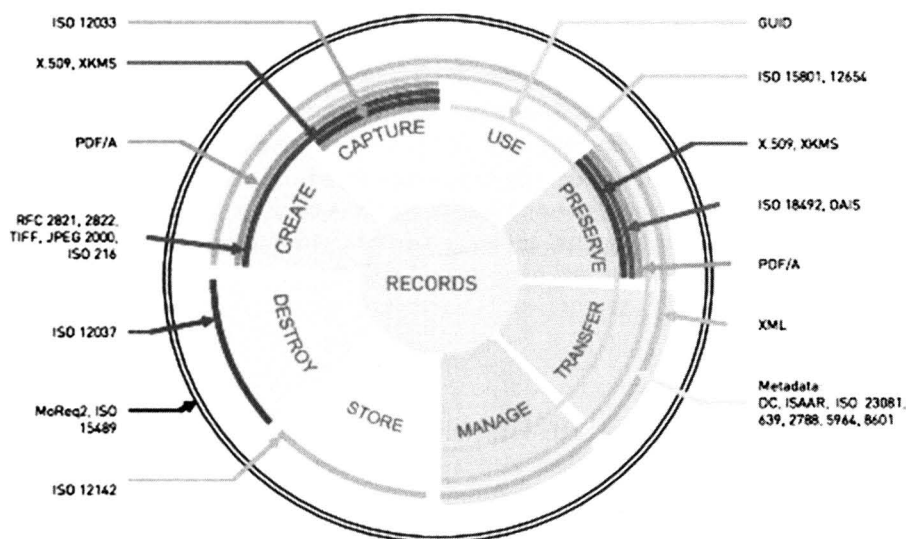
Figure 1.4.  
DataOne Data Life Cycle Model (<https://www.dataone.org>)



fields. For example, the PDF/A (Portable Document Format/Archiving) standard is a more open version of the (PDF) standard, which enables long-term archiving, hence the additional “/A” at the end of its abbreviation. This standard is the preferred document representation format for the creation process in the MoReq2<sup>13</sup> digital records management specification.

Figure 1.5.

MoReq2 Digital Records Management Model



## “ABOUTNESS” AND THE ROLE OF CLASSIFICATION IN INFORMATION

The information lifecycle is a useful framework for thinking the aspects of information and information services. Certain steps in the lifecycle, notably the process of description or discovery, require the creation of metadata that helps represent the “aboutness” of a resource in a discovery system. This process of assigning the “aboutness” of an item is also known as classification. Classification often involves the assignment of a series of subject headings or other categorizing data (e.g., place, date, format) to an item, in many cases also assigning a unique identifier (e.g., a call number in the Library of Congress schema) based on the content of the resource.

Classification is an important part of cataloging but is also a common activity in indexing and abstracting work. The process uses *conceptual analysis* of a work to identify the *aboutness* of a resource with the goal of enabling recall either through the search of information systems such as card catalogs or the accurate placement of a physical resource in context of other resources on the same topic. A sufficiently cataloged item should be cataloged down to the most *specific* level appropriate for a given classification, must be described thoroughly, or *exhaustively*, to ensure that all of the appropriate topics are captured in the classification, and often uses a *controlled vocabulary* to map topics in the resource onto a common classification schema. Classification often uses semantic and syntactic analysis as part of the content analysis process, extracting meaning from the actual content, or *semantics* of a work as well as from the content structure, or *syntactics* of a resource.

There are a number of classification systems commonly used in LAM communities including The Library of Congress system, the Colon Classification System, the Universal Decimal System, Bliss Bibliographic Classification system, and the Dewey Decimal system. Each of these systems focus on identifying the “aboutness” of a document and coding of that aboutness into a classification number. Broadly stated, these classification systems can be grouped into three broad types, Enumerative, Faceted, and Analytico-synthetic.

- Enumerative: Subjects are pre-defined and listed in a hierarchical notation. Application of the classification system involves finding the appropriate class in the classification system and applying the class without modification.
- Faceted: Faceted systems are non-hierarchical and involve the combination of multiple categorization areas (or facets) to create a classification. One of the most popular faceted classification systems is Ranganathan’s colon classification. Ranganathan’s system features five facets: Personality, Matter, Energy, Space, and Time (PMEST).
- Analytico-synthetic: Analytico-Synthetic systems are hierarchical but rather than relying on a completely pre-defined hierarchy, they allow the cataloger to add refining concepts to classification such as geographic, temporal, and topical refinements. In addition, an Analytico-Synthetic system allows the classifier to build a classification number using the combination of hierarchical and refining concepts.

These three types of classification systems (Enumerative, Faceted, and Analytico-synthetic) are the most common in LAM communities. In addition to these systems there are social classification systems known as folksonomies, which rely on the aggregation of tags assigned by users of information resources. Folksonomies are often represented in Tag Clouds, a visual representation of tags with emphasis based on tag occurrence.

Generally speaking, the Library of Congress Classification System and Dewey Decimal Systems are considered Analytico-Synthetic because they blend a hierarchical subject analysis (Enumerative) with refining classification schedules (quasi-faceted). For example, the LCC system allows you to assign geographic and time facet refinements to a subject classification and the Dewey system features 10 main divisions that are hierarchically arranged to create a classification.

Classification is not always a manual process or even a process that relies on human interpretation of meaning. Computational linguistic techniques including document clustering, topic modeling, and keyword clustering are each techniques that rely primarily on computer analysis of texts and the derivation of classification. While classification of bibliographic works in libraries is important in order to assign a unique call number for a resource, classification does not always require the assignment of a unique call number. In fact in the case of digital resources, the process of classification as well as the utility derived from it is a topic of research.

## **THE INTERNET AND ITS IMPACT ON INFORMATION USE**

The world of information is a very different place from what it was before the introduction of the Internet to the general public. The web has changed the way